# Political Game Theory

Nolan McCarty

Adam Meirowitz

To Liz, Janis, Lachlan, and Delaney.

# Contents

# Acknowledgements

The origin of this book lies in the utter inability of either of its authors to write legibly on a black board (or any other surface, for that matter). To save our students from what would have been the most severe form of pedagogical torture, we were forced to commit our lecture notes to the electronic form. This also compensated for our inability to spell without the aid of a spell checker.[1] We ultimately decided that all of the late nights spent typesetting game theory notes should not go in vain. So we undertook to turn them into this book, which of course, led to more late nights spent typing. We hope these weren't wasted either.

We are most grateful to the students at Columbia and Princeton on whom early versions of our notes and manuscript were inflicted. Puzzled looks and panicked office hours helped us to figure out how convey game theory to students of politics. We also benefited from early conversations with Chris Achen, Scott Ashworth, Larry Bartels, Keith Krehbiel, David Lewis and Thomas Romer on what a book on political game theory ought to look like. Along the way Stuart Jordan and Natasha Zharinova have provided valuable assistance and feedback. Finally, our greatest debts are to those who taught us political game theory: David Austen-Smith, Jeffrey Banks, David Baron, Bruce Bueno de Mesquito, Thomas Romer, and Howard Rosenthal.

Nolan McCarty
Adam Meirowitz

Princeton NJ

---

[1] Our mispelling styles are quite distinctive, however. For a given word, McCarty uses completely random spellings while Meirowitz consistently mispells the word in the exact same way.

CHAPTER 1

# Introduction

In a rather short period of time, game theory has become one of the most powerful analytical tools in the study of politics. From its earliest applications in electoral and legislative behavior, game theoretic models have proliferated in such diverse areas as international security, ethnic cooperation, to democratization. Indeed all of the major fields in political science have been the recipients of important contributions from political game theoretic models. Rarely does an issue of *the American Political Science Review,* the *American Journal of Political Science*, or *International Organization* appear without at least one article formulating a new game theoretic application to politics or providing an empirical test of an existing one.

Nevertheless, applications of game theory have not developed as fast as they have in economics. One of the consequences of this uneven development is that most political scientists who wish to learn game theory are forced to rely on textbooks written by and for economists. While there are many excellent economic game theory texts, their treatments of the subject are often not well-suited to the needs of many political scientists. First and perhaps most importantly, the applications and topics are generally those of interest to economists. For example, it is not always obvious to novice political scientists what duopoly or auction theory tells us about political phenomena. Alternatively, there are topics such as voting theory that are indispensable to political game theorists which receive scant coverage in economics texts. Finally, many economics treatments presume some level of exposure to ideas in classical price theory. Thus, the entry barriers to political scientists are not only the math, but also a knowledge of demand curves, marginal rates of substitution, and the like.

Certainly, there have been a few texts by and for political scientists such as those by Ordeshook and Morrow. However, we feel that each is dated both in terms of the applications but also in terms of the needs of modern political science. Ordeshook remains an outstanding treatment of social choice and spatial theory, yet it was written well before the emergence of non-cooperative theory as the dominant paradigm in

political game theory. Morrow does provide an accessible introduction to the tools of non-cooperative game theory. However, the analytical level of his presentation falls somewhat below the contemporary needs of students of political game theory. It has also been a decade since its publication – a decade in which there have been hundreds of important articles and books deploying the tools of game theory. We feel that there is a need to introduce today's students to today's literature.

So we kept several goals in mind while writing this book. First, we wanted to write a textbook on political game theory instead of a book on abstract or economic game theory. We wanted to focus on applications of interest to political scientists. We wanted to present topics that are unique to political analysis. Secondly, in writing a book for political scientists, we wanted to be cognizant of the diversity of backgrounds and interests in political science. We recognize that most doctoral students in the field enter graduate school with limited mathematical and modelling backgrounds. However, we felt that it would not serve even those students to ignore the role of mathematical rigor and the importance of theoretical concepts in contemporary political models. Thus, for those students we have included a detailed mathematical appendix covering necessary tools ranging from set theory to basic optimization. At the other end of the spectrum are students who come to graduate study in political science with strong backgrounds in mathematics and economics. We wanted to write a book that would be useful to that audience as well and have chosen to provide in depth coverage of some more difficult and subtle concepts that are of the first importance to political game theory. As a result, we have included a number of "starred" advanced sections which provide a bit more detail about the analytical and mathematical structure of the models that we encounter. All of these can be safely skipped upon first readings for those not quite ready for the more technical material.

## 1. Organization of the Book

In terms of the organization, our book departs from standard treatments by including a number of topics that are either specifically relevant for political science or designed for remediation in areas which students of political science often have limited background. Chapter 2 is reasonably self contained exposition of classical choice theory under certainty. In this chapter, we lay out the basic ideas of preferences and their relation to utility theory. We prove a few key results, but otherwise focus on providing the intuition and language of rational choice theory. We also provide a section on spatial or Euclidean preferences

which play a key role in voting theory as well as applications in electoral and legislative politics.

In chapter 3, we consider how agents make choices under uncertainty. We develop the standard von Neumann-Morgenstern expected utility model, but also consider some of the most serious criticisms levied against it. In addition to the standard treatment of preferences for risk, we discuss the special implications of risk when actors have spatial preferences.

Chapter 4 is a cursory review of social choice theory. The chapter is not intended to be a replacement for full-length texts such as those by Peter Ordeshook and David Austen-Smith and Jeff Banks, but primarily as a reference for those ideas and concepts in social choice theory that have become integral parts of formal political science, such as the impossibility theorem and the non-existence of the majority core.

Chapter 5 begins our treatment of non-cooperative game theory which lies at the heart of contemporary formal political theory. We examine normal form games with complete information and present the fundamental solution concept of Nash equilibrium. The theoretical development is fairly standard but we include a number of important political applications. We review the standard Downsian model of electoral competition as well as some more recent extensions by Donald Wittman and Randy Calvert. We also present several models of private contributions to public goods based on the work of Thomas Palfrey and Howard Rosenthal. In chapter 6, we extend the normal form model where agents are uncertain about the payoffs associated with different strategy combinations. After presenting the relevant solution concept, Bayesian Nash equilibrium, we look at incomplete information versions of many of the models reviewed in chapter 5. This allows the reader to get a good sense of the strategic implications of uncertainty. We present the Palfrey-Rosenthal model with complete and incomplete information to give the reader a sense of the implications (and often the lack of ) of different informational assumptions.

Chapter 7 considers dynamic, multi-stage games of complete information and develops the notions of subgame perfection. Here we focus on a number of applications drawn from legislative politics, democratic transitions, coalition formation, and international crisis bargaining. In chapter 8, we consider dynamic games where some players are imperfectly informed about the payoffs to different strategies. After developing the solution concepts relevant for such models, we explore a number of applications drawn from legislative politics, campaign finance, and international bargaining. Much of this chapter is focused on the important and broadly applied class of signaling games.

Chapter 9 reviews the theory of repeated games and their application in political science. The role of discounting and structure of folk theorems in repeated games is the primary focus of the chapter.

In chapter 12 we consider various applications of bargaining theory. We begin with the canonical model of Rubinstein and its majority rule version developed by Baron and Ferejohn. We then consider several examples of bargaining with incomplete information.

In chapter 11, we present the mechanism design approach to modeling institutions. Here the focus is on the selection of games that induce equilibrium behavior to meet certain ends. We discuss a number of recent applications to electoral politics and organizational design. We build on the work of chapter 8, drawing connections between signaling games and mechanism design.

Finally, to keep the book as self contained as possible, chapter 12 provides a review of all of the mathematics that are used in the book. Topics which are integral to the development of key theoretical results or tools for analyzing applications are drawn from the fields of set theory, real analysis, linear algebra, calculus, optimization and probability theory. Indeed this chapter may serve as a basis for review, self-study or a formal course in mathematics for students interested in working at the frontier of political game theory.

CHAPTER 2

# The Theory of Choice

The starting point for almost all of political game theory is the idea that individuals rationally pursue goals subject to constraints imposed by physical resources as well as the behavior of other actors. Such an assumption is often controversial. Indeed one of the most contentious debates in political science is the role of rationality and intentionality as a predictor of political behavior. However, we will defer debates between *homo economicus* and *homo sociologicus* and jump right into the classical model of rational choice.

For almost all of our purposes, it is sufficient to define rationality in terms of a few simple ideas:

(1) Confronted with any two options, which we might denote $x$ and $y$, an individual can determine whether he does not prefer option $x$ to option $y$ or whether he does not prefer $y$ to $x$, or both. When preferences satisfy this property, we say they are *complete*.

(2) Confronted with three options $x$, $y$, and $z$, if an individual does not prefer $y$ to $x$ and does not prefer $z$ to $y$ then it must be the case that she does not prefer $z$ to $x$. Preferences satisfying this property are *transitive*.

Roughly speaking, our definition of rational behavior is that consistent with complete and transitive preferences. Sometimes behavior dictated solely by properties 1 and 2 is called "thin" rationality. This is because properties 1 and 2 are not predicated in any way on assumptions about the substantive content of human desires. Thus, thin rationality contrasts with "thicker" notions of rationality where specific goals such as wealth, status, or fame are postulated. The thin characterization of rationality is consistent with a very large number of these substantive goals. In principal, thinly rational agents could be motivated by any number of factors including ideology, normative values, or even religion. As long as these belief systems produce complete and transitive orderings over personal and social outcomes, we can model the behavior they produce using the classical model of choice.

While it may be appealing to deal with models that are independent of assumptions about specific goals, it will often be desirable to make stronger assumptions about preferences. For example, we might assume that interest groups wish to maximize the wealth of its members or that politicians wish to maximize their reelection chances. In subsequent chapters, we will explore models that makes these types of assumptions about agent preferences. But rational models may be just as useful for models of activists who wish to minimize environmental degradation or the number of abortions for principled, non-material reasons.

In the following sections, we develop the classical theory of choice under *certainty*. By certainty, we mean simply that the agent has sufficient information about the choice environment that she can perfectly predict the consequences of each of her actions. Thus, certainty means only that there is no analytical difference between assuming that political actors choose actions based on the desired outcomes which result from those actions, or that they choose those outcomes directly. In later chapters, we will examine choice under uncertainty – the actor's lack of knowledge of some feature of the choice environment leads her to choose actions which have uncertain consequences.

## 1. Finite Sets of Actions and Outcomes

We begin by considering the simple case of an agent who faces a finite numbers of actions from which to choose. We denote these choices as a set $A = \{a_1, ..., a_k\}$. For example, a leader involved in an international crisis might face the following set of alternatives $A = \{$send in the troops, negotiate, do nothing$\}$ whereas an American voter might choose among $A = \{$vote Democrat, vote Republican, abstain$\}$.

As mentioned above in this chapter we assume that agents have *complete* information, that is they are sufficiently knowledge about the context of their choices that they can perfectly predict the consequences of each action. To capture this idea formally, we define outcome sets $X = \{x_1, ..., x_n\}$. Following one of our examples, these might be $X = \{$win major concessions and lose troops, win minor concessions, status quo$\}$. The assumption of certainty then implies that each action $a \in A$ maps directly on to one and only one $x \in X$. Formally, we assume that there exists a function $x : A \to X$ that maps each action into a specific outcome. We also assume that all of the outcomes in $X$ are feasible, that is each is the consequence of at least one action. Formally, $x_i$ is feasible if there exist an $a \in A$ such that $x(a) = x_i$. With the assumptions of certainty and feasibility, it makes

little difference whether we speak of an agent's preferences over actions or his preferences over outcomes. Thus, we concentrate on the agent's preferences over outcomes and are uninterested in the action that she must choose to attain the desired outcome. In chapter 3, the assumption of uncertainty or *incomplete* information makes the distinctions between actions and outcomes relevant.

We now turn to the concept of preferences and the types of restrictions that our two simple notions place on what outcome rational individuals may choose. Formally, preferences are modelled as a binary relation $R$ which represents "weak preference." The notation $x_i R x_j$ means that outcome $x_j$ is not preferred to policy $x_i$ or that $x_i$ is "weakly" preferred to $x_j$.[1] To help cement ideas, note that $R$ is similar to the binary relation $\geq$ (greater than or equal) which operates on real numbers.

Given the weak preference relation $R$, we define two other important binary relations: strict preference and indifference.

DEFINITION 2.1. *Given any $x, y \in X$ we say $xPy$ if and only if $xRy$ and not $yRx$. We say $xIy$ iff $xRy$ and $yRx$.*

Accordingly, $P$ denotes strict preference and $I$ denotes indifference. Returning to the example of $\geq$ on $X$, the strict preference relation derived from $\geq$ is equivalent to the relation $>$ and the indifference relation is equivalent to the relation $=$.

While preferences, in the form of binary relations, are a useful starting point, it is choices that we are ultimately interested in. Given a set of preferences, we could hardly call an agent's behavior rational unless she selected an outcome that she valued at least as much as any other. Consequently, we expect a rational agent to choose an $x^* \in X$ for which $x^* R y$ for every $y \in X$. However, without adding a little bit more structure on to her preferences, there is no guarantee that such a maximal outcome exists. Thus, we now turn to the conditions on $X$ and $R$ that insure such a "best" choice is meaningful and well-defined. We start with the following formal definition.

DEFINITION 2.2. *Given a set $X$ and weak preference relation $R$ on $X$, the maximal set $M(R, X) \subset X$ is defined as follows $M(R, X) = \{x \in X : xRy \ \forall \ y \in X\}$*

Thus, a fundamental tenant of rationality should be that **agents choose outcomes from the maximal set.** Of course, this requirement only makes sense if the maximal set is not empty i.e. $M(R, X) \neq$

---

[1]Formally, $R$ is a subset of $X \times X$ such that if $(x, y) \in R$ than $xRy$.

$\emptyset$. Thus, we are most interested in the properties of preferences that guarantee that $M(R, X)$ has at least one element.

The most obvious problem that might lead to an empty maximal set is that $R$ is silent between a pair of outcomes say $x$ and $y$. If neither $xRy$ or $yRx$ then it is not clear what a "rational choice" would be. Two conditions insure that all elements of $X$ are ordered.

DEFINITION 2.3. *A binary relation $R$ on $X$ is*
*(i) complete if for all $x, y \in X$ with $x \neq y$, either $xRy$ or $yRx$ or both.*
*(ii) reflexive if for all $x \in X$, $xRx$.*

Completeness simply means that the agent can compare any two outcomes. This is probably not a terribly controversial assumption (though we all know people who can't seem to make their minds up). Reflexivity is a technical condition and some authors choose to define completeness in a slightly different manner that also captures what we call reflexivity.

While completeness and reflexivity get us closer to a "rational" preference relation they are not sufficient. We need to rule out problems like $xRy$, $yRz$ and $zRx$. The problem with such preferences is that there is no reasonable choice–why choose $y$ when you can choose $x$, why choose $x$ when you can choose $z$, and why choose $z$ when you can choose $y$. Each of the following conditions on preferences resolve this problem.

DEFINITION 2.4. *A binary relation $R$ on $X$ is*
*(1) Transitive if for all $x, y, z \in X$  if $xRy$ and $yRz$ then $xRz$.*
*(2) Quasi-transitive if for all $x, y, z \in X$ if $xPy$ and $yPz$ then $xPz$.*
*(3) Acyclic if for all $\{x, y, z, ...., a, b\} \in X$  if $xPy$ and $yPz$ ... and $aPb$ then $xRb$*

Note the subtle differences among these definitions. Transitivity and quasi-transitivity may seem innocuous but they are reasonably strong assumptions which might be violated even by very reasonable behavior. For example, suppose $X$ is a set of 1000 different bottles of beer. Beer $b_1$ has had one drop of beer replaced with one drop of plain water while $b_2$ has had two drops replaced and so on to $b_{1000}$. Unless you are a master brewer, $b_1 I b_2$, and $b_2 I b_3, \ldots$, and $b_{999} I b_{1000}$. Since $xIy$ implies $xRy$ (by the definition of $I$) this would imply that $b_{1000} R b_{999}, \ldots$, $b_2 R b_1$ and if the relation is transitive we are left with $b_{1000} R b_1$. But clearly, $b_1 P b_{1000}$.[2] Note however that the assumption of

---

[2]This is approximately the difference between Guiness and Coors Light.

acyclicity does not suffer from this problem and is typically sufficient for our purposes. However, despite the problems associated with transitivity, we will maintain it as an assumption (rather than acyclicity) to simplify many of the results below.

DEFINITION 2.5. *Given a set $X$ a weak ordering is a binary relation that is complete, reflexive and transitive.*

It is not difficult to see that transitivity rules out exactly the cycle $xRy$, $yRz$ and $zRx$ considered above. Note that our recurring example of $\geq$ satisfies all of the conditions for a weak ordering. We can now state our first result.

THEOREM 2.1. *If $X$ is finite and $R$ is a weak ordering then $M(R,X) \neq \emptyset$.*

Thus, we can guarantee that there is a best choice so long as we are willing to assume that the choice set is finite and that $R$ is complete, reflexive, and transitive. The intuition behind this theorem is quite straightforward. Consider an outcome set $X$ with three elements say $x,y$, and $z$. If $M(R,X) = \emptyset$, then by definition $x$ and $y$ must be weakly preferred to $z$, $y$ and $z$ must be weakly preferred to $x$, and that $x$ and $z$ must be weakly preferred to $y$. Thus, the only possibility that does not violate transitivity is that $xIyIz$ which implies that $M(R,X) = X$. We can extend this logic to any size $X$.

PROOF. Assume that $X$ is finite and $R$ is complete, reflexive, and transitive. We establish the result by induction on the number of elements in $X$.

**Step 1:** If $X$ has 1 element (i.e. $X = \{x\}$), then by reflexivity $xRx$ and thus $M(R,X) = \{x\}$.

**Step 2:** We show that if it is true that for any $X'$ with $n$ elements $R'$ a weak ordering implies that $M(R',X') \neq \emptyset$ then for any $X$ with $n+1$ elements when $R$ is a weak ordering on $X$, $M(R,X) \neq \emptyset$.

-**Proof of step 2:** assume that for any $X'$ with $n$ elements $R'$ a weak ordering on $X'$ implies that $M(R',X') \neq \emptyset$. Now consider a set $X$ with $n+1$ elements. For arbitrary $x \in X$ it is true that $X = X' \cup \{x\}$ with $X'$ a set having $n$ elements. By assumption $M(R',X') \neq \emptyset$. So for an arbitrary $y \in M(R',X')$ either $yRx$ or $xRy$ or both by completeness.

–If $yRx$ then since $y \in M(R',X')$ we have $yRz$ for all $z \in X' \cup \{x\}$ and thus $y \in M(R,X)$ and the step 2 result is established.

–If it is not the case that $yRx$ then we have $xRy$. Since $y \in M(R',X')$ we have $yRz$ for any $z \in X'$. Thus for any $z \in X'$ we have

$xRy$ and $yRz$.Since $R$ is transitive this implies that we have $xRz$ for any $z \in X'$. This and $xRy$ imply that for any $w \in X'$ we have $xRw$ and thus $x \in M(R, X)$ and the step 2 result is established.

**Step 3:** By steps 1 and 2 for any finite sized $X$ if $R$ is a weak order on $X$ then $M(R, X) \neq \emptyset$.■                                           □

It turns out that a weak ordering is not needed for $M(R, X)$ but the proof is a bit more complicated so we leave it for an exercise.

THEOREM 2.2. *Assume $X$ is finite and $R$ is a complete and reflexive binary relation on $X$. $M(R, S) \neq \emptyset$ on any $S \subset X$ (except $S = \emptyset$) iff $R$ is acyclic.*

Even with a finite choice space and no uncertainty the theory of choice is fairly rich. Austen-Smith and Banks (1999) is a good first source for students interested in going further. Many economists and psychologists, have been concerned about the assumption of completeness and a theory of choice without this condition has been derived.

In the next, more technical, section, we consider rational choice when the outcome space is not finite, such as the real line. We derive an analog to theorem 1 for non-finite choice sets. While the results are conceptually similar, additional mathematical structure needs to be placed on the choice sets and preferences.

## 2. Continuous Outcome Spaces*

**2.1. Non-emptyness of $M(R, X)$.** Examination of the argument for theorem 1 demonstrates that the fact that the choice space was finite was useful. This allowed us to prove the result by induction for any number of outcomes. However, when there is an infinite number of choices, this approach is mathematically inappropriate. Thus, when agents choose from a choice space that is a continuum (e.g. the set of real numbers denoted $\mathbb{R}$ or the set $[0, 1] = \{x \in \mathbb{R} : x \geq 0 \text{ and } x \leq 1\}$) more structure on preferences is needed to insure that $M(R, S) \neq \emptyset$. Two simple examples demonstrate why things can go awry.

EXAMPLE 2.1. *Let $X = (0, 1)$ (or let $X = \mathbb{R}^1$) and let $R$ on $\mathbb{R}^1$ be equivalent to $\geq$ so that $xRy$ iff $x \geq y$   The set $M(\geq, X)$ is empty.*

To see why $M(\geq, (0, 1))$ is empty, note that for every $x \in X$ there exists some $y \in X$ for which $y > x$. Thus there can be no $x$ such that $xRy$ for all $y \in X$. In this example the fact that $(0, 1)$ has no biggest element results in the emptiness of the maximal set. Note however that if $X$ were a closed interval such as $[0, 1]$, we wouldn't have a problem as $M(\geq, [0, 1]) = 1$. This is a strong hint that general results about

the non-emptiness of the maximal set may depend on the choice set
being "closed."

EXAMPLE 2.2. *Let $X = [0, 1]$ and define $R$ on $\mathbb{R}^1$ as follows: $xRy$
if $x, y \leq \frac{1}{2}$ and $x \geq y$ or if $x, y > \frac{1}{2}$ and $x \leq y$ or if $x > \frac{1}{2}$ and $y \leq \frac{1}{2}$.
The set $M(R, X)$ is empty.*

To see this note that no element of $[0, \frac{1}{2}]$ can be in $M(R, X)$ as
any element of $(\frac{1}{2}, 1]$ is weakly preferred. However, elements of $(\frac{1}{2}, 1]$
cannot be part of $M(R, X)$ for reasons identical to that of the previous
example. Here the problem is not with $X$ – it is a closed interval as
we required to make the first example work. Instead, the problem is
with $R$. It jumps around at $\frac{1}{2}$. Outcomes slightly less than or equal
$\frac{1}{2}$ are among the least preferred while those slightly greater are among
the most preferred. It is this "discontinuity" in preferences that leads
to the empty maximal set in the example.

Before turning these examples and intuitions into general axioms,
we need to review a few mathematical concepts.[3] We begin with the
assumption that preferences are defined on $n$-dimensional Euclidean
space, and consider choice from subsets, $X \subset \mathbb{R}^n$. A point in such
a space can be written as a vector $x = (x^1, x^2, ...., x^n)$ where each
coordinate $x^i$ is a point in $\mathbb{R}^1$.

One of our main concerns is whether the set $X$ is open or closed.
These concepts can be grasped with the simplest example of $R^1$. In
this case a set $A \subset X$ is termed **open** if for every point $x \in A$ there
is some $\varepsilon > 0$ such that for any $y \in X$ satisfying $|x - y| < \varepsilon$ it is the
case that $y \in A$. Therefore, a set is open if given any point in the
set, all the points close to it are also in the set. Clearly, the set $(0, 1)$
is open. For each point in the set, there are some points higher and
some points lower which are also in the set. Thus, for any point $x$,
we can choose $\varepsilon$ so that $x - \varepsilon$ and $x + \varepsilon$ are also in the set. We say
that a set is closed if its complement is open. Therefore, since $(0, 1)$ is
open, $(-\infty, 0] \cup [1, \infty)$ is closed. Intervals such as $[0, 1]$ are also closed.
Some sets may be neither open or closed such as $[0, 1)$.

---

[3]More precisely we need a few Topological concepts. Students interested in
further study of choice theory would be well served examining the mathematical
appendix to this book or better, yet, a text on Real Analysis. An approach-
able introductory text is: Gaughan, Edward. 1993. *Introduction to Analysis, 4ed.*
Brooks/Cole Publishing Company. A more complete text is: Kolmogorov, A.N.
and S.V. Fomin. 1970. *Inroductory Real Analysis.* Dover.

To generalize these concepts to the $n$-dimensional Euclidean space, we begin with a measure of distance call the *norm*.

$$\|x - y\| = \sqrt{\sum_{i=1}^{n} (x^i - y^i)^2}.$$

The quantity $\|x - y\|$ is the distance between points $x$ and $y$ and generalizes the absolute value used in $\mathbb{R}^1$. Given this definition of distance, we can generalize the notion of an interval into that of a "ball."

DEFINITION 2.6. *An open ball of radius $\varepsilon > 0$ and center $x \in X$ , is denoted $B(x, \varepsilon) = \{y \in X : \|x - y\| < \varepsilon\}$.*

Now it is easy to generalize the concept of openness. A set is open if the set contains an open ball around each point for some $\varepsilon > 0$.

DEFINITION 2.7. *A set $A \subset \mathbb{R}^n$ is open if for every $x \in A$ there is some $\varepsilon > 0$ such that $B(x, \varepsilon) \subset A$.*

Just as before, a set closed if its complement is open. Thus, closed sets have the property that some points are on the boundary so that there does not exist an open ball that does not contain points outside the set.

DEFINITION 2.8. *A set $A \subset \mathbb{R}^n$ is closed if its complement $B = \mathbb{R}^n \backslash A$ is an open set.*

Recall our first example. Since $X$ is an open set, for each $x$ in $X$ there is an open ball around $x$ that is also in $X$. Since each of these balls contain points weakly preferred to $X$, no maximal set can exist. However, if $X = [0, 1]$, any open ball around 1 contains points outside of $X$. Since all of the points preferred to 1 lie outside of $[0, 1]$, $M(\geq, [0, 1]) = 1$. However, note that closedness itself is not sufficient. Recall that $(-\infty, 0] \cup [1, \infty)$ is a closed set, but $M(\geq, (-\infty, 0] \cup [1, \infty))$ is empty. The problem of course is that there is no upper bound on this set, so for any $x$ there is a $y > x$ so that $yRx$. Thus, another important condition is *boundedness*.

DEFINITION 2.9. *A set $A \subset \mathbb{R}^n$ is bounded if there exists a finite number $b$ such that for every $x \in A$ it is the case that $\|x - \mathbf{0}\| < b$ where $\mathbf{0}$ is the vector $(0, ..., 0)$.*

The set $(-\infty, 0] \cup [1, \infty)$ clearly fails this criteria so we can rule it out by requiring that choice sets be bounded. It is easy to see in example 1 so long as $X$ is closed and bounded $M(\geq, X)$ is non-empty. In $\mathbb{R}^n$, we call sets that are both closed and bounded *compact*.

DEFINITION 2.10. *A set $A \subset \mathbb{R}^n$ is compact if it is closed and bounded.*

Since all examples or problems and problems in this book will deal with subsets of Euclidean spaces, we could stop here. However, in arbitrary spaces, the equivalence between compactness and closed and bounded does not hold. It turns out that the proof of our main result is easier if we consider a more general definition of compactness (even though we lose some of the intuition of our examples). The more general definition of compactness is based on sets known as *open covers*. An open cover for a set $A$ is a collection of open sets whose union contains $A$.

DEFINITION 2.11. *Given a set $A$, an open covering of $A$ is a collection of sets $\{O_\theta\}_{\theta \in \Theta}$ where $\Theta$ is an arbitrary index set and $O_\theta$ is open for every $\theta \in \Theta$ such that $A \subset \{\cup_{\theta \in \Theta} O_\theta\}$ (in other words if $x \in A$ then there is some $\theta \in \Theta$ such that $x \in O_\theta$).*

Given this definition, we say that $A$ has a finite sub-cover if from every open cover we can select just a finite number of the open sets and be assured that $A$ is covered by this finite collection. The existence of such a sub-covering is equivalent to the compactness of $A$.

DEFINITION 2.12. *A set $A$ is compact if for any open covering $\{O_\theta\}_{\theta \in \Theta}$ of $A$ there exists some finite set $B \subset \Theta$, such that the finite covering $\{O_\theta\}_{\theta \in B}$ is a covering of $A$ i.e. $A \subset \cup_{\theta \in B} O_\theta$.*

Since the previous two definitions are subtle for those not familiar with analysis an example is warranted. Consider the space $\mathbb{R}^1$ and two subsets $[0, 1]$ and $(0, 1)$. We already know that $(0, 1)$ is not compact, as it is not closed and we have concluded that in Euclidean space compact sets are closed and bounded. To demonstrate that $(0, 1)$ is not compact using the open covering definition, consider the following open covering of $(0, 1)$. For each $\theta \in \Theta = \{3, 4, 5, .......\}$, let $O_\theta = (\frac{1}{\theta}, 1 - \frac{1}{\theta})$. This is a collection of open intervals centered at $\frac{1}{2}$, and the width of the intervals approaches 1 as $\theta$ gets big. Is $\{O_\theta\}_{\theta \in \Theta}$ and open covering of $(0, 1)$? Yes, for any element in $x \in (0, 1)$ you can pick a $\theta$ big enough so that $x \in O_\theta$. So we have constructed an open covering of $(0, 1)$. Now our definition of compactness says that if $(0, 1)$ is compact we need to be able to find a finite subset $B \subset \{3, 4, 5, .......\}$ and show that $(0, 1) \subset \cup_{\theta \in B} O_\theta$. But for a finite set $B$, there is a finite largest element $\theta^* \in B$.[4] This means that the value $\frac{1}{\theta^*}$ is strictly larger than 0 and

---

[4]You could actually prove this sentence by noting that $\geq$ is a weak ordering and applying our result about the non-emptyness of the maximal set for finite sets.

since $(0, 1)$ contains points arbitrarily close to 0, $\frac{1}{\theta^*}$ is strictly larger than some element of $(0, 1)$. Accordingly for any finite collection of subsets in the open covering (i.e. selection of $B$ that is finite), we can find an element of $(0, 1)$ that is not contained in any set $O_\theta$ for $\theta \in B$. Thus we have applied the open-covering definition to show what we already knew, $(0, 1)$ is not compact. The interested reader should try to proceed in the other direction, showing that $[0, 1]$ is compact according to the open-covering definition.[5]

Having elucidated properties that $X$ can satisfy (e.g. compactness), we turn to the properties that we would like to $R$ to satisfy. Not surprisingly, given example 2, we want $R$ to be "continuous" in a specific way. To define continuity, we use the concept of the upper contour set.[6] Given a binary relation $R$ on $\mathbb{R}^n$ the strict upper contour set of a point $x \in \mathbb{R}^n$ is $P(x) \equiv \{y \in \mathbb{R}^n : yPx\}$. The strict lower contour set of point $x$ is the set $P^{-1}(x) \equiv \{y \in \mathbb{R}^n : xPy\}$. So the upper contour set of $x$ contains the points that are preferred to $x$ and the lower contour set of $x$ contains the points that $x$ is preferred to. Similarly, the level set of $x$ is the set of points for which the agent is indifferent to $x$ or $I(x) \equiv \{y \in \mathbb{R}^n : yRx \text{ and } xRy\}$.

DEFINITION 2.13. *A binary relation $R$ on $\mathbb{R}^n$ is*
*(i) upper continuous if for all $x \in \mathbb{R}^n$, $P(x)$ is open*
*(ii) lower continuous if for all $x \in \mathbb{R}^n$, $P^{-1}(x)$ is open*
*(iii) continuous if is both lower and upper continuous.*

Consider the implications of these conditions. When preferences are complete, any point $y$ that is very close to $x$ is either in $P(x)$, $P^{-1}(x)$, or $I(x)$. When preferences are continuous, $y \in P(x)$ or $y \in P^{-1}(x)$ implies that points sufficiently close to $y$ will also be in the respective set. To see how continuity rules out anomalous behavior, recall example 2. There $P^{-1}(\frac{1}{2} + \varepsilon) = (-\infty, \frac{1}{2}] \cup (\frac{1}{2} + \varepsilon, 1]$ which is not an open set. Thus, the jump in preferences exhibited in the second example is ruled out when preferences are lower continuous.

We can now state the sufficient conditions for a non-empty maximal set.

---

[5]We direct the reader to the famed Heine-Borel theorem in any of the cited texts on Real-analysis which relates the topological open-covering and Euclidean closed-and bounded definitions of compactness for $\mathbb{R}^n$. Gaughan (1993) presents a particularly detailed proof of the result for $\mathbb{R}^1$.

[6]In political science, the upper contour set is often referred to as the "preferred to set". Keith Krehbiel has pointed out to both authors on numerous occassions that this terminology (along with many others) contains a redundancy. Thus, he and we implore all readers to use our preferred "preferred set."

THEOREM 2.3. *If $X \subset \mathbb{R}^n$ is non-empty and compact, and $R$ on $\mathbb{R}^n$ is complete, reflexive, transitive and lower continuous, then $M(R, X) \neq \emptyset$*

The proof of this result is more technical than most other sections of this book. The result also holds on for arbitrary topological spaces. This allows us to apply it to choice problems in which $x$ is a infinite sequence of outcomes, a function, or a probability distribution.

PROOF. Assume that $X$ is non-empty and compact, and $R$ on $X$ is complete, reflexive, transitive and lower continuous. To establish a contradiction, assume that $M(R, X) = \emptyset$. Thus, every point in $X$ is contained in $P^{-1}(\alpha)$ for some $\alpha \in X$. Since $R$ is lower continuous every such $P^{-1}(\alpha)$ is open. This means that $\{P^{-1}(\alpha)\}_{\alpha \in X}$ is an open covering of $X$. Since $X$ is compact, there exists a finite set of points $B \subset X$ for which the collection $\{P^{-1}(\alpha)\}_{\alpha \in B}$ is also a covering of $X$. That is if $x \in X$ then $x \in P^{-1}(\alpha)$ for some $\alpha \in B$. But we know from a previous result that $M(R, B) \neq \emptyset$, since $B$ is finite and $R$ is complete, reflexive and transitive. Thus a point $x^* \in M(R, B)$ exists. Now consider any arbitrary point $y \in X$. Either $y$ is an element of $M(R, B)$ or it is not. By definition, if $y \in M(R, B)$ then $x^*Ry$. If $y \notin M(R, B)$ since $\{P^{-1}(\alpha)\}_{\alpha \in B}$ covers $X$ there is some $\alpha \in B$ such that $y \in P^{-1}(\alpha)$. This means that $\alpha Ry$. However, since $x^* \in M(R, B)$ we know that $x^*R\alpha$. Since $R$ is transitive on $X$, this implies that $x^*Ry$. Thus, we have shown that for all $y \in X$, $x^*Ry$. This means that $x^* \in M(R, X)$, contradicting the assumption and establishing the non-emptiness of $M(R, X)$.∎ □

It is important to note that the theorem only establishes sufficient, not necessary conditions, for a non-empty maximal set. In particular, we will encounter situations where $X$ is either unbounded or not closed and $R$ is discontinuous. In each of these possibilities, the non-emptiness of $M(R, X)$ has to be established by other means. Violations of the compactness of $X$ will generally require stronger assumptions about $R$ while violations of continuity will require more structure on $X$.

**2.2. Uniqueness of** $M(R, X)$**.** Since we develop choice models to make predictions about behavior, it is certainly preferable that the model produce a single prediction, rather than a range of possible outcomes. Thus, it is valuable to know whether or not $M(R, X)$ has a unique element or whether a larger set of choices is consistent with rational behavior.

When the choice set is finite, we can typically only guarantee a unique element of $M(R, X)$ by assuming that all preferences are strict. Without indifference and a finite choice set, $M(R, X)$ must have only a single element if it exists.

When the choice space is not finite, we can impose sufficient structure to insure that $M(R, X)$ has only one element. We will need one condition on $X$ and one condition on $R$ to attain this result. The first condition is that $X$ be a convex set. This assumption requires that if $x$ and $y$ are in $X$ certain combinations of $x$ and $y$ must also be in $X$.

DEFINITION 2.14. $X \subset \mathbb{R}^n$ *is convex if for any* $x, y \in X$ *and the point* $\lambda x + (1 - \lambda)y$ *is in $X$ for every* $\lambda \in [0, 1]$.

The point $\lambda x + (1 - \lambda)y$ is often called the convex combination (or a weighted average) of $x$ and $y$. As an example the set $[0, 1]$ is convex because for any two points in the set, any point in between these two points is also in the set. Alternatively, $X = [0, \frac{1}{4}] \cup [\frac{3}{4}, 0]$ is not since $\frac{1}{4}\lambda + \frac{3}{4}(1 - \lambda) \notin X$ for any $\lambda \in (0, 1)$. Thus, convexity requires that there are no "holes" in the outcome set. When the outcome has more than dimension, convexity also requires that its surface not have any appendages. Think about your hand. Convex combinations of points on your thumb and index finger are not part of it.[7]

We will see that an important property is that preferences be convex as well.

DEFINITION 2.15. *Preference $R$ on the convex set $X$ is strictly convex if for any distinct points* $x, y \in X$ *if $xRy$ then* $[\lambda x + (1 - \lambda)y] \, Py$ *for any* $\lambda \in (0, 1)$.

Essentially, convex preferences have the property that if the agent prefers $x$ to $y$ she should also prefers convex combinations of $x$ and $y$ to $y$. Strict convexity goes a step further. Even if the agent is only indifferent between $x$ and $y$, she should still prefer the convex combination to either $x$ or $y$. We leave as an exercise that the strict convexity of $R$ implies that the lower contour sets $P^{-1}(x)$ are convex. Since the lower contour sets are convex, they cannot have holes, appendages, or flat spots.

The following result is easy to establish.

THEOREM 2.4. *If $X$ is convex and $R$ on $X$ is strictly convex, then if $M(R, X)$ is non-empty, it contains a single element.*

PROOF. By way of a contradiction assume that $X$ is convex, $R$ is strictly convex, and two distinct policies $x, y$ are both in $M(R, X)$.

---

[7]Game theorists spend a lot of time contemplating such ironies.

For arbitrary $\lambda \in (0,1)$ the point $[\lambda x + (1 - \lambda)y]$ is in $X$ since $X$ is convex. But since $R$ is strictly convex, $[\lambda x + (1 - \lambda)y] \, Py$. But this contradicts the assumptions that $y \in M(R,X)$. Thus the result is established.∎                                                                □

The importance of the last two theorems is clear. When the choice set is compact and a weak order is lower continuous, a "rational" choice exists. When the choice set is convex and the preference ordering is strictly convex, any optimal choice is unique.

## 3. Utility Theory

So far our model of choice and rationality is based on the use of binary preferences and the maximal set. However, binary operators are hard to work with except in the most trivial models. Numbers on the other hand are easy to work with. So if we can associate a number with each element of the outcome set, then we can just use the $\geq$ operator to compare alternatives. In this section we explore the conditions under which it is possible to represent outcome sets as sets of real numbers so that we can use $\geq$ as the preference operator. In other words, we would like to represent preferences using a utility function (a real valued function with domain $X$) such that

$$u(x) \geq u(y) \text{ implies } xRy$$
$$u(x) > u(y) \text{ implies } xPy$$
$$u(x) = u(y) \text{ implies } xIy$$

The idea of utility has been the subject of numerous philosophical and moral debates over the past 300 years, but again we will use a narrow definition. Utilities simply numerical representations of preferences for which $\geq$ is the appropriate preference operator – we imbue them with no additional normative content.

At our current level of generality, utility functions are ordinal as they are used only to rank alternatives. In particular, they are not used to tell is how much something is preferred to something else. The value $u(x) - u(y)$ has no meaning, because any function $w$ such that $w(x) \geq w(y)$ if and only if $u(x) \geq u(y)$ represents exactly the same preferences as $u$. This indicates that comparing utilities across agents is generally not a meaningful exercise. However, as we discuss in the next chapter, the standard model of choice under uncertainty presumes that utility functions contain more than simple ordinal information.

The following is a formal definition of a utility function.

DEFINITION 2.16. *Given $X$ and $R$ on $X$ we say the utility function $u : X \to \mathbb{R}^1$ represents $R$ if for all $x, y \in X$ $u(x) \geq u(y)$ iff $xRy$.*

From this definition it is quite easy to show that $u(x) > u(y)$ if and only if $xPy$ and $u(x) = u(y)$ if and only $xIy$. When $X$ is finite the existence of a utility representation of $R$ hinges only on $R$ being complete, reflexive,and transitive.

Just as we did in the last section, we wish to characterize the agent's optimal choice. Let $x$ be a maximizer of $u : X \to \mathbb{R}^1$ if $u(x) \geq u(y)$ for all $y \in X$. As the next result shows the existence of a maximizer and the non-emptiness of $M(R, X)$ are equivalent.

THEOREM 2.5. *If the function $u(\cdot)$ is a utility representation of $R$ on $X$ then $M(R, X) = \arg\max_{x \in X}\{u(x)\}$.*

PROOF. To show that $M(R, X) \subset \arg\max_{x \in X}\{u(x)\}$, assume that $u(\cdot)$ represents $R$ on $X$ and that $x' \in M(R, X)$. The latter assumption implies that $x'Ry$ for all $y \in X$. This and the former assumption imply that $u(x') \geq u(y)$ for all $y \in X$. Thus $x \in \arg\max_{x \in X}\{u(x)\}$. To show that $\arg\max_{x \in X}\{u(x)\} \subset M(R, X)$ assume that $u(\cdot)$ represents $R$ on $X$ and that $x' \in \arg\max_{x \in X}\{u(x)\}$. The latter assumption implies that $u(x') \geq u(y)$ for all $y \in X$. This and the former assumption imply that $x'Ry$ for all $y \in X$. Thus $x \in M(R, X)$.∎                               □

If $X$ is finite and $R$ is complete, reflexive and transitive, we know that $M(R, X)$ is non-empty (by theorem 1), thus a maximizer of $u(x)$ must exist. However, if $X$ is not finite, further conditions on $X$ and the function are needed to ensure the existence of maximizers. In the next advanced section we consider utility functions on non-finite outcome spaces.

## 4. Utility representations on Continuous Outcome Spaces*

We now review some basic properties of functions. The first desirable property of utility functions is continuity.

DEFINITION 2.17. *We say a function $f : X \to \mathbb{R}^1$ is continuous if for every $x \in X$ the following is true: For every $\varepsilon > 0$ there exists some $\delta > 0$ such that if $\|x - y\| < \delta$ $|f(x) - f(y)| < \varepsilon$.*

As is often taught to high school students, a continuous function is one that you can draw without lifting the pencil. Substantively, a continuous utility function is one that produces close utilities for outcomes that are close together.

The following sufficient conditions on preferences ensure that a continuous utility representation exists.

THEOREM 2.6. *(Debreu 1959) If $X \subset R^n$ and $R$ is complete, reflexive, transitive, and continuous, then there exists a continuous utility function $u : X \to \mathbb{R}^1$ that represents $R$.*

We will not undertake the proof of this claim. However the converse is not difficult to establish which we leave as an exercise. A result analogous to Theorem 2.3 is the following.

THEOREM 2.7. *If $X \subset \mathbb{R}^n$ is compact and $u : X \to \mathbb{R}^1$ is continuous then a maximizer exists.*

This result is sometimes known as the Weierstrass Theorem. We do not prove the result here (see Royden for a proof), since Theorem 2.3 is actually a result showing that only lower continuity and compactness are needed.

As we pointed out earlier, utility functions are arbitrary. Some texts call this an ordinal notion of utility as opposed to a cardinal notion of utility. There is nothing interesting about the particular value of a utility function at a specific point $x \in X$. All that matters is the ordering of $u(x)$ and $u(y)$ for any two $x, y \in X$. We say that $f : R^1 \to R^1$ is a strictly increasing function if for all $x, y \in X$ $x > y$ implies that $f(x) > f(y)$. Utility functions are defined only up to a strictly increasing transformation. This means that if $u : X \to \mathbb{R}^1$ represents $R$ on $X$ then $f \circ u : X \to \mathbb{R}^1$ represents $R$ on $X$ where $f \circ u : X \to \mathbb{R}^1$ is represents $f(u(x))$. Thus, scaling a utility function is of no consequence.

Thus, far we have not established any results that allow us to find the maximizer of a utility function. Fortunately, if we assume that utility functions are differentiable, the tools of calculus will allow us to characterize optimal choices. The mathematical appendix reviews key results from calculus.

## 5. Spatial Preferences

In most applications in economics, outcomes spaces are denominated in money (incomes, wealths, wages, profits etc.) or commodities (widgets, gizmos, chili burritos). It is sensible to assume that larger outcomes are preferred to smaller outcomes (except perhaps in the case of chili burritos). In other words, many of the preferences considered in economics are non-satiable in that agents either believe more is always better (i.e. money) or less is always better (air pollution). In political game theory, however, many of the outcomes we wish to study are policies such as taxes, welfare benefits, abortion restrictions where at least some agents have a most preferred outcome that is neither zero

or infinite. A voter's utility may be increasing in tax rates below some level and decreasing for higher levels. A voter may prefer restrictions on abortion only so stringent as outlawing them in third trimester but not more so. Thus, it often necessary to assume that political actors have satiable preferences. Formally, we can say an agent has such preferences when $M(X, R)$ contains elements that are interior to the outcome space $X$. Similarly, preferences are satiable when the maximizer of $u : X \to \mathbb{R}$ is in the interior of $X$. Figure 2.1 illustrates the differences between satiable and non-satiable preferences.

**Insert Figure 2.1 Here**

The most common application of satiable preferences is the *spatial model* where it is assumed that policy outcomes can be represented as points lying in a subset of $\mathbb{R}^d$. In principal, one could specify very general preferences over this space, but in practice (and the remainder of this book) it is generally assumed that voters have *single-peaked* and *symmetric* preferences. We will discuss single peakedness in more detail in the chapter on social choice, but for now we will simply note that it implies that the agent's maximal set has a single element and that the utility function has a single maximizer. This most preferred policy outcome is known as the agent's *ideal point.* The assumption of symmetry requires that the agent's utility declines at the same rate regardless of direction. This implies that preferences are a decreasing function of the distance between the policy outcome and the agent's ideal point.

If we assume that the policy space is one-dimensional, single-peaked, symmetric preferences are represented by utility functions of the form $u_i(x) = h(-|x - z_i|)$ where $z_i$ is agent $i$'s ideal point and $h : \mathbb{R}^1 \to \mathbb{R}^1$ is an increasing function. The two most popular examples are the linear, $u_i(x) = -|x - z_i|$ and quadratic utility functions $u_i(x) = -(x - z_i)^2$. These functions are plotted in Figure 2.2.

**Insert Figure 2.2 Here**

In higher dimensional applications $X \subset \mathbb{R}^d$ distances are measured Euclidean norm defined as

$$\|x - y\| = \sqrt{\sum_{j=1}^{n} \left(x^j - z_i^j\right)^2}.$$

Thus, symmetric, single-peaked preferences take the form of

$$u_i(x) = h(-\|x - z_i\|)$$

where again, $z_i \in \mathbb{R}^d$ is the ideal point of agent $i$, $h : \mathbb{R}^1 \to \mathbb{R}^1$ is an increasing function.

It is difficult to visualize utility functions over multidimensional spaces. However, for 2 dimensions graphical analysis is simplified by the fact that each agent's preferred sets i.e. $P(y) = \{x \in X | xRy\}$ form circular regions centered on the agent's ideal point. Similarly, the set of points for which the agent is indifferent to $y$ is a circle containing $y$ centered on the ideal point. These sets are illustrated in Figure 2.3. For any indifference curve, an agent prefers an outcome inside the circle to any that lies outside it.

### Insert Figure 2.3 Here

One of the reason that single-peaked, symmetric are so popular in applied political game theoretic models is the ease at which the predicted choices of agents can be characterized. As long as one is willing to make the appropriate assumptions, choice over a pair of outcomes can be characterized by an agent's ideal point and a "cutpoint" in $\mathbb{R}^1$ or a "cutting plane" in $\mathbb{R}^d$.

To see this, consider an agent with symmetric single peaked preferences over $\mathbb{R}^1$. Thus, agent $i$ prefers $x$ to $y$ if and only if $h(-|x - z_i|) > h(-|y - z_i|)$. Assuming that $x > y$, this condition becomes

$$z_i > c \equiv \frac{x + y}{2}$$

Conversely, $yPx$ if and only if $z_i < c$. Thus, given a set of agents and outcomes $x > y$, the model predicts that all agents with ideal points greater than the midpoint of $x$ and $y$ prefer $y$ and those with ideal points lower than the midpoint prefer $y$. Note that this prediction is completely independent of $h$.

This logic extends to $\mathbb{R}^d$ as well. Now agent $i$ prefers $x$ to $y$ if and only if $h(-\|x - z_i\|) > h(-\|y - z_i\|)$. Now we can define a *separating hyperplane* as follows. Let $C = \{c \mid \|x - c\| = \|y - c\|\}$. This hyperplane is equivalent to the cutpoint in $\mathbb{R}^1$. It divides the ideal points into those who prefer $x$ to $y$ and those who prefer $y$ to $x$. Again armed only with knowledge of ideal points and $C$, we can confidently characterize the choices of the agents.

## 6. Exercises

EXERCISE 2.1. *Prove the following: Assume $X$ is finite and $R$ is a complete and reflexive binary relation on $X$. Then $M(R, S) \neq \emptyset$ on any $S \subset X$ (except $S = \emptyset$) iff $R$ is acyclic.*

EXERCISE 2.2 (*). *Prove that $R$ is strict convex if and only if its lower contour sets $P^{-1}(x)$ are convex.*

EXERCISE 2.3 (*). *Show that if $X \subset \mathbb{R}^n$ and the continuous utility function $u : X \to \mathbb{R}^1$ represents the binary ordering $R$ on $X$ then $R$ is complete, reflexive, transitive, and continuous.*

EXERCISE 2.4 (*). *Use theorem 3 to prove the Weierstrass theorem.*

EXERCISE 2.5 (*). *Use Definition 2.12 to show that $[0, 1]$ is compact.*

EXERCISE 2.6. *For the following utility functions, describe the preferred set. You may do this either graphically or by formally characterizing $P(x) = \{y : yRx\}$ for all $x \in X$. Plot the utility curve if possible.*

(1) $u(x) = -|1 - x|$ for $x \in [0, 1]$
(2) $u(x) = -x^2$ for $x \in [0, 1]$
(3) $u(x) = \sqrt{x}$ for $x \in [0, 1]$
(4) $u(x) = -\alpha x_1^2 - (1 - \alpha)x_2^2$ for $x \in \mathbb{R}^2$

EXERCISE 2.7. *Let $x = (x_1, x_2)$ and $y = (y_1, y_2)$ be two outcomes from $\mathbb{R}^2$.*

(1) Assuming that all agents have single-peaked and symmetric preferences, compute the separating hyperplane $H$ as a function of $x_1, x_2, y_1$, and $y_2$. Verify that it is a straight line.
(2) Assume that each agent has non-symmetric preferences given by $-\alpha \left(x^1 - z_i^1\right)^2 - (1 - \alpha)\left(x^2 - z_i^2\right)^2$ for $x \in \mathbb{R}^2$. What does $C = \{c \mid \alpha \left(x^1 - c^1\right)^2 + (1 - \alpha)\left(x^2 - c^2\right)^2 = \alpha \left(y^1 - c^1\right)^2 + (1 - \alpha)\left(y^2 - c^2\right)^2\}$ look like now? Does it divide the ideal points of agents who prefer $x$ to $y$ from those who prefer $y$ to $x$?

EXERCISE 2.8. *Assume that an agent has spatial preferences $R$ over $\mathbb{R}^d$ represented by the utility function $u_i(x) = h(-\|x - z_i\|)$. Prove that for any $X \subset \mathbb{R}^d$, $M(R, X)$ is non-empty if $z_i$ is finite. Show that as long as $X$ is convex, $M(R, X)$ has a unique element.*

CHAPTER 3

# Choice Under Uncertainty

In this chapter we drop the assumption that individuals can perfectly predict the consequences of their actions. Instead we assume that outcomes arise probabilistically from the choice of actions i.e. that certain actions increase or decrease the likelihood of particular outcomes. Further, individuals are assumed to know which actions are most likely to produce which sorts of outcomes. Recall the example from the last chapter where $A = \{$send in the troops, try negotiating, do nothing$\}$ and $X = \{$win major concessions, win minor concessions, status quo$\}$. The agent may believe that major concessions are more likely when the troops are deployed than when negotiation is initiated. Thus, she would have to trade off this likelihood of generating a better outcome against her costs of taking each action. Deploying the troops would be rational if it is much more likely to lead to major concessions, the additional concessions are valuable to the agent, or if the costs of deployment are low. These are the basic trade-offs underlying the classical theory of choice under uncertainty.

There are two key elements of this approach. The first is the concept of *beliefs* which are modeled as probability distributions or "lotteries" over the outcomes associated with each action. The second is the specification of payoffs associated with each outcome. These payoffs are known as von Neumann-Morgenstern utility functions in honor of two of the pioneers of classical decision theory. As we shall see, the von Neumann-Morgenstern functions rely on a much stronger concept of utility than the ordinal functions discussed in chapter 2.

## 1. The Finite Case

Our presentation begins with the case of finite numbers of actions and outcomes. As in the previous chapter, we denote the feasible actions and outcomes as sets $A = \{a_1, ..., a_I\}$ and $X = \{x_1, ..., x_J\}$. However, we now assume that actions and outcomes are linked probabilistically. To formalize this assumption, we assume that the outcome depends both on the action taken and the "state of the world", $s$. From the point of view of the agent, $s$ is a random variable like rainfall on

election day or missile precision in a war. In decision-theoretic models (or in game theoretic models in which agents choose actions randomly), it can also represent the actions of other agents. We denote the set of states as $S = \{s_1, ..., s_K\}$. We assume that agents have beliefs about the likelihood of each of state represented by the probability function $\pi(s_k) \equiv \pi_k$. These probabilities have to satisfy the basic axioms of probability theory – they have to be between zero and one, inclusive, and they must sum to one. Formally, we require

$$0 \leq \pi_k \leq 1$$

$$\pi_1 + \pi_2 + ... + \pi_k = \sum_{k=1}^{K} \pi_k = 1$$

Given this setup, we can formalize the linkage between actions, states, and outcomes with an outcome function defined as $\chi(a, s) : A \times S \rightarrow X$.[1] As an example, consider Table 1 which specifies an outcome for each combination of states and actions:

| Table 3.1 | | | |
|---|---|---|---|
| $A\backslash S$ | $\mathbf{s}_1$ | $\mathbf{s}_2$ | $\mathbf{s}_3$ |
| $\mathbf{a}_1$ | $x_1$ | $x_1$ | $x_2$ |
| $\mathbf{a}_2$ | $x_1$ | $x_2$ | $x_3$ |

In this example, outcome $x_1$ occurs in state $s_1$ regardless of action. In states $s_2$ and $s_3$, the outcome depends on the action chosen by the agent. From the agent's perspective, however, it is not the state that matters so much as the likelihood of getting particular outcomes following each action. Since the agent does not know the state when she chooses $a_i$, the probability of receiving outcome $x$ is the probability that the state takes on a value $s$ such that $\chi(a_i, s) = x$. We let $p_{ij}$ be the probability that outcome $x_j$ occurs following action $a_i$.

Note that we can easily compute these probabilities from Table 1: $p_{11} = \pi_1 + \pi_2$, $p_{12} = \pi_3$, $p_{13} = 0$, $p_{21} = \pi_1$, $p_{22} = \pi_2$, and $p_{23} = \pi_3$. Thus, the general formula for these probabilities is

$$p_{ij} = \sum_{\{k:\chi(a_i,s_k)=x_j\}} \pi_k.$$

---

[1]Because we only focus on outcomes that can occur for some combination of $a$ and $s$, we require that the total number of action-state combinations be no less than the number of outcomes or $I \cdot K \geq J$.

The probabilities $p_{ij}$ inherit the following properties from the $\pi_k$:

$$0 \leq p_{ij} \leq 1$$

$$p_{i1} + p_{i2} + ... + p_{iJ} = \sum_{j=1}^{J} p_{ij} = 1 \ for \ each \ i.$$

In the remainder of this chapter, we simplify the notation by suppressing the dependence of the outcome probabilities on $s$ to focus solely on $p_{ij}$. However, in later chapters, we will use the action/state representation of the agent's problem more explicitly.

Also to keep notational clutter to a minimum, we define the *vector* $\mathbf{p}_i = (p_{i1}, \ldots, p_{iJ})$ and refer it as the *lottery* over the outcomes associated with action $a_i$. Because of the correspondence between the action and the lottery it generates, we will refer interchangeably to an agent choosing an action $a_i$ or simply choosing the lottery $\mathbf{p}_i$. Finally, let $\mathbf{P}$ be the set of all lotteries. Given that there are $J$ possible outcomes, the set $\mathbf{P}$ consists of the set of vectors of length $J$ that satisfy the above conditions (each element is between 0 and 1, and all coordinates sum to 1). This set is sometimes denoted $\Delta^J$ and termed the $J$ dimensional simplex.[2] For two dimensions, the simplex is simply the straight line from coordinate $(0, 1)$ to $(1, 0)$ as in Figure 3.1. For three dimensions, it is the triangular segment of the plane through $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$ as in the lower panel of Figure 3.11.

**Insert Figure 3.1 Here**

Another easy way to visualize lotteries is to model them as *trees* as in Figure 3.2. Beginning from the initial node, each branch corresponds to a particular outcome and is labeled with the probability of that outcome. Thus, we can see that lottery $\mathbf{p}$ generates a larger probability of $x_1$ and a lower probability of $x_3$ than the lottery $\mathbf{q}$, while both lotteries generate the same probability of $x_2$. To build some intuition for what follows, consider how an agent might choose between an action that generated $\mathbf{p}$ and one that generated $\mathbf{q}$. First, it would seem unreasonable for the agent to base her decision on her preferences for $x_2$ since both lotteries generate this outcome with identical probabilities. Since the difference between the lotteries is the relative likelihood of $x_1$ and $x_3$, it would also seem that a rational agent would choose $\mathbf{p}$ only if $x_1 R x_3$. Thus, using these two intuitive arguments (which we will formalize shortly), we can speculate that the agent would choose $\mathbf{p}$ if $x_1 P x_3$, $\mathbf{q}$ if $x_3 P x_1$, and either lottery if $x_1 I x_3$.

---

[2] We refer those readers who are unfamiliar with vectors and coordinate systems to the mathematical appendix.

**Insert Figure 3.2 Here**

One feature of the preceding example that facilitated our intuitive prediction is that the lotteries involved are *simple* in the sense that each outcome is associated with a single probability number. However, it is often necessary to consider more complicated situations where agents must choose between lotteries over lotteries over lotteries.... Such situations are known as *compound* lotteries. We can formalize this notion by defining a compound lottery over $\mathbf{P}$ by $\{\alpha_1, ..., \alpha_I\}$ where $\alpha_i$ represents the probability of playing lottery $\mathbf{p}_i$. As an example, consider how an agent might evaluate a lottery in which the agent gets to play $\mathbf{p}$ with probability $\frac{1}{4}$ and $\mathbf{q}$ with probability $\frac{3}{4}$. We'll call this lottery $\mathbf{r} = \frac{1}{4}\mathbf{p} + \frac{3}{4}\mathbf{q}$. Its tree representation is given in Figure 3.3. How should an agent choose between $\mathbf{p}, \mathbf{q}$, and $\mathbf{r}$? First, note that the availability of $\mathbf{r}$ should not change the preference ranking between $\mathbf{p}$ and $\mathbf{q}$ so we need only consider comparisons of $\mathbf{r}$ versus $\mathbf{p}$ and $\mathbf{r}$ versus $\mathbf{q}$. However, these comparisons would seem to be difficult because of $\mathbf{r}$'s compound structure. Fortunately, preferences over $\mathbf{r}$ are easy to analyze. We assume that our agents only care about the probabilities associated with each outcome, not the paths travelled to reach those outcomes. Thus, the agent can compute that the probability of receiving outcome $x_1$ is the probability of receiving lottery $\mathbf{p}$ $\left(\frac{1}{4}\right)$ times $\frac{1}{3}$ plus the probability of receiving $\mathbf{q}$ $\left(\frac{3}{4}\right)$ times $\frac{1}{4}$. The probabilities of $x_2$ and $x_3$ can be computed similarly. Since the agent can compute a single probability number for each outcome, $\mathbf{r}$ can be represented as a simple lottery as in the second panel of Figure 3.3. Formally, any compound lottery $\{\alpha_1, ..., \alpha_I\}$ over $\mathbf{P}$ can be represented as a simple lottery with the probability of $x_j$ given by $\sum_{i=1}^{I} \alpha_i p_{ij}$.

**Insert Figure 3.3 Here**

Given the reduction of $\mathbf{r}$ to a simple lottery, can we know say which lottery the agent will prefer? Note that in the reduction of $\mathbf{r}$, the probability of $x_2$ is again $\frac{1}{2}$ as it is in $\mathbf{p}$ and $\mathbf{q}$. So preferences over $x_2$ are again irrelevant and only the comparison of $x_1$ to $x_3$ matters. Since under $\mathbf{r}$ the outcome $x_1$ is more likely than $x_3$, any agent for whom $x_1 P x_3$ will prefer $\mathbf{r}$ over $\mathbf{q}$. However, it also seems intuitive that such an agent would also prefer $\mathbf{p}$ to $\mathbf{r}$ since $x_1$ is somewhat more likely under $\mathbf{p}$.

This discussion of how an agent with preferences over the outcomes $X$ might evaluate different lotteries demonstrates the key features of the theory of choice under uncertainty is standard in classical decision-theory and non-cooperative game theory. Our goal is to formalize

a notion of weak preferences on $\mathbf{P}$ that allows us to deal with more complicated problems involving choice under uncertainty. Just as utility functions greatly simply the analysis of choice under certainty, the concept of *expected utility* will simplify the analysis of choice under uncertainty. Now we are in a position to formalize some of the intuition in the previous example. All of the intuition can be succinctly summarized into four axioms about weak preferences $R$ on $\mathbf{P}$.

AXIOM 3.1. **Completeness and Transitivity:** *The weak preference relation $R$ over $\mathbf{P}$ is complete and transitive.*

AXIOM 3.2. **Reduction of Compound Lotteries:** *For any $\alpha \in [0,1]$ and $\mathbf{p} \in \mathbf{P}$, $\mathbf{p}I\left[\alpha\mathbf{p} + (1-\alpha)\mathbf{p}\right]$.*

AXIOM 3.3. **Continuity:** *Let $\mathbf{p}$, $\mathbf{q}$ and $\mathbf{r}$ be three lotteries in $\mathbf{P}$. The set of scalers $\alpha \in [0,1]$ such that $\left[\alpha\mathbf{p} + (1-\alpha)\mathbf{r}\right]R\mathbf{q}$ is a closed interval and the set of scalers $\beta \in [0,1]$ such that $\mathbf{q}R\left[\beta\mathbf{p} + (1-\beta)\mathbf{r}\right]$ is also a closed interval.*[3]

AXIOM 3.4. **Independence:** *Let $\mathbf{p}$, $\mathbf{q}$ and $\mathbf{r}$ be three lotteries in $\mathbf{P}$. For any scalar $\alpha \in (0,1)$, $\mathbf{p}R\mathbf{q}$ if and only if $\left[\alpha\mathbf{p} + (1-\alpha)\mathbf{r}\right]R\left[\alpha\mathbf{q} + (1-\alpha)\mathbf{r}\right]$.*

The substantive meaning of each of the axioms is pretty straightforward. First, we have to assume, just as in the case of outcomes, that any two lotteries can be compared and that preferences over lotteries do not cycle. This axiom is critical in our example as the predictions assume that the agent has well-behaved preferences over $x_1$, $x_2$, and $x_3$. As in the last chapter, transitivity can be extended to indifference and strict preference. The second axiom simply formalizes Figure 3.3 and guarantees that agents care only about the probabilities of the outcomes, and not the particular manner in which the probabilistic process of reaching those outcomes is represented.[4]

The 3.3 axiom is somewhat abstract but implies that small changes in the probabilities of outcomes should not lead to large changes in the preferences over lotteries. It requires that if $\mathbf{p}P\mathbf{q}$ then all lotteries sufficiently close to $\mathbf{p}$ should also be preferred to $\mathbf{q}$. If $\mathbf{p}P\mathbf{q}$, a modification of $\mathbf{p}$ that adds a very small probability of a really bad outcome will not reverse the preference ordering. The continuity axiom has the following straightforward useful implication.

---

[3]The mathematical appendix has a detailed discussion of closed sets. For the present purposes, however, it is sufficient to know that a closed interval $[a,b]$ is one that includes points $a$ and $b$ and all points in between.

[4]In some texts this assumption is implicit when authors define prefrences over lotteries. We choose to make the assumption explicit to highlight that this theory ignores details regarding the representation of lotteries.

LEMMA 3.1. *If* $\mathbf{p}R\mathbf{q}R\mathbf{r}$ *then there exists some* $\lambda \in [0,1]$ *such that* $[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]\,I\mathbf{q}$.

PROOF. Assume that $\mathbf{p}R\mathbf{q}R\mathbf{r}$ and that for any $\lambda \in [0,1]$, we have either $[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]\,P\mathbf{q}$ or $\mathbf{q}P\,[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]$. If $\mathbf{p}I\mathbf{q}$ or $\mathbf{r}I\mathbf{q}$, we obtain a contradiction at $\alpha = 1$ or $\alpha = 0$ so we must have $\mathbf{p}P\mathbf{q}P\mathbf{r}$. This implies that the sets $\{\alpha : [\alpha\mathbf{p} + (1 - \alpha)\,\mathbf{r}]\,R\mathbf{q}\}$ and $\{\beta : \mathbf{q}R\,[\beta\mathbf{p} + (1 - \beta)\,\mathbf{r}]\,\}$ are non-empty. By the continuity axiom, these sets are closed. This means that the first set contains a smallest element $\underline{\alpha}$ and the second set contains a largest element, $\overline{\beta}$. Since the strict preference is not reflexive, we cannot have $[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]\,P\mathbf{q}$ and $\mathbf{q}P\,[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]$ for any particular value of $\lambda$. Thus we must have $\overline{\beta} < \underline{\alpha}$. But then at $\lambda \in (\overline{\beta}, \underline{\alpha})$ we have neither $[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]\,P\mathbf{q}$ nor $\mathbf{q}P\,[\lambda\mathbf{p} + (1 - \lambda)\,\mathbf{r}]$ contradicting the original hypothesis. $\square$

Finally, consider the independence axiom which is perhaps the most controversial.[5] Suppose we have a preference ranking between two lotteries. If we mix each of those lotteries with a third (using the same probabilities), the independence axiom holds that the preference ordering will be the same as those over the original lotteries. As an example, consider two lotteries. The first pays \$100 with probability .5 and \$0 otherwise. The second pay \$75 for sure. If we compound each of these lotteries with a .5 chance of \$1,000,000 for sure and .5 chance of playing the original lottery, the independence axiom says that the preferences over the compound lotteries should correspond to the original lotteries. Using the "tree" metaphor for lotteries, the independence axiom says that the comparison of two lotteries is based only on the comparison of the outcome branches that are distinct across lotteries. Thus, this axiom formalizes the intuition behind ignoring $x_2$ in our first example. The independence axiom is easy to extend to the case of indifference and strict preference.

LEMMA 3.2. *For any scalar* $\alpha \in (0,1)$, $\mathbf{p}I\mathbf{q}$ *if and only if*

$$[\alpha\mathbf{p} + (1 - \alpha)\,\mathbf{r}]\,I\,[\alpha\mathbf{q} + (1 - \alpha)\,\mathbf{r}]\,.$$

PROOF. To show sufficiency, suppose that $\mathbf{p}I\mathbf{q}$. Then the independence axiom requires both $[\alpha\mathbf{p} + (1 - \alpha)\,\mathbf{r}]\,R\,[\alpha\mathbf{q} + (1 - \alpha)\,\mathbf{r}]$ and $[\alpha\mathbf{q} + (1 - \alpha)\,\mathbf{r}]\,R\,[\alpha\mathbf{p} + (1 - \alpha)\,\mathbf{r}]$. Thus, $[\alpha\mathbf{p} + (1 - \alpha)\,\mathbf{r}]\,I\,[\alpha\mathbf{q} + (1 - \alpha)\,\mathbf{r}]$ is the only possibility. The proof of necessity is very similar. $\square$

LEMMA 3.3. *For any scalar* $\alpha \in (0,1)$, $\mathbf{p}P\mathbf{q}$ *if and only if*

---

[5]In some texts and articles the independence axiom is called the substitution axiom.

$[\alpha \mathbf{p} + (1 - \alpha) \mathbf{r}] \, P \, [\alpha \mathbf{q} + (1 - \alpha) \mathbf{r}].$

PROOF. To show sufficiency, suppose that $\mathbf{p} P \mathbf{q}$. Then the independence axiom requires that $[\alpha \mathbf{p} + (1 - \alpha) \mathbf{r}] \, R \, [\alpha \mathbf{q} + (1 - \alpha) \mathbf{r}]$. To show indifference is inconsistent, assume that $[\alpha \mathbf{q} + (1 - \alpha) \mathbf{r}] \, I \, [\alpha \mathbf{p} + (1 - \alpha) \mathbf{r}]$. By the previous lemma this implies that $\mathbf{q} I \mathbf{p}$, contradicting the assumption that $\mathbf{p} P \mathbf{q}$. The proof of necessity is very similar. $\qquad \square$

An equally important but less direct implication of the independence axiom is the following lemma.

LEMMA 3.4. *If* $\mathbf{p} R \mathbf{q}$ *and* $\alpha \in (0, 1)$, *then* $\mathbf{p} R \, [\boldsymbol{\alpha} \mathbf{p} + (1 - \alpha) \mathbf{q}] \, R \mathbf{q}$.

PROOF. Since $\mathbf{p} P \mathbf{q}$, we can use the reduction of compound lotteries, lemma 3, and transitivity to show that
$\mathbf{p} I \, [\boldsymbol{\alpha} \mathbf{p} + (1 - \alpha) \mathbf{p}] \, R \, [\boldsymbol{\alpha} \mathbf{p} + (1 - \alpha) \mathbf{q}] \, R \, [\boldsymbol{\alpha} \mathbf{q} + (1 - \alpha) \mathbf{q}] \, I \mathbf{q}. \qquad \square$

This lemma just establishes that if we take a weighted average of two lotteries, the resulting lottery will have an intermediate preference ranking. Finally, the independence and continuity axioms have the following implication which is crucial for the existence of expected utility functions.

LEMMA 3.5. *Suppose the alternatives are indexed so that* $x_1 R x_j$ *for all* $j$ *and* $x_j R x_J$ *for all* $j$. *Then for all* $\alpha, \beta \in [0, 1]$,
$[\alpha x_1 + (1 - \alpha) x_J] \, R \, [\beta x_1 + (1 - \beta) x_J]$ *if and only if* $\alpha \geq \beta$.

PROOF. Suppose $\alpha \geq \beta$. We can then write $\alpha x_1 + (1 - \alpha) x_J$ as $\gamma x_1 + (1 - \gamma) [\beta x_1 + (1 - \beta) x_J]$ where $\gamma = \frac{\alpha - \beta}{1 - \beta} \in (0, 1]$. From lemma 4, $x_1 R [\beta x_1 + (1 - \beta) x_J]$. Applying lemma 4 again,
$[\gamma x_1 + (1 - \gamma) [\beta x_1 + (1 - \beta) x_J]] \, R \, [\beta x_1 + (1 - \beta) x_J]$.
The proof of necessity is identical to above where the roles of $\alpha$ and $\beta$ are reversed. $\qquad \square$

Sensibly when we compare lotteries over the best and worst outcome, we prefer the one with the greatest likelihood of producing the best outcome. With these axioms and lemma, we can now prove that preferences over lotteries can be represented by *expected utility functions*. We state the theorem in terms of preferences over lotteries. Recall that actions induce lotteries over outcomes and so an analogous statement can be made about preferences over actions.

THEOREM 3.1. *(von Neumann-Morgenstern) If axioms 3.1-3.4 hold, then there exists a function* $u(x_j)$ *(which assigns a number* $u_j$ *for each outcome) such that*

*i)* *the expected utility of lottery* $\mathbf{p}_i$ *(with is induced by action i ) is given by*

$$EU(\mathbf{p}_i) = p_{i1}u_1 + p_{i2}u_2 + \ldots + p_{iJ}u_J = \sum_{j=1}^{J} p_{ij}u_j$$

*ii)* $\mathbf{p}_i R \mathbf{p}_j$ *(i.e.* $a_i R a_j$*) if and only if* $EU(\mathbf{p}_i) \geq EU(\mathbf{p}_j)$.

The function $u(x_j)$ is sometimes called a Bernoulli utility function, to distinguish it from the expected utility function $EU(\mathbf{p})$. It is important to note that we are talking about two different types of utility functions. The Bernoulli, (or lower cased) functions are defined over outcomes. The expected utility (capitalized) functions are defined over lotteries. A subtle point, is that Theorem 1 starts with preferences over lotteries that satisfy the four axioms and states that we can construct an expected utility function over lotteries, that has a particular form –the expected utility of a lottery is nothing more that the expected *value* of lottery given the values of the outcomes specified by the Bernoulli utility functions.

The expected utility of a lottery is simply the average of the utilities over outcomes weighted by the probabilities of each outcome. For example, if we assigned utilities to outcomes $x_1$, $x_2$, and $x_3$ of $u(x_1)$, $u(x_2)$, and $u(x_3)$, then the expected utility of lottery $\mathbf{p}$ would be $EU(\mathbf{p}) = \frac{1}{3}u(x_1) + \frac{1}{2}u(x_2) + \frac{1}{6}u(x_3)$ while that of $\mathbf{q}$ would be $EU(\mathbf{q}) = \frac{1}{4}u(x_1) + \frac{1}{2}u(x_2) + \frac{1}{4}u(x_3)$. One of the most attractive properties of expected utility functions is that they are linear in the outcome utilities. Among other things this implies that $EU(\mathbf{r}) = EU\left(\frac{1}{4}\mathbf{p} + \frac{3}{4}\mathbf{q}\right) = \frac{1}{4}EU(\mathbf{p}) + \frac{3}{4}EU(\mathbf{q}) = \frac{13}{48}u(x_1) + \frac{1}{2}u(x_2) + \frac{11}{48}u(x_3)$ which is exactly what one gets from computing the expected utility of the reduced lottery.

We do not prove the von Neumann-Morgenstern theorem formally in this section, but we can sketch it in the case of three outcomes. The proof of the general result is very similar but uses mathematical induction to extend to an arbitrary number of alternatives. Let $X = \{x_1, x_2, x_3\}$ where $x_1 R x_2 R x_3$ and assume that at least one of the preferences in strict – $x_1 I x_2 I x_3$ is a trivial case. We will represent a lottery over $X$ as a vector $(p_1, p_2, p_3)$. From Lemma 1, we know that there exist $\alpha$ such that $x_2 I(\alpha, 0, 1 - \alpha)$. Similarly, we know that $x_1 I(1, 0, 0)$ and $x_3 I(0, 0, 1)$ Therefore, let $u_1 = 1$, $u_2 = \alpha$, and $u_3 = 0$. Now consider any lottery $\mathbf{p} = (p_1, p_2, p_3)$. From this lottery, we can form the compound lottery $(p_1 + u_2 p_2, 0, p_3 + p_2(1 - u_2))$ by substituting the lottery $(u_2, 0, 1 - u_2)$ for the degenerate lottery that

reaches $x_2$. Using Lemma 2, we know that the agent must be indifferent between this compound lottery and $\mathbf{p}$. Since we can make similar substitutions for $x_1$ and $x_2$, the agent is indifferent between $\mathbf{p}$ and $\{p_1 u_1 + p_2 u_2 + p_3 u_3, 0, p_1 (1 - u_1) + p_2 (1 - u_2) + p_3 (1 - u_3)\}$. Now consider an alternative lottery $\mathbf{q}$. By replicating the above arguments, we know that the agent is indifferent between $\mathbf{q}$ and $(q_1 u_1 + q_2 u_2 + q_3 u_3, 0, q_1 (1 - u_1) + q_2 (1 - u_2) + q_3 (1 - u_3))$.

For the grand finale, the application of Lemma 5 says that $\mathbf{p} R \mathbf{q}$ if and only if $p_1 u_1 + p_2 u_2 + p_3 u_3 \geq q_1 u_1 + q_2 u_2 + q_3 u_3$. Thus, we can represent preferences over lottery $\mathbf{p}$ by the scalar $p_1 u_1 + p_2 u_2 + p_3 u_3$. Note that the theorem does not say that any $\alpha \in [0, 1]$ will work. Rather the theorem ensures that there exists at least one such $\alpha$ so that the outcome utilities, $u_1 = 1, u_2 = \alpha$ and $u_3 = 0$ will work.


**1.1. Cardinal Utility.** In the previous chapter, we assuaged utility-skeptics with the argument "relax utility functions do nothing more than represent preference orderings." Once we move into the world of expected utilities, however, such a defense is no longer tenable. The utility functions over outcomes $u(x_j)$ are no longer simply ordinal, but *cardinal* in that they contain information about *relative* preferences over outcomes. Just as the Fahrenheit temperature scales ability to say that the difference between $212°$ and $32°$ is twice the difference between $122°$ and $32°$, cardinal utility functions allow us to say that "my preference for steak over chicken is 3.8 times my preference for chicken over fish." Thus, unlike the case of ordinal utilities, the value $u(x_j) - u(x_k)$ has a meaningful interpretation.

It is easy to see why expected utility theory depends on cardinal utility functions. Suppose that an agent were choosing between two lotteries over the three outcomes $x_1, x_2, x_3$ with $x_1 P x_2 P x_3$. Lottery 1 provides a .5 shot at $x_1$ and a .5 shot at $x_3$ while lottery 2 gives $x_2$ with certainty. Suppose we had an expected utility representation $u(x_1) = 1$, $u(x_2) = \alpha \in (0, 1)$, and $u(x_3) = 0$ which predicted that the agent would choose lottery 1. If this representation were simply ordinal, we could apply any order preserving transformation to the utility function, and the resulting function would represent the exact same preferences. But consider the following transformation of the utilities: $v(x_1) = 1$, $v(x_2) = 1 - \alpha$, and $v(x_3) = 0$. This transformation preserves the preference ordering, but now the agent would prefer lottery 2.

However, just as Fahrenheit is not the only temperature scale which produces identical relative information about heat, expected utility representations are not unique. To see this, consider two cardinal utility

functions $u(x_j)$ and $v(x_j) = a + bu(x_j)$. Given a lottery $\mathbf{p}$, these produce expected utility functions of $\sum_{j=1}^{J} p_j u(x_j)$ and $a + b\sum_{j=1}^{J} p_j u(x_j)$ respectively. Since the ordering of $\sum_{j=1}^{J} p_j u(x_j)$ and $\sum_{j=1}^{J} q_j u(x_j)$ will be the same as the ordering of $a + b\sum_{j=1}^{J} p_j u(x_j)$ and $a + b\sum_{j=1}^{J} q_j u(x_j)$ as long as $b > 0$, each cardinal utility function produces exactly the same behavior. This also implies that $u(x_j) - u(x_k)$ is unique up to a scale factor $b$ while relative differences $\frac{u(x_j) - u(x_k)}{u(x_l) - u(x_m)}$ are uniquely determined.

## 2. Risk Preferences

One aspect of choice under uncertainty that is not pinned down by the axioms of the previous section is the set of risks that a rational agent is willing to tolerate. Some agents may be willing to accept a substantial probability of a bad outcome in exchange for moderately higher probabilities of good outcomes, while others will seek to minimize the probability of bad outcomes by forgoing opportunities for high payoffs. Recall, that the continuity axiom says that given three outcomes $xPyPz$, an agent will prefer a lottery between $x$ and $z$ to a certainty of $y$ if the probability of $z$ is sufficiently small. However, the axiom is silent about how small this risk needs to be.

The way we characterize an agents's preference or toleration of risk is whether the agent is willing to accept a *fair bet*. A fair bet is one that pays its stake (or price) in expectation. We will suppose for now that the risky stakes and rewards (i.e. outcomes) are denominated in money or some other commodity which agents tend to prefer more of. In latter sections, we will extend these definitions to outcome spaces where agent's have satiable preferences. Let $w$ be the stake and $x_1 > x_2$ be distinct monetary outcomes and $p$ be the probability of $x_1$ and $1-p$ be the probability of $x_2$. Then we have the following definitions:

DEFINITION 3.1. *A bet is fair if* $w = px_1 + (1-p)x_2$. *A bet is favorable if* $w < px_1 + (1-p)x_2$. *A bet is unfair if* $w > px_1 + (1-p)x_2$

For example, a fair bet would be buying a \$1 lottery ticket that pays \$100 with probability $\frac{1}{100}$ and nothing with probability $\frac{99}{100}$ The bet would be favorable if the ticket cost less than a dollar and unfair if it cost more. Needless to say, all lotteries that are called "lotteries" in the real world (especially those run by state governments) are of the unfair variety. Using the notion of fair bets, we can characterize preferences for risk.

DEFINITION 3.2. *An agent is risk adverse if she will not accept any unfair bets, or* $u(px_1 + (1-p)x_2) > pu(x_1) + (1-p)u(x_2)$

DEFINITION 3.3. *An agent is risk acceptant if she will accept some unfair bet or* $u(px_1 + (1-p)x_2) < pu(x_1) + (1-p)u(x_2)$

DEFINITION 3.4. *An agent is risk neutral if she is indifferent between any fair bet and its stake or* $u(px_1 + (1-p)x_2) = pu(x_1) + (1-p)u(x_2)$

It turns out that an agent's preference for risk is closely related to the shape of her utility function for money. Consider Figure 3.4 which demonstrates the utility comparison for the fair bet $w = px_1 + (1-p)x_2$. Note that the line connecting the coordinates $(u(x_1), x_1)$ and $(u(x_2), x_2)$ must travel through the point

$$(pu(x_1) + (1-p)u(x_2), \ px_1 + (1-p)x_2).$$

Thus, we know that the value of $pu(x_1) + (1-p)u(x_2)$ lies at the intersection of the line between $u(x_1)$ and $u(x_2)$ and the vertical line beginning at $w$. Thus, we can see that $u(w) > pu(x_1) + (1-p)u(x_2)$ so that the agent is risk adverse and rejects the fair bet. Obviously, the property of the utility function responsible is the fact that the utility function always lies above any line connecting two utilities. This is the property of *concavity*.

**Insert Figure 3.4 Here**

In Figure 3.5, we can see that the utility function always lies below lines connecting two utility values. For a *convex* utility function such as this one, the agent will always accept the fair bet as $pu(x_1) + (1-p)u(x_2) > u(px_1 + (1-p)x_2)$. Figure 3.6 illustrates that linear utility functions produce risk-neutral behavior as they imply that the expected utility of a gamble is identical to the utility of the expected outcome.

**Insert Figures 3.5 and 3.6 Here**

**2.1. Risk Preferences and Stochastic Dominance\*.** It is not conceptually, difficult to extend these ideas to lotteries that assign positive probability to an arbitrary (finite) number of possible outcomes.

DEFINITION 3.5. *An agent (or its preference relation) exhibits risk aversion if for any non deterministic lottery* $\mathbf{p}$, $u(\sum_j p_j x_j) > \sum_j p_j u(x_j)$.

DEFINITION 3.6. *An agent (or its preference relation) exhibits risk acceptance if for any non deterministic lottery* $\mathbf{p}$, $u(\sum_j p_j x_j) < \sum_j p_j u(x_j)$

To relate the notion of risk aversion to potential behavior, we need to define a few additional concepts. Given a lottery $\mathbf{p}$ the *expected*

*value* is $E\left(\mathbf{p}\right) = \sum_j p_j x_j$ and the *variance* of the lottery is $V\left(\mathbf{p}\right) = \sum_j p_j(E\left(\mathbf{p}\right) - x_j)$ For a given lottery, $\mathbf{p}$ and expected utility representation $u = (u_1, ...., u_J)$, the *certainty equivalent*, $C\left(\mathbf{p}\right)$, is the amount of money that the agent values as much as the lottery. That is we have,

$$u(C\left(\mathbf{p}\right)) = EU(\mathbf{p}).$$

The behavior of risk adverse individuals is somewhat predictable. They are willing to sacrifice some expected value for reductions in variance so that certainty equivalent is less than the expected value. In addition if we consider two lotteries $\mathbf{p}$ and a lottery $\mathbf{q}$ with the same expected value as $\mathbf{p}$ but a greater variance, then the agent will prefer $\mathbf{p}$.

DEFINITION 3.7. *We say that $\mathbf{q}$ is a mean preserving spread of $\mathbf{p}$ if $\mathbf{q}$ is a compound lottery that first takes the realization of $\mathbf{p}$ and then adds to it a random term $\varepsilon$ with distribution $\mathbf{z}$ having $E\left(\mathbf{z}\right) = 0$ and $V\left(\mathbf{z}\right) > 0$.*

THEOREM 3.2. *Given a preference relation $R$ on lotteries and the Bernoulli utility function $u(x)$ that is used to represent $R$, the following statements are equivalent*

*1. The preference relation $R$ exhibits risk aversion*

*2. The utility function $u(x)$ is strictly concave*

*3. For any lottery $\mathbf{p}$, $C\left(\mathbf{p}\right) \leq E\left(\mathbf{p}\right)$ (and if $V\left(\mathbf{p}\right) > 0$ the inequality is strict).*

*4. For any two lotteries $\mathbf{q}$ and $\mathbf{p}$ in which $\mathbf{q}$ is a mean preserving spread of $\mathbf{p}$ we have $EU(\mathbf{q}) < EU(\mathbf{p})$*

PROOF. The equivalence between (1) and (2) is immediate. To show that (2) implies (3) assume that $R$ exhibits risk aversion. This implies that for any non-deterministic lottery $\mathbf{p}$ (so that $V\left(\mathbf{p}\right) > 0$), we have $u(E\left(\mathbf{p}\right)) > \sum_j p_j u(x_j)$. But since $u(C\left(\mathbf{p}\right)) = EU((\mathbf{p}))$ it must be the case that $u(E\left(\mathbf{p}\right)) > u(C\left(\mathbf{p}\right))$. Since $u(x)$ is in increasing function this implies that $E\left(\mathbf{p}\right) > C\left(\mathbf{p}\right)$.

To see that (3) implies (2), assume that (3) is true and consider any two outcomes $x_1$ and $x_2$ with $p_1 \in (0, 1)$. In this case (3) and the fact that $u(x)$ is increasing implies that $u(px_1 + (1 - p)x_2) > pu(x_1) + (1 - p)u(x_2)$ and thus the function $u(x)$ is concave. Since $V\left(\mathbf{p}\right) = 0$ and $\mathbf{p}$ deterministic result in equality this case is trivial.

To see that (2) implies (4) consider $\mathbf{p}$ and a mean preserving spread $\mathbf{q}$ that results in $x + \varepsilon$ with $x$ having the distribution $\mathbf{p}$ and $\varepsilon \in \{\varepsilon_1, ..., \varepsilon_T\}$ having the distribution $\mathbf{z}$. We have

$$EU(\mathbf{q}) = \sum_t \sum_j z_t p_j u(x_j + \varepsilon_t).$$

By (2), for each value of $x_j$ we have

$$\sum_t z_t p_j u(x_j + \varepsilon_t) < u(\sum_t z_t p_j (x_j + \varepsilon_t)).$$

Rearranging the right hand side yields

$$\sum_t z_t p_j u(x_j + \varepsilon_t) < u(p_j x_j + \sum_t \varepsilon_t) = u(p_j x_j).$$

Summing over $j$ yields

$$EU(\mathbf{q}) = \sum_t \sum_j z_t p_j u(x_j + \varepsilon_t) < \sum_j u(p_j x_j) = EU(\mathbf{p}).$$

To see that (4) implies (3) consider a lottery, $\mathbf{q}$. Note that the lottery $\mathbf{q}$, is a mean preserving spread of the lottery $\mathbf{p}$ which assigns probability 1 to $E\mathbf{q}$, (4) implies that $EU(\mathbf{p}) > EU(\mathbf{q})$. Since $p$ is deterministic, $u(E\mathbf{p}) = EU(\mathbf{p})$ so we have $u(E\mathbf{p}) > EU(\mathbf{q})$. Since $u(C(\mathbf{q})) = EU(\mathbf{q})$ monotonicity of $u(x)$ implies that $E(\mathbf{p}) = E(\mathbf{q}) > C(\mathbf{q})$.          □

If lottery $\mathbf{q}$ is a mean preserving spread of $\mathbf{p}$, lottery $\mathbf{q}$ is said to be second order stochastically dominated by $\mathbf{p}$. One result which can be proven by reapplying the logic of the last proof is gives a convenient characterization of second order stochastic dominance.

THEOREM 3.3. *Lottery* $\mathbf{q}$ *is second order stochastically dominated by* $\mathbf{p}$ *if and only if for any concave increasing function* $u(x)$,

$$\sum_j p_j u(x_j) \geq \sum_j q_j u(x_j).$$

Thus, risk averse individuals have preferences that conform with second order stochastic dominance –if $\mathbf{p}$ second order stochastically dominate $\mathbf{q}$, a risk averse agent will prefer $\mathbf{p}$ to $\mathbf{q}$. Now in some cases choice over lotteries is trivial. One common notion of "no brainer" decision problems involves choice between a lottery and another that is said to first-order stochastically dominate it.

DEFINITION 3.8. *Lottery* $\mathbf{q}$ *is first order stochastically dominated by* $\mathbf{p}$ *if for any non decreasing function* $u(x)$

$$\sum_j p_j u(x_j) \geq \sum_j q_j u(x_j).$$

So in choosing between lotteries which are ordered by first order stochastic dominance, risk attitudes are irrelevant.

**2.2. Risk Preferences with Satiable Preferences.** As we discussed in Chapter 2, many of the utility functions used in political science are satiable in that agents have most preferred outcomes. Assuming such preferences, however, entails implicit assumptions about risk. Consider the function in Figure 3.7 and a lottery over some $x_1$ less that the agent's ideal point and $x_2$ greater than the ideal point. As above, the expected utility of such a lottery is the intersection of the line between $u(x_1)$ and $u(x_2)$ and the vertical line beginning at $w$. However, note that since the ideal point lies between $x_1$ and $x_2$, there must be at least one outcome $w_1$ in this interval such that $u(w_1) > pu(x_1) + (1-p)u(x_2)$. Thus, satiable preferences must produce risk adverse behavior at least in regions near the ideal point. Satiable preferences need not produce global risk aversion, however. Gambles over a set of outcomes bounded away from the ideal point for which the agent's utility function is convex will produce risk acceptant behavior. The next section elaborates these points in more detail.

**Figure 3.7**

**2.3. Risk and Higher Dimension Euclidean preferences.** In this subsection we consider choice over lotteries on $\mathbb{R}^n$. In the case of preferences that are Euclidean, Bendor and Meirowitz (2004) have shown that we can extend the notions of risk aversion from the case of strictly increasing preferences. Recall, that preferences over $x \in \mathbb{R}^n$ are Euclidean if they are representable by a utility function of the form

$$u(x) = h(-\|x - x^*\|)$$

where $x^*$ is a point in $\mathbb{R}^n$ and $h : \mathbb{R}^1 \to \mathbb{R}^1$ is a strictly increasing function. If the function $h$ is strictly concave then it is not difficult to see that the utility function $u(x)$ is itself strictly concave. In this case the Bernoulli utility function will represent risk averse preferences. However, even if the function $h$ is not concave, the preferences exhibit a form of risk aversion. To extend the concept of a mean-preserving spread to $\mathbb{R}^n$, we simply apply Definition 3.7 with the relevant states in $\mathbb{R}^n$.

THEOREM 3.4. *If $u(x)$ is Euclidean then for any two lotteries $\mathbf{q}$ and $\mathbf{p}$ on $\mathbb{R}^n$ with expected value $x^*$ in which $\mathbf{q}$ is a mean preserving spread of $\mathbf{p}$ we have $EU(\mathbf{q}) < EU(\mathbf{p})$.*

PROOF. Assume $\mathbf{q}$ and $\mathbf{p}$ are lotteries on $R^n$ with expected value $x^*$ in which $\mathbf{q}$ is a mean preserving spread of $\mathbf{p}$.Define $d_j \equiv \|x_j - x^*\|$ to be the real variable that measures the distance between random variable $x$ and the point $x^*$. Now since preferences are Euclidean we

have $EU(\mathbf{q}) < EU(\mathbf{p})$ if and only if,

$$\sum_j q_j h(-\|x_j - x^*\|) < \sum_j p_j h(-\|x_j - x^*\|)$$

for some increasing function $h(\cdot)$. This condition is satisfied if and only if

$$\sum_j q_j h(-d_j) < \sum_j p_j h(-d_j).$$

Since $q$ is a mean preserving spread of $p$, it must be the case that

$$\sum_j q_j d_j > \sum_j p_j d_j.$$

But this means for any increasing function $g(\cdot)$ including $h(\cdot)$ we have

$$\sum_j q_j g(-d_j) < \sum_j p_j g(-d_j).$$

$\square$

Informally, for lotteries that are centered at the agent's ideal point second order stochastic dominance in outcomes corresponds to first order stochastic dominance in terms of disutility, and thus risk attitudes are not directly relevant to assessing how agents will compare lotteries. Alternatively put, all Euclidean agents are risk averse over lotteries that are centered at their ideal point.

## 3. Learning

Since agent's beliefs about the probabilities of various outcomes is the key to decision-making under uncertainty, it is important to analyze how a rational agent should respond to new information about the likelihood of various outcomes. Once again it is prudent to begin with an example. Consider Figure 3.8 where an agent believes that the incumbent politician is "good" with probability $\frac{3}{4}$ and "bad" with probability $\frac{1}{4}$. Suppose however that she could incorporate information about the incumbent's performance in office such as the inflation rate in the economy. How would this change her probability assessment of the incumbent's quality?

**Insert Figure 3.8 Here**

First, let's assume that the agent knows that good incumbents are more likely than bad incumbent to produce low inflation. In this example, we suppose that the agent knows that good incumbents produce low inflation with probability $\frac{2}{3}$ and that bad incumbent produces low inflation with only a $\frac{1}{5}$ probability. The first thing that our intuition

tells us is that when inflation is low the agent should increase her probability assessment that the incumbent is good higher than her original belief of $\frac{3}{4}$. Conversely, when inflation is high, the agent should lower the probability that the incumbent is good. Fortunately, we can take the analysis a step further and compute the exact probabilities that should be assigned a good incumbent after either realization of the inflation rate.

First, consider the case where inflation is low. A rational agent should know that the outcome is either the top node or the third node of the second panel of Figure 3.8. Further, she knows that and that there is a $\frac{3}{4} \cdot \frac{2}{3} = \frac{1}{2}$ of the top node and a $\frac{1}{4} \cdot \frac{1}{5} = \frac{1}{20}$ probability of reaching the third. Therefore, after observing low inflation, it is 10 times as likely that the incumbent is good than that he is bad. Let $p(l)$ be the probability of a good incumbent conditional on low inflation. Since probabilities must sum to one, $p(l) + \frac{p(l)}{10} = 1$ so that $p(l) = \frac{10}{11}$. We can use similar reasoning to show that $p(h) = \frac{4}{9}$. Note the confirmation of our intuition that a realization of a low inflation should raise the probability that the incumbent is good while high inflation lowers it.

To generalize this example, we will be a need to be a bit more precise about the underlying probability theory. Let $A$ and $B$ represent two events (such as the terminal nodes in Figure 3.8). Suppose we know that event $B$ has occurred and wish to compute the probability that event $A$ occurs. This is known as the *conditional probability of $A$ given event $B$*. We denote it as follows:

$$\Pr(A \mid B) = \frac{\Pr(A \& B)}{\Pr(B)} \text{ assuming } \Pr(B) > 0.$$

where $\Pr(A)$ is the probability of event $A$, $\Pr(B)$ is the probability of event $B$, and $\Pr(A\&B)$ is the probability that both events occur (known as the joint probability).

This formula which is often termed Bayes' law is attained by dividing both sides by $\Pr(B)$, an operation which is permissible as long as this value is non-zero. As a special case independent events have the property that $\Pr(A \& B) = Pr(A)Pr(B)$ so that

$$\Pr(A \mid B) = \frac{\Pr(A)\Pr(B)}{\Pr(B)} = \Pr(A)$$

To see this rule in action, note that the probability of low inflation and a good incumbent $\left(\frac{1}{2}\right)$ is the probability of low inflation conditional on a good incumbent $\left(\frac{2}{3}\right)$ times the probability of a good incumbent $\left(\frac{3}{4}\right)$. Given these definitions, we can state the main result.

THEOREM 3.5. *(Bayes' Law) Let $A_1$ ... $A_N$ be disjoint events (i.e. no two can occur simultaneously) such that $\sum Pr(A_n) = 1$ and $Pr(A_n) > 0$ for all $n$. Let $B$ be some other event (which may occur concurrently with any $A_n$ ).*

$$\Pr(A_n \mid B) = \frac{\Pr(B \mid A_n)\Pr(A_n)}{\sum \Pr(B \mid A_n)\Pr(A_n)}$$

Bayes' Law gives us an easy to use formula to compute how rational agents should update their probability assessments following new information. Note that we can easily apply it to the voter's problem from above. Let $A_1$ be the event that the incumbent is good and $A_2$ be the event that she is bad. Since the incumbent cannot be both good and bad, these events satisfy the requirement of disjointedness. Event $B$ is low inflation. For two events the formulas are:

$$\Pr(A_1 \mid B) = \frac{\Pr(B \mid A_1)\Pr(A_1)}{\Pr(B \mid A_1)\Pr(A_1) + \Pr(B \mid A_2)\Pr(A_2)}$$

and

$$\Pr(A_2 \mid B) = \frac{\Pr(B \mid A_2)\Pr(A_2)}{\Pr(B \mid A_1)\Pr(A_1) + \Pr(B \mid A_2)\Pr(A_2)}$$

We can obtain all of the following probabilities from Figure 3.8:

$$\Pr(A_1) = \frac{3}{4}$$

$$\Pr(A_2) = \frac{1}{4}$$

$$\Pr(B \mid A_1) = \frac{2}{3}$$

$$\Pr(B \mid A_2) = \frac{1}{5}$$

Thus, we can plug these numbers into Bayes' Law to get:

$$\Pr(A_1 \mid B) = \frac{\frac{2}{3} \cdot \frac{3}{4}}{\frac{2}{3} \cdot \frac{3}{4} + \frac{1}{5} \cdot \frac{1}{4}} = \frac{10}{11}$$

and

$$\Pr(A_2 \mid B) = \frac{\frac{1}{5} \cdot \frac{1}{4}}{\frac{2}{3} \cdot \frac{3}{4} + \frac{1}{5} \cdot \frac{1}{4}} = \frac{1}{11}.$$

Voi la!

While seemingly straightforward and logical, the application of Bayes' Law is often criticized as a poor model of learning. Not only can it be computationally challenging and exceed the typical individual's grasp of conditional probability, it can also produce counter-intuitive predictions. Consider the following scenario from the *Let's Make a Deal*

game show hosted by Monte Hall. Monte offers contestants the choice
of opening three doors. Behind one door is a luxury car while the
other doors hide prizes of little pecuniary value (goats seem to have
been a favorite). Once a door is selected but before it is opened,
Monte opens one of the remaining two doors to reveal a goat. He
then asks the contestant if he would like to switch his selection to the
remaining closed door. Should the rational contestant switch? Most
people would intuitively say there is nothing to gain from switching on
the grounds that getting the car from a subsequent switch is just as
likely as getting it on the original try. The probability of winning the
car is $\frac{1}{2}$ either way. Indeed a number of mathematicians and statisti-
cians took this position in response to the publication of this problem
in a popular newspaper column. However, this logic is incompatible
with Bayes' law.

To simplify, suppose the contestant chooses door 3. Since the doors
are *ex ante* the same, the analysis of the other cases is identical. First,
consider the probability of winning if the contestant does not switch
to the remaining door. Obviously, this is the same as the original
probability that there is a car behind door 3 or $\frac{1}{3}$. However, now
consider the probability of winning by switching. To formalize, let $A_1$,
$A_2, A_3$ correspond to the car being located behind doors 1, 2, and 3
respectively. Let $B_1$, $B_2$ corresponds to the event that Monte opens
door 1 or 2. Since we assume that $\Pr(A_1) = \Pr(A_2) = \Pr(A_3) = \frac{1}{3}$, we
simply need to compute for all of the events $\Pr(B_i|A_j)$. Since Monte
will never expose a car, $\Pr(B_1|A_1) = \Pr(B_2|A_2) = 0$. We also assume
that in the event $A_3$ Monte randomly selects which goat to expose.
Therefore, $\Pr(B_1|A_2) = \Pr(B_2|A_1) = 1$ and $\Pr(B_1|A_3) = \Pr(B_2|A_3) = \frac{1}{2}$.

Suppose Monte opens door 2, then the probability that a switching
contestant wins is equal to

$$\Pr(A_1 \mid B_2) = \frac{\Pr(B_2 \mid A_1)\Pr(A_1)}{\Pr(B_2 \mid A_1)\Pr(A_1) + \Pr(B_2 \mid A_2)\Pr(A_2) + \Pr(B_2 \mid A_3)\Pr(A_3)}$$

$$= \frac{1 \cdot \frac{1}{3}}{1 \cdot \frac{1}{3} + 0 \cdot \frac{1}{3} + \frac{1}{2} \cdot \frac{1}{3}} = \frac{2}{3}$$

Similarly, if Monte opens door 1, the probability of winning is $\Pr(A_2 \mid B_1) = \frac{2}{3}$. So a switching contestant wins with probability $\frac{2}{3}$ whereas a
sticking one only wins $\frac{1}{3}$ of the time.[6]

---

[6]The solution we present to this problem is somewhat convoluted in order to
provide an additional demonstration of Bayes' Rule. An easier proof is to note

So why does the intuition that switching doesn't pay fail us so badly? The reason is that most people do not appreciate the implication of the fact that Monte will never reveal a car. Thus, observing that he doesn't open a particular door is information that a switcher can use in his decision that a stand-patter cannot.

While the Monte Hall problem does expose a critical set of problems with Bayesian learning, such objections can be carried too far. Bayes' Law does tell us correctly that switchers will win $\frac{2}{3}$ of the time. Thus, a frequent viewer of the show can learn that one should switch even without ever doing a conditional probability calculation. So one can justify the use of Bayes' rule by appealing to the notion that agents are acting as if they had performed the calculation even if they are simply following rules that they have learned from experience.

## 4. Critiques of Expected Utility Theory

While most of the models used in this book will rely heavily on expected utility theory, it is worth pointing out that there is a large and influential body of work critical of expected utility theory. However, the application of these critical insights and alternative models to political game theory is still in its infancy.[7]

**4.1. Risk, Uncertainty, and Subjective Probability.** The economist Frank Knight argued that expected utility theory is a model of risk rather than uncertainty. He defines uncertainty as the situation where individuals lack sufficient statistical information to form estimates of the probabilities of various outcomes. In other words, in a situation of uncertainty, individuals do not know the true set of lotteries **P**. However, the statistician Leonard Savage responded that expected utility theory could be resuscitated by assuming that individuals have subjective beliefs about **P** which can be used to formulate probability distributions over outcomes.

However, Daniel Ellsberg formulated the following paradox which cast doubt over whether uncertainty was reducible to beliefs about beliefs. Suppose that there are two urns containing red and black balls. In urn 1, there are 100 red and black balls where the proportion of red balls is unknown. Urn 2, however, contains 50 red balls and 50 black balls.

---

that a switcher only loses if he picked the right door in the first place. Thus, a switcher loses $\frac{1}{3}$ of the time and wins $\frac{2}{3}$.

[7]"Behavioral" (as opposed to those based on expected utility theory) have become far more common in economics in recent years (see Camerer 200x).

Now suppose that subjects are given $100 for selecting a red ball. Most subjects choose to select from urn 2. But when offered $100 for selecting a black ball, the modal choice is again urn II. However, choosing urn 2 for both gambles violates the axioms of expected utility theory. According to expected utility theory, choosing urn II in search of a red ball indicates a belief that urn I has fewer than 50 red balls while selecting urn 2 for a black ball suggests that the subject believes that urn 1 has fewer than 50 black balls. Obviously, these beliefs are inconsistent with the knowledge that urn 1 contains 100 balls. Selecting urn 1 in both gambles similarly violates expected utility theory.

**4.2. The Allais Paradox.** The predictions of expected utility theory have been tested in a number of experimental settings. These studies have provided robust evidence for a number of decision-making anomalies inconsistent with expected utility theory. One of the earliest and most studied anomalies was first uncovered by the French economist Maurice Allais. This anomaly is based on the finding that subjects often make choices inconsistent with the independence axiom.

Initially, subjects are asked to choose between lotteries **a** and **b** where:

Lottery **a**: .33 chance of $2500, .66 chance of $2400, and .01 chance of 0

Lottery **b**: $2400 for sure

When given these choices, subjects overwhelmingly choose lottery **b**. For example, Kahneman and Tversky (1979) find that 82% choose lottery **b** when given this hypothetical choice.

Next the subjects are given the choice between lotteries **c** and **d**.

Lottery **c**: .33 chance of $2500, .67 chance of 0

Lottery **d**: .34 chance of $2400 and .66 chance of 0

Experimental subjects generally choose **c**. Kahneman and Tversky find that 83% choose this lottery. However, it can easily be shown that choosing **b** in the first experiment and **c** in the second violates the independence axiom and therefore expected utility theory. First, note that the choices of **b** and **c** imply that

$$u(2400) > .33u(2500) + .66u(2400) + .01u(0)$$

and

$$.33u(2500) + .67u(0) > .34u(2400) + .66u(0)$$

Rearranging the top inequality, we get $.34u(2400) > .33u(2500) + .01u(0)$ for the first inequality and for the second we get $.33u(2500) + .01u(0) > .34u(2400)$. Thus, we derive a contradiction. To see that the contradiction is attributable to a violation of the independence

axiom, note that lottery **a** can be written as the compound lottery $.34(33/34, 0, 1/34) + .66(0, 1, 0)$ over the outcomes $(2500, 2400, 0)$ while **b** is $.34(0, 1, 0) + .66(0, 1, 0)$. Thus, if **a**$P$**b**, then the independence axioms holds that $(33/34, 0, 1/34)\ P\ (0, 1, 0)$. But this in turn implies that $.34(33/34, 0, 1/34) + .66(0, 0, 1)\ P\ .34(0, 1, 0) + .66(0, 0, 1)$ which means that **c**$P$**d**.

**4.3. Prospect Theory.** In their classic article, Kahneman and Tversky (1979) propose an alternative model of decision-making to account for the Allais paradox and other experimental anomalies. Whereas many previous authors attributed the Allais paradox to a preference for certainty, Kahneman and Tversky note that the independence axiom is often violated when all of the lotteries are far from sure things. Consider the following pairs of lotteries.

Lottery **a**:  .45 chance of $6000, .55 chance of 0
Lottery **b**:  .90 chance of $3000, .10 chance of 0
Lottery **c**: .001 chance of $6000, .999 chance of 0
Lottery **d**: .002 chance of $3000 and .998 chance of 0

They find that the modal choices were **b** over **a** and **c** over **d**, choice which violate the independence axiom. Since lotteries **c** and **d** have miniscule probabilities, subjects seem inclined to go for the one with the bigger prize. However, when both probabilities are reasonably high, subjects are still inclined to take the one that is relatively more certain.

Kahneman and Tversky note, however, that this preference for certainty does not hold when gambles are over losses rather than gains. Consider the following pairs of lotteries:

Lottery **a**:  .80 chance of −$4000, .20 chance of 0
Lottery **b**:   −$3000 for sure
Lottery **c**: .20 chance of −$4000, .80 chance of 0
Lottery **d**: .25 chance of −$3000,.75 chance of 0

If the Allais paradox were simply due to preferences for certainty, the modal choices would again be for **b** and **c**. However, Kahneman and Tversky find that **a** and **d** are the modal choices. Their interpretation is that while individuals are risk adverse over gains, they are risk acceptant over losses.

Finally, Kahneman and Tversky argue that the presentation of the lotteries can affect the choices that individuals make. Suppose that an individual has been given $1000 and then offered:

Lottery **a**:  .5 chance of an additional $1000
Lottery **b**:   $500 for sure

Next consider an individual who has been given $2000 and offered the choice of:

Lottery **c**:   .5 chance of losing $1000

Lottery **d**:   Lose $500 for sure

Kahneman and Tversky find that **b** and **c** are the modal choices even though expected utility theory holds that **a** and **c** are identical lotteries as are **b** and **d**.

To account for these anomalies, Kahneman and Tversky propose prospect theory as an alternative to expected utility theory. According to their model, choice involves two distinct phases: editing and evaluation. In the editing phase, individuals "organize and reformulate the options so as to simplify subsequent evaluation and choice."

4.3.1. *The Editing Phase.* Kahneman and Tversky identify six distinct operations that occur during the editing phase.

(1) Coding: Since Kahneman and Tversky argue that individuals evaluate gains and losses separately. Thus, the first stage of editing involves determining a reference point and coding outcomes as either gains or losses.

(2) Combination: Individuals combine probabilities associated with identical outcomes.

(3) Segregation: Individuals identify and segregate the riskless components of a choice. For example, a lottery than produces $200 with probability .7 and $100 with .3 is interpreted as a riskless $100 gain and a lottery over an addition $100.

(4) Cancellation: When comparing two lotteries, individuals ignore the common elements of both lotteries. For example, the $2000 bonus in the last example "cancels out" and does not effect the choice between **c** and **d**.

(5) Simplification: Individuals may simplify the tasks by rounding probabilities such as recoding .49 to even odds or by dropping extremely unlikely outcomes from consideration.

(6) Detection of dominance: Individuals drop from consideration any lottery that is first-order stochastically dominated.

4.3.2. *The Evaluation Phase.* Kahneman and Tversky's model of evaluation is very similar in form to expected utility theory in that both models postulate that individuals evaluate gambles using a weighted average of the payoffs to the outcomes. However, in Kahneman and Tversky's model the weights used are not the subjective probabilities of the outcomes but rather functions of the probabilities. They also argue, contra expected utility theory, that the outcome value functions should treat gains and losses asymmetrically.

Let $x$ and $y$ be two distinct monetary outcomes where $p$ is the probability of $x$, $q$ is the probability of $y$. With probability $1 - p - q$, nothing happens or the payoff is $0$. Kahneman and Tversky define prospects as strictly positive if $x, y > 0$ and $p + q = 1$, strictly negative if $x, y < 0$ and $p + q = 1$, and regular in all other cases. For a regular prospect, individuals are assumed to maximize

$$V(x, p; y, q) = \pi(p) v(x) + \pi(q) v(y)$$

where $v(x)$ and $v(y)$ are the values of each outcome and $\pi(p)$ and $\pi(q)$ are weights based on the outcome probabilities. They assume that $v(0) = 0, \pi(0) = 0$, and $\pi(1) = 1$. Note this function would be an expected utility function if $v$ were a Bernoulli function and $\pi(p) = p$ for all $p$.

For strictly positive or strictly negative prospects such as $x > y > 0$ and $x < y < 0$ where $p + q = 1$, individuals maximize

$$V(x, p; y, q) = v(y) + \pi(p) [v(x) - v(y)].$$

This functional form captures the idea that individuals evaluate such lotteries as a risk-free component $v(y)$ plus a risky component $v(x) - v(y)$.

A key assumption of Prospect Theory is that $v(\cdot)$ is asymmetry with respect to gains and losses. Kahneman and Tversky make three specific assumptions.

(1) The value function is defined in terms of deviations from a reference point (no gains or losses).
(2) The value function is concave for gains and convex for losses.
(3) The value function is steeper for losses than for gains.

Figure 3.9 illustrates a function satisfying these properties.

### Insert Figure 3.9 Here

Additional, Kahneman and Tversky make several assumptions about the form of the decision weights $\pi(p)$.

(1) $\pi$ is an increasing function of $p$.
(2) $\pi(0) = 0$.
(3) $\pi(1) = 0$.
(4) For small values of $p$, $\pi(p) > p$.
(5) For small values of $p$, $\pi$ is subadditive i.e. $\pi(rp) > r\pi(p)$ for $0 < r < 1$.
(6) For all $p$, $\pi$ satisfies the property of subcertainty i.e. $\pi(p) + \pi(1 - p) < 1$.
(7) For all $0 < p, q, r < 1$, $\pi$ is subproportional i.e. $\frac{\pi(pq)}{\pi(p)} \leq \frac{\pi(pqr)}{\pi(pr)}$.

The first three assumptions are straight-forward. The fourth is simply the idea that individuals over-weight small probabilities. Subadditivity, which helps resolve the Allais paradox, implies in conjunction with Assumption 2 that $p$ is convex for small values of $p$ (see exercises 1 and 2). Subcertainty also helps resolve the Allais paradox. Recall that the modal choices require that $v(2400) > \pi(.33)v(2500) + \pi(.66)v(2400)$ and $\pi(.33)v(2500) > \pi(.34)v(2400)$. These two inequalities require that $1 > \pi(.66) + \pi(.34)$. Finally, subproportionality accounts for many of the violations of the independence axiom since it implies that for a fixed ratio of probabilities, the ratio decision weights will be closer to unity when the probabilities are high.

A function satisfying these assumptions is plotted in Figure 3.10.

### Insert Figure 3.10 Here

### 5. Time Preferences

In many of the dynamic models considered in this book, we need to characterize how individuals evaluate payoffs they receive now as opposed to those the receive in the future. We typically assume that individuals weight current utility more than future utility (if for no other reason, we could die tomorrow). We use the idea of the discount factor to capture this intuition. Let $0 < \delta < 1$ be relative weight that players put on utilities one period in the future. Utilities two periods in the future are weighted by $\delta^2$ and so on such that utilities $t$ periods in the future are discounted by $\delta^t$.

**5.1. Computing Payoff Streams.** Often we model games in which there is no determinate end date as infinite games where the number of periods goes to $\infty$. Clearly, in an infinite game, we can no longer simply add up the payoffs from each period in order to determine the utility from a sequence of actions. Fortunately, geometric discounting helps to facilitate these calculations.

We consider the easiest case first. Assume that an agent gets a payoff of $u_t = u$ over an infinite number of periods. There are two ways to calculate the *value function* $v^\infty$ of this stream of utilities.

Method 1: Note that we can write $v^\infty$ as $u + \delta u + \delta^2 u + ... = u \sum_{t=0}^{\infty} \delta^t$. Since $0 < \delta < 1$, $\sum_{t=0}^{\infty} \delta^t$ is a convergent power series. It is a well known result that $\sum_{t=0}^{\infty} \delta^t$ converges to $\frac{1}{1-\delta}$ so that $v^\infty = \frac{u}{1-\delta}$. We can easily derive the following facts about this particular power series:

(1) $\sum_{t=0}^{T} \delta^t = \frac{1-\delta^{T+1}}{1-\delta}$

(2) $\sum_{t=T}^{\infty} \delta^t = \frac{\delta^T}{1-\delta}$

(3) $\sum_{t=T}^{S} \delta^t = \frac{\delta^T - \delta^{S+1}}{1-\delta}$

Therefore, we can compute finite streams of utility easily as well. For example, the value of receiving $u$ for $T$ periods is $u \sum_{t=0}^{T} \delta^t = \frac{u(1-\delta^{T+1})}{1-\delta}$.

Method 2: Another way to derive $v$ is using Bellman's principle of optimality. Since $v$ is an infinite stream of utilities, we should be able to write it as a one period utility $u$ plus the discounted value of an infinite stream of utility beginning one period hence. Therefore,

$$v^{\infty} = u + \delta v^{\infty}$$

so that $v^{\infty} = \frac{u}{1-\delta}$. We can compute finite streams using this method as well. Again assume that the agent receives payoff $u$ for $T$ periods and we wishing to compute $v^T$. We know that $v^T = v^{\infty} - \delta^{T+1} v^{\infty}$ so that $v^T = \frac{u}{1-\delta} - \delta^{T+1} \frac{u}{1-\delta} = \frac{u(1-\delta^{T+1})}{1-\delta}$.

While the advantages of this method are small in the simple example of a constant stream of utility, they can be substantial in more complex settings. Now assume that there are $n$ states of the world $(s_1, ...., s_n)$. In each state, the agent receives $u_n$. We assume that the state evolves according to a *Markov process* such that $\Pr(S_t = s_i | S_{t-1} = s_j) = \pi_{ij}$.

Now suppose we want to compute the value $v_j$ of the stream of utilities beginning from state $j$. Using Bellman's principle, it is easy to see that

$$v_j = u_j + \delta \sum_{i=1}^{n} \pi_{ij} v_i$$

This creates a linear system of $n$ equations and $n$ unknowns (the $v_i$'s). Sometimes it will be easier to solve such a system by replacing one of the equations with the requirement that $\sum_{i=1}^{n} \pi_{ij} = 1$ for all $j$.

Let's consider an easy example. Suppose we wanted to compute the long term payoff to a political party who values holding a particular office $u_1$ per period and gets a payoff of $u_2$ in periods in which it does not hold the office. Suppose that there is an incumbent party effect so that if it holds office it wins with probability $p > \frac{1}{2}$ and remains in office (state 1). However, this also implies that when it is out of office (state 2), it will remain out of office in the next period with probability $p$. With probability $1 - p$, it transitions states either from office to out of office or vice versa. To compute, the party's payoffs from being

in states 1 or 2, we can set up the relevant Bellman equations. Note that $n = 2$, $\pi_{11} = \pi_{22} = p$, $\pi_{12} = \pi_{21} = 1 - p$, and $u_1 > u_2$ Thus, the Bellman's equations are

$$v_1 = u_1 + \delta(pv_1 + (1-p)v_2)$$
$$v_2 = u_2 + \delta((1-p)v_1 + pv_2)$$

It is straightforward to derive

$$v_1 = \frac{(1-\delta p)\, u_1 + \delta(1-p)u_2}{1 - 2\delta p + \delta^2(2p-1)}$$
$$v_2 = \frac{\delta(1-p)u_1 + (1-\delta p)\, u_2}{1 - 2\delta p + \delta^2(2p-1)}$$

In these examples, we have taken the utility streams as exogenous (either fixed or based on a fixed set of probabilities. This can be relax significantly. Suppose that the agent chooses $x_t \in X(s_t)$ in every time period to maximize the discounted value of the stream $u(x_t, s_t)$ where $s_t \in (s_1, ...., s_n)$ is the state of the world in time $t$. We may also allow the probability distribution of transitions from $s_t$ to depend on $x_t$ so let $\pi(s_{t+1}|x_t, s_t)$ be the probability of observing some state $s_{t+1}$ following state $s_t$ and choice $x_t$. We will only consider stationary plans i.e. those in which the prescription depends only on the state. Let $x(s)$ be a stationary plan specifying the action taken when the state is $s$.

We can then characterize the payoffs to implementing plan $x(s)$ in state $s$ as

$$v(x(s), s) = u(x(s), s) + \delta \sum_{s'} v(x(s'), s)\pi(s'|x(s), s)$$

Assuming that we can solve for $v(x(s), s)$ for all plans, we can compute the optimal one as

$$v^*(s) = \sup_x v(x(s), s)$$

Bellman's principle of optimality is that

$$v^*(s) = \sup_{x \in X(s)} \left[ u(x, s) + \delta \sum_{s'} v^*(s')\pi(s'|x, s) \right].$$

**5.2. Hyperbolic Discounting.** While most of the literature on repeated games uses the model of constant discounting, there is a growing literature in behavioral economics on alternatives more consistent with experimental evidence.[8] The most widely studied alternative is

---

[8] The reader should review optimization in the mathemtical appendix before proceeding to this section.

hyperbolic discounting which assumes that at time 0 agents discount the utility at time $t$ by

$$h(t) = (1 + \alpha t)^{-\frac{\gamma}{\alpha}}$$

for $\gamma > 0$ and $\alpha > 0$. Unless $\alpha$ is close to zero, hyperbolic discounting weighs the future much more heavily than constant discounting. It also implies that agents have a "time consistency" problem. What they consider optimal plan for time $t$ depends on how far time $t$ is in the future. Suppose that an agent has to decide how to allocate \$1 of consumption over three periods $0, 1, 2$. Assume that $U(x) = \sqrt{x}$ Using constant discounting the optimal plan solves

$$\sqrt{x_0} + \delta\sqrt{x_1} + \delta^2\sqrt{x_2}$$
$$\text{such that } \sum x_t = 1$$

The solution must satisfy $x_1 = \delta^2 x_0$ and $x_2 = \delta^4 x_0$. Substituting into the budget constraints, we get that

$$x_0 = \frac{1}{1 + \delta^2 + \delta^4}$$
$$x_1 = \frac{\delta^2}{1 + \delta^2 + \delta^4}$$
$$x_2 = \frac{\delta^4}{1 + \delta^2 + \delta^4}$$

Now consider what would happen if the agent re-optimized after consuming $x_0 = \frac{1}{1+\delta^2+\delta^4}$ in the first period. She would again optimally choose $x_2 = \delta^2 x_1$. Substituting this into the constraint $x_1 + x_2 = \frac{\delta^2+\delta^4}{1+\delta^2+\delta^4}$, we get

$$x_1 = \frac{1}{1 + \delta^2} \cdot \frac{\delta^2 + \delta^4}{1 + \delta^2 + \delta^4} = \frac{\delta^2}{1 + \delta^2 + \delta^4}$$
$$x_2 = \frac{\delta^4}{1 + \delta^2 + \delta^4}$$

Thus, she will wish to continue with her optimal consumption plan by consuming exactly as much as in period 2 as she had forecast.

Now consider the same allocation problem when the agent uses hyperbolic discounting. To keep the algebra simple, let $\alpha = \gamma = 1$. Thus, the agent solves

$$\sqrt{x_0} + \frac{1}{2}\sqrt{x_1} + \frac{1}{3}\sqrt{x_2}$$
$$\text{such that } \sum x_t = 1$$

The first order conditions for the optimum are $x_1 = \frac{1}{4}x_0$ and $x_2 = \frac{1}{9}x_0$. Therefore, the solution is

$$x_0 = \frac{36}{49}$$

$$x_1 = \frac{9}{49}$$

$$x_2 = \frac{4}{49}$$

Again consider what happens if the agent re-optimizes after consuming $x_0$. Now the first order condition is $x_2 = \frac{1}{4}x_1$. Substituting the constraint that $x_1 + x_2 = \frac{13}{49}$, we find that

$$x_1 = \frac{4}{5} \cdot \frac{13}{49} = \frac{52}{245} > \frac{9}{49}$$

$$x_2 = \frac{13}{245} < \frac{4}{49}$$

Thus, the agent will wish to change her optimal plan and shift more consumption to period 1. The reason for this anomaly is that she the relative weight of period 1 to period 2 consumption is higher in period 1 than it was in period 0.

While hyperbolic discounting has been useful in explaining experimental anomalies and temporal patterns in consumption (retirees consume less than a constant discounting model would predict), the applications in political science have been few.[9]

## 6. Exercises

EXERCISE 3.1. *Let Smith be a member of the House of Representatives. Smith is trying to decide whether or not to run for the Senate. He believes that he has a 50% chance of winning his party's nomination, and if he gets the nomination he has a 40% chance of winning the seat. Suppose that his utility from the Senate seat is $W$ while his utility of losing, returning home, and running his family used car lot is $L$. His utility of keeping his House seat is $H$.*

(1) Using a lottery tree, describe the lottery involved with running for the Senate.
(2) Compute the expected utility of running for the Senate.

---

[9]One conceptual obstacle is that utilities over infinite horizons may not be well defined. Suppose that an agent evaluated an infinite stream of constant utilities $u$. Evaluation requires that the series $\sum_{t=0}^{\infty} h(t)u$ converge. However, this will not be the case for large set of $\alpha$ and $\gamma$.

(3) How low must $H$ be relative to $W$ and $L$ before Smith will decides to run for the Senate?

EXERCISE 3.2. *Prove Theorem 3.1.*

EXERCISE 3.3. *Compute the expected payoff of the following lottery. There are 5 periods. In each period, the agent flips a coin and receives one dollar for each consecutive period for which she has obtained heads i.e. if she has received heads $x$ consecutive times, she receives $x.*

EXERCISE 3.4. *Suppose that instead of always revealing a goat, Monte Hall randomly selects a door to open and thus occasionally reveals the car. Clearly, a contestant should switch to the open door if the car is revealed, but should she switch to the closed door if a goat is revealed?*

EXERCISE 3.5. *Suppose that a country is fighting in a war. In each period, it cost $f > 0$ to fight a battle. The country wins each battle with probability $\pi$. The country wins the war and receives a payoff of $w > 0$ forever if it wins two consecutive battles. If it loses two consecutive battles, it loses and receives $l = 0$ forever. The county discounts future periods by $\delta$.*

*There are five states corresponding to the consecutive wins and losses in battle. Two of these are terminal states corresponding to victory of loss of the overall war. For each of the non-terminal states, compute the expected utility of continuing the war. Find a condition for $f$ in terms of $\pi$, $w$, $l$ for which the country chooses not to start the war. Find a condition for the country to surrender after losing one battle.*

EXERCISE 3.6. *Prove that $\pi(p) > p$ and sub-additivity imply that the decision weight function $\pi$ is convex for small values of $p$.*

EXERCISE 3.7. *Consider the following pair of lotteries:*
*Lottery $\mathbf{a}$: .45 chance of $6000, .55 chance of 0*
*Lottery $\mathbf{b}$: .90 chance of $3000, .10 chance of 0*
*Lottery $\mathbf{c}$: .001 chance of $6000, .999 chance of 0*
*Lottery $\mathbf{d}$: .002 chance of $3000 and .998 chance of 0*
*Which choices are predicted by Prospect Theory? Why?*

CHAPTER 4

# Social Choice Theory

## 1. The Open Search

In the pages that follow we consider a scenario that many readers of this book may soon encounter in their professional lives (if they have not already): the "open" faculty search. Consider a fictional political science department whose membership is spread evenly across five sub-fields: American ($A$), comparative ($C$), international relations ($I$), theory ($T$), and formal theory/methods ($F$). This year the fictional university is having a mediocre year financially so the dean only gives the department authorization for one additional hire. This dean, unwilling to alienate any of the department's various factions, does not specify which field the department should, but tells the department "since you study politics you should be able to settle this fairly." Those readers who have experienced a similar situation in their own departments should be smiling knowingly at the dean's folly.

Members of each sub-filed have homogeneous preferences as to which field the new hire should come. Indeed, each field has its own complete and transitive ordering over the field of the potential hire. These rankings are given by Table 4.1:

| Table 4.1 | | | | |
|---|---|---|---|---|
| **A** | **C** | **I** | **T** | **F** |
| $A$ | $C$ | $I$ | $T$ | $F$ |
| $F$ | $T$ | $C$ | $I$ | $A$ |
| $C$ | $I$ | $T$ | $C$ | $T$ |
| $I$ | $A$ | $F$ | $F$ | $C$ |
| $T$ | $F$ | $A$ | $A$ | $I$ |

So the department chair sets out to figure out how the department should decide. The first idea she entertains is to have the department vote based on *plurality rule*. Each member of the department is to cast a ballot for their favorite field and the one with the most votes wins. However, the chair quickly determines that the election would generate a five-way tie so she abandons that idea. Next she considers *pair-wise*

*majority* voting. Under this procedure, each field will be paired against each other field. If any field wins all of the pair-wise comparisons, the department would hire in that field. Sure that this is a fair way to decide, she implements this procedure in the next department meeting. The meeting begins with a vote between $A$ and $C$. Field $C$ wins with support from $T$ and $I$. However, in the vote between $C$ and $I$, $I$ wins 3-2, however $T$ beats $I$. While $T$ survives the vote with $F$, it loses to $C$. Thus, every field is defeated in at least one pairwise vote. The chair's procedure has failed to bring about any resolution. However, the chair does note that $A$ loses in every pairwise vote and $F$ loses to all fields except $A$. So at least she concludes that neither an Americanist nor a formal theorist should be hired.

Frustrated the chair decides that a scoring system such as the one used to rank college football teams might do the trick. Undeterred by previous failures, she proposes that each department member rank each field. A top ranking will give the field 5 points, a second ranking 4 points, and so on. If everyone voted according to his preferences, the chair calculated that the ranking would be $C$ (17 points), $T$ (16 points), $I$ (15 points), $F$(14 points), and $A$ (13 points). However, before the vote could take place, a theorist citing an obscure 16th century philosopher claimed it was inappropriate to weigh fourth and fifth rankings so heavily. He suggested that the vote be based solely on ranking the top three. Such a procedure would guarantee that the outcome was a tie between $C$ and $T$. Nevertheless, not feeling that the application of the dead philosopher's theory was appropriate in this circumstance, the chair moved forward with the original procedure.

However, she was taken aback by the results. The formal theorists, sensing the opportunity to be strategic, cast their ballots with $T$ in the first position and $C$ in the fifth position. This resulted in 18 points for $T$ and only 16 for $C$, an outcome preferred by $F$. Infuriated by the perceived duplicity, the $C$'s called for a revote. Their plan was to drop $T$ to the fifth position on their ballots and win 16 to 15. However, the chair quickly realized that $T$ would simply drop $C$ to the bottom in retaliation which might even lead to $I$ winning if they also cast their ballots strategically. She quickly adjourned the meeting. The next day she called the dean to have the line transferred to the economics department.

That the search ended in failure is not surprising. The fundamental result of social choice theory is that collective choice processes must either restrict the set alternatives or violate some desirable normative

properties. Furthermore, as we will soon see, all mechanisms for making collective decisions are subject to strategic manipulation by agents such as that perpetrated by $F$.

## 2. Preference Aggregation Rules

In this section, we lay out the basic notation and ideas for the formal analysis of preference aggregation rules. We will limit ourselves to the case of a finite set of agents $N = \{1, 2, ..., n\}$ ($n > 2$) who are to choose some outcome from the set $X$.[1] Our primary goal is to understand how the preferences of the individual agents map into the collective preferences so we assume that agent $i$ has preference ordering $R_i$ on $X$. As in chapter 2, we assume that these preference orderings are complete and transitive. We will denote the set of all possible complete and transitive preference orderings as $\mathcal{R}$. We denote a list of preference orderings for all $n$ agents as $\rho = \{R_1, R_2, ..., R_n\}$ which we call a preference profile. The set of profiles is therefore $\mathcal{R}^n$. By $\mathcal{B}$ we denote the set of complete orderings on $X$.

DEFINITION 4.1. *A preference aggregation rule is a function $f$ : $\mathcal{R}^n \to \mathcal{B}$.*

A preferences aggregation rule is simple a procedure that takes the set of individual preferences orderings and produces a social preference ordering. While we use subscripted $R_i$ to represent individual orderings, we denote the social ordering as $R$. As an example, consider the pairwise majority voting illustrated in the introduction to this chapter. We can formally define that function such that for two alternatives $x$ and $y$, $xRy$ if at least as many agents $xR_iy$ as $yR_ix$. Since a complete ordering can be produced for any set of preferences, this procedure satisfies our definition. Importantly, our definition does not restrict the outcomes of preference aggregation rules to be transitive. Indeed, in our fictional department pair-wise majority voting produces $T\ P\ I\ P\ C\ P\ T$.

What properties would we like our preference aggregation rule to satisfy? Perhaps the most important feature would be the ability to generate a best outcome so that the agents will actually have something to choose. In other words, we would like the social maximal set $M(R, X)$ to be non-empty, and we would like this to be true for all preference profiles. However, we know from chapter 2 that this requires that $R$ be transitive.

---

[1]Throughout this chapter, we will focus only on models of complete information so that we can speak interchangeably between choosing actions, policies, or outcomes.

DEFINITION 4.2. *A preference aggregation rule $f$ is transitive if for every $\rho \in \mathcal{R}^n$ the ordering $R$ is transitive.*

Secondly, it would be nice that the rule be at least minimally democratic so that the preferences of a single agent or dictator didn't completely determine the social ranking of the alternatives.

DEFINITION 4.3. *A preference aggregation rule $f$ is non-dictatorial if there does not exist an $i \in N$ such that for every $\rho \in \mathcal{R}^n$ for every $x, y \in X$, $xP_iy$ implies $xPy$.*

Next, we wouldn't look to favorably at aggregation rules that produced social rankings that all agents disagreed with. If all agents prefer $x$ to $y$, society's preferences should also reflect this ordering as well. This criteria was first expounded by the Italian economist Vilfredo Pareto is often referred to as *Pareto efficiency* or *optimality.*

DEFINITION 4.4. *A preference aggregation rule $f$ is Weakly Paretian if, for any $x, y \in X$, if $xP_iy$ for every $i \in N$ then $xPy$.*

Finally, it would be nice if the social preferences ordering for any two outcomes depended only on the individual preference orderings for those two outcomes. One of the reasons that the formal theorists were able to manipulate the outcome of the chair's counting procedure is that the social ranking between $C$ and $T$ depended on $F$'s relative preferences for $F$, $A$, and $I$.[2] From the perspective of a choice between $C$ and $T$ those preferences should be irrelevant. This property is known as the *independence of irrelevant alternatives* or IIA.

DEFINITION 4.5. *A preference aggregation rule $f$ is independent of irrelevant alternatives if, for any pair of policies $x, y \in X$ and any two profiles $\rho, \rho' \in \mathcal{R}^n$ with $xR_iy$ if and only if $xR'_iy$ for all $i \in N$, $xRy$ if and only if $xR'y$.*

These all seem like reasonable properties and each can be justified easily on normative or practical grounds (though the case for IIA is weaker). The properties of transitivity, Pareto optimality, and IIA can also be justified using the fact that they are all satisfied by an individual decision maker. Thus, these are the properties necessary so that social decision-making is as well behaved as that of a rational individual. Unfortunately, one of the fundamental results in all of social science

---

[2]This couting rule is known as the Borda count. For a formal definition, assume each agent has complete strict preferences and let $r_i(x) = |z : zP_ix|$ (read: number of outcomes prefered to $x$ by agent $i$). Then the Borda count rule is for all $x, y \in X$ $xPy$ if and only in $\sum_{i \in N} r_i(x) < \sum_{i \in N} r_i(y)$.

tells us that aggregation rules cannot satisfy all of these properties simultaneously . Arrow's Theorem says that the only aggregation function that produce transitive preferences while satisfying the Pareto principle and IIA is a dictator, not a happy result unless your last name happens to be Castro. Thus, the only way social preferences act like individual preferences is if they are in fact individual preferences.

We now state and prove Arrow's Theorem.

THEOREM 4.1. *If $X$ is finite and has at least three alternatives, then there is no preference aggregation rule $f : \mathcal{R}^n \to \mathcal{B}$ that is transitive, non dictatorial, weakly Paretian and independent of irrelevant alternatives.*

We need one more definition before turning to Arrow's Theorem.

DEFINITION 4.6. *Given a preference aggregation rule $f$ a set $L \subset N$ is semidecisive for $x$ against $y$ if for every $\rho \in \mathcal{R}^n$ with $xP_iy$ (all $i \in L$) and $yP_jx$ (all $j \in L^c = N \backslash L$) we have $xPy$. A set $L$ is decisive for $x$ against $y$ if for every $\rho \in \mathcal{R}^n$ with $xP_iy$ (all $i \in L$) we have $xPy$. A set $L$ is decisive if for every $x, y \in X$ it is decisive for $x$ against $y$.*

A convenient proof of Arrow's theorem rests on first establishing a property about decisiveness.

LEMMA 4.1. *Assume $f$ is a transitive preference aggregation rule that is independent of irrelevant alternatives and weakly Paretian. If $L \subset N$ is semidecisive for $x$ against $y$ for some $x, y \in X$ then $L$ is decisive.*

PROOF. Assume that $L \subset N$ is semidecisive for $x$ against $y$ and that under the profile $\rho \in \mathcal{R}^n$ $xP_iz$ for all $i \in L$. Consider a profile $\rho' \in \mathcal{R}^n$ such that for all $i \in L$ $xP_i'yP_i'z$ and for all $j \in L^c$ $yP_j'x$ and $yP_j'z$ with $z \notin \{x, y\}$ and for all $i \in L^c$ $xR_iz$ iff $xR_i'z$. Since $L$ is semidecisive for $x$ against $y$ $xP'y$. Since $f$ is weakly Paretian $yP'z$. Since $f$ is transitive $xP'z$. But since the preferences of $L^c$ on $x, z$ have not been specified in $\rho'$ and both $\rho$ and $\rho'$ agree on $x$ and $z$ (i.e. $xR_iz$ iff $xR_i'z$), the fact that $f$ is IIA implies $xPz$. Thus $L$ **is decisive for** $x$ **against** $z$. This of course implies that $L$ is semidecisive for $x$ against $z$ and an analogous argument demonstrates that $L$ **is decisive for** $x$ **against** $y$. We now verify that $L$ is decisive for $y$ against $z$. Consider a profile $\rho^0 \in \mathcal{R}^n$ with $yP_i^0z$ for all $i \in L$ and $\rho^+ \in \mathcal{R}^n$ such that for all $i \in L$ $yP_i^+xP_i^+z$ and for all $j \in L^c$ $zP_j^+x$ and $yP_j^+x$ and for all $i \in L^c$ $yR_i^0z$ iff $yR_i^+z$. Since we have already shown that $L$ is decisive for $x$ against $z$ we have $xP^+z$. Since $f$ is weakly Paretian we have $yP^+x$. Since $f$ is transitive we have $yP^+z$. Since the preferences of only members of $L$ have been specified on $\{y, z\}$ by $\rho^+$ and both

$\rho^0$ and $\rho^+$ agree on $y$ and $z$, IIA implies $yP^0z$.  Thus $L$ **is decisive for $y$ against $z$.**    This of course implies that $L$ is semidecisive for $y$ against $z$.  Relabeling the first step and using this fact implies that $L$ **is decisive for $y$ against $x$.**  Combining these (boldfaced) conclusions leads to the claim that $L$ is decisive.∎                                    □

Thus if any group is every semidecisive for some pairwise comparison than the group is decisive.  For preference aggregation rules satisfying IIA and the weak Pareto criterion if group its way once, it gets its way on all comparisons.  We now complete the proof of Arrow's theorem by showing that either an individual is decisive or the entire collective is not decisive.  The first finding violates the non-dictatorial condition and the second violates the weak Paretian condition.  Thus the proof demonstrates the incompatibility of Arrow's conditions.

PROOF OF ARROW'S THEOREM.  Assume that $X$ is finite and has at least three alternatives.  By way of a contradiction assume that we have a preference aggregation rule that is transitive, non dictatorial, weakly Paretian and independent of irrelevant alternatives.  Given the lemma, for any set $L \subset N$ either $L$ is decisive or there is no pair $x, y \in X$ for which $L$ is semidecisive for $x$ against $y$.  Consider two disjoint sets $A, B \subset N$ (disjoint means that $A \cap B = \emptyset$) which are not semidecisive for any $x$ and $y$ (and thus not decisive) Let $C = N\backslash\{A \cup B\}$.  Since $n > 2$, and no singleton set $\{i\}$ is decisive, three such sets $A, B, C$ exist.  Now consider the profile $\rho^- \in \mathcal{R}^n$ with $xP_i^- yP_i^- z$ for $i \in A$; $zP_j^- xP_j^- y$ for $j \in B$; and $yP_t^- zP_t^- x$ for $t \in C$.  Since $A$ and $B$ are not semidecisive for any pairs, we must have $zP^- x$ and $yP^- z$.  Since $f$ is transitive we must have $yP^- x$.  This implies that the set $A \cup B$ is not semidecisive for $x$ against $y$.  This means that the set is not decisive.  Thus the union of two disjoint sets which are not decisive is also not decisive.  Since $f$ is not dictatorial no singleton set is decisive.  This conclusion means that no (finite) union of individuals is decisive.  But this implies that $N$ is not decisive.  This contradicts the assumption that $f$ is weakly Paretian.  Thus the result is established.∎          □

The introduction to this chapter provides many examples of the implications of Arrow's theorem.  We have already pointed out that pairwise majority voting is not transitive and that the Borda count does not satisfy IIA.  As an additional example, consider unanimity rule which we define as $xPy$ if and only $xR_iy$ for all $i$ and $xP_iy$ for some $i$.  Clearly, this rule satisfies the weak Pareto criterion and it satisfies IIA since the rule operates only on pairs of alternatives.  But

it is not transitive. To see this suppose that we have the following individual preference orderings:

| Table 4.2 | | |
|---|---|---|
| 1 | 2 | 3 |
| $y$ | $z$ | $z$ |
| $z$ | $y$ | $x$ |
| $x$ | $x$ | $y$ |

Clearly, we have $xRy$ and $yRz$. However, $zPx$.

It is important to note that a preference aggregation rule has as its domain the set of all possible preference profiles. Thus, Arrow's theorem does not does not rule out the possibility that there is a satisfactory way to aggregate preferences for a given profile. One response to Arrow's theorem is to consider restrictions to the set of profiles and consider whether there are preference aggregation rules that satisfy the normative axioms on this smaller set.

One of the most common restriction is *single-peakedness*. Intuitively, single-peakedness means that there is some way of ordering the outcomes so that each agent's preferences rankings increase up to the most preferred outcome and then decline after that. Consider Figure 4.1 which plots preferences ordering for our fictional political science department. Given ordering of the outcomes $ACITF$, only fields $I$ and $T$ have preferences with a single peak as their preference rank is always increasing up to their ideal outcome and declining afterwards. The other fields have multiple peak preferences over the ordering $ACITF$. For example, field $A$ has peaks at $A$ and $F$. The motivated reader can verify that there is no way to order the outcomes so that all preferences have a single peak. Thus, the preference profile of our fictional department is no single-peaked. However, consider the following profile:

| Table 4.3 | | | | |
|---|---|---|---|---|
| **A** | **C** | **I** | **T** | **F** |
| $A$ | $C$ | $I$ | $T$ | $F$ |
| $F$ | $T$ | $C$ | $C$ | $A$ |
| $I$ | $I$ | $A$ | $I$ | $I$ |
| $C$ | $A$ | $F$ | $A$ | $C$ |
| $T$ | $F$ | $T$ | $F$ | $T$ |

**Insert Figure 4.1 Here**

Now we if we order the outcomes $TCIAF$ (or $FAICT$) then all fields have a single peak at the outcome associated with their own field as illustrated in Figure 4.2. To foreshadow our main result, consider the outcome of pairwise majority voting. Note that now $I$ defeats all of the other alternatives. Furthermore, pairwise majority voting produces the transitive ordering $I\ P\ A\ P\ C\ P\ F\ P\ T$, the identical to the preferences of $I$. Is it a coincidence that majority voting works well with our new single peaked preference profile? No, as we will see that single-peakedness is a sufficient (but not necessary) condition for the transitivity of majority rule.

### Insert Figure 4.2 Here

Before stating and proving this result, we need a bit more notation. Let $q$ be an ordering function which takes the set of outcomes and assigns each a unique rank. Formally, $q : X \to \{1, 2, .., |X|\}$ is a one-to-one and onto function (or bijection).[3] Now we can formally define single-peakedness.

DEFINITION 4.7. *Given a set $N$ and a choice space $X$ a preference profile $\rho \in \mathcal{R}^n$ is single-peaked if there exists some bijection $q : X \to \{1, 2, .., |X|\}$ such that for every $i \in N$ there is some $t_i \in X$ such that if $q(y) < q(t_i)$ then $t_i P_i y$ (and if $q(x) < q(y) < q(t_i)$ then $t_i P_i y P_i x$) and if $q(t_i) < q(b)$ then $t_i P_i b$ (and if $q(t_i) < q(b) < q(c)$ then $t_i P_i b P_i c$). The set of single-peaked profiles is denoted $\mathcal{S} \subset \mathcal{R}^n$.*

In the definition the policy $t_i$ is interpreted as $i$'s ideal policy, and the further the rank $q(y)$ is from $q(t_i)$ the less the agent prefers $y$. Thus, singlepeakedness implies that $x P_i y$ if $|q(t_i) - q(x)| < |q(t_i) - q(y)|$. This inequality implies that if $q(x) > q(y)$, $x P_i y$ if and only if $q(t_i) > \frac{q(x) + q(y)}{2}$. Conversely, if $q(y) > q(x)$, $x P_i y$ if and only if $q(t_i) < \frac{q(x) + q(y)}{2}$. We can now formally state the theorem.

THEOREM 4.2. *Given $\rho \in \mathcal{S}$ majority rule is transitive, weakly Paretian, IIA and non dictatorial.*

The proof is very straightforward, but we will simplify a little bit by assuming that there are an odd number of agents. We begin by showing the preference ordering majority rule produces the preference ordering of the agent with the "median ideal point." The median ideal point is defined as $t_m$ such that $q(t_i) > q(t_m)$ for exactly $\frac{N}{2}$ agents

---

[3]A function $q : X \to X$ is one-to-one if for every $y \in X$ the set $q^{-1}(y) = \{x \in X; q(x) = y\}$ is a singleton. The function is onto if for every $y \in X$ there is some $x \in X$ s.t. $q(x) = y$.

and $q(t_i) < q(t_m)$ for the remaining agents. The claim is that the social preference matches agent $m$'s preference so that $xPy$ if and only if $xP_my$. Let's prove necessity first. Suppose that $xP_my$. If $q(x) > q(y)$, then all agents with $q(t_i) < q(t_m)$ prefer $x$ to $y$. Since there are $\frac{N+1}{2}$ agents preferring $x$ to $y$, $xPy$. If $q(x) < q(y)$, then $m$ and the $\frac{N}{2}$ agents with $q(t_i) > q(t_m)$ prefer $x$ so that $xPy$. To show sufficiency, suppose that $xPy$. Since at least $\frac{N+1}{2}$ must prefer $x$ to $y$, at least one agent with $q(t_i) \leq q(t_m)$ and one agent with $q(t_i) \geq q(t_m)$ prefers $x$ to $y$. Denote these agents as $l$ and $h$ respectively. Suppose that $yP_mx$, but that $xP_ly$ and $xP_hy$. Suppose that $q(x) > q(y)$. Then single-peakedness and $yP_mx$ implies that $\frac{q(x)+q(y)}{2} > q(t_m)$ while $xP_hy$ and $xP_ly$ imply that $q(t_h) > \frac{q(x)+q(y)}{2}$ and $q(t_l) > \frac{q(x)+q(y)}{2}$. This contradicts $q(t_h) > q(t_m) > q(t_l)$. Finally, since the social preference profile matches that of the median agent, it has the properties of individual preferences which include transitivity and IIA. The social ordering is clearly weakly Paretian since the median's preferences never conflict with the social preferences. However, it might seem that $m$ is a dictator. However, recall that a dictator is one whose preferences determine social preferences for any preference profile. Since different profiles produce different median agents, $m$ is not a dictator.

## 3. Collective Choice

While it is useful to begin with the properties of aggregate preference orderings, we are ultimately interested in the set of policies that are maximal for a preference aggregation rule. As we did for individuals, we will assume that the social choices is the outcome that is maximal choice from the aggregate preference ordering. In social choice setting, we refer to the set of maximal choices as the core.

DEFINITION 4.8. *Given X, $\rho \in \mathcal{R}^n$ and a preference aggregation rule, the core is defined as $C_{f(\rho)}(X) = M(f(\rho), X)$.*

Applying theorem 1 of chapter 2, we know that if $X$ is finite and the collective preference is complete and transitive then the core is nonempty, and the social choice is well defined. However our analysis of Arrow's theorem indicates that transitivity of preference aggregation rules is not always satisfied. In such a case, we say that the core is empty or does not exist.

However, given the results of the last section, we know that majority rule is transitive under single peaked preferences is transitive so that a core does exist. Since the social preference under majority rule and single-peaked preferences is the preference ordering of the agent

with the median ideal point, it follows that the majority rule core is $t_m$. When a majority rule core exists, its outcome is known as a *Condorcet winner* after the Marquis de Condorcet who was among the first to formally study the properties of voting procedures. However, the result that the majority rule core is non-empty with single peaked preferences is generally attributed to Duncan Black.

THEOREM 4.3. *If $n > 2$ is odd and $\rho \in \mathcal{S}$ then letting $f(\cdot)$ be majority rule, $C_{f(\rho)}(X) = \{t_i : |j \in N\backslash i : t_j \leq t_i| = |k \in N\backslash i : t_k \geq t_i|\}$. That is the core is the median voter's ideal point.*

The proof is essentially the same as that of the transitivity of majority rule with single peaked preferences. This result indicates that if we are willing to assume that preferences satisfy the restrictive assumption of single-peakedness the majority rule core is well defined, and submitting to the "will of the majority" may be a reasonable way to make collective choices. However the restriction may not be appropriate for some settings. For example, suppose the set of policies is two dimensional. Then as we will see, the generalization of single peakedness is extraordinarily restrictive.

Consider Figure 4.3, which gives ideal points for five voters in two dimensions. We will assume that each agent has circular indifference curves so that it will vote for the alternative closest its ideal point. Our claim is that point 5 is a majority rule core point or the Condorcet winner as a majority prefers it to any other point in the policy space. To demonstrate this claim, we need to show that at least three voters will block any other policy. First, consider a move to any policy in the region marked $W$. Obviously, voter 5 will vote against any such move. as will voters 1 and 3. Thus, 5's ideal point is majority preferred to any policy in region $W$. Similarly, voters 1,2, and 5 will vote against moves to region $X$, voters 2,4, and 5 will vote against region $Y$, and 3,4, and 5 vote against region $Z$.

### Insert Figure 4.3 Here

The reason voter 5's ideal point is in the core is that voter 5 would be the median voter over any two alternatives in that if she prefers $x$ to $y$ at least two other voters will as well. Since this must also be true for comparisons of the ideal points of other players, voter 5's ideal point must lie on the lines connecting opposing pairs of ideal points (2-3 and 1-4). Thus, this condition (which we formalize in the next section) is fragile. A slight deviation from the intersection of these lines as in Figure 4.4 destroys the majority core. First, note that the intersection cannot be a core point because 3,4, and 5 prefer 5's ideal

point to the intersection. Secondly, note that 5's ideal point cannot be a Condorcet winner, because there are a set of points $Y'$ that are preferred by 1,2, and 3. Therefore, the existence of a "median in all directions" is a very restrictive assumption.

**Insert Figure 4.4 Here**

Given that the conditions on the existence of a majority rule core are so restrictive, an obvious question to ask is whether majority rule can at least reduce the set of possible outcomes by ruling some out as undesirable. Recall that in our introductory example, the department chair felt that it was reasonable to conclude that the department should not hire in American or Formal Theory since both of those fields were defeated by the three other fields. Thus, majority rule could eliminate two options even though it created a cycle among the three top alternatives. In this example, $C$, $I$, and $T$ represented a *top cycle:* a set of alternative that defeat all alternatives outside the set, but over which the aggregation rule is intransitive. Thus, perhaps while majority rule does not produce a core, it can produce a small top cycle. Unfortunately, this optimism is also unwarranted. Richard McKelvey showed that with under sincere voting (given any pair the of alternatives, the agent votes for the closest pair) and Euclidean preferences, the top cycle was either the core or the entire set of alternatives. While we treat this result more formally in the next section, the main intuition can be seen in the example presented in Figure 4.5. Here we have three voters with quadratic preferences. Assume that there is an initial status quo policy $a$. To illustrate McKelvey's result, it is sufficient to demonstrate that pairwise majority voting can lead from point $a$ to anywhere in the policy space and then return to $a$. First, note that point $b$ is majority preferred to $a$ since voters 1 and 2 prefer it. Continuing note that 2 and 3 prefer $c$ to $b$, and 1 and 3 prefer $d$ to $c$. Note that at each subsequent stage of the agenda, the set of policies that are preferred to the current status quo is getting larger. This allows us to reach points further and further away from the voters' ideal points. Finally, note that 1 and 2 prefer the very distant point $e$ to $d$. From $e$, we can either return to $a$ (the voter's unanimously prefer $a$) or we can leverage $e$ to get to points further and further away.

**Insert Figure 4.5 Here**

As a positive methodology the study of preference aggregation rules does not offer clear predictions. This is best exemplified by the result from McKelvey showing that it is generally the case that any policy can beat any other policy through a finite agenda. Some have interpreted this result as a prediction of chaos, whereby the theory predicts that

politics should be chaotic with observable cycles. This interpretation is naive, as it attributes a positive prediction to results that state that social choice does not generally offer predictions. A more reasonable interpretation is that the results demonstrate the need to investigate the political institutions within which collective choice is made. Under this interpretation the conclusion is that a model that takes as primitives only preferences and a preference aggregation rule may be underspecified. The tools of non-cooperative game theory will allow us to construct richer theories of collective choice.

### 3.1. Formal Analysis of the Plott and McKelvey Results**.

In this section, we present a much more formal analysis of Plott and McKelvey's results about majority rule when preferences are multi-dimensional. All of the following analysis is based on the following social choice environment.

CONDITION 1. *We assume $X \subset \mathbb{R}^d$ (d finite) is convex and agents have strictly convex, continuous preferences on $X$.*

If $X$ is compact then Theorems 3 and 4 of chapter 2 imply that each agent has a unique ideal point $y_i$ in $X$. Instead of assuming that $X$ is compact, we will assume directly that each agent has an ideal point. The assumption that preferences are strictly convex requires that the upper contour sets are strictly convex sets. The classic special case of Euclidean preferences, $u_i(x) = -\|x - y_i\|$ is convenient as in this case the upper contour sets are spherical. If we assume that preferences are Euclidean, we can specify exactly how utility changes as the policy alternatives are varied.

DEFINITION 4.9. *If preferences are Euclidean then for any $x \in X$ the gradient vector $\nabla u_i(x) = y_i - x$.*

The gradient vector is a directed vector or line segment that points away from the origin in the direction that agent $i$ would most like policy to move from point $x$.

The statement of Plott's result also uses the notion of a *pairing*. For a finite set $A$ we will call a mapping $p : A \to A$ a pairing if it is one-to-one. This means that each $i$ in $A$ is paired with exactly one $j$ in $A$. Now we can state Plott's conditions.

DEFINITION 4.10. *In the spatial model with Euclidean preferences the Plott conditions are satisfied at a policy $x \in X$ if there exists a pairing $p(\cdot)$ on the set $L = \{j \in N : y_j \neq x\}$ such that for every $i \in L$, $\nabla u_i(x) = -\lambda_i \nabla u_{p(i)}(x)$ for some $\lambda_i > 0$.*

The intuition is that when the Plott conditions are satisfied at $x$ the set of agents $L$ that do not have $x$ as their ideal point, can be paired so that each agent that wants to move in a particular direction is offset by a particular agent that wants to move in exactly the opposite direction. Since proponents of any change are paired with opponents, it is impossible to build a majority coalition to overturn $x$. The following result characterizes the relationship between the Plott conditions in the spatial model with Euclidean preferences and the majority rule core.

THEOREM 4.4. *In the spatial model with Euclidean preferences and $n$ odd the point $x$ in the interior of $X$ is in the core $C_{f(\rho)}(X)$ if and only if the Plott conditions are satisfied at $x$.*[4]

Even in the restrictive case of Euclidean preferences, it is clear that the Plott conditions will not in general be satisfied. Suppose that the conditions are satisfied for some $x$. This implies that for all $i$, $y_i - x = -\lambda_i \left( y_{p(i)} - x \right)$. If we perturb $y_{p(i)}$ so that it lies on different vector from the origin, this condition will not longer hold at $x$. More precisely if we think of the space $\mathbb{R}^{dn}$ as the space of possible ideal points of $n$ agents with Euclidean preferences on the choice space $\mathbb{R}^d$ then the subset of $\mathbb{R}^{dn}$ for which the Plott conditions are satisfied at some $x \in \mathbb{R}^d$ is incredibly small. Specifically it contains no open sets and thus has an empty-interior. Informally stated, if one imagined randomly picking an arbitrary profile from this space the probability of selecting one that satisfy the Plott conditions for some point would be 0.

While the set of profiles with a core point is very small, for any such profile there is another profile that is arbitrarily close and also has a core point. In the exercises, we ask the reader to show that if one consider small perturbations that also yield a core point, then the core point is only perturbed a little.[5]

We may conjecture that even though the core is generally empty, there is some other subset of the policy space which possesses normatively desirable properties and is therefore a reasonable prediction. One such concept is the following.

---

[4]A policy $x$ is in the interior of $X$ if there is an open ball $B(x, \varepsilon)$ that is contained in $X$.

[5]The assumption that preferences are Euclidean can be replaced by a differentiability condition to produce a more general result.

DEFINITION 4.11. *For a set $X$ a profile $\rho \in \mathcal{R}^n$ and a preference aggregation rule $f$ the top cycle set $T_{f(\rho)}$ is the set*

$$T_{f(\rho)} = \{x \in X : \forall y \in X \backslash x, \exists \{a_0, ..., a_t\} \subset X$$
$$s.t. \ a_0 = x, a_t = y \ t < \infty \ and \ \forall z < t \ a_{z-1} P a_z\}$$

The top cycle set is the set of points that can be reached from any other point via a finite chain of strict preferences. That is if $x \in T_{f(\rho)}$ then for every $y \in X \backslash x$ we can select a finite number of policies $\{a_1, a_2, ...., a_t\}$ for which $x P a_1 P a_2 P .... P a_t P y$. The following result indicates that either the Plott conditions are satisfied or the top cycle set covers the policy space.

THEOREM 4.5. *In the spatial model either $C_{f(\rho)}(X)$ is non-empty or $T_{f(\rho)} = X$.*

The implications of the last two theorems are striking. In the spatial model with Euclidean preferences unless a knife-edged condition holds (Plott conditions) any policy can be reached by any other policy in a finite chain of strict preferences.

## 4. Manipulation of Choice Functions

As the previous section illustrated, majority rule very often fails to provide sufficient guidance for making social choices. However, even in conditions in which a majority core exists, agents may not have the incentive to reveal their preferences truly so that the choice function can be implemented. As Gibbard (1973) and Satterwaite (1975) have shown, all social choice functions including majority rule are susceptible to manipulation of agents acting strategically. We have already seen one such example in the attempt of the formal theorists to manipulate the Borda count by misrepresenting their preferences over the fields. Such manipulation is also possible in voting. Consider a voting agenda where $x$ is first paired against $y$ and then against $z$. Assume that by majority vote, $xPy$, $zPx$, and $yPz$. Thus, if all voters voted according to their actual preferences, $x$ would defeat $y$ and then lose to $z$. However, voters who preferred $x$ to $y$ and $y$ to $z$ would have an incentive to vote strategically for $y$ (misrepresent their preferences between $x$ and $y$) in the first round so that $y$ might win round 1 and go on to defeat $z$.

To formalize the Gibbard-Satterwaite theorem, we need some additional notation and definitions.

DEFINITION 4.12. *The social decision function is an onto function $G : \mathcal{R}^n \to X$ that generates an outcome given a preference profile.*

Thus, a social decision function takes a preference profile and produces an outcomes (as opposed to a preference ordering). The requirement that $G$ be "onto" means that every profile produces an outcome and that every outcome can be supported by some profile. Now we can define manipulation.

DEFINITION 4.13. $G(\rho)$ *is manipulable at $\rho$ if and only if for some $i$ there exists $\rho' = \{R_1, ..., R_{i-1}, R'_i, R_{i+1}, ..., R_n\}$ such that $G(\rho')P_iG(\rho)$. $G$ is non-manipulable if it is not manipulable at any $\rho$.*

Thus, we say that a social decision function is manipulable if a single agent finds that it can change the outcome to one it prefers by reporting something other than her true preferences. In the definition, agent $i$ changes the outcome from $G(\rho)$ to $G(\rho')$ (an outcome she prefers) by claiming preferences $R'_i$ rather than $R_i$. We can now state Gibbard and Sattherwaite's result.

THEOREM 4.6. *If there are more than three alternatives and $G$ is non-manipulable, then there is a dictator (i.e. for some $i$ such that $G(\rho)P_ix$ for all $x$ and for all $\rho \in .\mathcal{R}^n$).*

The outline of the proof is as follows. First assume that $G$ is non-manipulable. Then we can construct a transitive and IIA preference ordering by applying $G$ to all of $X$ to get the most preferred outcome and then applying $G$ to the remaining elements of $X$ to get the second outcome, and so on. We can then show that this ordering satisfying the weak Pareto criterion. Thus, by Arrow's theorem, we must have a dictator.

Now consider the details. Suppose that $G$ is non-manipulable. For step 1, let $\rho$ and some $B \subseteq X$ be such that for all $i$, $x \in B$, and $y \in X/B$ then $xP_iy$. Thus, $B$ is a set of alternatives that all agents prefer to all alternatives not in the set. Our first claim is that the social decision should be from this "best" set or that $G(\rho) \in B$. Suppose that this were not true. Since $G$ is onto, we can pick an alternative profile $\rho'$ such that $G(\rho') \in B$. Then we can construct a series of alternatives:

$$y_0 = G(\rho)$$
$$y_1 = G(\rho|R'_1) = G(R'_1, ..., R_n)$$
$$y_i = G(\rho|R'_1..., R'_i) = G(R'_1, ...R'_i, ...R_n)$$
$$y_n = G(\rho')$$

Now let $k$ be the smallest integer such that $y_k \in B$. Since agent $k$ prefers everything inside $B$ to everything outside, she can get a better

alternative by reporting $R'_k$ instead of $R_k$. This contradicts $G$ being non-manipulable. Thus, it must be true that $G(\rho) \in B$.

The next step is to create a aggregation rule. Let the highest ranked element be $x_1 = G(\rho)$. Then move $x_1$ to the bottom of everyone's preference ranking to create $\rho_2$. Then let the second ranked choice be $x_2 = G(\rho_2)$. From step 1, we know that $x_2 \neq x_1$. We can continue this process until we have ranked all of the alternatives.

It is rather easy to see that the preference ordering will satisfy the weak Pareto criterion, since at every stage the decision rule has to choose an element of the "best" set for the constructed profile. So now we need to show that our aggregation rule is IIA. Suppose that it were not. Then there must be two profiles $\rho$ and $\rho'$ and alternatives $x, y$ $\in X$ such that $xR_iy$ if and only if $xR'_iy$ for all $i$ but that $x$ $f(\rho)$ $y$ and $y$ $f(\rho')$ $x$. Let $\rho(x, y)$ be the profile that agrees with $\rho$ everywhere except that $x$ and $y$ are moved to the top of everyone's ordering. We claim that $G(\rho(x, y)) = x$. Suppose it were not. Then let $\widehat{\rho}$ be the profile created by dropping alternatives to the bottom until $G(\widehat{\rho}\,) = x$. Then consider a sequence $y_i = G(\rho(x, y)|\widehat{R}_1, ..., \widehat{R}_i)$ so that $y_i$ is the social decision created by switching the first $i$ agents to the new profile. Note that $y_n = G(\widehat{\rho}) = x$ and $y_0 = G(\rho(x, y)) = y$ since step 1 implies $G(\rho(x, y)) \in \{x, y\}$. Similar to above, let $k$ be the smaller integer such that $y_k \neq y$. If $y_k = x$, then if $xP_ky$ $G$ can be manipulated by switching from $R_k$ to $\widehat{R}_k$. Alternatively, if $yP_kx$, a switch from $\widehat{R}_k$ to $R_k$ manipulates $G$. If $y_k \neq x$, then consider smallest $j > k$ such that $y_j \in \{x, y\}$. Using exactly the same logic as above, agent $j$ can manipulate $G$. So we have contradicted the assumption that $G$ is non-manipulable which implies that $G(\rho(x, y)) = x$.

Now consider a sequence $z_i = G(\rho(x, y)|R'_1(x, y), ..., R'_i(x, y))$. The assumption of no IIA implies that $z_n = G(\rho'(x, y)) = y$ and $z_0 = G(\rho(x, y)) = x$. Thus, as above, there must be some agent who prefers $y$ to $x$ and can switch their preference from $R(x, y)$ to $R'(x, y)$ and change the outcome from $x$ to $y$, or an agent preferring $x$ to $y$ who can switch from $R'(x, y)$ to $R(x, y)$ Thus, $G$ is manipulable. This contradiction implies that the preference ordering must be IIA. Thus, by Arrow's theorem, there must be a dictator for $f$ and therefore a dictator for $G$.

While essentially a negative result, this theorem does have important implications about the study of politics. Perhaps the most important is that strategic behavior is likely to be ubiquitous in politics in that mechanisms that are "strategy proof" typically do not exist. Thus, the next chapter begins our study of strategic models of politics.

## 5. Exercises

EXERCISE 4.1. *Suppose players have the following preferences:*

| 1 | 2 |
|---|---|
| a | e |
| b | b |
| c | d |
| d | a |
| e | c |

(1) What are the Borda counts for each of the alternatives?
(2) How can player 1 do better by misrepresenting her preferences?
(3) How can player 2 do better by misrepresenting his preferences?
(4) Is there any combination of statements (not necessarily truthful) for which the two players would not have an incentive to change, ex post?

EXERCISE 4.2. *Suppose there are three voters who are to decide on an alternative via pairwise majority rule. If there are three alternatives, all preferences are strict, and each voter has a different preference ordering from the other two, what percentage of the possible combinations of preferences result in a Condorcet winner? (Note that if two individuals share a common preference ordering, their most preferred must be a Condorcet winner. Why?)*

EXERCISE 4.3. *Assume that there are three voters with Euclidean preferences in two dimensions with ideal points at $(-1,0)$, $(0,1)$, and $(1,0)$ respectively. a. Construct an agenda to get from (0,0) to (2,2) b. Construct an agenda to get from (0,0) to (5,5). c. Construct an agenda to get from (0,0) to (-5,-5). Try to keep these agendas as short as possible.*

EXERCISE 4.4. *Show that if $\rho \in \mathbb{R}^{dn}$ is a profile of ideal points for which the Plott conditions are satisfied at some $x \in \mathbb{R}^d$ then for every $\varepsilon > 0$ there exists a profile $\rho^\varepsilon \in B(\rho, \varepsilon)$ for which the Plott conditions are not satisfied at any point for the profile $\rho^\varepsilon$.*

EXERCISE 4.5. *Show that if $\rho \in \mathbb{R}^{dn}$ (n odd) is a profile of ideal points for which the Plott conditions are satisfied at some $x \in \mathbb{R}^d$ then for every $\varepsilon > 0$ there exists a $\delta > 0$ such that if $\rho^\delta \in B(\rho, \delta)$ and the Plott conditions are satisfied for some point at the profile $\rho^\delta$ then the Plott conditions are satisfied for a point $x' \in B(x, \varepsilon)$ by the profile $\rho^\delta$.*

CHAPTER 5

# Games in the Normal Form

About 12.5 minutes into the broadcast, two murder suspects are arrested by Detectives Logan and Briscoe. District attorney Adam Schiff instructs assistant DA Jack McCoy to make the following offer to each separately:

- Confess and provide evidence of first degree murder by your accomplice. If she does not confess, you get a 1 year sentence on a weapons charge. If she does confess as well, you both get 8 years for murder II.
- Hold out. If your accomplice turns state's evidence, you will serve 25 to life for murder I. If she holds out, you will get 4 years for voluntary manslaughter.

Assuming each suspect loses a unit of utility for each year in prison, the following table shows the payoffs of each subject given all of the possible outcomes. The rows represent the choices of suspect 1 while the columns represent the actions of suspect 2. Each pair of numbers represents the payoffs for suspect 1 and suspect 2 respectively for each combination.

| Table 5.1: The Prisoner's Dilemma | | |
|---|---|---|
| 1\2 | *Hold Out* | *Confess* |
| *Hold Out* | -4,-4 | -25,-1 |
| *Confess* | -1,-25 | -8,-8 |

Note that the situation is strategic in the sense that the outcome of any action by suspect 1 depends on the choices of suspect 2 and vice versa. What should we expect the suspects to do? Collectively, they would like to hold out. If they both hold out, the total jail time would be only eight years, far less than any other outcome. However, unless they can reach some kind of binding agreement, its clear that the individual incentives of the suspects will preclude this outcome. Suppose that suspect 1 were to hold out, then suspect 2 would recognize that she could then do better by confessing, reducing jail time from 4 years to 1. Suspect 1 would have the same epiphany, making the

"socially optimal" agreement impossible. In fact, both suspects will recognize that she will do better by confessing regardless of the other's actions. Thus, they both confess, leading to 16 total years of jail. Thus, individual rationality leads to socially inferior outcomes (where society refers to the suspects, the DA and the police presumably prefer the outcome).[1]

In this game, the well-known "Prisoner's Dilemma", it is fairly straightforward to deduce what strategies rational actors will choose. However in other strategic situations, the predictions are more subtle. Consider the following game which we call the "Terrorist Hunt." Suppose that there are two agencies, the FBI and the CIA, which are responsible for investigating and apprehending terrorist suspects. We assume that there are two types of suspects, kingpins and operatives. Both agencies prefer capturing kingpins to capturing operatives to capturing no suspects. However, to capture a kingpin, the two agencies must cooperate by dedicating resources to a joint effort. Thus, if one agency fails to cooperate in the investigation, the other agency will fail to capture any suspects. On the other hand, each agency can capture an operative simply be acting on its own. Given this setting, we can now illustrate the strategic situation that each agency faces in deciding whether to go after the kingpin or the operative.

| Table 5.2: The Terrorist Hunt | | |
|---|---|---|
| FBI\CIA | *Kingpin* | *Operative* |
| *Kingpin* | 2,2 | 0,1 |
| *Operative* | 1,0 | 1,1 |

The rows of this matrix represent the possible strategies of the FBI (*hunt kingpin* or *hunt operative*) while the column represents those of the CIA. For both agencies, we have assigned a utility of 2 for the kingpin, 1 for the operative, and 0 for failing to capture either. Let's begin with the FBI's decision. Unlike the Prisoner's dilemma, the FBI's best choice depends on the CIA's choice. If the CIA hunts the kingpin, the FBI gets 2 from cooperating on that search versus 1 for hunting an operative by itself. However, if the CIA strikes out on its own, the FBI gets 0 from hunting the kingpin while it could generate 1 by pursuing an operative. Thus, the FBI's choice depends on what it believes the CIA will do and vice versa. So what is it reasonable

---

[1]Not to leave the reader in limbo, here is a quick summary of the rest of the episode. The confessions are thrown out on a technicality by an Upper Westside judge. The episode ends with a pithy piece of wisom by Shiff just as McCoy pours himself an 18 year old Scotch.

for each to believe? A key development in the study of strategic interaction was John Nash's characterization of rational behavior in such situations. In Nash's formulation, each agency should choose strategies that are "best responses" to the action of the other agency. If both agencies pursue such a course, the outcome is a best response to a best response. Since neither agency has a incentive to change it strategy, such an outcome is known as a Nash equilibrium.

To see whether an outcome is Nash equilibrium, it suffices to show that both agencies are doing their best given the actions of the other agencies actions. So consider whether the outcome where both agencies hunt the kingpin is a Nash equilibrium. Given that the CIA is hunting the kingpin, the best choice of the FBI is to hunt the kingpin. Similarly, the CIA's best response to the FBI's choice of the kingpin is to hunt the kingpin itself. Thus, both agencies pursuing the kingpin is a Nash equilibrium. However, this is not the only Nash equilibrium of the game. If the CIA decides to hunt the operative, the best that the FBI can do is to also settle for the operative. Since the CIA also prefers the hunt the operative when the FBI does, both agencies tracking an operative is also a Nash equilibrium.[2] While Nash's solution does not lead to a single prediction, it does rule out some outcomes. A situation where one agency hunts the kingpin while the other tracks an operative is not a Nash equilibrium since the agency hunting the kingpin would do better by switching to a search for an operative. Conversely, the agency hunting an operative would like to deviate from its strategy to search for the kingpin.

In the remainder of this chapter, we formalize and extend the concepts and issues raised by these examples.

## 1. The Normal Form

The first issue that one encounters when using game theory to model political phenomena is the question of how to represent the strategic situation. We begin with the simplest representation of a strategic situation: the *normal form* representation with complete and perfect information. Such a representation is based on the following elements:

(1) *Agents:* We let $N$ represent the set of agents. When we wish to refer to an arbitrary agent, we use the notation $i \in N$. We also use the symbol $-i \in N$ (read "not agent $i$") to refer to all agents other than $i$.

---

[2]While we defer the discussion of such possibilities, there is also a third equilibrium where each agency puruse the kingpin with probability .5 and the operative with probability .5.

(2) *Pure Strategies:* A pure strategy is an agent's plan of action such as "confess" or "hunt an operative" in our motivating examples. In a game with a single interaction, as our examples, a strategy is simply an action. However, in a game with multiple interactions, a strategy is specifies the action to be taken at each stage of the game. In a normal form representation, we must specify the set of pure strategies for each player which we denote $S_i$ for each $i \in N$. An arbitrary strategy by agent $i$ is given by $s_i \in S_i$. Given the strategy sets for each agent, we can generate the set of all possible strategy *profiles* $S$ by computing all possible combinations of strategies. Formally, $S \equiv \times_{i \in N} S_i$. A profile is then a vector $s = (s_1, ..., s_i, ..., s_n) \in S$. By $S_{-i} \equiv \times_{j \in N \setminus i} S_j$ we denote the space of strategies for every player except $i$. To economize on notation, we often represent $s$ as $(s_i, s_{-i})$. Below we discuss extensions of the definition of strategies which allows agents randomize over pure strategies.

(3) *Payoffs*: A normal form representation requires a specification of a von Neumann-Morgenstern utility function for each player over the set of strategy profiles or $u_i(s) : S \to \mathbb{R}^1$. Sometimes the utility function for $i$ is denoted $u_i(s_i, s_{-i})$. As in chapter 3, the functions $u_i(\cdot)$ are Bernoulli utility functions, and given any lottery over $S$ the agent calculates her expected utility under the lottery.

An interpretation of a normal form game is that at period 1 each player chooses their strategy $s_i \in S_i$ and in period 2 the agents receive $u_i(s)$ where $s = (s_1, ..., s_n)$. We will see that games with more periods can actually be reinterpreted as very large normal form games.

Accordingly a normal form game is completely defined by $\langle N, \{S_i, u(\cdot, ..., \cdot)\}_{i \in n} \rangle$. We sometimes use the shorthand $\langle N, S, u \rangle$ to represent a game where $u$ without a subscript represents the vector of utility functions $(u_1(\cdot), ..., u_n(\cdot))$. Some simple but quite interesting games involving two players can be represented as matrices.

To cement ideas, we know show that our two motivating examples can be fully described using the normal form. First consider the Prisoner's dilemma. Clearly $N = \{\text{player 1, player 2}\}$, $S_1 = S_2 = \{\text{hold out, confess}\}$. We can also write the payoff functions as

$$u_i(s_i, s_{-i}) = \begin{cases} -8 \text{ if } s_i = s_{-i} = \text{ confess} \\ -4 \text{ if } s_i = s_{-i} = \text{hold out} \\ -1 \text{ if } s_i = \text{confess } \& \ s_{-i} = \text{hold out} \\ -25 \text{ if } s_i = \text{hold out } \& \ s_{-i} = \text{confess} \end{cases}.$$

Similarly, the Terrorist Hunt can be represented as $N = \{CIA, FBI\}$, $S_1 = S_2 = \{\text{hunt kingpin, hunt operative}\}$ and

$$u_i(s_i, s_{-i}) = \begin{cases} 2 \text{ if } s_i = s_{-i} = \text{ hunt kingpin} \\ 1 \text{ if } s_i = s_{-i} = \text{hunt operative} \\ 1 \text{ if } s_i = \text{hunt operative } \& \ s_{-i} = \text{hunt kingpin} \\ 0 \text{ if } s_i = \text{hunt kingpin } \& \ s_{-i} = \text{hunt operative} \end{cases}.$$

Note that we were able to represent both of these normal forms with matrices. For two agent games, the relationship between the normal form and a game matrix generalizes to:

| Table 5.3:  Generic Normal Form Game | | | | |
|---|---|---|---|---|
| $1 \backslash 2$ | $s_{21}$ | $s_{22}$ | $\cdots$ | $s_{2k}$ |
| $s_{11}$ | $u(s_{11}, s_{21})$ | $u(s_{11}, s_{22})$ | $\cdots$ | $u(s_{11}, s_{2k})$ |
| $s_{12}$ | $u(s_{11}, s_{21})$ | $u(s_{12}, s_{22})$ | $\cdots$ | $u(s_{12}, s_{22})$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $s_{1l}$ | $u(s_{1l}, s_{21})$ | $u(s_{1l}, s_{22})$ | $\cdots$ | $u(s_{1l}, s_{2k})$ |

where $N = \{1, 2\}$, $S_1 = \{s_{11}, ..., s_{1l}\}$, and $S_2 = \{s_{21}, ..., s_{2k}\}$.

Representing a normal form with more than two players in a matrix is more difficult since it is difficult to represent the strategy combinations in two dimensions. However, sometimes it is useful to represent three player games using the following trick. Suppose we extended the terrorist hunt to include a third agency, the National Security Agency (NSA) whose strategy set is the same as the other two $S_3 = \{\text{hunt kingpin, hunt operative}\}$. We assume that capturing the kingpin requires cooperation by at least 2 agencies, but that each can capture an operative on its own so that the payoff function for the FBI is now:

$$u_1(s_1, s_{-1}) = \begin{cases} 2 \text{ if } s_2 = \text{kingpin or } s_3 = \text{kingpin} \\ 1 \text{ if } s_1 = \text{operative} \\ 0 \text{ if } s_2 = \text{operative } \& \ s_3 = \text{operative} \end{cases}.$$

Now consider the following pair of matrices:

| Table 5.4: Game if NSA Hunts Kingpin | | |
|---|---|---|
| FBI\CIA | Kingpin | Operative |
| Kingpin | 2,2,2 | 2,1,2 |
| Operative | 1,2,2 | 1,1,0 |

| Table 5.5:  Game if NSA Hunts Operative | | |
|---|---|---|
| FBI\CIA | Kingpin | Operative |
| Kingpin | 2,2,1 | 0,1,1 |
| Operative | 1,0,1 | 1,1,1 |

The top matrix shows the payoff triples corresponding the strategy combinations where the NSA hunts the kingpin while the lower matrix are those for which the NSA hunts the operative.   In general three player normal form games can be represented by matrices of payoff triples corresponding to each possible strategy for player 3.

## 2.  Solutions to Normal Form Games

The goal of game theory is to predict which element of $S$ will be chosen by the agents. In the the Prisoner's Dilemma we argued for the plausibility of {*confess, confess*} and that either {*kingpin, kingpin*} or {*operative, operative*} would be the result of the Terrorist Hunt.  Now we layout the general principles that lay behind these predictions.

**2.1.  Elimination of Dominated Strategies.**  A reasonable first principle for rational behavior in games is that agents should not play a strategy for which there exists an alternative strategy that raises her payoffs for all possible strategies by her opponent.  To make this criteria more concrete, recall the Prisoner's dilemma of Table 1.  The essence of the solution discussed in the introduction is that player 1 will never play "hold out" since it provides strictly less utility than "confess" for both possible choices by 2.  Thus, we say that "hold out" is *strictly dominated* for player 1 by confess and predict that she will not play it.   Similarly, "hold out" is strictly dominated for player 2 as well.   Thus, the only strategy combination that does not contain strictly dominated strategies is {*confess, confess*}.

DEFINITION 5.1. *(Strict dominance in pure strategies) A strategy $s_i$ is strictly dominated by $s_i'$ for player i iff $u_i(s_i, s_{-i}) < u_i(s_i', s_{-i})$  for all $s_{-i} \in S_{-i}$.*

DEFINITION 5.2. *(Elimination by Strict Dominance in Pure Strategies)  A strategy profile $s = (s_i, s_{-i})$ is a consistent with elimination by strict dominance if $s_i$ is not strictly dominated for any $i \in N$.*

It is also possible to tighten the predictions of elimination by strict dominance to note that it can be iterated as in the following example. Suppose that we have a normal form game represented by the following matrix:

| Table 5.6 | | | |
|---|---|---|---|
| 1\2 | *Left* | *Middle* | *Right* |
| *Up* | 1,0 | 1,2 | 0,1 |
| *Down* | 0,3 | 0,1 | 2,0 |

In this game agent 1 has no strictly dominated strategies. However, for agent 2, *Right* is dominated by *Middle* since *Middle* generates 2 versus 1 against *Up* and 1 versus 0 against *Down*. If agent 1 recognizes that agent 2 will not choose *Right*, he will perceive the game as

| Table 5.7 | | |
|---|---|---|
| 1\2 | *Left* | *Middle* |
| *Up* | 1,0 | 1,2 |
| *Down* | 0,3 | 0,1 |

In this reduced form, *Down* is now dominated for agent 1 by *Up* (payoff of 1 versus 0 for any strategy by agent 2). Since agent 2 knows that agent 1 will play *Up*, she prefers "Middle". Thus, {*Down*, *Middle*} is the solution consistent with iterated elimination of strictly dominated strategies.

DEFINITION 5.3. *(Iterated elimination of strictly dominated strategies) Given a normal form game* $\Gamma^0 = \langle N, S^0, u^0 \rangle$ *the process of iteratively deleting strictly dominated strategies is attained through the following algorithm:   for $t = 1, 2, ....$*

*In period t arbitrarily select a player $i^t \in N \backslash i^{t-1}$ and remove from $S_i^{t-1}$ each strategy that is strictly dominated in the game $\Gamma^{t-1}$ Call the set of strategies that survive $S_i^t$.  Let $S_j^t = S_j^{t-1}$ for $j \in N \backslash i^t$ and let $u_z^0$ be the restriction of $u_z$ to $S^t$ for each $z \in N$.*

*If at $\tau$ there is no $i^\tau \in N$ having a strictly dominated strategy in the game $\Gamma^{\tau-1}$ then call the set $S^{\tau-1}$ the set of outcomes that survive iterative deletion of strictly dominated strategies.*

It can be shown that regardless of what sequence of players is chosen the same set of actions will be reached. A justification for this procedure is to consider agents who reason in the following manner.

> I know that my opponents will not use strictly domi-
> nated strategies, and I know that my opponents know
> that I will not use strictly dominated strategies. Given
> this we are all really choosing from the smaller strategy
> space that survives the first $n$ iterations. But I know that
> my opponents will not use a strategy that is strictly dom-
> inated in this game, and I know that my opponents know
> that I won't play a strategy that is strictly dominated in
> this new game,.............*ad infinitum*.

While based on a much stronger premise, we can also use the idea
of weak dominance to generate predictions from games. A strategy
is weakly dominated if there is some other strategy that produces at
least as high payoffs against all opponents strategy profiles and a higher
outcome against at least one profile.

DEFINITION 5.4. *(Weak dominance in pure strategies) A strategy*
$s_i$ *is weakly dominated by* $s'_i$ *for player i iff* $u_i(s_i, s_{-i}) \leq u_i(s'_i, s_{-i})$ *for*
*all* $s_{-i} \in S_{-i}$ *and* $u_i(s_i, s_{-i}) < u_i(s'_i, s_{-i})$.*at least one* $s_{-i} \in S_{-i}$.

The definition of *elimination by weak dominance* is analogous to
that of elimination by strict dominance. An important application of
elimination by weak dominance in political science is in majority rule
voting games. Assume that a set of $n$ (odd) agents is voting between
two candidates $D$ and $R$. Each agent gets a payoff of 1 if her preferred
candidate wins and 0 otherwise. We will also define the strategy sets
so that $s_i = 1$ is a vote for $D$ and $s_i = 0$ is a vote for $R$. Since the
choice is by majority rule, the payoff for an agent who prefers $D$ is

$$u_D = \begin{cases} 1 \text{ if } \sum s_i > \frac{n+1}{2} \\ 0 \text{ otherwise} \end{cases}$$

and the payoff for an agent who prefers $R$ is $1 - u_D$. with this setup
it is easy to see that no strategies are strictly dominated. Unless
exactly $\frac{n-1}{2}$ agents choose $s_i = 1$ and exactly $\frac{n-1}{2}$ choose $s_i = 0$, an
agent's utility is not effected by her individual choice. Thus, under
most strategy profiles agents do not have strict prefers. However, since
agents do have a strict preference at opponent's profiles generating ties,
voting for the preferred outcome weakly dominates voting for the lesser
candidate. Thus, if we eliminate weakly dominated strategies, each
agent votes for her preferred candidate and the candidate preferred by
a majority wins.

As attractive as solutions based on dominance are, they suffer from
a number of problems. Perhaps the most important is that they have
little bite in a number of games. Neither version of the Terrorist Hunt

contain dominated strategies. Thus, all strategy profiles are plausible if we only impose a dominance criterion. Secondly, iterated dominance imposed strong rationality requirements on the agents. The solution to the game described in Table 6 requires that agent 2 know for certain that agent 1 will not play *Up* because he is certain that agent 2 will not play *Right*. Any slip in this logical chain leads to other outcomes. On these grounds, elimination by weak dominance is especially vulnerable since the agents may only have a strict preference against a very small set of profiles and be indifferent against the rest. Thus, the prediction of weak dominance arguments may be based agents basing decisions on very low probability occurrences.

**2.2. Nash Equilibrium.** As we discussed in the introduction, John Nash made a fundamental contribution to game theory by developing a solution for normal form games that can be applied to a very large class of models. Nash's solution requires that for all $i \in N$ agent $i$'s strategy $s_i$ be a best response to the the strategies played by the other players $s_{-i}$.

One of the most important concepts in game theory is the best response correspondence.[3] The best response to an opponent's profile $s_{-i}$ is simply the set of strategies that maximize an agent's utility when played against $s_{-i}$

DEFINITION 5.5. *The best response correspondence for agent $i \in N$ is a mapping $b_i(s_{-i}) : S_{-i} \longrightarrow\longrightarrow S_i$ defined as*
$$b_i(s_{-i}) = \{s_i \in S_i : u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i}) \text{ for every } s'_i \in S_i\} \text{for}$$
*every $s_{-i} \in S_{-i}$.*

To make these abstract definitions, consider the best response correspondences for our examples. In the Prisoner's dilemma, the best responses are:

$$b_1(confess) = \{confess\}$$
$$b_1(hold\ out) = \{confess\}$$
$$b_2(confess) = \{confess\}$$
$$b_2(hold\ ou\ t) = \{confess\}$$

---

[3]For a discussion of correspondences, see the Mathematical Appendix.

Similarly, the best response correspondences for the two agency version of Terrorist Hunt are

$$b_1(kingpin) = \{kingpin\}$$
$$b_1(operative) = \{operative\}$$
$$b_2(kingpin) = \{kingpin\}$$
$$b_2(operative) = \{operative\}$$

In these examples, the best response is a point, but in many cases it will be a set of strategies. Recall the majority voting game of the last section. Unless the opposing profile generates an exact tie, voting for either candidate is a best response. Thus, formally best response correspondence for agent $i$ preferring $D$ is

$$b_i\left(s_{-i} : \sum s_{-i} = \frac{n-1}{2}\right) = \{1\}$$

$$b_i\left(s_{-i} : \sum s_{-i} \neq \frac{n-1}{2}\right) = \{0, 1\}$$

Given the definition of the best response correspondence, we can define a Nash equilibrium as a strategy profile in which every agent is playing an element of her best response set against the strategies of the other agents.

DEFINITION 5.6. *A Nash equilibrium (in pure strategies) to a normal form game is a strategy profile $(s^*)$ satisfying the condition: for every $i \in N$*

$$s_i^* \in b_i(s_{-i}^*)$$

We can also state the definition without reference to the best response correspondence.

DEFINITION 5.7. *A Nash equilibrium (in pure strategies) to a normal form game is a strategy profile $(s^*)$ satisfying the condition: for every $i \in N$*

$$u_i(s_i^*, s_{-i}^*) \geq u_i(s_i', s_{-i}^*) \text{ for every } s_i' \in S_i.$$

The concept of a Nash equilibrium (NE) is deceptively simple. We require that agents correctly conjecture what the other players will do and that they play a best response to this conjecture. An alternative interpretation based on the second definition is that at a strategy profile which is a NE no player has an incentive to unilaterally change her strategy.

Now we can apply these definitions to our examples.

(1) The Prisoner's Dilemma: Since $\{confess\}$ is the sole element of the best response set for both agents against all outcomes, the unique Nash equilibrium is $\{confess, confess\}$.

(2) The Two Agency Terrorist Hunt: Since $b_i(kingpin) = \{kingpin\}$ for both agencies, $\{kingpin, kingpin\}$ is a Nash equilibrium. Similarly, the fact that $b_i(operative) = \{operative\}$ suggests that $\{operative, operative\}$ is also a NE.

(3) Three Agency Terrorist Hunt: First verify that the best response correspondence is

$$b_i(kingpin, kingpin) = \{kingpin\}$$
$$b_i(operative, kingpin) = \{kingpin\}$$
$$b_i(kingpin, operative) = \{kingpin\}$$
$$b_i(operative, operative) = \{operative\}$$

Using the first definition, it is easy to see that $b_i(kingpin, kingpin) = \{kingpin\}$ implies that $\{kingpin, kingpin, kingpin\}$ is a Nash equilibrium. $b_i(operative, operative) = \{operative\}$ implies that $\{operative, operative, operative\}$ is a NE. But the fact that $b_i(operative, kingpin) = \{kingpin\}$ and $b_i(kingpin, kingpin) = \{kingpin\}$ imply that there are no Nash equilibria where there are just two agencies pursuing the kingpin, even though cooperation of two agencies is sufficient to capture him.

4. Majority Voting Game: Here we claim that almost any strategy profile is a Nash equilibrium. First consider any profile such that $\sum s_i < \frac{n-1}{2}$ or $\sum s_i > \frac{n+1}{2}$. This implies that $b_i(s) = \{0, 1\}$ for all $i$. Thus, each such profile is a NE. Now consider $\sum s_i \in \left[\frac{n-1}{2}, \frac{n+1}{2}\right]$. Suppose that $\sum s_i = \frac{n+1}{2}$. This profile is a Nash equilibrium if and only if all agents choosing $s_i = 1$ prefer $D$. Suppose that this were not true and agent $i$ choose $s_i = 1$ but prefers $R$. However, since $\sum s_{\sim i} = \frac{n-1}{2}$, $b_i(s_{\sim i}) = \{0\}$. Thus, such an $s$ is not a NE. Similarly, if $\sum s_i = \frac{n-1}{2}$, $s$ is a NE if all agents choosing $s_i = 0$ prefer $R$. Thus, the set of NE includes every profile except those in which one candidate wins by a bare majority which includes a voter who prefers the losing candidate.

It is interesting to note the similarity and differences between the set of Nash equilibria and the predictions of iterated dominance arguments. In one case, the Prisoner's dilemma, the predictions are the same. In the Terrorist Hunt games, the set of Nash equilibria is smaller than the set of outcomes consistent with elimination of dominated strategies. However, in the case of majority rule voting, the set of Nash equilibria

is smaller than the set of strategies surviving strict dominance but much larger than the unique prediction of the elimination of weakly dominated strategies.

The following theorems establish the exact link between NE and profiles surviving iterated dominance.

THEOREM 5.1. *If a strategy profile* $(s_1^*, ..., s_N^*)$ *is a Nash equilibrium, then none of its strategies can be eliminated through iterated dominance.*

PROOF. Suppose the theorem is false. Let $s_i^*$ be the first of the strategies in the equilibrium profile to be eliminated. That it is eliminated requires that $u_i(s_i^*, s_{-i}) < u_i(s_i, s_{-i})$ for some $s_i$ and for all $s_{-i}$ that have not been eliminated. Since by assumption $s_{-i}^*$ has not been eliminated, $u_i(s_i^*, s_{-i}^*) < u_i(s_i, s_{-i}^*)$. This violates the assumption that $(s_1^*, ..., s_N^*)$ is a Nash equilibrium. □

THEOREM 5.2. *If only the strategy profile* $(s_1^*, ..., s_N^*)$ *survives iterated elimination of strictly dominated strategies then it is the unique pure strategy Nash equilibrium.*

PROOF. Uniqueness follows from the pervious theorem as any Nash equilibrium must survived iterated elimination. Now we show that the remaining profile must be a Nash equilibrium. Suppose $(s_1^*, ..., s_N^*)$ is not a Nash equilibrium. Then there exists $s_i'$ such that $u_i(s_i^*, s_{-i}^*) < u_i(s_i', s_{-i}^*)$ but $s_i'$ is eliminated by weak dominance. Let $s_i$ be the strategy that eliminates $s_i'$. If $s_i = s_i^*$, we have a contradiction. Assume that $s_i \neq s_i^*$. Then $u_i(s_i', s_{-i}^*) < u_i(s_i, s_{-i}^*)$ since $s_i$ eliminates $s_i^*$. We can continue this process, no more than a finite number of times until $s_i = s_i^*$. □

In all of our examples, there is at least one Nash equilibrium. However, it is quite possible that there will be no strategy profiles satisfying the requirements of Nash equilibria. To see this consider the following game. Suppose that there are two armies, an defending army $(D)$ and an invader $(I)$. The invading army must decide whether to invade through the mountains $M$ or to come through the plains $P$. Similarly, the defenders must decide where to fortify its defenses in the mountains or those in the plains. If the invader attacks an undefended area it wins a payoff of 1, however it loses 1 if it attacks a fortification. Similarly, the defenders get 1 by correctly predicting the direction of the attack and loses 1 otherwise. Thus, this normal form game, known as Colonel Blotto, can be represented by the following matrix.

| Table 5.8: Colonel Blotto | | |
|---|---|---|
| $D \backslash I$ | $M$ | $P$ |
| $M$ | 1,-1 | -1,1 |
| $P$ | -1,1 | 1,-1 |

Note that $b_D(M) = \{M\}$, $b_D(P) = \{P\}$, $b_I(M) = \{P\}$, and $b_I(P) = \{M\}$. A Nash equilibrium requires that there exist $\{s_D, s_I\}$ such that $b_I(s_D) = s_I$ and $b_D(s_I) = s_D$. Note that it is impossible for the best response correspondences of this game to satisfy these conditions. For any pair of strategies, one agent will have an incentive to choose a different one. Absent a Nash equilibrium (or any restrictions imposed by dominance), we lack a prediction for how this game will be played.

In applications one generally defines a game and seeks to characterize the set of NE (or some other set of strategy profiles). Accordingly in characterizing the equilibrium set, one is interested in existence and uniqueness. From the perspective of applied researchers, it is generally the case that a unique NE is most desirable as it means that we are analyzing a well specified model that makes clean predictions. The case of multiple equilibria is less desirable as it may mean that the model yields ambiguous predictions. The case of no equilibria may be very unsatisfactory as the model makes no predictions.

## 3. Application: The Hotelling Model of Political Competition

To provide an additional example of Nash equilibria, we turn one of the most widely used models in political game theory: the Hotelling (1927) model of political competition which was extended by Downs (1957). Suppose that a small town wants to decide where to build a school. Its citizenry are stretched out evenly along a one mile stretch of road and would like the school to be built as close to their homes as possible. Thus, we assume that the voter's ideal points are distributed uniformly over $[0, 1]$. The decision will be made following an election where two candidates compete for office by campaigning on promises of the school's location. The winning candidate builds the school at the promised location and receives a payoff of 1. The losing candidate gets $-1$. In the case of a tie, we assume that the election is decided by a coin toss. To keep the discussion simple, we will assume that the voters are not strategic agents in the game, but vote for the closest candidate.

Given our setup, the candidate's strategy sets are $S_1 = S_2 = [0, 1]$. Given that the voters vote for the closest candidate, we can compute the vote shares for both candidates for any strategy profile $(s_1, s_2)$. Since the voters are distributed uniformly, the number of voters in any interval is equal to width of that interval. So if $s_2 > s_1$, all voters to the left of $\frac{s_1 + s_2}{2}$ vote for candidate 1 and her share is thus $\frac{s_1 + s_2}{2}$. The remaining $1 - \frac{s_1 + s_2}{2}$ voters vote for candidate 2. Conversely, if $s_1 > s_2$, candidate 2 receives a vote share $\frac{s_1 + s_2}{2}$ while candidate 1 gets the rest. Given the candidate's utilities for winning, their payoff functions over strategies are:

$$u_1(s_1, s_2) = \begin{cases} 1 \text{ if } s_1 < s_2 \text{ and } \frac{s_1 + s_2}{2} > .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1 + s_2}{2} < .5 \\ 0 \text{ if } s_1 = s_2 \text{ and } \frac{s_1 + s_2}{2} = .5 \\ -1 \text{ if } s_1 < s_2 \text{ and } \frac{s_1 + s_2}{2} < .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1 + s_2}{2} > .5 \end{cases}$$

and

$$u_2(s_1, s_2) = -u_1(s_1, s_2)$$

Our claim is that the unique Nash equilibrium is $s_1 = s_2 = .5$. To demonstrate, we begin by computing the best response functions. We start with candidate 1. Suppose that $s_2 < .5$, then candidate 1 can win for sure by choosing any platform generate a vote share of .5 or more. Thus, $b_1(s_2) = (s_2, 1 - s_2)$. For $s_2 > .5$, similar calculations produce $b_1(s_2) = (1 - s_2, s_2)$. If $s_2 = .5$, candidate 1 can at best generate a tie by choosing .5 as well. Thus, $b_1(.5) = .5$. Since candidate 2's situation is entirely symmetric, her best response correspondence is

$$b_2(s_1) = (s_1, 1 - s_1) \quad \text{if } s_1 < .5$$
$$b_2(s_1) = (1 - s_1, s_1) \quad \text{if } s_1 > .5$$
$$b_2(s_1) = s_1 \qquad\qquad \text{if } s_1 = .5$$

That $s_1^* = s_2^* = .5$ is a NE follows trivially since $b_1(.5) = .5$ and $b_2(.5) = .5$. The trick is to show that it is the only one. Suppose to the contrary that $s_1^* = s_2^* \neq .5$. However, this cannot be a NE since $b_1(s_2^*)$ does not include $s_2^* = s_1^*$. Now suppose $s_1^* < s_2^*$. Now $b_1(s_2^*) = (1 - s_2^*, s_2^*)$ while $b_2(s_1^*) = (s_1^*, 1 - s_1^*)$. Putting all of these conditions together, a NE requires that $1 - s_2^* < s_1^*$ and $s_2^* < 1 - s_1^*$ which imply the contradictory inequalities that $s_2^* > 1 - s_1^*$ and $s_2^* < 1 - s_1^*$. Since the $s_1^* > s_2^*$ is analogous, we have established the uniqueness of $s_1^* = s_2^* = .5$ as a Nash equilibrium.

A more intuitive proof of our claim follows from our second definition of Nash equilibrium (a strategy profile for which no agent has a strict preference to deviate). Clearly, no candidate will defect from $s_1^* = s_2^* = .5$ payoffs would fall from .5 to $-1$. Now consider any other possible equilibrium $(s_1^*, s_2^*)$. If one candidate wins in this equilibrium,

the other candidate move to .5 and at least generate a tie. Thus, any other equilibrium must generate a tie. However, unless $s_1^* = s_2^* = .5$, either candidate can move to .5 and win for sure. Thus, $s_1^* = s_2^* = .5$ is the only NE.

Finally, we note that the same outcome is generated by applying elimination of weakly dominated strategies. Each candidate generates no worse than a tie by choosing .5. The tie occurs only if the opponent chooses .5 as well. Against this profile anything other platform will lose. Thus, $s_i = .5$ weakly dominates all other strategies.

**3.1. Vote Maximizing Candidates.** Suppose instead of maximizing their chance of winning, each candidate maximizes his vote share. Thus, the payoffs are

$$u_1(s_1, s_2) = \begin{cases} \frac{s_1+s_2}{2} & \text{if } s_1 < s_2 \\ .5 & \text{if } s_1 = s_2 \\ 1 - \frac{s_1+s_2}{2} & \text{if } s_1 > s_2 \end{cases}$$

and

$$u_2(s_1, s_2) = 1 - u_1(s_1, s_2)$$

Our claim is that $s_1^* = s_2^* = .5$ is again the unique Nash equilibrium. Again, let's start with the best response correspondences. Suppose that $s_2 < .5$, candidate 1 would like to choose the smallest platform greater than $s_2$. However, since the strategy sets are continuums, no such platform exists. Similarly, if $s_2 > .5$, candidate 1 would like to choose the smallest platform less than $s_2$ which does not exist either. Candidate 2 faces the same situation so that $b_i(s_{-i} \neq .5) = \phi$. Now consider the best response to $s_2 = .5$. candidate 1 gets .5 for proposing .5 generating the tie and a strictly lower vote share for any other platform. Thus, .5 is the best response. Since candidate 2 faces the same incentives, $b_i(s_{-i} = .5) = .5$. So clearly $s_1^* = s_2^* = .5$ is a Nash equilibrium. Since the best response sets for any other strategy pair are empty, this is also the unique NE.

**3.2. Ideological Candidates.** We now consider one last version of the Hotelling model. In this version, candidate are ideological in that they care about the policy implemented by the winning candidate. We will assume that candidate 1 wants the school to be as close to 0 as possible so that her utility from the winning outcome $x$ is $-|x|$. Similarly, candidate 2 would like the outcome to be as close to 1 as possible and so gets a utility of $-|1-x|$ from a school located at $x$. In the case of a tied election, the voters flip a coin and the winner gets to implement her platform. Given that the candidates' incentives to move policy to the extremes, it might seem that the outcomes would

no longer be located at the median voter. However, this is not the
case as we show once again that $s_1^* = s_2^* = .5$ is the unique NE.

First let us specify the payoff functions for each candidate.

$$u_1(s_1, s_2) = \begin{cases} -|s_1| & \text{if } s_1 < s_2 \text{ and } \frac{s_1+s_2}{2} > .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1+s_2}{2} < .5 \\ -.5 \cdot |s_1| - .5 \cdot |s_2| & \text{if } s_1 = s_2 \text{ and } \frac{s_1+s_2}{2} = .5 \\ -|s_2| & \text{if } s_1 < s_2 \text{ and } \frac{s_1+s_2}{2} < .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1+s_2}{2} > .5 \end{cases}$$

and

$$u_2(s_1, s_2) = \begin{cases} -|1 - s_1| & \text{if } s_1 < s_2 \text{ and } \frac{s_1+s_2}{2} > .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1+s_2}{2} < .5 \\ -.5 \cdot |1 - s_1| - .5 \cdot |1 - s_2| & \text{if } s_1 = s_2 \text{ and } \frac{s_1+s_2}{2} = .5 \\ -|1 - s_2| & \text{if } s_1 < s_2 \text{ and } \frac{s_1+s_2}{2} < .5 \text{ or if } s_1 > s_2 \text{ and } \frac{s_1+s_2}{2} > .5 \end{cases}$$

Now consider the best response functions. We begin with candidate 1.
If $s_2 < .5$, no proposal less than $s_2$ defeats $s_2$. Thus, candidate 1's best
response is to choose $s_2$ or a proposal that loses to $s_2$. This implies that
$b_1(s_2 < .5) = \{s_2\} \cup (1 - s_2, 1]$. Alternatively, if $s_2 > .5$, candidate 1
will choose the smallest platform that defeats $s_2$. However, just as in
the last section, such a platform does not exist so that $b_1(s_2 > .5) = \phi$.
By similar arguments, $b_2(s_1 > .5) = \{s_1\} \cup [0, 1 - s_1)$ and $b_2(s_1 < .5) =$
$\phi$. Finally, consider the best response to $s_i = .5$. Any proposal by the
other candidate loses for sure generating a utility of .5. Responding
with .5 leads to a lottery over $s_1 = s_2 = .5$ which has an expected value
of .5 for both candidates. Thus, $b_i(.5) = [0, 1]$.

Given these correspondences, it is easy to see that $s_1^* = s_2^* = .5$ is a
Nash equilibrium since $.5 \in b_i(.5)$ for both candidates. Now we show
uniqueness. It's clear that since $b_1(s_2 > .5) = \phi$ and $b_2(s_1 < .5) =$
$\phi$ the only possible candidates for NE are $s_1^* > .5 > s_2^*$. However,
this condition in conjunction with $b_1(s_2 < .5) = \{s_2\} \cup (1 - s_2, 1]$ and
$b_2(s_1 > .5) = \{s_1\} \cup [0, 1 - s_1)$ implies that $s_1^* > 1 - s_2^*$ and $s_2^* < 1 - s_1^*$.
Since these last two inequalities are inconsistent, $s_1^* = s_2^* = .5$ is the
only possible NE.

## 4. Existence of Nash Equilibria

As we saw in our Colonel Blotto example, there is no guarantee
that a Nash equilibrium in *pure strategies* will exist. In this section,
we explore the conditions necessary for the existence of Nash equilibria.

First we consider a set of sufficient conditions for equilibria to exist
in pure strategies. Before doing so we need an additional definition to
describe an important property payoff functions: quasi-concavity.

DEFINITION 5.8. *A function $f(x) : X \to \mathbb{R}^1$ with $X$ a convex set is
strictly quasi-concave if for any $t \in \mathbb{R}^1$, $x \neq y \in X$ and $\lambda \in (0, 1)$ with
$f(x) \geq t$ and $f(y) \geq t$ it is the case that $f(\lambda x + (1 - \lambda)y) > t$.*

Alternatively, a function is strictly quasi-concave if its upper contour sets are convex. Since we have already shown strictly convex preferences have singleton maximal sets (when the sets are non empty). This means that if a game has convex $S$ and utility functions $u_i(s_i, s_{-i})$ that are strictly quasi concave in $s_i$ for each $s_{-i} \in S_{-i}$ the best response correspondence will be a function (i.e. a single valued correspondence). This feature of the best response correspondences guarantees the existence of a NE. In the next subsection we prove the following result.

THEOREM 5.3. *If the normal form game* $\langle N, S, u \rangle$ *satisfies the following conditions:*

*(1) $S_i$ is a convex and compact subset of a Euclidean space for each $i \in N$.*

*(2) $u_i(s_i, s_{-i}) : S \to \mathbb{R}^1$ is a continuous function for each $i \in N$*

*(3) for every $i \in N$ and every $s'_{-i} \in S_{-i}$ the function $u_i(s_i, s'_{-i}) : S_i \to \mathbb{R}^1$ is strictly quasi concave*

*a Nash Equilibrium exists.*

As useful as the previous result, it is obviously restrictive in that many games will not satisfy its assumptions. An alternative approach is to consider *mixed strategies*. A mixed strategies is a randomization over the pure strategies. We will denote a mixed strategy as $\sigma_i$ and use $\sigma_i(s_i)$ to denote the probability that agent $i$ chooses strategy $s_i$. The set of mixed strategies for player $i$ will be the set of probability distributions over $S_i$ which we denote. $\Delta_i = \Delta(S_i)$.

Thus, a game in mixed strategies is very much like a game in pure strategies except that each agent chooses $\sigma_i \in \Delta_i$ rather than $s_i \in S_i$ and the players evaluate strategies according to expected utility over the lotteries induced by the mixed strategies. The following is a formal definition of a game in mixed strategies.

DEFINITION 5.9. *Given a normal form game* $\Gamma = \langle N, S, u \rangle$ *the mixed extension game* $\Gamma^m = \langle N, \Delta, u^m \rangle$ *is constructed as follows:* $\Delta_i = \Delta(S_i)$ *with an arbitrary strategy* $\sigma_i \in \Delta_i$ *for all $i \in N$ and* $\Delta = \times_{i \in N} \Delta_i$ *with $\sigma_i(s_i)$ denoting the probability that mixed strategy $\sigma_i$ assigns to pure strategy $s_i$. The expected utility function, $U_i(\sigma_i, \sigma_{-i}) : \Delta \to \mathbb{R}^1$ is defined as*

$$U_i(\sigma_i, \sigma_{-i}) = \sum_{s_{-i} \in S_{-i}} \sum_{s_i \in S_i} u_i(s_i, s_{-i}) \sigma_i(s_i), \sigma_{-i}(s_{-i}) \text{ for all } i \in N.$$

Since the mixed extension of a normal form game is itself a normal form game, our definition of NE applies directly to mixed extensions.

To understand the mechanics of mixed strategy games, recall the Colonel Blotto game. Suppose now that each side can choose lotteries

over its strategies.    So let $\sigma_1 = \sigma_1(M)$ be the probability that the defender protects the mountains and let $\sigma_2 = \sigma_2(M)$ be the probability that the invader attacks the mountains.  Given these strategies, we can compute the expected utility for each player for each action.

$$u_1(M, \sigma_2) = \sigma_2 - (1 - \sigma_2) = 2\sigma_2 - 1$$
$$u_1(P, \sigma_2) = -\sigma_2 + (1 - \sigma_2) = 1 - 2\sigma_2$$
$$u_2(\sigma_1, M) = -\sigma_1 + (1 - \sigma_1) = 1 - 2\sigma_1$$
$$u_2(\sigma_1, P) = \sigma_1 - (1 - \sigma_1) = 2\sigma_1 - 1$$

From these payoffs, it is straightforward to compute $b_1(\sigma_2)$ and $b_2(\sigma_1)$. Note that $u_1(M, \sigma_2) > u_1(P, \sigma_2)$ if $\sigma_2 > \frac{1}{2}$, $u_1(M, \sigma_2) = u_1(P, \sigma_2)$ if $\sigma_2 = \frac{1}{2}$, and $u_1(M, \sigma_2) < u_1(P, \sigma_2)$ if $\sigma_2 < \frac{1}{2}$.    Thus, we can write the defender's best response to all possible values of invader's mixed strategy of $\sigma_2$.[4]

$$b_1(\sigma_2) = \begin{cases} M \text{ if } \sigma_2 > \frac{1}{2} \\ P \text{ if } \sigma_2 < \frac{1}{2} \\ \{M, P\} \text{ if } \sigma_2 = \frac{1}{2} \end{cases}$$

Clearly, since $b_1(\sigma_2) = \{M, P\}$ when $\sigma_2 = \frac{1}{2}$, any randomization over this set is also a best response.    So we may also write $b_1(\sigma_2) = \sigma_1$ if $\sigma_2 = \frac{1}{2}$.  Now we can use exactly the same process to compute the invader's best response function.  Note that

$$u_2(\sigma_1, M) > u_2(\sigma_1, P) \text{ if } \sigma_1 < \frac{1}{2}$$
$$u_2(\sigma_1, M) < u_2(\sigma_1, P) \text{ if } \sigma_1 > \frac{1}{2}$$
$$u_2(\sigma_1, M) = u_2(\sigma_1, P) \text{ if } \sigma_1 = \frac{1}{2}$$

Thus, the invader's best response correspondence is

$$b_2(\sigma_1) = \begin{cases} M \text{ if } \sigma_1 < \frac{1}{2} \\ P \text{ if } \sigma_1 > \frac{1}{2} \\ \{M, P\} \text{ or } \sigma_2 \text{ if } \sigma_1 = \frac{1}{2} \end{cases}$$

Now we can compute the Nash equilibrium.  We know from our previous discussion that there can be no Nash equilibrium where one agent plays any of its pure strategies.  So the remaining combination to check is whether there is a combination of $\sigma_1$ and $\sigma_2$ that satisfies Nash's criteria.    Clearly, $\sigma_1$ is only a best response by the defender if $\sigma_2 = \frac{1}{2}$

---

[4]Note that the pure strategies of mountain or plains are just the special cases of $\sigma_2 = 0$ and $\sigma_2 = 1$.

and $\sigma_2$ is only a best response if $\sigma_1 = \frac{1}{2}$. Thus, $\sigma_1 = \frac{1}{2}$ and $\sigma_2 = \frac{1}{2}$ is a Nash equilibrium in mixed strategies.

For a two player, two strategy game plotting best response functions is often useful for finding mixed strategy equilibria. Consider Figure 5.1. The horizontal axis plots the defender's mixed strategy raging from $\sigma_1 = 0$ (plains) to $\sigma_1 = 1$ (mountains). The vertical axis plots the invader's mixed strategy from $\sigma_2 = 0$ (plains) to $\sigma_2 = 1$ (mountains). The solid line represents the defenders best response to $\sigma_2$ and the dotted line represents the invader's best response. Note that the only intersection of these best response curves is at the Nash equilibrium $\sigma_1 = \frac{1}{2}$ and $\sigma_2 = \frac{1}{2}$.

**Insert Figure 5.1 Here**

Mixed strategy equilibrium do not only exist in games without pure strategy equilibria. To see this, let's return to the two agency Terrorist Hunt game represented in Table 5.2. Now let $\sigma_1$ be the probability that the FBI hunts the kingpin and $\sigma_2$ be the probability the CIA does. Thus, we can compute the expected utilities for each agency for each action.

$$u_1(kingpin, \sigma_2) = 2\sigma_2$$
$$u_1(operative, \sigma_2) = 1$$
$$u_2(\sigma_1, kingpin) = 2\sigma_1$$
$$u_2(\sigma_1, operative) = 1$$

As before we can compare these utilities to generate the best response functions.

$$b_1(\sigma_2) = \begin{cases} \text{kingpin if } \sigma_2 > \frac{1}{2} \\ \text{operative if } \sigma_2 < \frac{1}{2} \\ \sigma_1 \in [0,1] \text{ if } \sigma_2 = \frac{1}{2} \end{cases}$$

$$b_2(\sigma_1) = \begin{cases} \text{kingpin if } \sigma_1 > \frac{1}{2} \\ \text{operative if } \sigma_1 < \frac{1}{2} \\ \sigma_2 \in [0,1] \text{ if } \sigma_1 = \frac{1}{2} \end{cases}$$

Figure 5.2 plots these best response functions as we did for Colonel Blotto. Now note that there are three intersections of $\{\sigma_1, \sigma_2\}$: $\{0,0\}$, $\{1,1\}$, and $\left\{\frac{1}{2}, \frac{1}{2}\right\}$. Of course the first two correspond to the pure strategy equilibria we have already computed, but the third is an additional mixed strategy equilibrium.

**Insert Figure 5.2 Here**

A feature of mixed strategy equilibria is that the probability that an agent plays a particular strategy is not a function of her preferences but those of her opponents. This is because an agent will only play

mixed strategies is she is indifferent between a set of pure strategies. Thus, the opponent must be choosing its own mixed strategy to insure that the agent is indeed indifferent. Thus, its the preferences of the opponent that determine the mixing probabilities. Occasionally, this leads to counter-intuitive predictions. For example, suppose that we modified the Terrorist Hunt so that the CIA received a much higher payoff that the FBI for capturing the kingpin. We modify the payoffs as follows

| Table 5.9: Modified Terrorist Hunt | | |
|---|---|---|
| FBI\CIA | *Kingpin* | *Operative* |
| *Kingpin* | 2,4 | 0,1 |
| *Operative* | 1,0 | 1,1 |

A naive prediction would be that since the CIA gets a higher payoff from the kingpin that it will be more likely to hunt him. This prediction is wrong. The reader can check that best response functions are now:

$$b_1(\sigma_2) = \begin{cases} kingpin \text{ if } \sigma_2 > \frac{1}{2} \\ operative \text{ if } \sigma_2 < \frac{1}{2} \\ \sigma_1 \in [0,1] \text{ if } \sigma_2 = \frac{1}{2} \end{cases}$$

$$b_2(\sigma_1) = \begin{cases} kingpin \text{ if } \sigma_1 > \frac{1}{4} \\ operative \text{ if } \sigma_1 < \frac{1}{4} \\ \sigma_2 \in [0,1] \text{ if } \sigma_1 = \frac{1}{4} \end{cases}$$

Figure 5.3 plots these best responses and reveals that the mixed strategy equilibrium is now $\left\{\frac{1}{4}, \frac{1}{2}\right\}$. Thus, the change in the CIA's preferences did not lead it to hunt the kingpin with a higher probability, it still hunts him $\frac{1}{2}$ of the time. However, the change actually *decreases* the likelihood that the FBI will hunt the kingpin. Thus, the probability that the the operative will get caught in the mixed strategy equilibrium goes down. The logic behind this result is straightforward. In the mixed strategy equilibrium, the CIA must choose $\sigma_2$ to make the FBI indifferent between hunting the kingpin and the operative. Since the FBI's preferences did not change, $\sigma_2$ does not change. Since the FBI choose $\sigma_1$ to make the CIA indifferent, an increased utility for the kingpin means that the FBI must lower the probability of searching for the kingpin to maintain this indifference.

## Insert Figure 5.3 Here

While mixed strategy equilibria do have some undesirable properties, they are guaranteed to exist in games with finite strategy sets as per the following theorem.

THEOREM 5.4. *(Nash) Given a normal form game* $\Gamma = \langle N, S, u \rangle$ *in which $S$ is finite, the mixed extension $\Gamma^m = \langle N, \Delta, u^m \rangle$ has at least one NE. In other words every finite game has a mixed strategy NE.*

The proof of this theorem (just as the last) utilizes some advanced mathematics (fixed point theorems) and are thus relegated to the advanced section below.

**4.1. Dominance and Mixed Strategies.** Now that we have defined mixed strategies, we can extend our definition of dominance to include mixed strategies

DEFINITION 5.10. *A pure strategy $s_i \in S$ is strictly dominated if there exists a $\sigma'_i \in \Delta(S_i)$ such that*

$$U_i(\sigma'_i, s_{-i}) > u_i(s_i, s_{-i}) \text{ for every } s_{-i} \in S_{-i}.$$

*The strategy $s_i$ is weakly dominated if there is a $\sigma'_i$ for which the inequality holds weakly for every $s_{-i} \in S_{-i}$ and strictly for some $s'_{-i} \in S_{-i}$.*[5]

The extension is straightforward in that we know allow strategies to be dominated by mixed strategies. As the following example show, this extended definition of dominance is much stronger that dominance in pure strategies as some strategies may be dominated by mixed strategies which are not dominated by pure strategies.[6] To see how a mixture may dominated a strategy undominated by pure strategies, consider the following game:

| Table 5.10 | | | |
|---|---|---|---|
| $1\backslash 2$ | $L$ | $M$ | $R$ |
| $U$ | $3,1$ | $4,2$ | $1,4$ |
| $D$ | $2,4$ | $1,2$ | $3,1$ |

Note that neither player has any strategies which are dominated by pure strategies. However, consider a mixed strategy by player 2 of $\sigma_2(L) = \frac{1}{2}$ and $\sigma_2(R) = \frac{1}{2}$. This mixture has an expected value of 2.5 when played against both $U$ and $D$. This is a higher utility than player 2's pure strategy of $M$. Having extended the definition of dominance, we can now state the following relationship between mixed strategy equilibria and iterated dominance.

---

[5] Recall that $U_i(\sigma'_i, s_{-i}) = \sum_{s'_i \in S_i} u_i(s'_i, s_{-i}) \sigma_i(s_i)$.

[6] Obviously, the converse cannot be true as any strategy dominated by a pure strategy is dominated by a mixed strategy placing probability one on the dominating strategy.

THEOREM 5.5. *Given a finite normal form game with a mixed strategy NE $\sigma^*$ if the strategy $s_i$ played with positive probability under $\sigma_i^*$ then it survives iterated deletion of strictly dominated strategies.*

We leave the proof as an exercise. This theorem can be quite useful in computing mixed strategy equilibria. Instead of computing best responses to lotteries over all strategies, we need only compute those for mixtures of strategies surviving iterated dominance.

**4.2. Calculating Nash Equilibria.** Even though we can specify sufficient conditions for the existence of NE, actually computing the equilibrium of a game is sometimes more art than science. In fact, there are games which we know have equilibria like Chess for which the actual equilibrium strategies have never been calculated. Nevertheless, there are a few tricks and algorithms which can facilitate computation

4.2.1. *Pure Strategy Nash Equilibria in Finite Games.* We begin by outlining the process for checking whether a given profile is an equilibrium. Given a finite game, one way to characterize all of the pure strategy NE is to test whether each profile $s' \in S$ is a Nash equilibrium. To do this one starts with a profile $s' = (s'_1, ..., s'_n)$ and asks the following sequence of questions:

1. Holding $s'_2, ..., s'_n$ fixed is there a strategy $s''_1$ for which $u_1(s''_1, s'_{-1}) > u_1(s'_1, s'_{-1})$. If so then $s'$ is not a Nash equilibrium. If not then continue

2. Holding $s'_1, s'_3..., s'_n$ fixed is there a strategy $s''_2$ for which $u_2(s''_2, s'_{-2}) > u_2(s'_2, s'_{-2})$. If so then $s'$ is not a Nash equilibrium. If not then continue
.
.
.

$i$. Holding $s'_1, ...s'_{i-1}, s'_{i+1}.., s'_n$ fixed is there a strategy $s''_i$ for which $u_i(s''_i, s'_{-i}) > u_i(s'_i, s'_{-i})$. If so then $s'$ is not a Nash equilibrium. If not then continue
.
.
.

$n$. Holding $s'_1, ..., s'_{n-1}$ fixed is there a strategy $s''_n$ for which $u_n(s''_n, s'_{-n}) > u_n(s'_n, s'_{-n})$. If so then $s'$ is not a Nash equilibrium. If not then $s'$ is a Nash equilibrium.

This algorithm is then repeated for each profile in $S$.

In two player finite games –which are representable by matrices– the algorithm is particularly straightforward. Start with a profile i.e. matrix entry and see whether there is an entry in the same column that makes the row player better off. If so then the original profile is not

a Nash equilibrium. If not then repeat the exercise interchanging row and column. Consider the following example:

| Table 5.11 | | | |
|---|---|---|---|
| 1/2 | $l$ | $c$ | $r$ |
| $t$ | 5,4 | 2,3 | 6,2 |
| $m$ | 2,5 | 3,6 | 5,5 |
| $b$ | 5,2 | 0,3 | 7,4 |

We begin by conjecturing that $(t, l)$ is a pure strategy Nash equilibrium. Note that player 1 can only affect the row choice. Given that $s_2 = l$ player 1 chooses between $u_1(t, l) = 5$, $u_1(m, l) = 2$ and $u_1(b, l) = 5$. Accordingly $b_1(l) = \{t, b\}$. Given $s_1 = t$ player 2 chooses between utilities of 4, 3 and 2 so $b_2(t) = \{l\}$. Accordingly $(t, l)$ is a NE. Since $b_2(l)$ has a single element, we know that this is the only pure strategy NE in which $t$ is played. Recall that $b_1(l) = \{t, b\}$ so $(m, l)$ is not a NE as player 1 would deviate to either $t$ or $b$ if she anticipated player 2 selecting $l$. Now we conjecture that $(m, c)$ is a NE and note that $b_1(c) = \{m\}$ and $b_2(m) = \{c\}$ and thus our conjecture is correct. Since $r \notin b_2(m)$ we note that the only pure strategy NE in which $m$ is played is $(m, c)$. Now if we conjecture that $(b, l)$ is a NE we will note that $b_2(b) = \{r\}$ and thus we see that our conjecture is incorrect. If we conjecture that $(b, c)$ is a NE we note that $b_1(c) = \{m\}$ and so our conjecture is incorrect. Finally conjecturing that $(b, r)$ is a NE we observe that neither player has an incentive to deviate and that the conjecture is correct. We thus conclude that the set of pure strategy NE is $\{(t, l), (m, c), (b, r)\}$.

## 5. Pure Strategy Nash Equilibria in Non-Finite Games*

In games where the strategy space is not finite, the algorithm exhibited above will not work. However, sometimes we can use the techniques of optimization to compute equilibria. If we assume that the utility functions $u_i(s)$ are twice differentiable, we can use simple calculus to compute best responses. Since a best response to $s_{-i}$ is the maximizer of $u(s_i, s_{-i})$ over $S_i$, a sufficient condition for $s_i \in b_i(s_{-i})$ is that $\frac{\partial u_i(s_i, s_{-i})}{\partial s_i} = 0$ and $\frac{\partial^2 u_i(s_i, s_{-i})}{\partial s_i^2} < 0$. The first of these conditions is termed the first order condition (FOC) while the second is known as the second order condition (SOC).[7] Thus, when the FOC and SOC hold, we can use the solutions to $\frac{\partial u_i(s_i, s_{-i})}{\partial s_i} = 0$ to solve for each of the

---

[7] If $S_i$ has more than one dimension then the term $\frac{\partial u_i(s_i, s_{-i})}{\partial s_i}$ is a vector where each coordinate is the partial derivative with respect to one coordinate of $s_i$, and the

best response functions. Then a Nash equilibrium is the solution to
the system

$$s_1^* = b_1(s_{-1}^*)$$
$$.$$
$$s_i^* = b_i(s_{-i}^*)$$
$$.$$
$$s_n^* = b_n(s_{-n}^*)$$

However, this procedure can be unnecessarily cumbersome as it requires solving for each best response function as well as the system of
equations. Often it is more convenient to solve the system of first
order conditions directly

$$\frac{\partial u_1(s_1, s_{-1})}{\partial s_1} = 0$$
$$.$$
$$\frac{\partial u_i(s_i, s_{-i})}{\partial s_i} = 0$$
$$.$$
$$\frac{\partial u_n(s_n, s_{-n})}{\partial s_n} = 0$$

To guarantee that the solution is an equilibrium, we must check the
second order conditions for each agent at the solution.

If either system of equations has multiple solutions, there are more
than one equilibrium in pure strategies. However, the absence of a solution does not necessarily imply that there are no pure strategy Nash
equilibria. Since the FOC and SOC are sufficient but not necessary,
there are many games in which the best responses do not satisfy them.
Such situations can arise for a variety of reasons. In cases where payoffs against $s_{-i}$ are either not quasi-concave or they are monotonically
increasing or decreasing over $S_i$, agent $i$'s best response will be at the
boundary of $S_i$. For such best responses, known as "corner solutions",
the FOC will typically not hold. A second situation in which the FOC
approach will not work is when the payoff functions are discontinuous
in $s_{-i}$. Since the payoffs will not be differentiable at the discontinuities,
we cannot compute the FOC or the SOC.

---

quantity 0 denotes the vector of 0's. The second order condition is the requirement
that the matrix of second derivatives be negative definite.

## 6. Application: Interest Group Contributions

To demonstrate this procedure for computing Nash equilibria, we consider a simple interest group game. Two interest groups $N = \{1, 2\}$ want to influence a government policy. Both groups know that the final policy will be a function of how much support they give to the government. The first group's most preferred policy is $0$ and the second group's most preferred policy is $1$. The government favors the policy $\frac{1}{2}$ but may be influenced by the contributions. Each group can contribute an amount $s_i \in [0, 1]$ and the final policy is given by $x(s_1, s_2) = \frac{1}{2} - s_1 + s_2$. So both groups contribute to the government simultaneously and then the government enacts the policy given by the function $x(s_1, s_2)$. The government keeps all of the contributions to buy advertisements for the next election.[8] We assume that the interest groups each have utility functions over their contribution and the final policy of the form:

$$u_1(s_1, s_2) = -(x(s_1, s_2))^2 - s_1$$
$$u_2(s_1, s_2) = -(1 - x(s_1, s_2))^2 - s_2$$

Substituting the policy function into the utility functions we attain:

$$u_1(s_1, s_2) = -(\frac{1}{2} - s_1 + s_2)^2 - s_1$$
$$u_2(s_1, s_2) = -(1 - (\frac{1}{2} - s_1 + s_2))^2 - s_2$$

The first order conditions are given by differentiation:

$$FOC_1 : 2(\frac{1}{2} - s_1 + s_2) - 1 = 0$$
$$FOC_2 : 2(1 - (\frac{1}{2} - s_1 + s_2)) - 1 = 0$$

Solving $FOC_1$ yields the best response correspondence

$$b_1(s_2) = s_2.$$

Solving $FOC_2$ yields the best response correspondence

$$b_2(s_1) = s_1.$$

Now the set of pure strategy NE to this game is infinite and it is given by:

$$\{(s_1, s_2) \in [0, 1]^2 : s_1 = s_2\}$$

---

[8]This game is closely related to the all-pay auction in economics.

This result has a very straightforward interpretation. Any pair of equivalent contributions is a Nash equilibrium. The resulting policy from such a profile of contributions is $\frac{1}{2}$. No contributor wants to unilaterally deviate because the marginal gain of an additional unit of contribution (in terms of pulling policy in the desirable direction) is exactly offset by the marginal cost of loosing another unit of resources. As in the Prisoner's dilemma, the Nash equilibria of this game are inefficient. That is, the contributors would rather commit to not giving any money to the government. However, since no such commitment is possible in the game, the fact that each contributor has an incentive to deviate from such an agreement means that this outcome is not supportable.

## 7. Application: International Externalities

Suppose that we have two countries who must decide how much to invest in pollution abatement. The strategies are a levels of investment $s_1$ and $s_2 > 0$. Each country pays a cost $c(s_i) = k_i s_i$ We will assume that $k_1 < k_2$ so that country 1 can lower pollution a given amount at a lower costs than country 2. The payoffs to the investments are based on the total investment made so that $u_i(s_1, s_2) = \sqrt{s_1 + s_2}$. Therefore, the payoffs for each country are

$$\sqrt{s_1 + s_2} - k_i s_i$$

Given these payoffs it is straightforward to get the FOC conditions:

$$\frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}} - k_1 = 0$$

$$\frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}} - k_2 = 0$$

while the SOCs are both $-\frac{1}{4}(s_1 + s_2)^{-\frac{3}{2}} < 0$. However, note that there can be no solution to the system of FOCs since $k_1 < k_2$. Suppose that $s_1 > (2k_2)^{-2}$, the left hand side of country 2's FOC is always negative. This implies that country 2's payoff are always decreasing in its investment. This is shown graphically in Figure 5.4 which plots country 2's payoffs as a function of country 1's investment. Since all investments are assumed to be non-negative, country 2's best response to this level of country 1 investment is to invest zero. Similarly, if $s_2 > (2k_1)^{-2}$, country 1's best response is zero investment. Since these "corner solutions" are part of each countries best response functions, we may not use the FOC approach.

**Insert Figure 5.4 Here**

It is easy to see the actual solution graphically.    Figure 5.5 plots the best response functions for both countries.    The vertical axis represents both country 2's strategy and its best response to country 1. Similarly, the horizontal axis represents country 1's strategy and best response.    The dotted line plots country 2's best response function which is declining in $s_1$ until $s_1 > (2k_2)^{-2}$ and then it is zero.    The solid line is the country 1's best response which also declines to zero at $s_2 = (2k_1)^{-2}$.    Clearly, the only intersection of the best responses is at the point where $s_1 = (2k_1)^{-2}$ and $s_2 = 0$.

### Insert Figure 5.5 Here

If we assumed that instead of $k_1 < k_2$ that country 2 is the low cost country, it is easy to check that the Nash equilibrium strategies would be $s_1 = 0$ and $s_2 = (2k_2)^{-2}$.    Thus, the unique Nash equilibrium of this game is one where the high cost country free rides completely on the low cost country.    Thus, the low cost country chooses optimal investment knowing that the high cost will invest nothing.

## 8.  Computing Equilibria with Constrained Optimization*

A useful alternative when a game may have corner best response functions is to use the techniques of constrained maximization such as Kuhn-Tucker programing to compute necessary conditions for Nash equilibria.    To illustrate, consider the externality game where we impose the constraint that $s_1 \geq 0$ and $s_2 \geq 0$.  We can formally incorporate these constraints with Lagrange multipliers $\lambda_1$ and $\lambda_2$ so that agent $i$ chooses $s_i$ to maximize

$$\sqrt{s_1 + s_2} - k_i s_i + \lambda_i s_i$$

Now the necessary conditions for a Nash equilibria are the first- order conditions

$$\frac{1}{2} (s_1 + s_2)^{-\frac{1}{2}} - k_1 + \lambda_1 = 0$$

$$\frac{1}{2} (s_1 + s_2)^{-\frac{1}{2}} - k_2 + \lambda_2 = 0$$

along with the "slackness" conditions

$$\lambda_1 s_1 = 0$$
$$\lambda_2 s_2 = 0$$

and the constraints:

$$\lambda_i \geq 0 \text{ for all } i$$
$$s_i \geq 0 \text{ for all } i.$$

Thus, a Nash equilibrium is a solution to the four equations that satisfy the constraints. Note that the first two equations imply that $\lambda_i = k_i - \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}}$. We can use rewrite the slack conditions as

$$s_1 \left[ k_1 - \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}} \right] = 0$$

$$s_2 \left[ k_2 - \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}} \right] = 0$$

and the constraints on $\lambda$ as

$$k_1 \geq \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}}$$

$$k_2 \geq \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}}$$

It is easy to see that there can be no solution in which both countries are making positive investments because this would require the bracketed terms of each FOC to be zero which they cannot be. It is also easy to see that $s_1 = s_2 = 0$ is not an equilibrium because the requirement that $k_i \geq \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}}$ would be violated. Suppose that $s_1 = 0$ and $s_2 > 0$. Then the second first order condition implies that $k_2 = \frac{1}{2}(s_1 + s_2)^{-\frac{1}{2}} > k_1$ which violates the non-negativity constraints on $\lambda_1$. Thus, the only possible NE involves $s_1 > 0$ and $s_2 = 0$. The first FOC implies that $s_1^* = (2k_1)^{-2}$ which exactly the result we derived in the last section.

## 9. Proving the Existence of Nash Equilibria**

In this section, we provide a rigorous proof of Nash's Theorem. But first we need to develop some additional mathematical ideas. The most important is that of a *fixed point*. Intuitively, a fixed point of a function (correspondence) is a point (set) in the domain that maps into itself in the range. Formally, given a correspondence. $c : A \to A$ a fixed point $x^* \in A$ is a point such that $x^* \in c(x^*)$. If $c(\cdot)$ is a function then a fixed point is a point $x^*$ such that $x^* = c(x^*)$. Note that since a Nash equilibrium is a strategy profile $s^*$ for which $s_i^* \in b_i(s_{-i}^*)$ for every $i \in N$, it is a fixed point of the the best response correspondence

$$b(s) = (b_1(s_{-1}), ..., b_i(s_{-i}), ..., b_n(s_{-n})).$$

Thus, proving the existence of a Nash equilibrium is simply a matter of determining whether or not a fixed point exists for a given best response correspondence. Fortunately, there is a large body of mathematics dedicated to determining the sufficient conditions for the existence of a fixed-point. A specification of such conditions is known as a *fixed point theorem*. Thus, if we can demonstrate that the properties of

the best response correspondence for a game match the conditions of an established fixed point theorem, we have proven the existence of a Nash equilibrium.

Before proceeding, we need to elaborate some of the conditions that will be required for the existence of a fixed point. The first is that we require that the correspondence be convex valued.

DEFINITION 5.11. *A correspondence $c : A \to\to A$ is convex valued if for every $a \in A$, $c(a)$ is a convex subset of $A$.*

In other words a correspondence $c(\cdot)$ is convex valued if for every $x \in A$ if $y, z \in c(x)$ then for any $\lambda \in [0, 1]$ the point $\lambda y + (1-\lambda)z \in c(x)$. Now we define the upper inverse of a set $B$ which is the set of points in the domain which the correspondence maps into subsets of $B$.

DEFINITION 5.12. *For a correspondence $c : A \to\to A$, the upper inverse of a set $B \subseteq A$, is $c^+(B) = \{x \in A : c(x) \subset B\}$.*

We can define *upper hemi-continuity,* an important property for establishing the existence of a fixed point.

DEFINITION 5.13. *A correspondence $c : A \to\to A$ is upper hemi-continuous if for every open set $O \subseteq A$ the set $c^+(O)$ is open.*

Definition 5.13 is often hard to verify. The existence of a closed graph is easier to check.

DEFINITION 5.14. *A correspondence $c : A \to\to A$ has a closed graph if for any two sequences $x^n \to x \in A$ and $y^n \to y \in A$ with $x^n \in A$ and $y^n \in c(x^n)$ for every $n$ we have $y \in c(x)$.*

The following theorem (whose proof we leave as an exercise) establishes that correspondences with closed graphs are upper hemi-continuous.

THEOREM 5.6. *If $A$ is compact a correspondence $c : A \to\to A$ is upper hemi-continuous if it has a closed graph.*

The intuition behind the closed graph condition is not difficult to see. When a correspondence has a closed graph, if we have two sequences $x^n$ and $y^n$ of points each in $A$ that converges to $x$ and $y$ both in $A$ with $y^n \in c(x^n)$ it must be the case that $y \in c(x)$. In other words for any sequence in the domain converging to $x$ and any selection of points $y^n$ that are in the image converging to a point $y$, it must be the case that the limit $y$ is in the image of the correspondence evaluated at the limit of $x$. Graphically for a correspondence that has a closed graph, the set $\{(x, y) \in A^2 : y \in c(x)\}$ is closed in the space $A^2$. For a more complete treatment of these concepts see Border (1989).

It is not difficult to see that if the correspondence $c : A \longrightarrow A$ is single-valued for every $a \in A$ then $c(\cdot)$ is a function. If a single valued correspondence is upper hemi-continuous then it is also a continuous function.

To establish Theorem 5.3, we will use the following fixed point theorem

THEOREM 5.7. *(Brouwer) Suppose $A \subset R^d$ is a compact and convex set. If $f : A \to A$ is a continuous function then $f(\cdot)$ has a fixed point in $A$.*

To establish Theorem 5.4, we will use the following fixed point theorem.

THEOREM 5.8. *(Kakutani) Suppose that $A \subset R^d$ is a compact and convex set with $c : A \longrightarrow A$ a correspondence satisfying the conditions:*
*(1) $c(x)$ is non-empty for every $x \in A$*
*(2) $c(\cdot)$ is convex valued*
*(3) $c(\cdot)$ is upper hemi-continuous*
*then $c(\cdot)$ has a fixed point in $A$.*

Several proofs of these results appear in Border. In order to establish the existence of Nash equilibria in either mixed strategies or pure strategies when the appropriate assumptions are satisfied, we will need to show that in the case of Theorem 5.7 $b(s)$ is a continuous function and in the case of Theorem 5.8 $b(s)$ is a correspondence that satisfies the conditions 1-3. A result that is useful in its own right, as well as helpful in demonstrating that $b(s)$ is non-empty and upper hemi-continuous is the Theorem of the Maximum.

THEOREM 5.9. *(Theorem of the Maximum) Let $X \subset \mathbb{R}^d, M \subset \mathbb{R}^z$ be compact and convex sets. Let $f(x, m) : X \times M \to R^1$ be continuous in $x$ and $m$ then the correspondence $c : M \longrightarrow X$ defined as*

$$c(m) = \arg \max_{x \in X}\{f(x, m)\}$$

*is non-empty and upper hemi-continuous.*

The fact that the set of optimal choices is non-empty is interesting on its own. This result was stated in a previous section. The fact that the correspondence $c(\cdot)$ defined in the theorem is upper hemi-continuous has the following interpretation. Calling the vector $m$ a parameter vector of the optimization problem, if we consider a sequence of parameter vectors $m^n$ converging to $m$ then for any selection of optimal policies $x^n \in c(m^n)$ that converge to $x$ it will be the case that $x \in c(m)$.

We can now prove Theorem 5.3.

PROOF OF THEOREM 5.3. Assume that: $S_i$ is a convex subset of $\mathbb{R}^d$ for each $i \in N$ and for each $i \in N$, $u_i(s) : S \to \mathbb{R}^1$ is continuous and for each $s'_{-i} \in S_{-i}$, $u_i(s_i, s'_{-i})$ is strictly quasiconcave in $s_i$. By the Theorem of the Maximum the correspondence $b_i(s_{-i}) : S_{-i} \to\to S$ defined as

$$b_i(s_{-i}) = \arg \max_{s_i \in S_i} \{u_i(s_i, s_{-i})\}$$

is non-empty and upper hemi-continuous. By theorem 4 of chapter 2, $b_i(s_{-i})$ is a singleton for every $s_{-i} \in S_{-i}$. The fact that an upper hemi-continuous correspondence that is single valued is a continuous function (see Exercise 5.7) implies that $b_i(s_{-i})$ is a continuous function from $S_{-i}$ into $S_i$ for each $i \in N$. We now construct the function

$$b(s) : S \to S$$

by defining $b(s_1, ..., s_n) = (b_1(s_{-1}), ...., b_n(s_{-n}))$. Since $s_{-i}(s)$ is a projection it is continuous. Since $b_i(\cdot)$ is continuous and the composition of continuous functions is continuous, and the product of continuous functions is continuos in the product space the function $b(s)$ is continuous. By Brouwer's fixed point theorem this mapping has a fixed point, $s^* = b(s^*)$. But this means that for every $i \in N$, $b_i(s^*_{-i}) = s^*_i$ and thus $s^*$ is a NE.∎  □

The proof of Theorem 5.4 is similar. We will show that for any finite game $\Gamma$ in the mixed extension game $\Gamma^m$ the best response correspondence satisfies the conditions of Kakutani's fixed point theorem.

PROOF OF THEOREM 5.4. Assume that in the game $\Gamma$, $S_i$ is finite for each $i \in N$. This implies that in the mixed extension $\Gamma^m$ $\Delta_i$ is a compact and convex subset of a finite dimensional Euclidean space. By definition 5.9 we can see that $U(\sigma_i, \sigma_{-i})$ is linear and therefore continuous in $\sigma$. Letting $b_i(\sigma_{-i}) : \Delta_{-i} \to\to \Delta_i$ be defined as

$$b_i(\sigma_{-i}) = \arg \max_{\sigma_i \in \Delta_i} \{U(\sigma_i, \sigma_{-i})\}$$

the Theorem of the Maximum implies that this correspondence is non-empty for every $\sigma_{-i} \in \Delta_{-i}$ and upper hemi-continuous. Since $U(\sigma_i, \sigma_{-i})$ is linear for any $\sigma_{-i}$ and any two $\sigma'_i, \sigma''_i$ if $U(\sigma'_i, \sigma_{-i}) = U(\sigma''_i, \sigma_{-i})$ we have $U(\lambda \sigma'_i + (1 - \lambda)\sigma''_i, \sigma_{-i}) = U(\sigma'_i, \sigma_{-i})$ so $b_i(\sigma_{-i})$ is convex valued. Combining these facts we see that the correspondence

$$b(\sigma) : \Delta \to \Delta$$

defined as $b(\sigma_1, ..., \sigma_n) = (b_1(\sigma_{-1}), ...., b_n(\sigma_{-n}))$ satisfies the requirements of Kakutani's fixed point theorem. Thus there is a mixed strategy profile satisfying the condition $\sigma^* \in b(\sigma^*)$. Such a profile is a NE to $\Gamma^m$ and thus a mixed strategy NE to $\Gamma$.∎                                    □

## 10. Strategic Complementarity

Consider two candidates 1 and 2 running for office. Each candidate selects a level $s_i > 0$ of campaign effort. The opportunity cost of this effort is $\beta_i s_i$ where $\beta_i > 0$. We assume that the election outcome is a probabilistic function of $s_1$ and $s_2$. Let $\pi(s_1, s_2)$ denoting the probability that candidate 1 wins so that candidate 2 wins with probability $1 - \pi(s_1, s_2)$. It is reasonable to think that this probability function is increasing in $s_1$ and decreasing in $s_2$. In formulating a reasonable model of campaigns, we confront an important question: should candidate 1's best response be increasing or decreasing in candidate 2's effort? It seems natural to think that candidate 1 will select higher levels of $s_1$ in response to higher levels of $s_2$ and vice versa. If this is the case, we say that the two choice variables are *strategic complements.* Games with strategic complementarity are among a class of games known as *supermodular* and are particularly amenable to equilibrium and comparative static analysis. In this section we sketch the intuition behind the analysis of such games using an example. We leave the technical details for the subsequent section.

Kahn and Kenney's (1999) study of Senate elections posits that competitiveness is an important factor in determining how much media coverage campaign activities will generate. The competitiveness shapes the way that voters respond to campaigns while campaigns influence the competitiveness of a race. Absent high competitiveness, the media and voters tune-out and thus the marginal value of advertisements or speeches is low. On the other hand when the competitiveness is high these messages have large effects. Kahn and Kenney's theory portrays the media as an mechanism that determines competitiveness as a function of campaigning. Consistent with this interpretation, higher levels of campaigning are likely to result in higher states of competitiveness. In highly competitive races advertisements and speeches are more influential. The feedback loop that Kahn and Kenney discuss suggests that candidate effort levels may indeed be strategic complements. If so, each of the best responses are strictly increasing functions. Figure 5.6 depicts a game of this form with a unique equilibrium point $(s_1^*, s_2^*)$.

**Insert Figure 5.6**

While pictures with multiple equilibria could be drawn, it is not difficult to see that the equilibria will all be completely ordered such that if $(s_1^*, s_2^*)$ and $(s_1^{**}, s_2^{**})$ are two Nash equilibria and $s_1^* > s_1^{**}$ then $s_2^* > s_2^{**}$ (and if $s_2^* > s_2^{**}$ then $s_1^* > s_1^{**}$). To see that this must be true, try drawing a counterexample while maintaining the assumption that the best responses are strictly increasing.

We can take Kahn and Kenney's hypothesis about competitiveness further. Suppose at the beginning of a campaign, there is an exogenous level of competitiveness. Sources of variation might include the importance of office, the media environment, and the attentiveness of the voters. So consider two electoral environments where one is more competitive than the other. For a fixed level of $s_i$, we would expect that $b_{-i}(s_i)$ is higher in the more competitive race. In Figure 5.6, $b_i'(s_{-i})$ denotes the best responses for the more competitive while $b_i(s_{-i})$ represents those of the less competitive election. The figure demonstrates that the equilibrium campaign efforts of the more competitive race $(s_1'^*, s_2'^*)$ are higher than the equilibrium efforts of the less competitive race, $(s_1^*, s_2^*)$. Thus, we can conclude that competitiveness generates more campaigning.

It is crucial to note that this comparative static result was generated by nothing more than the assumption of strategic complementarity. As long as both best responses are increasing, a common upward shift in the best responses must lead to a higher intersection point. It is easy to see that the effect of competitiveness would be ambiguous if one of the players had a downward sloping best response Furthermore, as long as candidate efforts are strategic complements, the result does not hinge on any further assumptions about functional forms or player preferences. It also obviates the need for the continuous and differentiable best response functions required for use of the implicit function theorem.

The trick, of course, is to determine what model primitives are consistent with complementarity. In the following technical section, we present the underlying theory behind supermodular games and present some sufficiency conditions for the existence of equilibria exhibiting monotone comparative statics.

## 11. Supermodularity and Monotone Comparative Statics*

Recall that Nash's theorem relies on the continuity of best responses and compactness of strategy sets so that Kakutani's fixed point theorem ensures the existence of equilibria. Unfortunately, these existence

requirements often require very strong assumptions about the primitives of the model. However, we can use a different fixed point theorem to establish existence so long as the best responses satisfy certain monotonicity conditions, such as the ones discussed in the last section. As a bonus, the comparative statics analysis of the equilibria set is immediate when the best responses satisfy these conditions. These analyses are much simpler and straightforward than those based on the implicit function theorem. In this section, we summarize several results on a class of games known as supermodular games. Of particular interest is the existence of equilibria with discontinuous best response correspondences and direct comparative static results.

In principal the concepts of this section can be developed for any partial orderings, with the goal being the attainment of results about monotonicity relative to the partial order. To keep things concrete we will consider only the natural partial ordering $\geq$ on sets contained in $\mathbb{R}^n$. For any two numbers $x, y \in \mathbb{R}^1$ the *join* is $x \vee y = \max\{x, y\}$ and the *meet* is $x \wedge y = \min\{x, y\}$. For any two vectors in $x, y$ in $\mathbb{R}^n$ we have $x \vee y = (\max\{x_1, y_1\}, ..., \max\{x_n, y_n\})$ and $x \wedge y = (\min\{x_1, y_1\}, ...., \min\{x_n, y_n\})$. A set which contains all of the joins and meets of its elements is called a lattice.

DEFINITION 5.15. *A set $A$ is a lattice if for each $x, y \in A$ we have* $x \vee y \in A$ *and* $x \wedge y \in A$

Note that intervals and the products of intervals are lattices, but sets like $\{x \in \mathbb{R}^3 : x_1 + x_2 + x_3 \leq 1\}$ are not lattices.[9] Intuitively, squares (products of intervals) are lattices but triangles (simplices) are not.[10] Typically we are interested in single agent or multi-agent decision theory problems in which there is a set of choice variables $X$ and a set of exogenous parameters $P$, both of which are lattices. The objective function depends on both types of variables, $f(x, p) : X \times P \to \mathbb{R}^1$. The key condition that we now focus on is called supermodularity.

DEFINITION 5.16. *The function $f(\cdot, \cdot) : X \times P \to \mathbb{R}^1$is supermodular in $(x, p)$ if for all $z, z' \in X \times P$, $f(z) + f(z') \leq f(z \vee z') + f(z \wedge z')$*

While verification of this condition may be difficult, a more intuitive condition, *increasing differences,* is often easily verified.

---

[9]But see page 14 of Topkis 1998 for a dicusssion of translations that convert sets like the simplex into lattices.

[10]Although, we do not consider examples where the sets are discrete (or just not convex) the definition of a lattice and subsequent results can be readily applied to non convex sets.

DEFINITION 5.17. *The function $f(\cdot,\cdot) : X \times P \to \mathbb{R}^1$ has increasing differences in $(x,p)$ if for all $p \le p'$ and $x \le x'$, $f(x',p') - f(x,p') \ge f(x',p) - f(x,p)$.*

A function with increasing differences exhibits a greater marginal effect of $x$ when $p$ is higher. In other words increasing difference formalizes the idea of complementarity. Since increasing differences is easier to interpret in terms of the substance of models than is supermodularity, it is convenient that the following equivalence holds.

PROPOSITION 5.1. *(Topkis) The function $f(\cdot,\cdot) : X \times P \to \mathbb{R}^1$ has increasing differences in $(x,p)$ if and only if it is supermodular in $(x,p)$.*

In the case of a twice differentiable function $f(\cdot,\cdot) : X \times P \to \mathbb{R}^1$, supermodularity coincides with the condition

$$\frac{\partial^2 f}{\partial z_1 \partial z_2} \ge 0$$

for any coordinates $z_1$ and $z_2$ of $X$ or $P$ or both. Supermodularity is preserved under several types of operations.

PROPOSITION 5.2. *Suppose that $X$ is a lattice then*
*(a) if $f(x)$ is supermodular on $X$ and $\alpha > 0$ then $\alpha f(x)$ is supermodular on $X$.*
*(b) if $f(x)$ and $g(x)$ are supermodular on $X$ then $f(x) + g(x)$ is supermodular on $X$.*
*(c) if $f_t(x)$ is supermodular on $X$ for $t = 1,2,...$ and $f(x) = \lim_{k \to \theta} f_k(x)$ for each $x \in X$ then $f(x)$ is supermodular in $X$.*
*(d) if $F(\omega)$ is a distribution function on a set $\Omega$ and $g(x,\omega)$ is supermodular on $X$ for each $\omega \in \Omega$ then $f(x) = \int_\Omega g(x,\omega)dF(\omega)$ is supermodular on $X$.*

We first present the implications of supermodularity in the simple case of an individual optimization problem. Supermodularity implies that the optimal solution has monotone comparative statics. By this we mean that each element of $\arg\max_{x \in X} f(x,p)$ is an increasing function of $p$. Specifically we have the following result.

PROPOSITION 5.3. *(Topkis) If $f(\cdot,\cdot) : X \times P \to \mathbb{R}^1$ is supermodular in $(x,p)$ and $g(p) = \arg\max_{x \in X} f(x,p)$ then each component of $g(p)$ is increasing in $p$ where $g(p)$ exists.*

Of course nothing in the theorem says that an optimal solution exist. Recall that 2.3 establishes the non-emptiness of the maximal set for lower continuous orderings on compact sets. Lower continuity

of preferences corresponds to what is called upper semicontinuity of objective functions.

DEFINITION 5.18. *A function $f(x) : X \to \mathbb{R}^1$ is upper semicontinuous at $x_0 \in X$ if for all $\varepsilon > f(x_0)$ there is some $\delta > 0$ such that for all $x \in B(x_0, \delta)$ (an open ball around $x_0$ with radius $\delta$) $\varepsilon \geq f(x)$. A function is upper semicontinuous on $X$ if it is upper semicontinuous at every $x_0 \in X$.*

Upper semicontinuity at $x_0$ requires that $\liminf_{x \to x_0} f(x) \leq f(x_0)$. Using the result of Exercise 5.11, Theorem 2.3 and restating Topkis' result leads to the conclusion.

PROPOSITION 5.4. *If $f(\cdot, \cdot) : X \times P \to \mathbb{R}^1$ is supermodular in $(x, p)$ and upper semicontinuous then $g(p) = \arg\max_{x \in X} f(x, p)$ is non-empty for every $p \in P$ and each component of $g(p)$ is increasing in $p$.*

In our argument for existence with continuity, we used properties of the objective functions to establish properties of the best responses and then applied a fixed point theorem to the best response correspondence. Analogously Tarsky's fixed point theorem allows us to use monotonicity of best responses to establish existence of Nash equilibria. This and the previous proposition then lead to the conclusion that a game with supermodular and upper semicontinuous utility functions has a Nash equilibria. We first present Tarsky's theorem.

PROPOSITION 5.5. *(Tarsky) If $f(x)$ is an increasing mapping from a compact lattice $X$ into itself then there exists at least one fixed point $x^*$ such that $f(x^*) = x^*$.*

### Insert Figure 5.7

In Figure 5.7 we demonstrate the case where $X$ is one dimensional. As long as the function $f(x)$ is increasing the presence of discontinuities does not enable us to skip all of the crossings of $f(x)$ and the dotted $45°$ line. Just as Kakutani's theorem generalized Brouwer's to the case of correspondences, Zhou (1994) represents a generalization of Tarsky to the case of correspondences. We do not present this result as much additional notation is needed, but instead present directly the result for supermodular games. A regular supermodular game is a normal form game $\langle N, \{S_i, u(\cdot, ..., \cdot)\}_{i \in n} \rangle$ in which $S_i$ is a compact lattice for each $i \in N$ and $u(\cdot, ..., \cdot)$ is supermodular and upper semicontinuous for each $i \in N$.

PROPOSITION 5.6. *A regular supermodular game has at least one pure strategy Nash equilibrium and the set of such equilibria is a lattice with a smallest and biggest equilibrium profile.*

Why do we care about the existence of a smallest and biggest equilibrium? Our second interesting result for supermodular games deals with comparative statics of these particular equilibria. To motivate the finding we return to the case of a function $f(x)$ on $\mathbb{R}^1$.

**Insert Figure 5.8**

In Figure 5.8 we depict an increasing function with four fixed points $f(x) = x$. The dotted curve $f'(x)$ represents the result of shifting $f(x)$ up. Note that the biggest and smallest fixed points shift to the right (get larger) when the function is shifted up. In contrast some of the fixed points move the other way. The smallest and biggest fixed points tend to be the result of the function crossing the 45° line from above and behave the same way. This intuition forms the basis for the following comparative static result.

PROPOSITION 5.7. *In a regular supermodular game the smallest and biggest equilibrium profiles are increasing in p.*

As Figure 5.8 suggests other equilibria may behave differently. Echenique (2002) shows that in games like this with complementarities any selection of equilibria that does not also exhibit the monotone comparative static of the biggest and smallest equilibria will be unstable under a large range of adaptive dynamics. In other words there are clear reasons to select equilibria which exhibit the comparative static described in the proposition.

Returning to the discussion of Kahn and Kenney's competitiveness hypothesis, suppose that candidate 1 wishes to maximize $\pi(s_1, s_2, c) - \beta_1 s_1$ and candidate 2 wishes to maximize $1 - \pi(s_1, s_2, c) - \beta_2 s_2$ where $\pi(s_1, s_2, c)$ depends on an exogenous level of competitiveness $c \in \mathbb{R}^1$. The fact that campaign activity influences competitiveness and that voters pay closer attention in more competitive races suggests that the incremental effect of $s_1$ on $\pi(s_1, s_2, c)$ should be higher when $s_2$ and/or $c$ are higher. A symmetric argument follows for the incremental effect of $s_2$. Accordingly the assumption that payoffs are supermodular is consistent with (if not implied by) Kahn and Kenney's explanation. Assuming compactness of the choice sets and upper semicontinuity of $\pi(s_1, s_2, c)$ in $s_1$ and $-\pi(s_1, s_2, c)$ in $s_2$ allows us to conclude (from Propositions 5.6 and 5.7) that equilibria exist and in the biggest and smallest equilibria $s_i$ is increasing in $c$. This may seem like a pretty strong conclusion to reach without having to specify very much about the function $\pi(\cdot, \cdot, \cdot)$. The power of monotone comparative statics is that they tend to isolate the minimal structure that is needed for a particular comparative static result.

For a thorough review of results for supermodular games and monotone comparative statics see Topkis (1998). A focused summary of results and applications to political science appear in Ashworth and Bueno de Mesquita (2004).

## 12. Refining Nash Equilibria

Recall the majority rule voting game of section where we showed that any profile which no single voter is pivotal is a Nash equilibrium. We justified ignoring these equilibria on the basis that they involved playing weakly dominated strategies. However, the elimination of weakly dominated strategies is not the only way that we can justify "refining" the set of Nash equilibria in this game. Suppose instead of assuming that every agent is capable of playing her best response with probability one, we assume that there some small probability that each player will "tremble" and play another strategy. To keep things very simple, lets look at a three player version of the game where 2 voters prefer $D$ and one prefers $R$. Formally, we will assume that each player must play each pure strategy with a small probability $\varepsilon$ which must be less than $\frac{1}{2}$.[11] This captures the idea that mistakes ensure that all strategies are played with a least a minimal probability.

Clearly it is a best response for each agent to attempt to maximize the probability of their preferred candidate winning which is the same as minimizing the probability that the least desired candidate wins. So we need to compute the probability that each candidate wins under various combinations of strategies. The probabilities of a $R$ victory are:

$$\Pr\left(R|3 \text{ attempted votes for } R\right) = (1-\varepsilon)^3 + 3\left(1-\varepsilon\right)^2 \varepsilon$$

$$\Pr\left(R|2 \text{ attempted votes for } R\right) = (1-\varepsilon)^3 + (1-\varepsilon)^2 \varepsilon + 2\left(1-\varepsilon\right)\varepsilon^2$$

$$\Pr\left(R|1 \text{ attempted votes for } R\right) = \varepsilon^2\left(1-\varepsilon\right) + 2\left(1-\varepsilon\right)^2 \varepsilon + \varepsilon^3$$

$$\Pr\left(R|0 \text{ attempted votes for } R\right) = \varepsilon^3 + 3\varepsilon^2(1-\varepsilon)$$

The reader should verify that since $\varepsilon < \frac{1}{2}$, the probability that a Republican wins is strictly increasing in the number of intended votes.

First, we consider the "bad" equilibrium where $R$ wins with unanimously. Does this outcome survive in the presence of trembles i.e. does each voter vote $R$ with the maximal probability $1 - \varepsilon$? Clearly, the $R$ preferring voter will maximize the probability of an $R$ victory by voting $R$ with probability $1 - \varepsilon$. However, consider the choice of a $D$

---

[11]The idea that voters might not vote for the candidate that they intended has a renewed substative importantance after the 2000 presidential election.

voter? She has the choice of conforming the equilibrium by intending to vote $R$ in which case $R$ wins with probability $(1 - \varepsilon)^3 + 3(1 - \varepsilon)^2 \varepsilon$ or defecting to an intended vote of $D$ in which case $R$ wins $(1 - \varepsilon)^3 + (1 - \varepsilon)^2 \varepsilon + 2(1 - \varepsilon)\varepsilon^2$. This defection reduces the probability that $R$ wins by $2(1 - \varepsilon)^2 \varepsilon - 2(1 - \varepsilon)\varepsilon^2 > 0$. Thus, the $D$ voter will prefer the defection. It is easy to show that the equilibrium corresponding where all voters vote $D$ also does not survive trembles.

This concept of refining the set of equilibria by focusing on those that are robust to small mistakes by the agents was first developed by Selten. He named such equilibria *perfect*. A formal definition of perfect equilibria follows:

DEFINITION 5.19. *An "$\varepsilon$-constrained" equilibrium is a totally mixed strategy profile $\sigma^\varepsilon$ such that for each player $i$, $\sigma_i^\varepsilon$ solves $max_{\sigma_i} u_i(\sigma_i, \sigma_{-i}^\varepsilon)$ subject to $\sigma_i(s_i) \geq \varepsilon$. A perfect equilibria is the limit of an $\varepsilon$-constrained equilibrium as $\varepsilon$ goes to 0.*

It is to see how the unanimous voting equilibria fail to meet this definition since in the $\varepsilon$-constrained equilibria all players vote for their least preferred candidate with probability $\varepsilon$. Thus, the limit of these equilibria are strategy profiles which place zero probability on voting for the lesser candidate. Indeed, the only perfect equilibrium is the one where all voters vote for their favorite candidate.

Perfect equilibria have some desirable properties. First, all perfect equilibria are Nash equilibria of the game without the $\varepsilon$ constraints. Thus, the set of perfect equilibria are a proper subset of the set of Nash equilibria. Secondly, it can be shown using arguments similar to the existence of Nash equilibria that all finite game normal form games have at least one perfect equilibria.[12]

## 13. Application: Private Provision of Public Goods

Since the publication of Mancur Olson's *Logic of Collective Action*, a central question in political science has been the conditions under which individual rational agents would be willing to incur personal costs to contribute to the public good. In this section, we present a game theoretic model of such contributions. This model is based on the work of Palfrey and Rosenthal (1984).

Assume that there are $n$ agents who must decide whether to provide a public good. Provision requires the contribution of a single individual and produces a utility of 1 unit for each individual. However, any

---

[12]The proof follws from the fact that the $\varepsilon$-constrained mixed strategy space is compact, convex, and non-empty.

contributor pays a cost $c < 1$. Thus, the strategy set for each agent is $\{contribute, \, don't \, contribute\}$. To simplify notation, we define "contribute" as $s_i = 1$ and "don't contribute" as $s_i = 0$. Since we will also consider mixed strategy equilibria, let $\sigma_i$ be the probability that agent $i$ contributes. Given this setup, the payoff for each agent is $1 - c$ if she contributes, 1 is she doesn't contribute but some other agent does, and 0 otherwise.

First, consider the set of pure strategy equilibria. It is easy to check that for each agent $i$ there is an equilibrium where $s_i = 1$ and $s_{-i} = 0$. Agent $i$ receives $1 - c$ while all other agent receive a utility of 1. Clearly, agent $i$ will not defect since failing to contribute will lower her payoff to 0. Similarly, no other agent with defect to contributing since it will lower his payoff from 1 to $1 - c$. Now we can check to see that no other combination of strategies is a pure strategy Nash equilibrium. First, consider $s_i = 0$ for all $i$. In this case, any agent would do better by contributing so that this profile cannot be an equilibrium. Next consider a strategy combination where more than one agent contributed to the public good. Clearly, any contributor would increase her utility by withholding the contribution since it would be provided by the contribution from another agent.

It is important to note that the equilibria to this game are quite different from the decision theoretic predictions of Olson. In fact, all of the equilibria to this game are Pareto efficient since that the public good is provided by just enough contributions. However, there are many reasons to think that this equilibrium is not a very valid description of how this game would actually be played. Most importantly, since there are so many Nash equilibria, how would the agents ever coordinate on one of them? Secondly, each of the pure strategy Nash equilibria involves ex ante identical agents playing different strategies. An equilibrium where identical agents played identical strategies would seem more plausible. For these reasons, a more plausible equilibrium is a mixed strategy equilibrium where every agent contributes with a probability determined in equilibrium. Since there may be many such equilibria, we will only consider the symmetric mixed strategy equilibrium where $\sigma_i = \sigma$ for all $i$. This restriction is consistent with our criticism of the asymmetry inherent in the pure strategy Nash equilibria.

Recall that for a mixed strategy to be a best response for agent $i$, she must be indifferent among the pure strategies that she mixes over. Thus, if agent $i$ is willing to play $\sigma_i$ against $n - 1$ players contributing with probability $\sigma$, it must be true that $u_i \left( contribute, \sigma \right) =$

$u_i \left( not\ contribute, \sigma \right)$ or

$$1 - c = 1 - \left( 1 - \sigma \right)^{n-1}$$

Using this condition, we can solve for the equilibrium value: $\sigma = 1 - c^{\frac{1}{n-1}}$. Note that in some ways this equilibrium is more consistent with the predictions of Olsen. As $n$ gets large, $\sigma$ goes to 0 while if $c$ goes to zero $\sigma$ goes to 1.

**13.1. Multiple Contributions.** Now we will consider an extension of this model where the public good is provided only if $k$ out of the $n$. It is easy to see from the logic presented in the previous section that there are many pure strategy equilibria where exactly $k$ agents make contributions. Clearly, in such an equilibrium, any agent not contributing has no incentive to defect and make a contribution as it will cost $c$ and not change the probability of obtaining the public good. Conversely, any contributor who defected would cause the good to not be provided. Since saving the contribution cost is less valuable than losing the public good, such a defection will not occur. Therefore, there is an equilibrium corresponding to contributions by every possible combination of $k$ agents. From a basic result in combinatorics, we know that there are exactly

$$\left( \begin{array}{c} n \\ k \end{array} \right) \equiv \frac{n!}{k!\,(n-k)!}$$

distinct Nash equilibria where $n! = 1 \cdot 2 \cdot 3... \cdot n$. The notation $\left( \begin{array}{c} n \\ k \end{array} \right)$, known as the binomial coefficient, represents the number of combinations of $k$ elements drawn from $n$ objects.

The plausibility of these pure strategy equilibria are perhaps even less compelling than the pure strategy equilibria of the one contribution game. So again we will compute the symmetric mixed strategy equilibria for this game.

Let $x_{-i}$ represent a realization of $\sigma_{-i}$ so that it is number of contributions made agents other than $i$. The payoff to $i$ from contributing is

$$\Pr \left( x_{-i} < k - 1 \right) \cdot 0 + \Pr \left( x_{-i} \geq k - 1 \right) \cdot 1 - c$$

while the payoff from abstaining is

$$\Pr \left( x_{-i} < k \right) \cdot 0 + \Pr \left( x_{-i} \geq k \right) \cdot 1$$

As before, playing mixed strategies requires that the agents be indifferent among the pure strategies in the mixture. Equating these payoffs

and doing a bit of algebra, the necessary condition for mixing is

$$\Pr\left(x_{-i} = k - 1\right) = c$$

This condition has a nice intuitive interpretation. Clearly, contributing only has a positive benefit if exactly $k - 1$ other agents contribute. Thus, the payoff from the contributing to the public good is discounted by the probability that a contribution is pivotal. Because agents will be mixing, this expected benefit has to be equated to the contribution costs $c$.

Since we are assuming that all agents are independently playing the same mixed strategy $\sigma$, we can compute the exact value of $\Pr\left(x_{-i} = k - 1\right)$ as it equals the probability of obtaining exactly $k-1$ successes in $n-1$ trials with a success probability of $\sigma$. Thus, a standard result in probability theory implies that

$$\Pr\left(x_{-i} = k - 1\right) = \binom{n-1}{k-1} \sigma^{k-1} \left(1 - \sigma\right)^{n-k}$$

Therefore, computing the symmetric mixed strategy equilibrium involves find the set of $\sigma$ that solve:

$$\binom{n-1}{k-1} \sigma^{k-1} \left(1 - \sigma\right)^{n-k} = c$$

or

$$\sigma^{k-1} \left(1 - \sigma\right)^{n-k} = \frac{(k-1)!\,(n-k)!}{(n-1)!} c$$

Before characterizing exactly what the set of solutions looks like, it is worthwhile to look at a couple of examples. First, suppose that $n = 5$ and $k = 3$. Then the equation reduces to $\sigma^2 \left(1 - \sigma\right)^2 = \frac{1}{6}c$. The solid lines of Figure 5.9 plots the left and right sides of this equation. Note that as long as $c$ is sufficiently low, there are two solutions to the equation which represent distinct mixed strategy equilibria, $\sigma_L^* < \frac{1}{2} < \sigma_H^*$. It is easy to see how the equilibrium mixtures change as a function of $c$. The effect of increasing $c$ is to raise $\sigma_L^*$ and lower $\sigma_H^*$.

### Insert Figure 5.9

Figure 5.9 also plots the conditions for $k = 4$ and $k = 5$ which are given by $\sigma^3 \left(1 - \sigma\right) = \frac{1}{4}c$ and $\sigma^4 = c$ respectively. The case of $k = 4$ is similar to $k = 3$ in that it also has two mixed strategy equilibria. However, $\sigma^3 \left(1 - \sigma\right) > \sigma^2 \left(1 - \sigma\right)^2$ if $\sigma > \frac{1}{2}$. This effect plus the fact that $\frac{1}{4}c > \frac{1}{6}c$ implies that $\sigma_H^*$ is increases in $k$. Since $\sigma^3 \left(1 - \sigma\right) < \sigma^2 \left(1 - \sigma\right)^2$ if $\sigma < \frac{1}{2}$, $\sigma_L^*$ also increases. In the case of $k = 5$, the fact that $\sigma^4$ is an increasing function for $0 \leq \sigma \leq 1$ implies that there can

only be one mixed equilibrium. It has a higher contribution probability than $\sigma_L^*$ when $k = 4$ which increases in $c$.

Many of the implications of the examples generalize. First regardless of $n$ and $k$, there can be at most two mixed strategy equilibria. To see this, let $c(\sigma)$ be the level of costs that supports $\sigma$ as the equilibrium mixing strategy or

$$c(\sigma) = \binom{n-1}{k-1} \sigma^{k-1} (1-\sigma)^{n-k}$$

Differentiating with respect to $\sigma$, we find that

$$c'(\sigma) = \binom{n-1}{k-1} [(k-1)(1-\sigma) - (n-k)\sigma] \sigma^{k-2} (1-\sigma)^{n-k-1}$$

It is easy to see that if $k < n$ the function $c(\sigma)$ is single peaked since

$$c' \gtreqqless 0 \text{ if } \sigma \lesseqqgtr \frac{k-1}{n-1}$$

Thus, if $k < n$ and $c < c\left(\frac{k-1}{n-1}\right) \equiv c_{\max}$, there will be two mixed strategy equilibria, but none if $c > c_{\max}$.[13] If $k = n$ or $k = 1$, there will be one mixed strategy equilibrium so long as $c < 1$.

We can also use this result to say something about what happens as $n$ gets very large. Since $c(0) = 0$, we know that that $c_{\max}$ goes to zero as $n$ gets large. Thus, in the limit there are no mixed strategy equilibria.

For $\sigma_H^* > \frac{k-1}{n-1}$, contribution probabilities fall in $c$ and $n$ and increase in $k$. For $\sigma_L^* < \frac{k-1}{n-1}$, the contribution probabilities fall increase in $c$ and $n$. For $\sigma_L^* < \frac{k-1}{n}$, contributions are falling in $k$ while they increase in $k$ if $\sigma_L^* \in \left(\frac{k-1}{n}, \frac{k-1}{n-1}\right)$.[14]

## 14. Exercises

EXERCISE 5.1. *In the Hotelling model, show that there is no equilibrium in pure strategies if there are three parties for any specification of the parties objectives.. What is the Nash equilibrium with four parties if parties maximize vote share?*

EXERCISE 5.2. *Show that the prisoner's dilemma has no mixed strategy equilibrium.*

---

[13]Of course in the unlikely event that $c = c_{\max}$, there will be a single symmeteric mixed stratgey equilibrium.

[14]Some derivations useful in proving these claims is available in the appendix to Palfrey and Rosenthal (1988).

EXERCISE 5.3. *Find the mixed strategy equilibrium of the following game:*

$$
\begin{array}{c|cc}
1\backslash 2 & L & R \\
T & 2,1 & 0,2 \\
B & 1,3 & 3,0
\end{array}
$$

EXERCISE 5.4. *Verify that the two definitions of Nash Equilibrium are equivalent. Hint show that if a strategy profile satisfies the first then it must satisfy the second, and then that if it satisfies the second it must satisfy the first.*

EXERCISE 5.5. *Prove Theorem 5.5.*

EXERCISE 5.6. * Prove Theorem 5.6.*

EXERCISE 5.7. *Show that an upper hemi-continuous correspondence that is single valued is a continuous function.*

EXERCISE 5.8. *Characterize the pure strategy Nash equilibria to the International Externality game when $k_1 = k_2$.*

EXERCISE 5.9. *Verify that if $f(\cdot, \cdot)$ has increasing differences in $(x, p)$ then for all $p \leq p'$ and $x \leq x'$ $f(x, p') - f(x, p) \geq f(x', p') - f(x', p)$.*

EXERCISE 5.10. *Prove parts a-c of Proposition 5.2.*

EXERCISE 5.11. *Assume that $X$ is a compact subset of $\mathbb{R}^n$ and $R$ is a lower continuous partial order on $X$. Show that if $u(x)$ is a utility function that represents $R$ on $X$ then $u(\cdot)$ is upper semi-continuous on $X$. Now show that if $u(x)$ is upper semi-continuous on $X$ then any preference relation that it represents is lower continuous.*

EXERCISE 5.12. *Consider the Palfrey-Rosenthal contribution game. Construct an asymmetric Nash equilibrium where $l$ agents contribute ($\sigma_i = 1$), $m$ agents do not contribute ($\sigma_i = 0$), and $n - m - l$ agents choose a mixed strategy $\sigma_i = q \in (0, 1)$. Show that if $l > 0$ or $m > 0$, $q^*$ is unique. Is this a stable equilibrium?*

EXERCISE 5.13. *Consider an extension of the Palfrey-Rosenthal model where contributions are refunded if the public good is not provided (i.e. fewer than $k$ contributions are made). Characterize the pure strategy and mixed strategy equilibria of this game.*

## CHAPTER 6

# Bayesian Games in the Normal Form

In the normal form games considered above there was no uncertainty. The agents know their own payoffs and those of their opponents. More precisely, the entire structure of the game involves what is called *common knowledge* –each player knows all the details of the game and each player knows that each player knows the details of the game, and each player knows that each player knows that each player knows ..... *ad infinitum*. However, in many settings this assumption is inappropriate and we might suspect that interesting strategic incentives are created by uncertainty. Returning to the "Terrorist Hunt" game,

| Table 6.1: The Terrorist Hunt | | |
|---|---|---|
| FBI\CIA | *Kingpin* | *Operative* |
| *Kingpin* | 2,2 | 0,1 |
| *Operative* | 1,0 | 1,1 |

We might think that the CIA is not sure if the FBI prefers arresting Operatives to not arresting anyone. The loss of an operative may not lower terrorism risks but involve dramatic administrative headaches. So the CIA may think that it is possible that the game looks like

| Table 6.2: Modified Hunt 1 | | |
|---|---|---|
| FBI'\CIA | *Kingpin* | *Operative* |
| *Kingpin* | 2,2 | 0,1 |
| *Operative* | 0,0 | 0,1 |

But suppose that the CIA fears that the FBI may have yet a different preference ordering, preferring to arrest Operatives more than Kingpins, because Operatives usually fold under the pressure providing lots of information, whereas Kingpins remain silent. In this case the game might look like

| Table 6.3: Modified Hunt 2 | | |
|---|---|---|
| FBI″\CIA | *Kingpin* | *Operative* |
| *Kingpin* | 1,2 | 0,1 |
| *Operative* | 2,0 | 2,1 |

In this setting the CIA's calculation of which strategy to play may differ. We saw that in the original game the pure strategy Nash equilibria are (*Kingpin, Kingpin*) and (*Operative, Operative*). But if the CIA thinks that the FBI is possibly of one of these alternative types then it is less clear which strategy the CIA should play. In this case we cannot apply our notion of Nash equilibria or dominance to the game.

In this section we develop tools to analyze richer models involving agents that do not know the payoffs of the other players. This feature is termed **incomplete information**. The standard practice (originated by Harsanyi 1967-68) is to convert such a game into one where a fictional player (nature) moves first drawing the utility functions of the agents from a probability distribution that is known to the players. Following this draw, agents simultaneously select their actions. This approach is called one of **imperfect information.** So in the modified "Terrorist Hunt" we might think that Nature selects the preferences (or type) of the FBI by tossing a fair three sided die –making each type equally likely. The FBI knows its type and chooses an action. The CIA does not know the FBI's type and chooses its action. In this case, specifying strategies for the CIA is somewhat more complicated. In assessing the desirability of a strategy, the CIA needs to form conjecture about the strategy that each of the three possible FBI types will use. Given such a conjecture, the CIA can compare the expected utility of choosing *Kingpin* or *Operative*. Conversely, given a conjecture about the CIA's strategy, each of the possible types of FBI should choose a strategy that maximizes its utility. In essence we have translated a problem in which the CIA does not know the preferences of the FBI to a new game in which the CIA is playing one of three possible FBI players (which are drawn from a known distribution) and each player (three FBI types plus the CIA) are all playing strategies.

Aside from the possibility that agents may not know each others preferences, it is possible that agents do not know their own preferences. Returning to the original "Terrorist Hunt" problem, a more realistic model involves a probability of catching a terrorist conditional on strategies by the players. We might think that the original matrix is justified by the following matrix.

| Table 6.4: Terrorist Hunt with Uncertainty | | |
|---|---|---|
| FBI\CIA | $Kingpin$ | $Operative$ |
| $Kingpin$ | $10 \times \frac{1}{5} + 0 \times \frac{4}{5}, 10 \times \frac{1}{5} + 0 \times \frac{4}{5}$ | $0, 6 \times \frac{1}{6} + 0 \times \frac{5}{6}$ |
| $Operative$ | $6 \times \frac{1}{6} + 0 \times \frac{5}{6}, 0$ | $6 \times \frac{1}{6} + 0 \times \frac{5}{6}, 6 \times \frac{1}{6} + 0 \times \frac{5}{6}$ |

Here 10 is the utility payoff to catching a Kingpin and $\frac{1}{5}$ is the probability of catching a $Kingpin$ if both agencies cooperate on Kingpin searching. Alternatively 0 is the payoff to failed $Kingpin$ searching which occurs with probability $\frac{4}{5}$ when the agencies cooperate on $Kingpin$ searching. While the overall payoffs are the same in this matrix and the earlier one, this representation explicitly captures the fact that payoffs may depend on both strategies and some additional randomness in the world.

A comment is in order about the role of common knowledge in games of incomplete information. As we pointed out at the beginning of the chapter, we assume in games of complete information that all elements of the game – players, strategies, and payoffs – are known to all players and all players know this. In games of incomplete information, we still maintain the common knowledge assumption. In particular, we assume that all players known the probability distribution that Nature uses in selecting player types and that all players know that all players know and so on.

## 1. Formal Definitions

We now modify our basic normal form structure $\Gamma$ to account for both imperfect information about player types and payoffs that depend on additional randomness.[1] We begin by adding to $\Gamma$ player types and a known lottery over these types and an additional random state variable.

(1) Types: We assume that for each player, $i \in N$ there is a finite set $\Theta_i$ of possible types. In the first example of this section, this set is a singleton for the CIA and it includes three elements for the FBI. By $\theta_{-i}$ and $\Theta_{-i}$ we denote a profile (and the set of such profiles) of types for all players other than $i$.

(2) Random state: In addition there is a random variable $\omega \in \Omega$ (which is assumed to be a finite state). In the second example of this section $\omega$ could be thought of as the realization of searching given different strategies. So $\omega$ specifies whether

---

[1]Some readers may want to review probability theory in the mathemtical appendix before proceeding.

a Kingpin or Operative will be found given every strategy in $S$.

(3) Natures randomization: At the beginning of the game nature selects the vector of player types $\theta = (\theta_1, ..., \theta_n) \in \Theta = \prod_{i \in N} \Theta_i$ and $\omega \in \Omega$ from a joint lottery assigning each pair $(\theta, \omega)$ probability $p(\theta, \omega)$. By $p(\theta_{-i}, \omega \mid \theta_i)$ we denote the conditional probability of $\theta, \omega$ given $\theta_i$.

(4) Strategies: Each player selects an action $s_i$ in the set $S_i$.

Expected utilities: For each possible strategy profile $s$, type $\theta_i$ and state $\omega$ agent $i$ has utility $u_i(s_i, s_{-i}, \theta_i, \omega)$. Given her type $\theta_i$ agent $i$'s conditional expected utility from strategy profile $s$ is

$$Eu_i(s; \theta_i) = \sum_{\omega \in \Omega} p(\theta_{-i}, \omega \mid \theta_i) u_i(s, \theta_i, \omega).$$

Accordingly a normal form Bayesian game is defined by the collection: $\langle N, \Omega, \{S_i, \Theta_i, u(\cdot, ..., \cdot)\}_{i \in n}, p(\cdot, \cdot) \rangle$. We sometimes use the shorthand $\langle N, \Omega, S, \Theta, u, p \rangle$. Just as normal form games may be defined with non finite strategy spaces, Bayesian games may be defined with infinite type and action spaces. We provide several such examples below.

In a normal form a strategy profile is just a list $s \in S$. In a Bayesian game we need to record the strategy that each possible type of player will use. Accordingly in a Bayesian game a strategy for player $i$ is a function $\phi_i(\theta_i) : \Theta_i \to S_i$ that selects a strategy $s_i \in S_i$ for each possible type $\theta_i \in \Theta_i$. In the version of "Terrorist Hunt" in which the CIA is not sure of the FBI's preferences, if we consider the FBI types as $\Theta_{FBI} = \{standard, pro - Kingpin, pro - Operative\}$ each occurring with equal probability then an example of a strategy for the FBI is $\phi_{FBI}(standard) = Kingpin$, $\phi_{FBI}(pro - Kingpin) = Kingpin$, $\phi_{FBI}(pro - Operative) = Operative$.

In a Bayesian normal form game we can extend the idea of Nash equilibria to the concept of Bayesian Nash equilibria. Since discussion of best responses when a strategy itself is a function is a little complicated, the easiest way to define a Bayesian Nash equilibrium is to extend our second definition of Nash equilibrium.

DEFINITION 6.1. *Given a normal form Bayesian game $\langle N, \Omega, S, \Theta, u, p \rangle$ a Bayesian Nash equilibrium is a profile of strategies, $(\phi_1^*(\cdot), ... \phi_n^*(\cdot))$ such that for every $i \in N$, for each $\theta_i \in \Theta_i$*

$$(6.1) \quad EU_i(\phi_i^*(\theta_i), \phi_{-i}^*(\cdot); \theta_i) \geq EU_i(s_i', \phi_{-i}^*(\cdot); \theta_i) \text{ for every } s_i' \in S_i.$$

Thus, in a Bayesian Nash equilibrium, every type of each player chooses strategies that maximize their expected utility given strategies of all other player types and the probability distribution of those types.

Returning to the multiple-type FBI game, we can solve for a Bayesian Nash equilibria. Suppose that the FBI strategy is as specified above, $\phi_{FBI}(standard) = Kingpin$, $\phi_{FBI}(pro{-}Kingpin) = Kingpin$, $\phi_{FBI}(pro{-}Operative) = Operative$. In this game the CIA has only one possible type and there is no uncertainty other than FBI types so we can suppress $\omega$. Given this strategy, we have

$$EU_{CIA}(Kingpin, \phi_{FBI}(\cdot)) = \frac{2}{3}2 + \frac{1}{3}0 = \frac{4}{3}$$
$$EU_{CIA}(Operative, \phi_{FBI}(\cdot)) = \frac{2}{3}0 + \frac{1}{3}1 = \frac{1}{3}$$

Thus, the CIA's best response to $\phi_{FBI}(\cdot)$ is to select Kingpin. Given this we must verify if any of the potential FBI types wish to deviate. For the standard type FBI we know that matching the CIA is a best response so $\phi_{FBI}(standard) = Kingpin$ is a best response. For the pro-Kingpin type selecting Kingpin when the CIA selects Kingpin results in the highest possible payoff (2) and thus no desirable deviation exists. Finally, for the pro-Operative FBI, selecting $Operative$ when the CIA selects $Kingpin$ results in utility of 2 while a deviation to $Kingpin$ results in utility 1, thus no desirable deviation exists. This means that we have characterized a Bayesian Nash equilibrium to the game.

## 2. Application: Trade restrictions

We consider a simple setting involving two nations that are contemplating restrictive trade policies. Let $N = \{1, 2\}$ and suppose that each country has two possible types $\Theta_i = \{u, b\}$ A type $u$ country wishes to limit its imports from the other country *unilaterally*, while a type $b$ country wishes to pursue a *bilateral* policy of limiting trade only if the other country does so. We assume the country types are independently drawn with type $u$ occurring with probability $p \in (0, 1)$. The strategy space for each country is $S = \{l, f\}$ where $l$ denotes a enacting an import limit and $f$ denotes a free-trade policy. Country $i$ has the following payoffs

$$u_i(s_i, s_{-i}; \theta_i) = \begin{cases} 3 \text{ if } s_i = l, s_{-i} = f \text{ and } \theta_i = u \\ 2 \text{ if } s_i = f, s_{-i} = f \text{ and } \theta_i = u \\ 1 \text{ if } s_i = l, s_{-i} = l \text{ and } \theta_i = u \\ 0 \text{ if } s_i = f, s_{-i} = l \text{ and } \theta_i = u \\ 3 \text{ if } s_i = f, s_{-i} = f \text{ and } \theta_i = b \\ 2 \text{ if } s_i = l, s_{-i} = f \text{ and } \theta_i = b \\ 1 \text{ if } s_i = l, s_{-i} = l \text{ and } \theta_i = b \\ 0 \text{ if } s_i = f, s_{-i} = l \text{ and } \theta_i = b \end{cases}$$

A strategy is a mapping $s_i(\theta_i) : \{u, b\} \to \{l, f\}$. Note that a $u$ country always gets a higher payoff from $l$ regardless of the actions of the other country. Thus, if there were common knowledge that both countries were type $u$ then the game would be a Prisoners' Dilemma and each country would have a dominant strategy to choose $l$. Alternatively, if there was common knowledge that both countries were type $b$ then there would be two pure strategy Nash equilibria $(f, f)$ and $(l, l)$.

In computing, Bayesian Nash equilibria, it is often useful to make conjectures about equilibrium strategies and check to see if they satisfy the required conditions. Since we know that a type $u$ country has a dominant strategy of selecting $l$ every equilibrium will involve $s_i(u) = l$. Thus, the only possible equilibria are $(s_i(u), s_i(b)) = (l, l)$ and $(l, f)$. Thus, we first investigate the possibility that $s_i(u) = l$ and $s_i(b) = f$ for both $i = 1, 2$. If country 2 uses this strategy, then $s_2 = l$ with probability $p$ and $s_2 = f$ with probability $(1 - p)$. Thus, country 1's expected utility is

$$Eu_1(s_1, \theta_i = u) = \begin{cases} p + (1 - p)3 \text{ if } s_1 = l \\ 2(1 - p) \text{ if } s_1 = f \end{cases}$$

$$Eu_1(s_1, \theta_i = b) = \begin{cases} 2(1 - p) + p \text{ if } s_1 = l \\ 3(1 - p) \text{ if } s_1 = f \end{cases}$$

When is the conjectured strategy a best response? The strategy $s_i(u) = l$ is a best response since type $u$ has a dominant strategy to erect trade barriers. Alternatively, $s_i(b) = f$ is a best response if

$$3(1 - p) \geq 2(1 - p) + p$$

which is true as long as $p \leq \frac{1}{2}$. Since the calculations for country 2 are identical, we know that for $p \leq \frac{1}{2}$ the profile $s_i(u) = l$ and $s_i(b) = f$ is a Bayesian Nash equilibrium.

Now we lets check to see if $s_i(b) = l$ for both countries can be a best response. We note that if country 2 uses the strategy $s_i(\theta_i) = l$ regardless of $\theta_i$ then country 1 with type $b$ has the expected utility

$$Eu_1(s_1, \theta_i = b) = \begin{cases} 1 \text{ if } s_1 = l \\ 0 \text{ if } s_1 = f. \end{cases}$$

This means that regardless of $p$, the strategies $s_i(b) = l, s_i(u) = l$ for both countries is a Bayesian Nash equilibrium. So we have seen that there is always a Bayesian Nash equilibrium in which bilateral limits $(l, l)$ occur. Moreover, if $p \leq \frac{1}{2}$ there is a second equilibrium with $s_i(u) = l$ and $s_i(b) = f$. In this equilibrium free trade occurs if both countries are of type $b$, which happens with probability $(1 - p)^2$. So the possibility of bilateral policies requires that a country is not the

unilateral type and the expectation that her opponent is also unlikely to be the unilateral type.

## 3. Application: Jury Voting

As another example of Bayesian Nash Equilibrium, we now consider a simple example of the jury model developed by Austen-Smith and Banks (1997). Suppose that three jurors $N = \{1, 2, 3\}$ are responsible for deciding whether to convict or acquit a defendant. Thus, they collectively must choose an outcome $x \in \{c, a\}$. The jurors simultaneously cast ballots $v_i \in S_i = \{c, a\}$ and the outcome is chosen by majority rule. Each player faces uncertainty about whether or not the defendant is guilty, $G$, or innocent, $I$. So the set of state variables can be denoted as $\Omega = \{G, I\}$. In the guilty state the jurors receive utility 1 from convicting and 0 from acquitting. Alternatively, in innocent state the jurors receive utility 1 from acquitting and 0 from convicting. Each player assigns prior probability $\pi > \frac{1}{2}$ to the guilty state.

If we assumed that each of the jurors had identical information, each juror would receive an expected utility of $\pi$ from a guilty verdict and $1 - \pi$ from an acquittal. Since $\pi > 1 - \pi$, the Nash equilibrium that survives the elimination of weakly dominated strategies calls for each juror to vote guilty.

However, assume that before voting each player receives a private signal concerning the defendants guilt $\theta_i \in \{0, 1\}$. We assume that this signal is informative in the sense that a juror is more likely to receive the signal $\theta_i = 1$ when the defendant is guilty than she is when the defendant is innocent. Furthermore, to keep things very simple, we assume that the probability of receiving the guilty signal ($\theta_i = 1$) when the defendant is guilty is the same as that of receiving the innocent signal ($\theta_i = 0$) when the defendant is innocent. Formally, we assume that $\Pr(\theta_i = 1 \mid \omega = G) = \Pr(\theta_i = 0 \mid \omega = I) = p > \frac{1}{2}$ which obviously requires that $\Pr(\theta_i = 0 \mid \omega = G) = \Pr(\theta_i = 0 \mid \omega = I) = 1 - p$.

After receiving her signal, voter $i$ will select her vote $v(\theta_i)$ to maximize the probability that guilty defendant are convicted and innocent defendants are acquitted. Suppose that each voter uses a straightforward strategy $v_i(0) = a$ and $v_i(1) = c$ of voting to convict when they get a signal of 1 and voting to acquit when the get a signal of 0. To verify whether this strategy combination constitutes a Bayesian Nash equilibrium, we need to verify that voter 1 is willing to use this strategy if she conjectures that voters 2 and 3 are using this strategy. Given

these conjectures, the expected utility of voting to convict is

$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = G \mid \theta_1)+$$
$$\Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{and } \omega = G \mid \theta_1)+$$
$$\Pr(\theta_2 = 1 \text{ and } \theta_2 = 1 \text{ and } \omega = G \mid \theta_1)+$$
$$\Pr(\theta_2 = 0 \text{ and } \theta_2 = 0 \text{ and } \omega = I \mid \theta_1)$$

while the expected utility of voting to acquit is

$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = I \mid \theta_1)+$$
$$\Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{and } \omega = I \mid \theta_1)+$$
$$\Pr(\theta_2 = 0 \text{ and } \theta_2 = 0 \text{ and } \omega = I \mid \theta_1)+$$
$$\Pr(\theta_2 = 1 \text{ and } \theta_2 = 1 \text{ and } \omega = G \mid \theta_1).$$

Note that the last two terms of each sum are the same and thus cancel out in comparing these two expected utilities. Accordingly voting to convict is a best response if

$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = G \mid \theta_1) + \Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{and } \omega = G \mid \theta_1) \geq$$
$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = I \mid \theta_1) + \Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{ and } \omega = I \mid \theta_1)$$

while voting to acquit is a best response if the reverse weak inequality holds. Since these expressions depend on conditional probabilities of observing combinations of the state variable and the signals of the remaining jurors, juror 1 must use Bayes' Rule to evaluate each term. Suppose that juror 1 receives $\theta_1 = 1$. If can be easily shown that

$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = G \mid \theta_1 = 1)$$

$$= \Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{ and } \omega = G \mid \theta_1 = 1) = \frac{\pi p^2 (1-p)}{\pi p + (1-\pi)(1-p)}$$

and

$$\Pr(\theta_2 = 1 \text{ and } \theta_3 = 0 \text{ and } \omega = I \mid \theta_1 = 1)$$

$$= \Pr(\theta_3 = 1 \text{ and } \theta_2 = 0 \text{ and } \omega = I \mid \theta_1 = 1) = \frac{(1-\pi)\, p \,(1-p)^2}{\pi p + (1-\pi)(1-p)}$$

Thus, juror 1 will choose $v_i(1) = c$ if

$$2\frac{\pi p^2(1-p)}{\pi p + (1-\pi)(1-p)} \geq 2\frac{(1-\pi)p(1-p)^2}{\pi p + (1-\pi)(1-p)}$$

or after simplifying

$$\pi p(1-p)p \geq (1-\pi)p(1-p)(1-p).$$

We can rearrange this expression to

$$\frac{\pi p^2(1-p)}{\pi p^2(1-p) + (1-\pi)p(1-p)^2} \geq \frac{1}{2}.$$

The left hand side is just the conditional probability of guilt given two signals of 1 and one signal of 0. In other words agent 1 is willing to vote to convict if she thinks that the defendant is more likely to be guilty than innocent when she conditions on her signal (one of the 1 signals) and the assumption that she is pivotal so that the remaining two agents have received different signals. Similarly, we can express the requirement for a vote of innocence conditional on a signal of 0 as

$$\frac{\pi p(1-p)^2}{\pi p(1-p)^2 + (1-\pi)p^2(1-p)} \leq \frac{1}{2}.$$

So the conclusion is that in a Bayesian equilibrium to this game a voter must vote for the action which she prefers in the state she believes is more likely given that she conditions on both her own signal and the profile of other signals that must occur for her to be pivotal. Austen-Smith and Banks note that in many cases this simple strategy (voting to convict if $\theta_i = 1$ and voting to acquit if $\theta_i = 0$) is not equilibrium behavior. In this example this point can be demonstrated by choosing parameters $\pi$ and $p$ for which one of the last two inequalities does not hold. Of course alternative types of strategies might be played. Voters can randomize for some types, or they may choose to vote the same way regardless of their type, or different voters may use different strategies. Fedderson and Pessendorfer (1998) consider the properties of equilibria to this game as one varies the voting rule and population size.

## 4. Application: Jury Voting with a Continuum of Signals*

One extension of the jury model involves a larger type space.[2] Suppose instead of receiving a binary signal, each juror receives a signal $\theta_i \in [0, 1]$. One way to model the case of an informative signal taking on a continuum of possible values is to assume that $\theta_i$ is drawn from a state conditional distribution $F(\theta_i|\omega)$ with a differentiable density function $f(\theta_i|\omega)$ that satisfies the *monotone likelihood ratio* condition.

DEFINITION 6.2. *The conditional densities satisfy the strict monotone likelihood ratio condition (SMLR) if $\frac{f(\theta_i|G)}{f(\theta_i|I)}$ is a strictly monotone function of $\theta_i$ on $[0, 1]$.*

---

[2]Duggan and Martinelli (2001) and Meirowitz (2002) consider this extension.

For the convenience, we assume that $\frac{f(\theta_i|G)}{f(\theta_i|I)}$ is strictly increasing. To see why this assumption is important, note that Bayes' rule implies that

$$\Pr(G|\theta_i) = \frac{f(\theta_i \mid G)\pi}{f(\theta_i \mid G)\pi + f(\theta_i \mid I)(1 - \pi)}$$

$$= \frac{\frac{f(\theta_i|G)}{f(\theta_i|I)}\pi}{\frac{f(\theta_i|G)}{f(\theta_i|I)}\pi + (1 - \pi)}$$

It can easily be verifies that $\Pr(G|\theta_i)$ is increasing in $\theta_i$ if and only if $\frac{f(\theta_i|G)}{f(\theta_i|I)}$ is increasing in $\theta_i$. Thus, the SMLR condition implies that the higher the signal agent $i$ receives the higher her posterior belief about the probability that the state is $\omega = G$.

To keep things simple, we will focus exclusively on symmetric strategies where voters who receive the same signal choose the same strategy. Thus, a symmetric strategy profile is characterized by a mapping $v_i(\theta_i) : [0,1] \rightarrow \{c,a\}$. Following the logic from above, a Bayesian Nash equilibrium must involve a strategy that is optimal when the agent conditions on her private information and the conjecture that she is pivotal. An agent will vote to convict if she thinks the probability of guilt is no less than $\frac{1}{2}$ and she will vote to acquit if she thinks the probability of guilt is no more than $\frac{1}{2}$. Given that higher signals are better indicators of guilt, one natural conjecture is that the strategy must be weakly increasing. For low values of $\theta_i$ an acquittal vote is cast and for high values of $\theta_i$ a conviction vote is cast. Let $\widehat{\theta} \in [0,1]$ denote a cutpoint and lets assume that agents $i \in N\backslash i$ use the strategy

$$v_i(\theta_i) = \begin{cases} c \text{ if } \theta_i \geq \widehat{\theta} \\ a \text{ if } \theta_i < \widehat{\theta} \end{cases}$$

Suppose that the juror decision rule is a $q$-rule, requiring at least $q \geq \frac{n+1}{2}$ votes to convict. Thus, if players $N\backslash i$ use the cutpoint strategy then the posterior probability of $\{\omega = G\}$ given $\theta_i$ and that $i$ is pivotal is given by

(6.2) $\quad pr(G \mid piv, \theta_i; \widehat{\theta}) =$

$$\frac{\pi f_G \widehat{F}_G^{n-q-1} \left[1 - \widehat{F}_G\right]^{q-1}}{\pi f_G \widehat{F}_G^{n-q-1} \left[1 - \widehat{F}_G\right]^{q-1} + (1 - \pi) f_I \widehat{F}_I^{n-q-1} \left[1 - \widehat{F}_I\right]^{q-1}}$$

where $f_\omega = f(\theta_i \mid \omega)$ and $\widehat{F}_\omega = F(\widehat{\theta} \mid \omega)$. We leave the derivation of this expression as an exercise. This probability is a function of the

parameter $\widehat{\theta}$ and we make this point explicit in the left hand side. In this model the existence of a symmetric equilibrium in which voters use a cutpoint hinges on finding a value of $\widehat{\theta}$ such that

$$pr(G \mid piv, \widehat{\theta}; \widehat{\theta}) = \frac{1}{2}$$

and demonstrating that if $\theta_i < \widehat{\theta}$ $pr(G \mid piv, \theta_i; \widehat{\theta}) \leq \frac{1}{2}$ and if $\theta_i > \widehat{\theta}$ $pr(G \mid piv, \theta_i; \widehat{\theta}) \geq \frac{1}{2}$. While analysis of examples is cumbersome, it is easy to come up with conditions on the primitives of the game to insure that such a $\widehat{\theta} \in (0,1)$ exists. First, since $pr(G \mid piv, \theta_i; \widehat{\theta}) \geq \frac{1}{2}$ if and only if

$$\frac{\pi f(\theta_i \mid G) F(\widehat{\theta} \mid G)^{n-q-1} \left[1 - F(\widehat{\theta} \mid G)\right]^{q-1}}{(1-\pi) f(\theta_i \mid I) F(\widehat{\theta} \mid I)^{n-q-1} \left[1 - F(\widehat{\theta} \mid I)\right]^{q-1}} =$$

$$\frac{f(\theta_i \mid G)}{f(\theta_i \mid I)} \frac{\pi F(\widehat{\theta} \mid G)^{n-q-1} \left[1 - F(\widehat{\theta} \mid G)\right]^{q-1}}{(1-\pi) F(\widehat{\theta} \mid I)^{n-q-1} \left[1 - F(\widehat{\theta} \mid I)\right]^{q-1}} \geq 1$$

the strict monotone likelihood ratio conditions implies that if $pr(G \mid piv, \widehat{\theta}; \widehat{\theta}) = \frac{1}{2}$ then $\theta_i < \widehat{\theta}$ implies $pr(G \mid piv, \theta_i; \widehat{\theta}) \leq \frac{1}{2}$ and $\theta_i > \widehat{\theta}$ implies $pr(G \mid piv, \theta_i; \widehat{\theta}) \geq \frac{1}{2}$. This means that existence of an equilibrium hinges on establishing the existence of a solution to the equation

$$pr(G \mid piv, \widehat{\theta}; \widehat{\theta}) = \frac{1}{2}$$

If $pr(G \mid piv, 0; 0) \leq \frac{1}{2} \leq pr(G \mid piv, 1; 1)$ then the intermediate value theorem implies that such a cutpoint will exist since the function $pr(G \mid piv, \cdot; \cdot)$ is continuous. For a large class of games these boundary conditions will be satisfied.

So while the simple binary type model demonstrates that equilibria where everyone uses the same rule and voting is determined by private information will not exist, equilibria of this type generally exist in the continuum model. In the case of unanimity rule, the modeling technology is consequential for the conclusions that can be reached about the desirability of particular political institutions. Using the binary model, Fedderson and Pesendorfer (1998) show that unanimity rule is an uniquely poor way to aggregate information for large populations because in equilibrium voters condition on the assumption that everyone else is voting to convict. In the continuum model (Meirowitz 2002)

unanimity rule turns out to be as good as the other voting rules in some cases.

## 5. Application: Public Goods and Incomplete Information

In this section, we consider a version of the Palfrey-Rosenthal contribution game where potential contributors are uncertain about the contribution costs of other players. To keep the model as close to the one we have already analyzed, we assume that every agent receives a utility of 1 if at least $k$ agents contribute and 0 otherwise. Agent $i$ pays a cost $c_i$ to contribute where we assume that $c_i$ is distributed uniformly on $[0, 1]$. Each agent learns their own costs, but remain uncertain about the costs of others.

**5.1.** $k = 1$. First we consider the case where a single contribution is necessary for the provision of the good. Since $c_i < 1$ with certainty for all contributors, there are always $n$ Bayesian Nash equilibrium corresponding to agent $i$ contributing with certainty. However, as we did before, we will concentrate on symmetric equilibrium where all player types with the same cost play the same strategy. We therefore focus on equilibria in cutpoint strategies. In such equilibria, agent $i$ contributes if and only if $c_i < \widehat{c}$. If we assume that all players other than $i$ choose a cutpoint strategy, agent $i$ utility from contributing is $1 - c_i$. To compute her utility for not contributing, note that she receives 1 as long as there as at least one contributor. Since $c$ is distributed uniformly on $[0, 1]$, other contributors contribute with probability $\widehat{c}$ so that the probability of no contributions is $[1 - \widehat{c}]^{n-1}$. Thus, agent $i$'s utility from not contributing is $1 - [1 - \widehat{c}]^{n-1}$. Thus agent $i$ contributes so long as

$$[1 - \widehat{c}]^{n-1} \geq c_i$$

We can use this expression to solve implicitly for $\widehat{c}$ since agent $i$ must be indifferent at the cutpoint. Thus, we have $\widehat{c}^{n-1} + \widehat{c} = 1$. It is very easy to show that $\widehat{c}$ is declining in $n$ and goes to zero for very large $n$. In turn, this means that the probability that any individual will contribute goes to zero as the group expands. Thus, the incomplete information version also predicts that there will be more free-riding in large groups.

**5.2.** $k > 1*$. Now we turn to the case where multiple contributions are required for the provision of the good. We again assume that agents use cutpoint strategies and contribute only if $c_i \leq \widehat{c}$.

Let $x_{-i}$ be the realized number of contributions from agents other than $i$. From arguments identical to those of the last chapter, we know

that agent $i$'s net utility from contributing is

$$\Pr\left(x_{-i} = k - 1\right) - c_i$$

Since contribution probabilities are $\widehat{c}$, we can also show that

$$\Pr\left(x_{-i} = k - 1\right) = \left(\begin{array}{c} n - 1 \\ k - 1 \end{array}\right) \widehat{c}^{\,k-1} \left(1 - \widehat{c}\right)^{n-k}$$

Since the agent $i$ must be indifferent at the cutpoint cost, we again get an implicit solution for $\widehat{c}$:

$$\left(\begin{array}{c} n - 1 \\ k - 1 \end{array}\right) \widehat{c}^{\,k-2} \left(1 - \widehat{c}\right)^{n-k} = 1$$

Note that this solution is very similar to that of the mixed strategy equilibrium with complete information. The main difference is that $\widehat{c}$ plays the role of $\sigma^*$. Thus, many of the implications of our previous analysis carry over.

To reduce notation let $\Pi(\widehat{c}) = \left(\begin{array}{c} n - 1 \\ k - 1 \end{array}\right) \widehat{c}^{\,k-2} \left(1 - \widehat{c}\right)^{n-k}$ so that our equilibrium condition is $\Pi(\widehat{c}) = 1$. First note that

$$\Pi' \gtreqless 0 \text{ if } \widehat{c} \lesseqgtr \frac{k - 2}{n - 2}$$

Thus, so long as $2 < k < n$, $\Pi(\widehat{c})$ is single peaked for $\widehat{c} \in [0, 1]$. Thus, if $2 < k < n$ and $\max \Pi > 1$, there will be two equilibrium cutpoints $\widehat{c}_H$ and $\widehat{c}_L$ just as there were two symmetric equilibria in mixed strategies. If $\max \Pi < 1$, there will be no cutpoint equilibria.[3]

As before, it is easy to demonstrate that the cutpoint equilibria disappear as $n$ gets very large. Note that the $\widehat{c}$ that maximizes $\Pi$ is bounded by $\frac{k-2}{n-2}$. Thus, as $n$ gets large, this maximizer goes to zero. Since $\Pi(0) = 0$, $\lim_{n \to \infty}\left[\max \Pi\right] = 0$. Thus for sufficiently large $n$, $\max \Pi < 1$. We leave it to the readers to verify that the effects of $n$ and $k$ on $\widehat{c}_H$ and $\widehat{c}_L$ are essentially the same as the effects of $n$ and $k$ on $\sigma_L^*$ and $\sigma_H^*$ in the symmetric mixed strategy equilibrium of the previous chapter.

## 6. Application: Electoral Competition under Uncertainty

We now return to the Hotelling model of candidate competition with policy motivated candidates and consider two extensions. We first assume that candidate preferences are known (candidates have ideal points of 0 and 1 as before) but that instead of knowing that the

---

[3]In the unlikely event that $\max \Pi = 1$, there will be a unique cutpoint equilibrium.

voters are arranged uniformly on the unit interval, we assume that the candidates believe that the median voter's location is randomly drawn from the uniform distribution on $[0, 1]$. This model is an example of models treated by Wittman (1977) and Calvert (1985). To formulate this model as a Bayesian game, we let $\Omega = [0, 1]$ denote the set of possible locations of the median voter. Since candidate preferences are common knowledge in this example the type spaces are singletons. So all of the uncertainty in the game is captured by the probability distribution $F(\omega) = \omega$ on $[0, 1]$. Here, we assume that candidate preferences over policy are quadratic so that $u_1(x) = -x^2$ and $u_2(x) = -(1 - x)^2$. Given two platforms $s_1 < s_2$ candidate 1 wins if the median is closer to $s_1$ than $s_2$. This is true if the median is less than $\frac{s_1+s_2}{2}$. Since the median is uniformly distributed, candidate 1 wins with probability $\frac{s_1+s_2}{2}$. In this case candidate expected utilities are

$$Eu_1(s_1, s_2) = \begin{cases} -s_1^2 \frac{s_1+s_2}{2} - s_2^2(1 - \frac{s_1+s_2}{2}) & \text{if } s_1 < s_2 \\ -s_2^2 \frac{s_1+s_2}{2} - s_1^2(1 - \frac{s_1+s_2}{2}) & \text{if } s_1 > s_2 \end{cases}$$

and

$$Eu_2(s_1, s_2) = \begin{cases} -(1 - s_1)^2 \frac{s_1+s_2}{2} - (1 - s_2)^2(1 - \frac{s_1+s_2}{2}) & \text{if } s_1 < s_2 \\ -(1 - s_2)^2 \frac{s_1+s_2}{2} - (1 - s_1)^2(1 - \frac{s_1+s_2}{2}) & \text{if } s_1 > s_2 \end{cases}.$$

To construct an equilibrium, suppose that candidate 1 knows that candidate 2 will locate at $z \geq \frac{1}{2}$. Then candidate 1 chooses $s_1 \in [0, z]$ to optimize

$$\max_{s_1}\{-s_1^2 \frac{s_1 + z}{2} - z^2(1 - \frac{s_1 + z}{2})\}$$

Note that we can ignore the possibility of choosing $s_1 > z$ because this strategy is always dominated by $s_1 = z$.[4] To find the optimal choice of $s_1$, we can differentiate the objective function with respect to $s_1$ and set this derivative to 0. This first order condition is

$$-\frac{3}{2}s_1^2 - zs_1 + \frac{z^2}{2} = 0.$$

Solving for $s_1$ yields two solutions, but only one is in the appropriate range $[0, 1]$. This solution then gives us the best response function

$$s_1(s_2) = \frac{s_2}{3}.$$

To be sure that this solution characterizes a local maxima (as opposed to a local minima or saddle point) we need to check that the second derivative of the objective function is negative when evaluated at this

---

[4]Indeed if candidate chose $s_1 > z$, she would prefer candidate 2 to win.

value. This second order condition simplifies to $-2s_2$ and is negative for any value of $s_2 \in (0, 1]$.

Similarly, treating $s_1$ as a fixed parameter $z \leq \frac{1}{2}$ we can differentiate candidate 2's objective function

$$\max_{s_2}\{-(1-z)^2 \frac{z+s_2}{2} - (1-s_2)^2 \left[1 - \frac{z+s_2}{2}\right]\}$$

to find an optimal $s_2 \in [z, 1]$. The solution is

$$s_2(s_1) = \frac{2}{3} + \frac{1}{3}s_1.$$

We leave verification of the second order condition to the reader. A Bayesian Nash equilibrium is then a strategy combination $(s_1^*, s_2^*)$ that solves the system

$$s_1^* = \frac{1}{3}s_2^*$$

$$s_2^* = \frac{2}{3} + \frac{1}{3}s_1^*.$$

The unique solution to this system is

$$s_1^* = \frac{1}{4}, s_2^* = \frac{3}{4}.$$

So with policy motivated candidates and uncertainty about voter preferences candidate divergence is predicted.

**6.1. Private Information about Candidate Preferences.** Now we consider a model where in addition to uncertainty about the location of the median voter, candidates have private information about their policy preferences. One simple example involves candidate 1 having ideal point $\theta_1 \in \{0, \frac{1}{2}\}$ and candidate 2 having ideal point $\theta_2 \in \{\frac{1}{2}, 1\}$. So candidate the utility to candidate $i$ of policy location $x$ is $u(x) = -(\theta_i - x)^2$. For simplicity we assume that each type is drawn with equal probability and that choice of types across candidates are independent. As before we assume that the median voter's ideal point is randomly drawn from a uniform distribution over $[0, 1]$. In this case a strategy for candidate 1 is a mapping $s_1(\theta_1) : \{0, \frac{1}{2}\} \to [0, \frac{1}{2}]$ and a strategy for candidate 2 is a mapping $s_2(\theta_2) : \{\frac{1}{2}, 1\} \to [\frac{1}{2}, 1]$. For simplicity, we ignore the possibility of a candidate selecting a policy that is further from her ideal point than the value $\frac{1}{2}$. One justification for this assumption might be that parties constrain candidates from crossing over into the far side of the ideological spectrum. In this case suppose that candidate 2 uses the strategy $s_2(\frac{1}{2}) = a$ and $s_2(1) = b$. Now, in considering the optimal location for a candidate 1 with type

$\theta_1 = \frac{1}{2}$ note that any location $s_1 < \frac{1}{2}$ is dominated by the location $\frac{1}{2}$. This is true because for fixed $a, b$ the location $\frac{1}{2}$ will win with probability strictly higher than the probability that a location of $s_1 < \frac{1}{2}$ wins. Moreover, conditional on winning a location of $s_1 < \frac{1}{2}$ is less desirable to candidate 1 with type $\theta_1 = \frac{1}{2}$ than a location of $\frac{1}{2}$. Accordingly, we know that $s_1(\frac{1}{2}) = 0$ is strictly dominant. Similarly, $s_2(\frac{1}{2}) = \frac{1}{2}$ is strictly dominant. Now to find the best response for candidate 1 of type $\theta_1 = 0$ we need to solve an optimization problem. If $\theta_1 = 0$ then candidate 1's objective function is

$$\max_{s_1}\{-s_1^2\left(\frac{s_1 + \frac{1}{2}}{4} + \frac{s_1 + b}{4}\right) - \frac{(\frac{1}{2})^2}{2}\left(1 - \frac{s_1 + \frac{1}{2}}{2}\right) - \frac{b^2}{2}\left(1 - \frac{s_1 + b}{4}\right)\}$$

Differentiating with respect to $s_1$ and setting this term equal to 0, yields the first order condition.

$$\frac{1}{8}b^2 - \frac{1}{2}bs_1 - \frac{1}{4}s_1 - \frac{3}{2}s_1^2 + \frac{1}{16} = 0.$$

The solution (in the appropriate range) is

$$s_1(0; b) = \frac{1}{12}\sqrt{4b + 16b^2 + 7} - \frac{1}{6}b - \frac{1}{12}.$$

Now instead of solving candidate 2's problem in an analogous manner, we can notice that given $s_1(\frac{1}{2}) = \frac{1}{2}$ and $s_2(0) = a$ candidate 2's problem is the mirror image of candidate 1's. This means that we can find the equilibrium values of $s_2(1) = b$ and $s_1(0) = 1 - b$ that solve the relevant first order conditions by solving for $b$ that satisfies the equality

$$1 - b = \frac{1}{12}\sqrt{4b + 16b^2 + 7} - \frac{1}{6}b - \frac{1}{12}.$$

The solution is $b = \frac{11}{7} - \frac{1}{14}\sqrt{106} \simeq 0.836\,03$. Thus, the Bayesian Nash equilibrium is $s_1(0) = 0.164$, $s_1(\frac{1}{2}) = \frac{1}{2}$, $s_2(\frac{1}{2}) = \frac{1}{2}$, $s_2(1) = 0.836$. It is instructive to compare the platforms of $\theta_1 = 0$ and $\theta_2 = 1$ with the outcomes of the game where candidate preferences are known. One might expect that the platforms would be more convergent given that each candidate think that they might be playing against a moderate candidate. This intuition is misleading though since we observe platforms in the candidate uncertainty game that are even more divergent. The rationale for the seemingly anomalous outcome is that each extreme candidate type knows that they will lose with certainty against a moderate type unless they locate at the at the expected median. However, they are indifferent between wining in losing if both candidate locate at .5. So all the candidate uncertainty does is make the election outcome

more random.   This additional randomness mitigate the penalty for taking extreme positions.   Therefore, the candidates polarize.

## 7.  Application: Campaigns, Contests and Auctions*

An alternative perspective on electoral competition frames the competition as a contest.   Candidates each select a level of costly effort, and a winner is chosen.   More effort increases the likelihood that one wins.   This approach allows us to focus on the role of money in campaigns.   The model that we present here shares many features with models of auctions – a topic we take up in some detail in Chapter 11. Consider a set $N = \{1, ..., n\}$ of candidates that are running for office. Candidates compete by raising money and spending it on advertisements.   Let $a_i \in R_+^1$ denote the level of fundraising by candidate $i$. Given the accumulations $a = (a_1, ..., a_n)$, the winner is determined by the function $p(a) : R_+^n \to N$.   This function should be weakly increasing.   A reasonable example includes the mapping $p(a) = \arg\max_{i \in N} a_i$ which awards the office to the candidate that raises the most money.[5] Candidate $i$'s utility depends on the identity of the winner, the level of accumulation $a_i$ and the candidates value of winning office, $\theta_i \in [0, 1]$. For simplicity we assume that the values to winning office (types) of each candidate are private information independently drawn from a uniform distribution on $[0, 1]$.   Specifically, candidate $i$'s utility takes the form

$$u_i(a) = \theta_i 1_{\{p(a)=i\}} - a_i$$

where $1_{\{p(a)=i\}}$ is an indicator function that takes the value 1 if $p(a) = i$ and 0 otherwise.   In this Bayesian game each candidate simultaneously selects their level of $a_i$ and then the payoffs are realized.   In the language of auction theory, this is a first price all play auction with independent types.   A Bayesian Nash equilibrium is a function (for each candidate) that $\theta_i$ into $a_i$.   Once again, we focus on symmetric equilibria in which each candidate uses the same strategy.

Directly solving the model for continuous strategy functions is often quite difficult.   A trick is to assume that the strategy function have a specific functional form, solve for any free parameters, and verify that the solutions constitute equilibria.   Here we conjecture that players $j \neq i$ use a strategy of the form, $a_j(\theta_j) = b\theta_j^c$ where $b$ and $c$ are parameters to be determined.   If players $2, ..., n$ use the conjectured

---

[5]An alternative interpretation of this model is to treat $a_i$ as the level of effort or time that the candidate spends running for office.

strategy, and candidate 1 selects $a_1$ then the probability that 1 wins is $pr\{\max_{j\neq i} b\theta_j^c < a_1\}$. This probability is

$$pr\left\{\max_{j\neq i}\theta_j < \left(\frac{a_1}{b}\right)^{\frac{1}{c}}\right\} = \left(\frac{a_1}{b}\right)^{\frac{n-1}{c}}.$$

Accordingly the expected utility to player 1 with type $\theta_1$ from action $a_1$ is

$$\left(\frac{a_1}{b}\right)^{\frac{n-1}{c}}\theta_1 - a_1.$$

Differentiating with respect to $a_1$ yields the first order condition

$$\theta_1\frac{n-1}{cb}\left(\frac{a_1}{b}\right)^{\frac{n-1-c}{c}} = 1,$$

and solving for $a_1$ yields

(6.3)
$$a_1 = b\left(\frac{cb}{(n-1)\theta_1}\right)^{\frac{c}{n-1-c}}.$$

We began by conjecturing that players $j = 2, ..., n$ used a strategy of the form $a_j(\theta_j) = b\theta_j^c$ and found that player 1's best response was to use a strategy of the form given above. An equilibrium can then be found by solving for values of $b$ and $c$ such that

(6.4)
$$b\theta_1^c = b\left(\frac{cb}{(n-1)\theta_1}\right)^{\frac{c}{n-1-c}}.$$

Suppose that $b = \frac{n-1}{n}$ and $c = n$ then substitution into the right hand side yields

(6.5)
$$a_1 = \frac{n-1}{n}\left(\frac{n\left(\frac{n-1}{n}\right)}{(n-1)\theta_1}\right)^{\frac{n}{n-1-n}}$$

This simplifies to

$$a_1 = \frac{n-1}{n}\theta_1^n$$

confirming that $b = \frac{n-1}{n}$ and $c = n$ correspond to an equilibrium. Thus, a symmetric Bayesian Nash equilibrium is for each candidate to accumulate

$$a_i(\theta_i) = \frac{n-1}{n}\theta_i^n.$$

The relationship between this model and other auctions can easily be seen. In this game, a candidate suffers disutility $a_i$ regardless of whether or not she wins. An alternative model might involve each agent announcing promises to pay if they win. Another example of this form would involve interest groups that make promises to bribe a

committee chairman if their preferred nominee is conferred an appointment. In the chapter on mechanism design we will consider auctions of this type and trace out the relationships between equilibria to different types of auction.

## 8. Existence of Bayesian Nash equilibria*

When will a Bayesian Nash equilibrium exist?. In considering this question it is useful to note that the Bayesian normal form games considered are really just special cases of Normal form games, in which each player simultaneously selects a strategy (where a strategy is a mapping from $\Theta_i$ into $S_i$) and the payoffs are defined as the agents' expected utility over strategy profiles. An alternative statement is also true, a Bayesian game is equivalent to a normal form game in which every possible agent-type pair is a player of the normal form game. The requirement that a Bayesian Nash equilibrium involves strategies that are best responses for every possible type is equivalent to the requirement that in this larger normal form game every player (a player-type pairing) select a best response. This observation allows us to apply our previous results to establish the existence of Bayesian Nash equilibria in Mixed strategies to Bayesian normal form games with finite type and action spaces.

Specifically, let us start with a Bayesian game $\langle N, S, \Theta, u, p \rangle$ with $N, S, \Theta$ all finite sets. Without loss of generality we denote types in the following manner, $\Theta_i = \{\theta_i^1, ...., \theta_i^{k_i}\}$. We can define a new normal form game $\Gamma'$ as follows: Let $N' = \{\theta_1^1, .., \theta_1^{k_1}, .\theta_2^1, ..., \theta_n^{k_n}\}$. In this normal form game all agents with subscript $i$ have strategy space $S_i' = S_i$. Let $\Theta_{-i} = \times_{j=N\backslash i}\Theta_j$ denote the set of possible type profiles for the agents $N\backslash i$ in the original Bayesian game. Given a strategy profile $s^+ = (s_1^1, ...., s_i^j, .....s_n^{k_n}) \in \times_{i=1}^n S_i'$ to the game $\Gamma'$ we can identify this strategy with one in $\Gamma$ by letting $s_i^+(\theta_i = \theta_i^j) = s_i^j$. The utility to agent $\theta_i^j$ is then defined by using the notion of expected utility in $\Gamma$

$$v_i^j(s^+) = EU_i(s_i^+(\theta_i), s_{-i}^+(\cdot); \theta_i^j).$$

The new normal form game $\Gamma' = \langle N', S', v \rangle$ is then well defined, which leads to the result.

PROPOSITION 6.1. *Given the Bayesian game* $\langle N, S, \Theta, u, p \rangle$ *with* $N, S, \Theta$ *all finite sets a Bayesian Nash equilibrium in mixed strategies exists.*

PROOF. Given the Nash's Theorem (Theorem 5.4), the finite game $\langle N', S', v \rangle$ has a Nash equilibrium in mixed strategies. Let $\sigma_i^j$ denote

the lottery over $S_i$ that such a mixed strategy equilibrium specifies. Since the profile $\sigma$ satisfies the condition for a Nash equilibrium, the strategy $\sigma_i(\theta_i = \theta_i^j) \equiv \phi_i(\theta_i)$ will satisfy the condition 6.1.          $\square$

## 9. Exercises

EXERCISE 6.1. *Consider the jury voting game where $p = \frac{3}{4}$ and $\pi = \frac{2}{3}$. Characterize the set of BNE to the game. Now assume that instead of majority rule, a version of unanimity rule is used – if all agents vote to convict the defendant is convicted, if at least one agent votes to acquit the defendant is acquitted. Characterize the BNE to this game (again assuming that $p = \frac{3}{4}$ and $\pi = \frac{2}{3}$).*

EXERCISE 6.2. *Consider the Jury Voting Game with a continuum of types. Prove equation 6.2.*

EXERCISE 6.3. *Consider a version of the Palfrey-Rosenthal model where $k$ contributions are required for the provision of the public good. However, assume that contributions are refunded if there are fewer than $k$. How does this modification effect the value of the cutpoint $\widehat{c}$? Now suppose that contributions in excess of $k$ are returned randomly to the agents. What happens?*

CHAPTER 7

# Extensive Form Games

Normal form representations of games are static in that all players choose their strategies simultaneously. However, many applications in political science involve players choosing strategies in sequence or in multiple stages. As we will see, it is possible to model these games in the normal form. However, it is often more convenient to model these games in the *extensive* form.

To motivate the extensive form, consider the following application from international relations. Assume that there are two countries $A$ and $B$ who are involved in a dispute over territory. We assume that $B$ controls the territory in question, thus the first stage of the game involves $A$'s decision about whether to *initiate* conflict by moving troops into the disputed region. After observing whether $A$ initiates, $B$ then decides whether to *acquiesce* and let $A$ maintain control or *escalate* in an attempt to expel $A$'s army from the territory. If country decides to escalate, it is successful in repelling $A$ with probability $p$.

We assume that, apart from the resources of the disputed region, each country has a national wealth of $a_0$ and $b_0$ and that the territory is worth an additional 4 units of national wealth to each country. $A$'s invasion of the region costs one unit as long as $B$ does not attack. However, an escalation by $B$ costs each country 6 units. The following table gives the payoffs from each of the possible outcomes.

| Table 7.1:  Escalation Game |
| --- |
| If $A$ does not initiate and $B$ acquiesces, $(a_0, b_0 + 4)$ |
| If $A$ does not initiate and $B$ escalates, $(a_0 - 6, b_0 - 2)$ |
| If $A$ initiates and $B$ acquiesces, $(a_0 + 3, b_0)$ |
| If $A$ initiates and $B$ escalates $(a_0 - 6 + 4p, b_0 - 6 + 4p)$ |

Suppose we were to model this game as a normal form where $A$ chooses whether or not to initiate and $B$ decides whether to acquiesce or escalate. Then we would be ignoring the fact that $B$ knows $A$'s choice when it makes its decision.

A better way of representing this game is using a game tree as in Figure 7.1. A game tree consists of *nodes* representing all of the

decisions made to date. Alternatively, we say that the nodes represent the *histories* There is an initial node at the beginning of the game (representing that nothing has happened). At each node there are *branches* representing the actions available to the player who chooses at that node. Each of these branches lead to the nodes of the next stage. The end of the game is represented by terminal nodes which specify the payoffs for each play of the game.

### Insert Figure 7.1 Here

At the initial node of our war game, $A$ makes its decision at the initial node from which there are two branches corresponding to *initiate* and *not initiate*. At each of the subsequent nodes, $B$ makes its choice between *acquiesce* or *escalate*. At each of the four terminal nodes, the corresponding payoffs are denoted.

Just as a matrix is used to represent normal form games, the game tree is a representation of the extensive form. The elements of the extensive game are:

(1) The set of agents $N$

(2) A set of histories of the game $H$. The elements of $H$ correspond to nodes of the game tree. $H^T$ is the set of terminal histories. By convention, the initial node is represented as $H^0 = \{\phi\}$.

(3) A mapping $p(h) : H \backslash H^T \rightarrow N$ assigns to each non-terminal history $h$ an agent who must make a decision at $h$.

(4) A set of actions $A(h)$ that $p(h)$ may take following history $h$. These may involve randomizations over actions.

(5) Information sets $I \subseteq H \backslash H^T$ which form a partition of the set of histories. If $h \in I$, $p(h)$ is uncertain whether she is at node $h$ or some other node $h' \in I$. In the game above, each player knows the history when it is called upon to play so that each information set contains a single element. We call such situations games of *complete and perfect* information (or simply perfect information). This assumption will be relaxed in later sections so that players may not observe all actions preceding their moves so that some information set $I$ contains multiple elements. These are games of *complete but imperfect* information (or simply imperfect information).

(6) Payoffs $U$, a list of Bernoulli utility functions $u_i(h) : H^T \rightarrow \mathbb{R}^1$ for each $i \in N$.

Thus, a finite extensive from game $\Gamma^E$ is the collection $\langle N, H, p(\cdot), U \rangle$. In the extensive form, a strategy is a complete plan of action. Therefore, it specifies a feasible action for the player in every history that

the player might be called upon to act. A formal definition of strategy follows.

DEFINITION 7.1. *Given an extensive form game* $\Gamma^E$, *a* **strategy profile** *for player* $i \in N$ *is a mapping* $s_i(h) : H_i \to A(h)$. *A strategy profile is a mapping* $s(h) : H \backslash H^T \to A(h)$.

Given this definition, we can specify the strategy sets for both players. Since $A$ only moves at the initial node, it strategy sets is simply $\{initiate, don't\ initiate\}$. However, $B$ strategies must be conditioned on each history. Thus, $B$ must choose from the following strategies: $\langle always\ acquiesce \rangle$, $\langle always\ escalate \rangle$, $\langle escalate\ if\ initiate,\ acquiesce\ otherwise \rangle$, $\langle acquiesce\ if\ initiate,\ escalate\ otherwise \rangle$.

Now that we have specified the strategies, it is easy to see that we can also represent this game as a normal form.

| Table 7.2: Escalation Game in Normal Form | | |
|---|---|---|
| B\A | *Initiate* | *Don't Initiate* |
| $\langle$*always acquiesce*$\rangle$ | $a_0 + 3, b_0$ | $a_0, b_0 + 4$ |
| $\langle$*always escalate*$\rangle$ | $a_0 - 7 + 4p, b_0 - 6 + 4p$ | $a_0 - 6, b_0 - 2$ |
| $\langle$*escalate if initiate, acquiesce otherwise*$\rangle$ | $a_0 - 7 + 4p, b_0 - 6 + 4p$ | $a_0, b_0 + 4$ |
| $\langle$*acquiesce if initiate, escalate otherwise*$\rangle$ | $a_0 + 3, b_0$ | $a_0 - 6, b_0 - 2$ |

With this representation, it is easy to see that there are three Nash Equilibria. The first two are the strategy profiles:

$$(Initiate,\ \langle always\ acquiesce \rangle)$$

and

$$(Initiate,\ \langle acquiesce\ if\ initiate,\ escalate\ otherwise \rangle)$$

Each of these predict that $A$ will invade the disputed region and $B$ will not respond. The third Nash equilibrium is the profile

$$(Don't\ Initiate,\ \langle escalate\ if\ initiate,\ acquiesce\ otherwise \rangle)$$

Of course, the third equilibrium predicts that $A$ will be deterred from entering the disputed region by the threat that $B$ will escalate.

This example shows some of the limitations of the Nash equilibrium concept in dynamic games. In particular, the predictions in the second and third equilibria are somewhat implausible. First, consider the second equilibrium. Suppose that Country $A$ defected from its equilibrium strategy and decided not to initiate. The equilibrium then calls for $B$ to escalate the conflict even though $A$ did not initiate.

Clearly, $B$ is worse off by carrying out its strategy. Thus, Nash equilibria allows for behavior that is not rational at histories that are "off the equilibrium path." One might discount the problem in the second equilibrium since it predicts the same behavior as equilibrium 1 which does not suffer from the problem – $B$ is happy to acquiesce if $A$ does not initiate. However, consider how the problem emerges in the third equilibrium. Again suppose that $A$ defected by choosing to initiate. If this happens $B$ is clearly better off acquiescing for any value of $p$. Thus, $B$'s threat to escalate is not credible. It would never rationally carry it out if it were called on to do so. Thus, the "peaceful" outcome is built on behavior which is not *sequentially rational* i.e. rational at every possible information set.

In the next couple of sections, we discuss refinements of Nash equilibria appropriate for dynamic games which eliminate strategies which are not sequentially rational. Next we discuss the concept of backward induction which eliminates non-credible threats and sequentially irrational behavior in games of perfect information. Then we will introduce games of imperfect information and the refinement of subgame perfection.

## 1. Backward Induction

The most common way of solving dynamic games of perfect information is through backward induction. In backward induction, we assume that the last player to act chooses the action at each node that maximizes her utility. The second to last player then chooses his actions optimally knowing that the last player will choose optimal actions at each node. This process is continued until each player has chosen optimally under the assumption that all future players will make optimal choices at each history.

It is easy to apply backward induction to our conflict game. First, we require that $B$ make optimal choices at each node. At the *initiate* node, $B$ clearly gets a higher utility from acquiescing The same is true at the *don't initiate* node so that $A$ knows that $B$ will always acquiesce. Given this knowledge, $A$ optimally chooses to initiate. Thus, the solution from backward induction is (*Initiate*, ⟨*always acquiesce*⟩) – the Nash equilibrium that did not involve sequentially irrational behavior.

Let's consider some other examples before formalizing the procedure.

**1.1. Application: The Centipede Game.** Figure 7.2 presents a game tree that has been studied extensively in experimental economics, known as the Centipede game. Two players take turns choosing

between $Down$ and $Left$. The choice of $D$ ends the game, but $L$ continues it until stage 5. One of the reasons that this game is of such interest to experimentalists is that a naive player 1 may attempt to continue to play $L$ in order to get the large payoff of 10 at stage 5. However, we will see such a strategy is not sequentially rational and does not survive backwards induction. We begin at stage 5 where player 1 will clearly play $L$. Backing up to stage 4, player 2 knows that player 1 will play left in the last stage which would give her a payoff of $-10$ so she does better playing $D$. Backing up one more stage, player knows that $L$ generates $-4$ while $D$ guarantees 3. Thus, he chooses $D$. Clearly, if we continue this process back to the first stage, we see that in fact player 1 will rationally choose $D$. Indeed, the only strategy profile that survives backward induction is $\{D, D, D, D, L\}$.

**Insert Figure 7.2**

**1.2. Application: Sequential Bargaining.** The application of bargaining models has become increasingly important in political game theory. Indeed, we dedicate an entire chapter to it later in the book. Here we consider one of the simplest versions. Assume that there are two players, 1 and 2, who are bargaining over how to allocate \$1. In the first period, player 1 proposes a division of the dollar where she keeps $x_1$ and gives $x_2 = 1 - x_1$ to player 2. If player 2 accepts, the dollar is divided accordingly and the game ends. However, if player 2 rejects, the value of the dollar decreases to $\delta$ where $1 > \delta > 0$. This is intended to capture the fact that the players are impatient in that they prefer to settle sooner than later. In round 2, player 2 may make an offer such that she keeps $x_2$ and gives $x_1 = \delta - x_2$ to player 1. If player 1 accepts, remaining $\delta$ is divided. However, if she rejects, the dollar disappears and both players get 0. For simplicity, we assume that the payoffs to each player are $u_i(x_i) = x_i$.

It turns out that there are lots of Nash equilibria to this game. In fact any allocation can be supported with Nash equilibrium strategies. To see this consider, the following strategy combination:

Player 1: Propose $x_2 = z$. If it is rejected, reject any offer in round 2.

Player 2: Accept in round 1 if $x_2 \geq z$, reject otherwise and then propose $x_2 = \delta$ in round 2.

Clearly, the best response of player 1 is to propose $x_2 = z$ in round 1 for any $z \leq 1$. Otherwise, player 1 would receive 0. Similarly, player 2's best response is to accept $z$. However, these strategies are clearly not sequentially rational. Player 1 does not profit by rejecting all second period proposals. He should accept any proposal that gives

at least as high a utility as the 0 received from rejecting. Thus, player 2 will get to keep all $\delta$ in round 2. Thus, she will accept in round 1 if and only $x_2 \geq \delta$. Knowing this, player 1's best proposal in round 1 is $x_1 = 1 - \delta$ and $x_2 = \delta$. This outcome is the only Nash equilibrium that survives backward induction.

**1.3. Solving Games via Backward.** Having seen a few examples, we can to generalize the notion of backward induction. Let $H^{T-1}$ be the set of histories of play that can be reached at stage $T - 1$. At each of these histories $h \in H^{T-1}$, backward induction the player who acts, $p(h)$, to choose its action optimally to maximize her utility. Thus, for each $h \in H^{T-1}$, $p(h)$ selects $a^*(h) = \arg\max_{a \in A(h)} u_{p(h)}((h, a))$.

Next consider the set of histories that immediately proceed $H^{T-1}$ and only lead to histories in $H^{T-1}$ (we denote this set by $H^{T-2}$). For each $h \in H^{T-2}$, $p(h)$ select the action $a \in A(h)$ which is optimal for $p(h)$ given the choices made from $H^{T-1}$ or $a^*(h) = \arg\max_{a \in A(h)} u_{p(h)}((h, a, a^*(h, a)))$. This process can be iterated to stage $k$ where we solve for $a^*(h) = \arg\max_{a \in A(h)} u_{p(h)}((h, a, a^*(h, a)))$ for each $h \in H^{T-k}$. This process continues to the initial node $H^0 = \{\emptyset\}$.

## 2. Dynamic Games of Complete but Imperfect Information

So far we have only considered models in which at every stage, the player who moves knows all of the previous moves and so knows exactly which game node she is at. Using the terminology of the last section, all information sets contain a single element. Now we consider models in which information sets contain multiple histories. Games of this form are said to have **imperfect information**. This can occur either because not all moves are observed or because some moves are taken simultaneously.

Let's first look at a simple game where some actions may not be observable. Consider a game between a bureaucrat $B$ and a politician $P$. The bureaucrat has to choose a regulatory enforcement level from $\{H, L\}$ which represent high and low respectively. High enforcement is costs $c > 0$ to $B$ but low enforcement is assumed to be costless. To keep things simple, we will assume that $B$ gets no utility from its enforcement. Therefore, it gets $-c$ for $H$ and 0 for $L$. However, assume that $P$ prefers $H$ to $L$ and that $u_P(H) = 1$ and $u_P(L) = 0$. $P$ cannot observe $B$'s enforcement level unless it conducts oversight of $B$ at a cost $1 > k > 0$. If $B$ is found to have chosen the lax enforcement, it suffers a penalty $f$ which we assume is greater than $c$ and is forced to choose $H$.

The main difference between this game and ones we have seen before is that $P$ does not know whether the history is $H$ or $L$ at the point at which she has to decide whether to conduct oversight. Consequently, $I(H) = \{H, L\}$ and $I(L) = \{H, L\}$. In extensive form given in Figure 7.3, we denote that $H$ and $L$ are in the same information set by connect the nodes with dotted lines.

### Insert Figure 7.3 Here

Since $P$ does not observe $B$'s action, she must play the same action at each node. Since a strategy by $B$ is simply a choice at the first node, we can write this game in the normal form.

| Table 7.3: Oversight Game | | |
|---|---|---|
| $B \backslash P$ | *Oversight* | *No Oversight* |
| $H$ | $-c, 1-k$ | $-c, 1$ |
| $L$ | $-f, 1-k$ | $0, 0$ |

Given the assumption that $f > c$, there are no pure strategy Nash equilibria in this game. If $B$ chooses $H$, $P$'s best response is to not conduct oversight, but the best response to no oversight is low enforcement in which case $P$ would prefer oversight. The mixed strategy equilibrium of this game involves $B$ choosing the high enforcement level with probability $1 - k$ and $P$ conducting oversight with probability $\frac{c}{f}$.

Now we turn to a familiar example to illustrate how the extensive form can accommodate simultaneous actions. The trick is to treat simultaneous moves as sequential ones in which subsequent players do not observe the action taken. Consider the prisoner's dilemma where two crooks have to decide whether or not to confess. We can model this is extensive form by letting player one move first and then placing both *confess* and *don't confess* in the same information set for player 2 as we have done in Figure 7.4.

### Insert Figure 7.4 Here

Finally to show how flexible the extensive form can be, consider the abstract game with three stages in Figure 7.5. Each player has three moves $Left$, $Middle$, $Right$. When player 1 plays left it is observed, but when player 2 plays right it is observed. Player 2 therefore has two information sets $\{L\}$, $\{M, R\}$. Player 3 has four $\{LL, LM\}$, $\{LR\}$, $\{ML, MM, RL, RM\}$, and $\{MR, RR\}$.

### Insert Figure 7.5 Here

The technical and conceptual difficulty with games of imperfect information is that we no longer can apply backward induction since players do not know which node they are on. We need a more general

notion of sequential rationality. Rather than assume that each player takes the best action given the node they are on, we assume that at each stage that can be conceptualized as a distinct game all the players play Nash equilibrium strategies. This concept is based on the idea of the *subgame.* A subgame is a subset of an extensive form that satisfies the following criteria:

(1) It begins at a node that is a singleton information set.
(2) It includes all nodes following this initial node, but only nodes that follow the initial node.
(3) It does not cut any information sets. If a histories $h$ and $h'$ are in the same information set, there are part of the same subgame.

The example in Figure 7.5 has three subgames: the original game, a subgame following $L$, and a subgame game following the history $LR$.

A formal definition of the set of subgames follows.

DEFINITION 7.2. *Given an extensive form game* $\Gamma^E$, *the set of* **subgames** *are all of the extensive form games constructed by selecting all* $h \in H \backslash H^T$ *which are singleton information sets and restricting* $H$, $p(\cdot)$ *and* $u_i(\cdot)$ *to histories that can be reached from* $h$.

Given the definition of subgames, the new requirement of sequential rational is that all agents play Nash equilibria in all subgames. Thus requirement is known as subgame perfect Nash equilibrium or SPNE.

DEFINITION 7.3. *Given an extensive form game* $\Gamma^E$, *a strategy profile* $s(\cdot)$ *is a subgame perfect Nash equilibrium (SPNE) if in every subgame to* $\Gamma^E$ *the restriction of the strategy profile* $s(\cdot)$ *to the subgame is a NE of the subgame.*

An important result establishes the existence of SPNE for finite games.

THEOREM 7.1. *Every finite extensive form game has a SPNE. Moreover, if no player is indifferent between any two terminal histories then the SPNE is unique.*

As an example we will consider a problem of sequential voting by 3 players $N = \{1, 2, 3\}$. Suppose that the choices $x, y, z$ are to be voted on with the agenda: choose between $x$ and $y$ first, and then compare the winner with $z$, enacting either the winner from the first vote or $z$ depending on which proposal gets the most votes. We assume that at each stage of voting ballots are cast simultaneously. Figure 7.6 depicts the game tree

**Insert Figure 7.6 Here**

We assume that the player have the following preferences over the enacted policy $xP_1yP_1z$; $yP_2zP_2x$; $zP_3xP_3y$. Applying subgame perfection and requiring that strategies are not weakly dominated (so voting is sincere) we see that if the final vote is between $x$ and $z$ then players 2 and 3 will vote for $z$. In contrast if the final vote is between $y$ and $z$ then players 1 and 2 will vote for $y$ Accordingly, in voting over $x$ and $y$ in the first period, strategic agents will anticipate that the real choice is between the **sophisticated equivalents**, $z$ and $y$. Accordingly players 1 and 2 will vote for $y$ over $x$. Note that player 1 prefers $x$ to $y$, but in a SPNE she casts a strategic vote for $y$ over $x$ because she realizes that a vote for $x$ is really a vote for $z$ which she finds very unappealing.

It is important to note that if we drop the requirement that voters do not use weakly dominated strategies, the set of SPNE can be quite large. Recall that any unanimous vote is a Nash equilibrium in any of the subgames at the second stage of the agenda. Thus, a large number of SPNE can be constructed by specifying different Nash equilibrium strategies for each second stage subgame.

As a second example, consider a model similar to one used by Weingast (1997) to explain the development of the rule of law. This game consists of a ruler $R$ who can choose whether or not to expropriate wealth $x$ from one of two social groups $A$ or $B$. After observing which group the ruler attempts to expropriate, $A$ and $B$ decide simultaneously whether or not to challenge him. Each incurs cost $c$ from challenging. If both challenge, the attempted expropriation fails and each receives a benefit $b$. A successful challenge also costs the ruler $k$. If one or zero groups challenge, the expropriation succeeds. The extensive form is shown in Figure 7.7.

**Insert Figure 7.7 Here**

We begin our analysis by computing the Nash equilibria of the subgame following the decision to expropriate from $A$. The normal form for this subgame can be represented as:

| Figure 7.4: | Expropriation Subgame | |
|---|---|---|
| $B \backslash A$ | Challenge | Don't Challenge |
| Challenge | $b, b$ | $-x, -c$ |
| Don't Challenge | $-x - c, 0$ | $-x, 0$ |

Clearly, there are two pure strategy Nash equilibria to this subgame corresponding to both groups challenging and to both groups

not challenging. There is also a mixed strategy equilibrium but we will ignore it to keep the example simple. Since the subgame following an attempted expropriation of $B$ is symmetric, there are also two pure strategy equilibria corresponding to both challenging and neither challenging.

Now we can back up to the first stage of the game where $R$ anticipates that the Nash equilibrium that is played in each of the subgames of the second stage. Suppose he anticipates that both groups will challenge in both subgames, then $R$'s best response is not to expropriate. Suppose that the groups challenge in only one of the subgames but not the other. Then $R$ will expropriate appropriate group. If there is no challenge in either subgame, the ruler might expropriate either. Thus, there are five SPNE in pure strategies. Weingast argues that the key to establishing the rule of law is that $A$ and $B$ coordinate on the Nash equilibrium where they both challenge attempted expropriations by the ruler.

It is important to note that solution via backward induction is just a special case of SPNE. Since in a game of perfect information, all information sets are singletons, each node begins a new subgame of the extensive form. Clearly, optimization at every node constitutes a Nash equilibrium of all subgames. Therefore, any solution using backward induction is a SPNE.

## 3. Subgame Perfection and Perfect Equilibria

One of the justifications of SPNE as a solution concept is its close relationship to Selten's perfect equilibrium concept. To see the close link, consider the normal form representation of our crisis game and the Nash equilibrium profile

$$(Don't\ Initiate,\ \langle escalate\ if\ initiate,\ acquiesce\ otherwise\rangle).$$

Suppose that we computed a completely mixed equilibrium where each strategy had to be played with at least probability $\varepsilon$. In particular, assume that country $A$ initiates with probability $\varepsilon$. Then country $B$'s expected utility from $\langle$escalate if initiate, acquiesce otherwise$\rangle$ is $(1 - \varepsilon)(b_0 + 4) + \varepsilon(b_0 - 6 + 4p)$ while its utility from always acquiescing is $(1 - \varepsilon)(b_0 + 4) + \varepsilon b_0$. Since the expected utility from $\langle$always acquiesce$\rangle$ is larger for any $\varepsilon > 0$, $B$ will want to play it with the maximum probability in the mixed equilibrium. As a result $\langle$escalate if initiate, acquiesce otherwise$\rangle$ cannot be the limit of completely mixed Nash equilibria. Thus, it not only fails the requirements of SPNE, but it is not a perfect equilibrium either.

However, while all SPNE are perfect, there are perfect equilibria that are not subgame perfect. This problem arises because extensive form games represented in the normal form often generate correlation in the trembles when the same player moves more than once in the extensive form.

## 4. Applications

### 4.1. Agenda Control.

4.1.1. *The Romer-Rosenthal Model.* Suppose, as is the case in many localities in the U.S., that local school budgets have to be approved by the voters. Further suppose that only the school board can place the measure on a referendum ballot. Thus, formally, the board has monopoly agenda control over proposals for school spending $s \in [0, \infty)$. We assume that the school board would like to maximize the amount of spending so that $u_B(s)$ is always increasing in $s$.

Once the referendum has been placed on the ballot, voters decide via majority rule whether or not to approval it. We assume that all voters will turnout so that their only possible strategies are $\{Y, N\}$. If a majority chooses $Y$, then $s$ becomes the new level of spending. If a majority chooses $N$, then some reversion (or status quo) spending level $q$ is adopted. We assume that voters have single peaked and symmetric preferences over school spending $u_i(s)$. Let $v_i$ be the ideal point of voter $i$. As we showed in chapter 2, such preferences take the form $u_i(s) = h(-|s - v_i|)$.

We propose to solve this game using subgame perfect equilibrium. However, note that since the last stage of the game is a majority rule voting game, there are always Nash equilibria where the proposal is accepted and equilibria where it is rejected for any $s$. Thus, we will assume that voters do not use weakly dominated strategies in the voting subgame. Thus, given any proposal $s$, each voter will vote $Y$ if $u_i(s) \geq u_i(q)$. Note that single-peakedness implies that if voter $i$ prefers $q$ to $s$, all voters with ideal points to the left of $i$ do so as well. Further, if voter $i$ prefers $s$ to $q$, all voters to the right of $i$ do so as well. Thus, if the median voter votes $Y$ then the proposal will pass.

Given the equilibrium of the voting subgame, the school board's best response is to choose the largest $s$ that is acceptable to the median voter. Let $v_m$ be the ideal point of the median voter. Given $v_m$ and $q$, we can compute which policies that the median prefers to $q$ or $u_m(s) \geq u_m(q)$. Note that since $h$ is a non-decreasing function, this inequality requires that

$$-|s - v_m| \geq -|q - v_m|$$

Therefore, if $q < v_m$, this inequality holds for $s \in [q, 2v_m - q]$. Conversely, if $q > v_m$, a successful proposal requires $s \in [2v_m - q, q]$. Thus, the highest obtainable policy the school board can get is the maximum of $2v_m - q$ and $q$. Since the board wants to maximize $s$, it will choose $\max\{q, 2v_m - q\}$. Figure 7.8 plots the equilibrium value of $s^*$ as functions of $v_m$ and $q$.

## Insert Figure 7.8 Here

This simple model produces some clear predictions about the relationship between voter preferences, statutory reversions, and policy outcomes. First, note that the board can use its agenda control to generate higher spending outcomes when the statutory reversion is low so long as the median prefers to spend more than the reversion amount. This is because the voter's threat to reject large spending proposals is not credible when the reversion is bad. A second important implication is that while spending outcomes are responsive to changes in voter preferences (at least when $q < v_m$), spending grows twice as fast as the median voter's preferred spending level.

4.1.2. *The Presidential Veto.* In the United States and many other presidential systems, the executive has a veto power over legislative enactments. We can use a version of the Romer-Rosenthal model to explore how the veto enhances the executive's influence over legislation.

To keep things as simple as possible, we will model the legislature as a single actor $L$ who has single-peaked symmetric preferences on a single dimension with an ideal point of $l$. Thus, we denote the legislature's preferences as $u_l(x) = h(-|x - l|)$ for policy outcomes $x \in \mathbb{R}$. Extending the model to the case where the legislature is a collectivity is straightforward. Similarly, we assume that the president's has an ideal point $p$ and preferences given by $u_p(x) = h(-|x - p|)$.

The game form is very simple. In the first stage, $L$ proposes a bill $b$ to change the status quo policy $q$. Subsequently, the president $P$ decides whether to accept $b$ or to veto it the bill which results in maintaining the status quo $q$. Thus, we ignore the legislature's ability to override vetoes, an issue we take up in the next section.

We can solve this game very easily using backward induction. Clearly, in the last stage, the president's best response is to accept any bill for which $u_p(b) \geq u_p(q)$ or $-|b - p| \geq -|q - p|$. Thus, if $p > q$, she will accept any $b \in [q, 2p - q]$. Alternatively, if $p < q$, she will accept $b \in [2p - q, q]$. Let $P(q)$ denote the set of bills that the president will accept over the status quo. Now we back up to the legislature's decision node. Since the legislature knows which policies will be accepted, it will choose its most preferred policy from $P(q)$. If $l \in P(q)$, then

clearly $b^* = l$. If $c$ is below min $P(q)$, then $b^* =$min $P(q)$. If $l >$ max $P(q)$ , then $b^* =$max $P(q)$.

Suppose that $l > p$. Then, given our derivations of $P(q)$, the equilibrium policy outcome is

$$b^* = \begin{cases} 2p - q \text{ if } p > q \text{ and } l > 2p - q \\ l \text{ if } p > q \text{ and } l < 2p - q \\ l \text{ if } l < q \\ q \text{ if } l > q > p \end{cases}$$

If $p > l$, the equilibrium outcomes are

$$b^* = \begin{cases} 2p - q \text{ if } p < q \text{ and } l < 2p - q \\ l \text{ if } p < q \text{ and } l > 2p - q \\ l \text{ if } l > q \\ q \text{ if } l > q > p \end{cases}.$$

Figure 7.9 plots the equilibrium outcomes as a function of $l$, $p$, and $q$. The comparative statics results are quite similar to the original Romer-Rosenthal model. In particular, the legislature does better off when the status quo is far from the president's ideal point. Another important implication is that the influence conferred by the veto is not very large. In all of the cases that the veto has some impact i.e. $b^* \neq l$, the president is indifferent between the equilibrium proposal and the status quo. Finally, since the model is one of perfect information, the legislator perfectly predicts the president's behavior and no vetoes occur in equilibrium. In later chapters, we consider models in which vetoes may occur as part of equilibrium strategies.

### Insert Figure 7.9 Here

4.1.3. *The Veto Override.* Now we consider a simple extension to the model of the previous section. Instead of assuming that $q$ is the outcome following any veto, we consider a model where the legislature can override the veto with a supermajority. Assume that the legislature has $n$ members and that $k > \frac{n+1}{2}$ votes are need to override the executive veto. Further, assume that each legislator has single peaked preferences of the form $u_i (x) = h (- |x - l_i|)$ and that the ideal points $l_i$ are ordered such that $l_i > l_j$ if and only in $i > j$. Motivated by a model in which legislative proposals are made according to an open rule agenda process, we assume that the legislative proposer is the median with ideal point $m \equiv l_{(n+1)/2}$.

Given these assumptions (most importantly single-peakedness), a successful override requires that $u_k (b) \geq u_k (q)$ and $u_{n-k-1} (b) \geq u_{n-k-1} (q)$. To see that this is true, consider the case where $u_k (b) \geq u_k (q)$ and $u_{n-k-1} (b) < u_{n-k-1} (q)$. Since preferences are single-peaked, this means

that there is some $i \in [n - k - 1, k]$ such that $u_i(b) < u_i(q)$ for all legislators with ideal points lower that $l_i$. Therefore, the number who support the override must be strictly less than $k$. The logic of the other possibilities is similar. Because their support is necessary and sufficient, legislators $n - k - 1$ and $k$ are commonly referred to as the *override pivots.*

Since an override is only necessary in case of a presidential veto, only one of the override pivots is strategically relevant. To see this consider a hypothetical vetoed bill where $u_p(b) < u_p(q)$ and $u_m(b) > u_m(q)$. If $p < m$, single peakedness and the fact $l_k > m$ imply that $u_k(b) > u_k(q)$. Thus, the override depends solely on $n - k - 1$'s preferences. Similarly, if $p > m$, single peakedness and $l_{n-k-1} < m$ imply that $u_{n-k-1}(b) > u_{n-k-1}(q)$. The implication is that only the preferences of the pivot that lies on the same side of the median as the president matter for a successful override.

Thus far we have established that it is necessary for the proposer to attract the support of either the president or the override pivot on his side median. Now we consider how the proposer chooses her optimal proposal. First suppose that $l_{n-k-1} < p < m$. Again using singlepeakedness, we know that any bill that $l_{n-k-1}$ and $m$ prefer to $q$, $p$ also prefers. Thus, the proposer does not require the support of $l_{n-k-1}$. A similar argument establishes that when $p < l_{n-k-1} < m$, the proposer need only attract $n - k - 1$'s support. The corresponding cases where $p > m$ are symmetric so that we know that the proposer need only attract the closer of the president and the override pivot on the president's side of the median.

Thus, we can define the pivotal actor and her ideal point as $v = \max(l_{n-k-1}, p)$ if $p < m$ and $v = \max(l_k, p)$ otherwise. Furthermore, the game can be treated as a direct application of the Romer-Rosenthal model where the ideal point of the veto player is $v$. Thus, the SPNE proposal is given by:

$$
b^* = \begin{cases}
2v - q & \text{if } v > q \text{ and } m > 2v - q \\
m & \text{if } v > q \text{ and } m < 2v - q \\
m & \text{if } m < q \\
q & \text{if } m > q > v
\end{cases}
$$

if $v > m$, and

$$
b^* = \begin{cases}
2p - q & \text{if } p < q \text{ and } m < 2p - q \\
m & \text{if } p < q \text{ and } m > 2p - q \\
m & \text{if } m > q \\
q & \text{if } m > q > p
\end{cases}
$$

otherwise. Figure 7.10 illustrates the equilibrium outcomes for differ-
ent values of $k$. Not surprisingly, when the number of votes need to
override goes down, the effect of the veto power is diminished.

### Insert Figure 7.10 Here

4.1.4. *Structure Induced Equilibrium.* Recall from chapter 2 that
the core of multi-dimensional majority rule voting models is typically
empty. Primarily to explain how legislatures overcame this "chaos"
problem, Shepsle (1979) developed the idea of "structure-induced equi-
librium." The basic idea is that legislative institutions such as commit-
tee systems restrict the types of legislative proposals and amendments
that may be considered, and that such restrictions lead to non-empty
legislative cores.

While Shepsle's initial work is developed within the paradigm of
social choice, we present his various models as extensive form games
which we solve for subgame perfect Nash equilibria.

To keep things simple, we focus on a legislature with $n$ members
and two committees with ideal points $C_1$, and $C_2$. The policy space is
assumed to be a subset of $\mathbb{R}^2$. Player $C_1$ represents a committee with
jurisdiction over dimension 1 while $C_2$ has jurisdiction over dimension
2. In the absence of legislative activity, we assume that the status quo
$q = (q^1, q^2)$ remains in force. Each legislator has quadratic policies
$(x^1, x^2)$ given by

$$u_i\left(x^1, x^2\right) = -\left(x^1 - l_i^1\right)^2 - \left(x^2 - l_i^2\right)^2$$

where $(l_i^1, l_i^2)$ is the ideal points of legislator $i$. We also assume that
each committee has quadratic preferences with ideal points $(c_1^1, c_1^2)$,
and $(c_2^1, c_2^2)$ respectively. An important feature of these preferences is
that they are separable across dimensions. Each players preferences
on dimension 2 are independent of outcomes on dimension 1 and vice
versa.

We know consider several extensive forms representing various leg-
islative institutions.

The Open Rule with Germaneness. The first extensive form we con-
sider is the open rule with a germaneness requirement. In this game,
each committee sequentially reports a bill to change the status quo on
its dimension. Thus, $C_1$ makes a proposal $b_1$ to change $q^1$ and $C_2$
makes a proposal $b_2$ to change $q^2$. Each bill may be amended, but only
on the germane dimension. Thus, amendments to $b_j$ can only move it
along dimension $j$.

Let's begin in the last stage of the game where $c_2$ proposes $b_2$. Since
amendments to move the bill along dimension 2 can be freely made,

the median voter theorem suggests that final outcome will be the ideal point of the median voter on dimension 2. We denote this ideal point as $m^2 = median\{l_i^2\}$. Thus, $c_2$ has a weakly dominant strategy to propose $b_2 = m^2$. Note that this results depends on the fact that preferences are separable so that the outcome on dimension 1 does not effect dimension 2 preferences. Clearly, by the same logic, $c^1$ will propose $b_1 = m^1 = median\{l_i^1\}$. Finally, note that the separability of preferences implies that the order in which the committees make proposals does not matter.

The policy $(m^1, m^2)$ (or the dimension by dimension median) is the structure induced equilibrium under the open rule with germaneness. It is precisely because amendments can only be made on one dimension at a time that we get this generalization of the median voter result rather than a majority rule cycle.

Gatekeeping. While in the last section, we assumed that each committee had to make a proposal, now we assume that committees have discretion over whether to make a proposal at all. If committee $j$ "keeps the gates closed" on dimension $j$, the policy outcome remains $q^j$. However, if the committee does report a a bill it can be freely amended subject to germaneness.

First consider committee 2's decision. If it opens the gates, the policy on its dimension will be $m^2$. Thus, it exercises gatekeeping if $-\left(b_1^* - c_2^1\right)^2 - \left(q^2 - c_2^2\right)^2 > -\left(b_1^* - c_2^1\right)^2 - \left(m^2 - c_2^2\right)^2$. Simple algebra reveals that the committee will close the gates if and only if $c_2^2 > q^2 > m^2$ or $m^2 > q^2 > c_2^2$. When one of these conditions hold, there are no policy revisions to the status quo on dimension 2 that the committee and the 2nd dimension median simultaneously prefer. Thus, the committee reports no bill. Since the case of committee 1 is symmetric, we obtain gatekeeping when $c_1^1 > q^1 > m^1$ or $m^1 > q^1 > c_1^1$. Thus, the policy outcomes $(x^{1*}, x^{2*})$ of the SPNE are

$$x^{j*} = \begin{cases} q^j \text{ if } c_j^j > q^j > m^j \text{ or } m^j > q^j > c_j^j \\ m^j \text{ otherwise} \end{cases}$$

The Closed Rule. Finally, suppose that committees make proposals under closed rules so that amendments are not allowed. Therefore, they are Romer-Rosenthal agenda setters within their jurisdictions. We can directly apply the results of section 4.1.2 where the proposer is $c^j$ and the vetoer is $m^j$. Note that gatekeeping powers are irrelevant in this game since committee $j$ can do at least as well as the gatekeeping outcome by proposing $q^j$. In fact, as Groseclose and Krehbiel (2002) point out, both $c^j$ and the $jth$ dimension median are better off under the closed rule than the open rule with gatekeeping.

**4.2. A Model of Power Transitions.** Our next application is based on the work of Powell (1999) who uses similar model to study how dramatic shifts of power in the international system might lead to violent conflict. Suppose that there are two countries $A$ and $B$. Country $A$ is making a claim against a region controlled by $B$. The total value of the region is normalized to \$1 per period. We focus on a two-period version of this game and assume that each country weigh the outcome of each period equally.

First consider country $B$'s options. It can appease $A$ in each period $t$ by offering it a share of the region's output $0 \leq x_t \leq 1$ or it can attempt to settle the dispute militarily by attacking $A$. If $B$ chooses to attack and wins the war, $A$ drops its claim and the game ends. If $B$ loses the war, $A$ takes undisputed control of the region. Country $A$'s available choices in period $t$ is either to accept $x_t$ or to refuse it and go to war. Fighting a war costs $c$ to both sides.

An important feature of the environment that Powell studies is that country $A$'s military capability is increasing relative to $B$'s over time. Assume that in the first period, A wins a war with probability $p_1$ where as in the second period $A$ wins with probability $p_2 > p_1$. To keep things interesting, we assume that $p_2 > c$.

With respect to the incidence of violent conflict, there are two types of equilibria that we might observe: one in which $B$ appeases $A$ in both periods and one in which $B$ attacks $A$ in the first period.

First, suppose that $2c > p_2 - 2p_1$. Our claim is that in this case, the SPNE equilibrium is one where $B$ gives $p_2 - c$ to country $A$ in the second period and $max\{0, 2p_1 - p_2\}$ in the first period. Since this is a game of perfect information, we can verify that these strategies are part of a SPNE using backward induction. In the second period, $A$'s expected utility of fighting is $p_2 - c$ so that $B$ must offer at least as much to avoid a conflict, leaving it with a payoff of $1 - p_2 + c$. Since $B$'s expected utility of fighting is $1 - p_2 - c$, it strictly prefers appeasing. Now consider period 1. $A$ receives \$1 in each period if it wins a war and 0 if it loses. Therefore, the expected utility of fighting is $2p_1 - c$. Therefore, to appease $A$, $B$ must choose $x_1$ so that $x_1 + p_2 - c \geq 2p_1 - c$ or $x_1 \geq 2p_1 - p_2$. Since $B$ will rationally offer the minimum amount $x_1^* = \max\{0, 2p_1 - p_2\}$. Now we need only check to see that $B$ would prefer to pay $x_1^* + x_2^*$. Since $B$'s expected utility of war is $2(1 - p_1) - c$, it prefers the payment if and only if $2c \geq p_2 - 2p_1 + x_1^*$ which always holds when $2c > p_2 - 2p_1$.

Now suppose that this condition does not hold so that $p_2 - 2p_1 > 2c$. Since $c > 0$, this requires that $x_1^* = 0$. Thus, $B$ would prefer to fight

than to make the payments since $p_2 - 2p_1 + x_1^* > 2c$. Thus, when the condition fails, the SPNE has $B$ attacking $A$ in the first round.

To generate intuition, note that the necessary condition for a peaceful resolution only fails when $p_2$ is much greater than $p_1$ so that $A$ is much weaker in the first stage than it will be in the second stage. Thus, $B$ prefers to attack when $A$ is weak to avoid making large concession when $A$ becomes more powerful. If the distribution of power were stable, their would be no war in equilibrium.

### 4.3. A Model of Transitions to Democracy.

Acemoglu and Robinson (forthcoming) develop a number of models designed to explore the conditions under which authoritarian polities will adopt democratic institutions. In this section, we provide a simple sketch of their framework and one of their models.

Suppose that there are two types of agents: rich and poor. Let $\lambda > \frac{1}{2}$ be the proportion of citizens who are poor while $1 - \lambda$ is the proportion of rich citizens. Since these agent differ in their incomes, they have different preferences over tax rates. Rich citizens each receive income $y^r$ and poor citizens have income $y^p$. The average income in the society is $\lambda y^p + (1 - \lambda)y^r$. Clearly, $y^r > \overline{y} > y^p$. An important parameter in Acemoglu and Robinson's analysis is $\theta$ which represents the share of income held by the poor so that

$$y^p = \frac{\theta \overline{y}}{\lambda} \text{ and } y^r = \frac{(1 - \theta)\overline{y}}{1 - \lambda}.$$

Thus, an increase in $\theta$ represents a decrease in inequality.

The primary policy instrument in this political economy is a linear tax and transfer scheme where the government sets a proportional tax rate $\tau$ and then transfers the tax revenue back to the citizens in each lump sum. Given a tax rate $\tau$, the per capita tax levy is $\tau \overline{y}$. However, as a simple way of capturing the distortionary effects of income taxation, Acemoglu and Robinson assume that revenues are lower than the levy by a function $C(\tau)\overline{y}$. To keep things simple and get a closed form solution, we assume that $C(\tau) = \frac{1}{2}\tau^2$. Thus, the distortion is an increasing convex function of the tax rate. After deducting this deadweight loss, the amount of money available for transfers is given by $T = (\tau - C(\tau))\overline{y}$ and the after-tax and transfer income of each agent is

$$V^i(\tau) = (1 - \tau)y^i + \left(\tau - \frac{1}{2}\tau^2\right)\overline{y}$$

Now we consider the preferred tax rates by rich and poor voters. The first order condition for the optimal tax rate choice by poor voters is

$$\overline{y} - y^p - \tau\overline{y} = 0$$

Using the fact that $\frac{\theta\overline{y}}{\lambda}$, we can write the poor's most preferred tax rate as

$$\tau^p = \frac{\lambda - \theta}{\lambda}$$

Since the poor have lower incomes than the rich, their income share is lower than their share in the population so that $\lambda > \theta$. Thus, $0 < \tau^p < 1$. Also note that $\tau^p$ is decreasing in $\theta$ so that the poor's preferred tax rate is increasing in inequality.

Now consider the preferences of the rich. Their first order condition is

$$\overline{y} - y^r - \tau\overline{y} = 0$$

This produces an infeasible negative tax rate. Thus, the rich's most preferred feasible tax rate is $\tau^r = 0$.

In Acemoglu and Robinson's model, there is a political shock in each period which determines the consequences of overthrowing the regime and replacing it by a dictatorship of the left. They assume that when the shock is $S$, $1 - \mu^S$ of the economy's income is destroyed where $S = H, L$ and $\mu^H > \mu^L$. Thus, in state $H$ the costs of overthrowing the regime are low compared to state $L$. During a revolution the income of the rich is confiscated and evenly divided among the poor, thus is state $S$ the payoff to the poor is given by

$$V^p\left(R, \mu^S\right) = \frac{\mu^S\overline{y}}{\lambda}$$

For simplicity, they assume that following the revolution the rich get no income so that $V^r\left(R, \mu^S\right) = 0$.

To parameterize the outcomes of Acemoglu and Robinson's model, we need to consider the extent to which revolution is a threat. We say that the revolution constraint binds in state $S$ if the poor prefer revolution to an authoritarian outcome at the rich's ideal tax rate of zero or that $V^p\left(R, \mu^S\right) > V^p\left(0\right)$. From substituting the above expressions, we find that this constraint binds when $\mu^S > \theta$.

Given this specification of the economy and the costs of revolution, we turn to one of Acemoglu and Robinson's extensive form games, illustrated in Figure 7.11. To simplify the figure, we show only the extensive form following the realization of the state $S$.

**Insert Figure 7.11 Here**

First, the state, $H$ or $L$, is revealed. Then the rich move first and decide whether to move to a democracy $D$ or to maintain control in an authoritarian system $N$. If the rich choose $N$, they also choose a tax rate $\widehat{\tau}$.

After the rich make their decisions, the poor move next and decide whether to initiate a revolution $R$ or to accept the rich's decision $(NR)$. If they revolt, the payoffs are $V^p\left(R, \mu^S\right) = \frac{\mu^S \overline{y}}{\lambda}$ and $V^r\left(R, \mu^S\right) = 0$. If they do not revolt against $D$, the tax rate is chosen by majority rule. Since the median voter is poor, the equilibrium tax rate is $\tau^p$ and the payoffs to $D$ are therefore $V^r\left(\tau^p\right)$ and $V^r\left(\tau^p\right)$ for the rich and poor respectively. By comparing the payoffs of $R$ and $D$, it is easy to establish that in state $S$ the poor will prefer to revolt rather than accept democracy if and only if $\mu^S > \theta + \tau^p\left(\lambda - \theta\right) - \frac{1}{2}\tau^{p2}\lambda$ or $\mu^S > \theta + \frac{(\lambda - \theta)^2}{2\lambda}$.

Now suppose the rich chose $N$ and the poor prefer not to revolt. Acemoglu and Robinson assume that the rich may not be able to commit to maintaining $\widehat{\tau} > 0$ after the revolutionary threat has passed. To model this commitment problem, they assume that with probability $p$ rich maintain the tax rate $\widehat{\tau}$ but with probability $1 - p$ they have the opportunity to renege and choose $\tau^r = 0$. Given the rich's initial choice of tax rate and the possibility of reneging, we can compute that the utilities from $N$ are

$$V^p\left(N, \widehat{\tau}\right) = (1 - p)\, y^p + p\left[(1 - \tau)\, y^p + \left(\tau - \frac{1}{2}\tau^2\right)\overline{y}\right]$$

$$= y^p + p\left[\tau\left(\overline{y} - y^p\right) - \frac{1}{2}\tau^2\overline{y}\right]$$

and

$$V^r\left(N, \widehat{\tau}\right) = (1 - p)\, y^r + p\left[(1 - \tau)\, y^r + \left(\tau - \frac{1}{2}\tau^2\right)\overline{y}\right]$$

$$= y^r + p\left[\tau\left(\overline{y} - y^r\right) - \frac{1}{2}\tau^2\overline{y}\right]$$

if the poor do not revolt. Given these payoffs, it is easy to see that the poor will prefer to revolt against $N$ if

$$\mu^S > \theta + p\left[\widehat{\tau}\left(\lambda - \theta\right) - \frac{1}{2}\widehat{\tau}^2\lambda\right]$$

In order to reduce the number of cases, we follow Acemoglu and Robinson and assume $\mu^L < \theta$ so that the poor never revolt in state $L$. This leaves us with three cases:

(1) Suppose that $\mu^H < \theta$. Then the revolution constraint binds in neither case. Thus, the unique SPNE consists of $N$, a tax rate of zero, and no revolution.

(2) Suppose that $\mu^H > \theta + \frac{(\lambda-\theta)^2}{2\lambda}$. Then even democracy does not deter the poor from revolting, so a revolution occurs.

(3) Suppose that $\theta + \frac{(\lambda-\theta)^2}{2\lambda} > \mu^H > \theta$. In this case, it may be possible to the rich to avoid a revolution by accommodating the poor with a tax rate $\widehat{\tau}$. From above we know that doing so requires that the Rich set the tax rate so that

$$p < \frac{\mu^H - \theta}{\widehat{\tau}\,(\lambda - \theta) - \frac{1}{2}\widehat{\tau}^2 \lambda}$$

However, if $p < \frac{\mu^H - \theta}{\tau^p(\lambda-\theta) - \frac{1}{2}\tau^{p2}\lambda} = \frac{2\lambda\left(\mu^H - \theta\right)}{(\lambda-\theta)^2}$ the rich will prefer to choose $D$ than to set the tax rate higher than $\tau^p$. Thus, there is a critical value of $p^* = \frac{2\lambda\left(\mu^H - \theta\right)}{(\lambda-\theta)^2}$ such that democracy is the outcome if $p^* > p$. Thus, when the rich have difficulty committing to a high tax rate, they can avoid revolution by transitioning to democracy.

To generate some predictions about when democratic transitions are likely to occur, we can look at how $p^*$ is affected by changes in the parameters. Not surprisingly, $p^*$ is increasing in $\mu^H$ suggesting that when the costs of revolution are low, the rich is more likely to support democratization. Secondly, $p^*$ and the likelihood of democracy are decreasing in $\theta$. This occurs because greater inequality makes a revolution a more attractive option for the poor. In turn, the rich have to make more concessions to prevent it. If committing to these concessions is sufficiently difficult, a democratic transition will occur.

**4.4. A Model of Coalition Formation.** One of the earliest applications of political game theory is the study of coalition formation (Riker 1962). While the earliest models were developed within the cooperative game theoretic and social choice traditions, there have been a number of recent applications using non-cooperative bargaining models.

In this section, we look at one such model of coalition governments developed by Banks and Austen-Smith (1989). Assume that there are three parties $\alpha$, $\beta$, and $\gamma$ where $\Omega = \{\alpha, \beta, \gamma\}$, who have known policy positions $p_\alpha$, $p_\beta$, and $p_\gamma$ on a single dimension policy-space $P \subset \mathbb{R}$ where $p_\alpha > p_\beta > p_\gamma$. Let $w = \{\omega_\alpha, \omega_\beta, \omega_\gamma\}$ be the vector of votes shares for party in the last election where we assume that all vote shares are

less than $\frac{1}{2}$ so that the government has to be a coalition. To simplify matters, we will assume that these vote shares are exogenous parameters, whereas Banks and Austen-Smith derive them endogenously from a model of voting. We will also be interested in the vote shares for parties in various coalitions. Let $C \subset \Omega$ be a coalition then the vote share of each coalition is given by

$$\omega_C = \sum_{k \in C} \omega_k$$

We say that $C$ is a winning coalition if $\omega_C > \frac{1}{2}$. So let $D(w) = \left\{ C \subset \Omega : \omega_C > \frac{1}{2} \right\}$ be the set of winning coalitions and $D_k(w) = \{C \subset D(w) : k \in C\}$ be the set of winning coalitions that include party $k$.

The three parties will bargain over the formation of a new government. In doing so, they will choose a policy $y \in P$ and allocate a fixed set of portfolios $G$. To keep things simple, we follow Austen-Smith and Banks and assume that $G$ is infinitely divisible. We assume that the allocations $\mathbf{g} = \{g_\alpha, g_\beta, g_\gamma\}$ satisfy $\sum_{k \in C} g_k = G$.

We assume that each party has quadratic preferences over policy and additive linear preferences in portfolios. Therefore, the payoff to party $k$ to policy $y$ and allocation $\mathbf{g}$ is given by

$$- (y - p_k)^2 + g_k$$

The protocol for bargaining is as follows. First, the party with the largest vote share, say $k$, is selected as formateur and chooses a coalition from $D_k(w)$. The formatuer then proposes a policy $y_k$ and an allocation $\mathbf{g}$. If its coalition partners accept, $y_k$ and $\mathbf{g}_k$ are implemented and the game ends. However, if one of the coalition partners vetoes, the second largest party becomes the formatuer, selects a coalition from $D_l(w)$, and proposes $y_l$ and $\mathbf{g}_l$. If this is defeated, the smallest party becomes the formatuer. If the smallest party is unsuccessful, a caretaker government takes office and maintains a status quo policy $p_q$ and chooses $\mathbf{g} = \{0, 0, 0\}$.

We can solve this game via backward induction. The payoffs to party $k$ from a caretaker government are $v_k^c = - (p_q - p_k)^2$. To simplify, we will assume that $v_k^c < - (p_j - p_k)^2$ for all $j$ and $k$ so that any party $k$ prefers party $j$'s ideal point and a zero share of the portfolios to a caretaker government. Formally, this requires that $p_q \notin [2p_\gamma - p_\alpha, 2p_\alpha - p_\gamma]$.

We first consider the example of vote shares such that $\omega_\alpha > \omega_\beta > \omega_\gamma$. So consider the third stage where party $\gamma$ is the formateur. By assumption, all parties prefer $y_\gamma = p_\gamma$ and $\mathbf{g}_\gamma = (0, 0, G)$ to the caretaker government so this must be party $\gamma$'s optimal choice.

Now consider party $\beta$'s choice. The utilities of defeating any proposal by party $\beta$ and moving to party $\gamma$'s proposal stage are $v_\gamma^\gamma = G$ and $v_\alpha^\gamma = -\left(p_\gamma - p_\alpha\right)^2$. Since $\gamma$ receives the highest possible utility from voting against $\beta$'s offer, $\beta$ must make an offer to party $\alpha$. However, note that $\alpha$ prefers $y_\beta = p_\beta$ and $\mathbf{g}_\beta = (0, G, 0)$ to a government formed by $\gamma$, so $\alpha$ will accept $\beta$'s ideal point and no portfolios.

Now we back up to the first stage of the game where $\alpha$ makes a proposal. The utilities from defeating party $\alpha$'s proposal and moving to party $\beta$'s stage are now $v_\beta^\beta = G$ and $v_\gamma^\beta = -\left(p_\beta - p_\gamma\right)^2$. Clearly, $\alpha$ has nothing to offer $\beta$ and will thus try to form a coalition with $\gamma$. Thus, $\alpha$ will choose $y_\alpha$ and $g_\gamma$ to maximize $-\left(y_\alpha - p_\alpha\right)^2 + G - g_\gamma$ subject to $-\left(y_\alpha - p_\gamma\right)^2 + g_\gamma \geq -\left(p_\beta - p_\gamma\right)^2$ and $G \geq g_\gamma \geq 0$. There are three cases depending on whether there are corner solutions $g_\gamma^* = G$ or $g_\gamma^* = 0$.

(1) If $p_\alpha - p_\beta \geq p_\beta - p_\gamma$ and $G \geq \frac{1}{4}\left(p_\alpha - p_\gamma\right)^2 - \left(p_\beta - p_\gamma\right)^2$, $y_\alpha^* = \frac{p_\alpha + p_\gamma}{2}$ and $g_\gamma^* = \frac{1}{4}\left(p_\alpha - p_\gamma\right)^2 - \left(p_\beta - p_\gamma\right)^2$.

(2) If $p_\alpha - p_\beta \geq p_\beta - p_\gamma$ and $G < \frac{1}{4}\left(p_\alpha - p_\gamma\right)^2 - \left(p_\beta - p_\gamma\right)^2$, $y_\alpha^* = p_\gamma + \sqrt{G + \left(p_\beta - p_\gamma\right)^2}$ and $g_\gamma^* = G$.

(3) If $p_\alpha - p_\beta < p_\beta - p_\gamma$, $y_\alpha^* = p_\beta$ and $g_\gamma^* = 0$.

In the first two cases, the distance from $\alpha$'s ideal point and $p_\beta$ is greater than the distance from $p_\beta$ to $\gamma$'s ideal point. Thus, $\alpha$ is more willing to give up portfolios in favor of a policy better than $p_\beta$ than party $\gamma$ requires as compensation. Thus, party $\alpha$'s offers a compromise policy. When $G$ is sufficiently large, $\alpha$ offers the compromise policy $y_\alpha^* = \frac{p_\alpha + p_\gamma}{2}$ which reflects an optimal trade-off in its policy goals and its desire to hold portfolios. When $G$ is small, however, $\alpha$ is willing to give up all of the portfolios in order to move policy in the direction of its ideal point. Finally, in the last case, $\alpha$ is sufficiently well off under $p_\beta$ compared to $\gamma$ that $\alpha$ is unwilling to compensate $\gamma$ for moving policy towards its ideal point. An interesting feature of this outcome is that the coalition is a *non-connected* one of the extreme parties. This is in contrast to arguments stressing that policy motivated parties will seek form coalitions with ideological allies (Axelrod 1970).

Now consider the case where $\omega_\beta > \omega_\alpha > \omega_\gamma$. Once again $\gamma$ will choose $y_\gamma^* = p_\gamma$ and $\mathbf{g}_\gamma = (0, 0, G)$ in the last period. Now consider $\alpha$'s choice in the second period. Clearly, it has nothing it can offer $\gamma$ and will therefore try to build a coalition with $\beta$. Thus, $\alpha$ will choose $y_\alpha$ and $g_\beta$ to maximize $-\left(y_\alpha - p_\alpha\right)^2 + G - g_\beta$ subject to $-\left(y_\alpha - p_\beta\right)^2 + g_\beta \geq -\left(p_\beta - p_\gamma\right)^2$ and $G \geq g_\beta \geq 0$. The solutions has four distinct cases.

(1) If $\frac{p_\alpha + p_\beta}{2} \geq 2p_\beta - p_\gamma$ and $G \geq \frac{1}{4}\left(p_\alpha - p_\beta\right)^2 - \left(p_\beta - p_\gamma\right)^2$, $y_\alpha^* = \frac{p_\alpha + p_\beta}{2}$ and $g_\beta^* = \frac{1}{4}\left(p_\alpha - p_\beta\right)^2 - \left(p_\beta - p_\gamma\right)^2$.

(2) If $\frac{p_\alpha + p_\beta}{2} \geq 2p_\beta - p_\gamma$ and $G < \frac{1}{4}\left(p_\alpha - p_\beta\right)^2 - \left(p_\beta - p_\gamma\right)^2$, $y_\alpha^* = p_\beta + \sqrt{G + \left(p_\beta - p_\gamma\right)^2}$ and $g_\beta^* = G$.

(3) If $p_\alpha > 2p_\beta - p_\gamma > \frac{p_\alpha + p_\beta}{2}$, $y_\alpha^* = 2p_\beta - p_\gamma$ and $g_\beta^* = 0$.

(4) If $p_\alpha < 2p_\beta - p_\gamma$, $y_\alpha^* = p_\alpha$ and $g_\beta^* = 0$.

In the last two cases, $\alpha$ finds it optimal to play the Romer-Rosenthal setter game with $\beta$ with a reversion of $p_\gamma$ over policy and offer no portfolios. In cases 1 and 2, party $\alpha$ makes portfolio concessions to move policy further than the Romer-Rosenthal inflection point of $2p_\beta - p_\gamma$.

Despite the complexity of these cases, the important thing to note is that all predict $y_\alpha^* > p_\beta$ and $g_\gamma^* = 0$. Thus, when party $\beta$ makes its offer in the first period, it knows that party $\gamma$ will accept $y_\beta^* = p_\beta$ and $g_\gamma^* = 0$. Thus, the subgame perfect Nash equilibrium outcomes are $y^* = p_\beta$ and $\mathbf{g}^* = (0, G, 0)$. Thus, this case generates a connected coalition that implements the ideal point of the median party.

We leave proofs for the remaining cases as exercises. However, the following table summarizes the outcomes for all cases.

| Table 7.5: Outcomes of Austen-Smith and Banks' Model | | |
|---|---|---|
| **Case** | **Governing Coalition** | **Policy** |
| $\omega_\alpha > \omega_\beta > \omega_\gamma$ | $\alpha$ and $\gamma$ | $\frac{p_\alpha + p_\gamma}{2}$ |
| $\omega_\beta > \omega_\alpha > \omega_\gamma$ | $\beta$ and $\gamma$ | $p_\beta$ |
| $\omega_\beta > \omega_\gamma > \omega_\alpha$ | $\beta$ and $\alpha$ | $p_\beta$ |
| $\omega_\alpha > \omega_\gamma > \omega_\beta$ | $\alpha$ and $\beta$ | $\frac{p_\alpha + p_\beta}{2}$ |
| $\omega_\gamma > \omega_\alpha > \omega_\beta$ | $\gamma$ and $\beta$ | $\frac{p_\gamma + p_\beta}{2}$ |
| $\omega_\gamma > \omega_\beta > \omega_\alpha$ | $\gamma$ and $\alpha$ | $\frac{p_\alpha + p_\gamma}{2}$ |

A key point of Austen-Smith and Banks model is that composition of the government and the policies it implements are driven by the voting weights which determine the sequence of proposals. Since these weights are determined voting behavior, the key to making predictions is an understanding of how voters behave in anticipation of the parliamentary bargaining. We refer the readers to the original for an analysis of the voting game.

## 5. Exercises

EXERCISE 7.1. *Diane is collecting money for the Center for the Study of Democratic Politics coffee fund. She needs to collect $2 from*

*at least three faculty members to operate the fund for the month. No member can contribute more that $2 and Diane cannot exclude non-contributors from drinking coffee. Each center member has an estimated $10 benefit from coffee service. If less than $6 is contributed, Diane keeps the money and no coffee is provided. If more than $6 is contributed, Diane provides the coffee and pockets the difference.*

  a. Assume that Diane decides to ask the faculty members in the following order: Arnold, Bartels, Lewis, Prior, and Romer. Assume that each faculty member can observe who has contributed and who hasn't.. What are the set of Nash equilibria to this game? What is the unique Nash equilibrium that survives backward induction?

  b. Now modify the game somewhat so that Lewis, Prior, and Romer do not know whether or not Arnold and Bartels contributed. Further suppose that Lewis, Prior, and Romer must decide simultaneously. Draw a game tree for this game in extensive form. Pay particular attention to the information sets. What are the sub-game perfect equilibria to this game? Does Diane do better or worse in this game according to her personal payoffs?

  c. Now let Diane choose the information structure of the game i.e. she can choose which contribution decisions are revealed at which stage. Suppose she wants to maximize her payoffs. Which game should she choose?

EXERCISE 7.2. *Vote Buying (This exercise is based on Groseclose (1996)). Assume that there are $N$ legislators with policy preferences $u_i(x)$ for $x \in \mathbb{R}$. They must vote for bill $x_B$ against the status quo $x_0$ where $x_B > x_0$. So let $\alpha_i = u_i(x_B) - u_i(x_0)$ the degree of preference for $x_B$ over $x_0$ for legislator $i$. We assume that each legislators policy payoff for voting for the bill is $\alpha_i$ whether or not the bill passes. We also assume that there are two vote buyers $L$ and $R$ with net preference parameters $\alpha_L$ and $\alpha_R$ respectively. Vote buyer $L$ wants to defeat $x_B$ so that $\alpha_L < 0$ while $R$ wants to pass it as $\alpha_R > 0$. Consider the following model. $R$ moves first and offers $z_i^R$ to each legislator who agrees to vote for the bill. $L$ moves second and offers $z_i^L$ to each legislator in exchange for voting against $x_B$. Thus, the payoff for voting in favor of $x_B$ is $\alpha_i + z_i^R$ while the payoff for voting against is $-\alpha_i + z_i^L$.*

  a. Assume that $N = 5$ and that $\alpha_i = -3 + i$ for $i = 1, 5$. Characterize the subgame perfect Nash equilibria to this game for various levels of $\alpha_L$ and $\alpha_R$.

  b. Can there be greater than minimum winning coalitions?

EXERCISE 7.3. *Derive the policy outcome and governing coalition for the remaining cases of the Austen-Smith and Banks model.*

## CHAPTER 8

# Dynamic Games of Incomplete Information

In chapter 6, we saw that uncertainty about the preferences (pay-offs) of others fundamentally alters the strategic situation players face in static normal form games. The implications of this in dynamic, multi-stage games lead to even more interesting strategic possibilities. Consider the deterrence game of chapter 7.

**Insert Figure 8.1 Here**

Recall that the unique subgame perfect equilibrium is $\{Initiate, Acquiesce\}$. Suppose however the game is changed to the following:

**Insert Figure 8.2 Here**

Now the subgame perfect Nash Equilibrium is $\{Do\ Not\ Initiate, Escalate\}$. But what happens if the players do not know which game they are play-ing? We now consider games in which players face uncertainty about qualities of some of the other players. Games of this form are called games of **incomplete information.** Just as in chapter 6, we can model such uncertainty with the Harsanyi maneuver. The perspective is that uncertainty about the payoffs of other players can be interpreted as playing a game in which players are not certain about what history they are at. This trick involves the use of a fictitious player–Nature–that randomly selects players types from a distribution which is known to the players. If we want to model a setting where player $i$ does not know player $j$'s preferences at a particular time then we assume that nature chooses player $j$'s payoffs (type) prior to agent $i$'s decision and we treat player $i$ as facing an information set with multiple modes since she did not see which player $j$ type was drawn by nature. This trick converts games of incomplete information – I don't know the game– to games of imperfect information –I know the game but I don't know what Nature did.

In principle, analysis of games of incomplete information strains our ability to design satisfactory notions of equilibrium. But Mertens and Zamir (1985) have shown, subject to some very technical condi-tions, any description of incomplete information can be characterized as a Bayesian game, though one with a potentially very large type space. Following this conclusion, applied game theoretic models of

incomplete information begin with a description of the incomplete information game as a game of imperfect information, with uncertainty about players showing up as uncertainty about the realization of nature's randomization.[1]

To make sense of this discussion, we return to our example. Suppose that Nature chooses game $I$ with probability $p$ and game $II$ otherwise. Now consider the following information structures.

(1) Suppose neither player observes Nature's move as in Figure 8.3. Then we say that information is imperfect but symmetric in that both players are in the same situation. This game is easy to analyze since all we need to do is compute country $B$'s expected utility of escalation and modify the game accordingly. Since $B$'s expected utility of escalation is $-p8 - 3(1 - p) = -3 - 5p$, it prefers escalation whenever $p < \frac{1}{5}$. Thus, if $p < \frac{1}{5}$, the outcome will be $\{Do\ Not\ Initiate, Escalate\}$ otherwise it will be $\{Initiate, Acquiesce\}$.

**Insert Figure 8.3 Here**

2. Suppose that only $B$ observes Nature's choice. This situation is depicted in Figure 8.4. This information structure implies that $A$ will be uncertain of $B$'s choice. Since $B$ only escalates in game II, $A$'s expected utility from initiating is $4p - (1-p)8 = -8 + 12p$. Thus, $A$ prefers initiating only if $p > \frac{2}{3}$.

**Insert Figure 8.4 Here**

3. Suppose that only player $A$ observes Nature's move as in Figure 8.5. Now the game has asymmetric information. This changes the strategic situation dramatically. While $B$ doesn't know Nature's choice, it knows that $A$ knows it. Thus, $A$ must consider what information her choices provide about Nature's draw. To see how these informational incentives effect behavior, consider the seemingly natural way of playing the game where $A$ initiates in game I, but not in game II. If $A$ played these strategies $B$ would be able to infer from $A$'s initiation that they are playing game I and should acquiesce. However, if $B$ responded in this way, $A$ would have a strong incentive defect by initiating even in game II.

---

[1]In this book and nearly all applied game theory, multiperson decision problems with incomplete information are converted into games of imperfect information using Harsanyi's trick of letting nature select types from a well defined set with common beliefs. This means that uncertainty about the preferences of players is no different (theoretically speaking) than uncertainty introduced by simultaneous moves or hidden actions.

**Insert Figure 8.5 Here**

In this chapter, we focus on the strategic use of information in dynamic settings. As we will see, incomplete information raises a number of important issues.

- *Strategic Use of Information*: Do any of the players have a strategic advantage based on how information is allocated? In many games, informed player have important advantages. However, we will see that often the uninformed player is advantaged.
- *Learning*: Can the uniformed players get more information from observing the actions of the informed players? How do these possibilities effect the strategies of the informed players?
- *Signaling*: Can the informed players credibly communicate information about the game to the uninformed players? Can informed players mislead uninformed players?

## 1. Perfect Bayesian Equilibria

In dynamic games with imperfect information, players are often uncertain as to which histories (including Nature's move) have been reached at the point in which they move. While subgame perfection can rule out some unreasonable Nash equilibrium, in many extensive form games with imperfectly observed actions a stronger equilibrium concept is needed. Consider the extensive form game depicted in Figure 8.6. Player 1 chooses whether to secretly deploy military capability to attack an island. She can either not deploy any ships, $ND$, or she send a small fleet of ships ($S$) or a big line of ships ($B$). Player 2 can only observe whether there was a deployment, as she can see the ships coming, but cannot determine how many ships are coming. If no deployment occurs then the payoffs are $(0, 5)$ as player 2 keeps the island. If there is an deployment, then player 2 must decide whether to respond to the attack ($R$). If there is no response ($NR$) then player 1 wins the island. If there is a response, then player 2 wins the island but the casualties for player 2 are much higher under $S$ then under $B$. The casualties for player 1 are higher under $B$ then under $S$.

**Insert Figure 8.6 Here**

There are three Nash Equilibria to this game. The first is $(ND, R)$. This means that player 1 does not deploy, but if she did player 2 would respond. The second Nash equilibrium is $(B, NR)$. Player 1 deploys a big line of ships, and player 2 does not respond. The profile $(S, NR)$ is also an equilibrium.

There is something very perplexing about first Nash equilibrium. Regardless of whether $B$ or $S$ is played, player 2 is better off playing $NR$. Shouldn't player 1 recognize this and send the ships? In the last chapter, we used subgame perfection to get us out of such conundrums, but that want help us here. Since this game has no proper subgames, $(NS, R)$ is also subgame perfect Nash equilibrium. .

Our argument that this profile is not reasonable is based on the idea that player 1 should anticipate a rational response from player 2 at player 2's information set. We incorporate this type of **sequential rationality** into an equilibrium concept, by requiring that at each information set agents form **beliefs** about which history they have reached and select best responses given these beliefs. These equilibria are called Perfect Bayesian Equilibria or PBE for short.

Returning to the example, we can see that no belief about the history of play at Player 2's information set justifies the selection of $R$ as a best response. Player 2 has to believe that either $S$ or $B$ ships have been deployed. In either case, she is better off choosing $NR$.

Our example leans toward the trivial side of the spectrum so consider a slight modification. In this game player 1 can only win the island if she selects $B$. Moreover, player 2 would rather defend the island if player 1 has selected $S$. Figure 8.7 depicts the relevant payoffs.

**Insert Figure 8.7 Here**

In this version whether $R$ or $NR$ is sequentially rational depends on what beliefs player 2 assigns to the two possible histories in her information set. If she believes that $S$ was played then $R$ is sequentially rational. Conversely if she believes that $B$ was played then $NR$ is sequentially rational. What should she believe? Clearly, her beliefs are based on expectations about what player 1 does. But player 1's choice will depend on what she expects player 2 to believe. How do we close this loop?

**1.1. Formal Definitions.** In this section we present the techniques needed to analyze games of this form. We now define the concepts needed to characterize PBE. We start with beliefs over histories.

DEFINITION 8.1. *Given an extensive form game with imperfectly observed actions,* $\Gamma^{EI}$ *a **belief on information set** $I_j \in I$ is a probability distribution on $I_j$. A **belief** profile is a mapping $b : H \to [0, 1]$ such that for every $I_j \in I$ $b(\cdot)$ is a belief on $I_j$ (that is for every $I_j \in I$ $\sum_{h \in I_j} b(h) = 1$ if $I_j$ is finite and $\int_{h \in I_j} db(h) = 1$ if $I_j$ is not finite).*

So in the examples above, a belief on player 2's information set, is a probability distribution over $\{S, B\}$. We use the term *belief profile* to describe a complete list of beliefs for all information sets. Since only one player makes a decision at each information set, there is no ambiguity about whose beliefs are relevant on each portion of the belief profile. If player $i$ is called to make a choice at information set $I_j$ then the portion of the belief profile which describes the belief at information set $I_j$ describes player $i$'s belief at information set $I_j$.

Given a belief profile, we can define a condition on strategies known as *sequential rationality*. Loosely speaking, sequential rationality requires that all strategies be optimal at each information set given a belief profile. To formalize this notion, let $p(I_j)$ denote the player and $s(I_j)$ denote the action called for at information set $I_j$. These terms are equivalent to $p(h)$ and $s(h)$ when $h \in I_j$. For a fixed strategy profile $s(\cdot)$ we denote the expected utility to player $p(I_j)$ associated with the choice $a$ at history $h$ by $Eu_{p(h)}(a, h, s(\cdot))$. This is an expected utility (as opposed to a utility) because players other than $p(h)$ may play mixed strategies. When player $p(I_j)$ assigns probability $b(h)$ to being at history $h \in I_j$ conditional upon being at the information set $I_j$, the expected utility to taking action $a$ at information set $I_j$ is

$$Eu_{p(I_j)}(a, I_j, s(\cdot), b(\cdot)) = \sum_{h \in I_j} b(h) Eu_{p(h)}(a, h, s(\cdot)).$$

A strategy profile is sequentially rational relative to a belief if it involves optimal actions at each information set, when players evaluate the desirability of action $a$ using $Eu_{p(I_j)}(a, I_j, s(\cdot), b(\cdot))$. Note that when there is infinite set of possible histories, the summation is replaced with integration.

DEFINITION 8.2. *Given an extensive form game with imperfectly observed actions, $\Gamma^{EI}$ and a belief $b(h)$ on each information set, the strategy profile $s(\cdot)$ **is sequentially rational (relative to the beliefs) at information set** $I_j$ if given any available action $s'$ we have*

$$Eu_{p(I_j)}(s(I_j), I_j, s(\cdot), b(\cdot)) \geq Eu_{p(I_j)}(s', I_j, s(\cdot), b(\cdot)).$$

*If the strategy profile is sequentially rational (relative to the beliefs) at every information set, then it is **sequentially rational (relative to the beliefs)**.*

Returning to the example in Figure 8.7 above, if the beliefs assign a probability close to 1 on $S$ then $R$ is sequentially rational at the information set. Similarly if player 2 believes $B$ then $NR$ is a sequentially rational response.

We now consider a condition on beliefs known as *consistency*. Consistency essentially requires that agents use Bayes' Rule to formulate their beliefs at $I_j$ whenever possible.

Recall that Bayes' rule provides us with probability that event $A$ occurs conditional the occurrence of $B$. It is given by

$$\Pr(A \mid B) = \frac{\Pr(A\&B)}{\Pr(B)}.$$

Consequently, consistency requires that agents compute the probability of particular history $h \in I_j$ conditional on reaching $I_j$. Of course, the probability of reaching $I_j$ depends on the strategy profile that the players are using. So we use $\Pr(I_j|s(\cdot))$ to denote the probability that $I_j$ is reached conditional the strategy profile $s(\cdot)$. Secondly, note that since $h$ is assumed to be an element of $I_j$ and therefore by definition no other information set, the probability that $h$ and $I_j$ are both reached under strategy $s(\cdot)$ is simply the probability that $h$ is reached. We denote this probability as $\Pr(h|s(\cdot))$.

Therefore, Bayes Rule implies that the probability of reaching history $h \in I_j$ conditional on reaching information set $I_j$ under strategy profile $s(\cdot)$ is

$$\Pr(h \mid I_j, s(\cdot)) = \frac{\Pr(h \mid s(\cdot))}{\Pr(I_j \mid s(\cdot))}.$$

Thus, weak consistency is the requirement that $b(h) = \Pr(h \mid I_j, s(\cdot))$ whenever possible.

DEFINITION 8.3. *Given an extensive form game with imperfectly observed actions, $\Gamma^{EI}$ and a strategy profile $s(\cdot)$ we say that the beliefs $b(\cdot)$ are **weakly consistent** relative to strategy $s(\cdot)$ if $b(h) = \Pr(h \mid I_j, s(\cdot))$ whenever $\Pr(I_j \mid s(\cdot)) > 0$.*

Of course, if $I_j$ is not reached under $s(\cdot)$ then $\Pr(I_j \mid s(\cdot)) = 0$ so that Bayes rule is undefined. Thus, consistency places no requirements on beliefs on "out of equilibrium" information sets. This weakness is sometimes problematic, as we will see. Combining weak consistency of beliefs and sequential rationality of strategies yields the equilibrium concept PBE.

DEFINITION 8.4. *Given an extensive form game with imperfectly observed actions, $\Gamma^{EI}$ a **perfect Bayesian equilibrium** (PBE) is a pair $(s(\cdot), b(\cdot))$ such that: (1) the strategy profile $s(\cdot)$ is sequentially rational relative to the belief $b(\cdot)$, and (2) the belief $b(\cdot)$ is weakly consistent relative to the strategy profile $s(\cdot)$.*

Thus a PBE requires the construction of beliefs. The existence of beliefs allows us to define a notion of sequential rationality (optimality of choices at histories). Moreover, the beliefs that players entertain are related to the equilibrium strategies, in that histories which are relatively more likely to be reached under a strategy profile, are believed to occur with a higher probability. It should not be surprising that this notion equilibrium has a certain circularity to it. Recall, that Nash equilibrium requires that strategies are individually best responses given a conjecture of other players strategies and that the conjecture turn out to be correct. Similarly PBE requires that strategies be best responses given beliefs which depend on the conjectured strategies of other, that the beliefs are reasonable given the strategies and that the conjecture about strategies be correct.

Returning to the game in Figure 8.6, we can now consider what strategy profiles occur in a PBE. Clearly the Nash equilibrium $(ND, R)$ is not supportable as a PBE, because for any beliefs about which history $S$ or $B$ player 2 is at when her information set is reached, $NR$ is the unique response that is sequentially rational for 2 at the information set. Now given that player 2 is choosing $NR$, player 1's optimal choice is to play either $S$ or $B$. Now if player 1 chooses $B$ then consistent beliefs must assign probability 1 to player 2 being at history $B$. Thus, one PBE is $(B, NR)$, $\Pr(B) = 1$, where $\Pr(B)$ is the posterior probability of $B$ given that player 2's information set is reached under player 2's beliefs. Similarly there is a PBE of the form $(S, NR)$, $\Pr(B) = 0$.

Now consider the game in Figure 8.7. If player 2 believes that $\Pr(B) = 1$ then $NR$ is the best response. On the other hand if player 2 believes that $\Pr(B) = 0$ then $R$ is the best response. One candidate for a PBE is $(ND, R)$, $\Pr(B) = 0$. Note that since no constraint is imposed on beliefs over the histories $B$ and $S$ when player 1 plays $ND$, the belief $\Pr(B) = 0$ is consistent relative to the strategy $ND$ But, the strategy profile $(ND, R)$ is not sequentially rational as player 1 would prefer to play $B$ than $ND$ when she conjectures that player 2 is playing $R$. It is also clear that $ND$ cannot be a best response to $NR$.

Alternatively we can try to characterize a pure strategy PBE in which $ND$ is not played. If $B$ is played and beliefs are consistent the only sequentially rational strategy by 2 will involve $NR$. But if player 1 conjectures that player 2 is playing $NR$ she will want to play $S$. So we cannot have $B$ played in a pure strategy PBE. On the other hand if $S$ is played then consistent beliefs must assign probability 1 to player 2 being at this history. Thus, the only sequentially rational action will involve playing $R$. But if player 1 conjectures that player 2 is playing

$R$ then she will want to play $B$. Thus we cannot have a pure strategy PBE in which $S$ is played. We have thus shown that there is no pure strategy PBE to the game.

It is not difficult to characterize the mixed strategy PBE to the game. Suppose that player 1 plays $B$ with probability $q$ and $b$ with probability $(1 - q)$. Further suppose that player 2 plays $R$ with probability $z$ and $NR$ with probability $(1 - z)$. Consistency of beliefs requires that $\Pr(B) = q$. Now for player 2 to be indifferent between $R$ and $NR$ it must be the case that

$$q(-5) + (1 - q)2 = q(0) + (1 - q)0$$

This requires that $q = \frac{2}{7}$. Now in order for player 1 to be indifferent between playing $B$ and $S$ it must be the case that

$$(1 - z)5 + z(-2) = (1 - z)4 + z3$$

This requires that $z = \frac{1}{6}$. Accordingly the strategy profile, $S$ with probability $\frac{5}{7}$, $B$ with probability $\frac{2}{7}$, $R$ with probability $\frac{1}{6}$ and $NR$ with probability $\frac{5}{6}$ is supportable as a PBE, with the beliefs $\Pr(B) = \frac{2}{7}$.

**1.2. Signaling Games.** An important class of games of imperfect information involve asymmetric information with the more informed agent, the *sender*, moving first followed by the less informed *receiver*. These games take their name from the possibility that the sender's action will convey information about her type to the receiver. We begin with the simplest possible signaling game to demonstrate some of the potential incentives faced by the sender.

Let Nature draws a type $\theta \in \{a, b\}$ for player 1. Player 1 observes her type and chooses a "message" $m \in \{a, b\}$. Player 2 observes the message but does not observe player 1's type. Following the message, player 2 chooses a "policy" $p \in \{a, b\}$. The payoffs to each player from a type, action pair are denoted $u_i(p, \theta)$. Figure 8.8 depicts the game form.

### Insert Figure 8.8 Here

Here the non-trivial information sets involve moves by player 2 and are represented by the dotted lines. When player 2 makes a policy selection, she knows what message was sent by player 1, but she does not know which state has been chosen by nature.

Assume that $u_2(p, \theta) > u_2(p', \theta)$ and $u_2(p', \theta') > u_2(p, \theta')$ so that player 2 would like to know $\theta$ before selecting $p$. We can motivate this assumption with a simple interpretation: The value of $\theta$ effects the desirability of each policy and player 1 wants to match $\theta$ and $p$ while player 2 wants the pair to be unmatched $(\theta \neq p)$. An example

is a legislature (agent 2) choosing between two policy alternatives $a$ and $b$ that are both risky but desirable in expectation. An informed expert (agent 1) gives unverifiable testimony before Congress about the relative risks of the two alternatives and $\theta$ denotes the identity of the policy which is actually more risky. If the legislature is more risk averse than the expert then the preference profile described above is appropriate. In this game player 1's action is called a *cheap talk* speech because agent 2 cannot verify the accuracy of 1's speech, and there is no explicit cost to lying. Our assumption that 1 wants to match $p$ and $\theta$ while 2 does not is consistent with the following ordering of payoffs.

$$u_1(b, a) < u_1(a, a)$$
$$u_1(a, b) < u_1(b, b)$$
$$u_2(b, a) > u_2(a, a)$$
$$u_2(a, b) > u_2(b, b)$$

To make the description one of imperfect information, we need to further specify a pair of prior beliefs over the state $\theta$. Suppose that $\theta = a$ with probability $\pi > \frac{1}{2}$. One natural question to ask is whether there is a PBE in which player 1, the informed player, reveals her private information. This would require that she use one of the following strategies

$$m(\theta) = \begin{cases} a \text{ if } \theta = a \\ b \text{ if } \theta = b \end{cases}$$

or

$$m(\theta) = \begin{cases} b \text{ if } \theta = a \\ a \text{ if } \theta = b \end{cases}.$$

We begin by focusing on the first message strategy. If player 1 uses this strategy profile, consistency of beliefs requires that

$$b(\theta = a \mid m = a) = \frac{\pi \cdot 1}{\pi \cdot 1 + (1 - \pi) \cdot 0} = 1$$

and

$$b(\theta = a \mid m = b) = \frac{\pi \cdot 0}{\pi \cdot 0 + (1 - \pi) \cdot 1} = 0.$$

Given these beliefs sequential rationality requires that agent 2 select policy according to the following mapping

$$p(m) = \begin{cases} b \text{ if } m = a \\ a \text{ if } m = b \end{cases}.$$

The last thing to check is whether the specified $m(\cdot)$ strategy is sequentially rational. Here the critical question is whether it represents a best response to the mapping $p(\cdot)$. Note that if $m = a$ the policy is $b$ and if

$m = b$ the policy is $a$. Since $u_1(b, a) < u_1(a, a)$ and $u_1(a, b) < u_1(b, b)$ a player 1 that has observed $\theta = a$ can deviate from the strategy and announce $m = b$ which will result in the outcome $a$. This outcome is more desirable than the outcome of playing the conjectured strategy. Similarly if player 1 observed $\theta = b$, she can gain be deviating and announcing $m = a$. This argument demonstrates that there is no PBE in which the "truthful" message strategy $m(\theta) = \theta$ is deployed. It is left as an exercise to show that there cannot be a PBE in which the second message strategy listed above is used.

The message strategies defined above are called separating because if they are used by the sender then the receiver learns the type $\theta$. We showed that if the truthful message strategy is used then consistency of beliefs required that the beliefs are deterministic.

While the game does not possess a separating equilibrium, there are other possible PBE. Suppose that the sender sends the same message (say $a$) regardless of $\theta$. Thus $m(\theta) = a$. Given this message strategy, the receiver learns nothing from observing $m$. In this case consistency of beliefs requires that

$$b(\theta = a \mid m = a) = \frac{\pi \cdot 1}{\pi \cdot 1 + (1 - \pi) \cdot 1} = \pi.$$

What beliefs should the receiver form following an $m = b$? This question is tricky. If we try to use Bayes' rule we get

$$b(\theta = a \mid m = b) = \frac{\pi \cdot 0}{\pi \cdot 0 + (1 - \pi) \cdot 0} = \frac{0}{0}.$$

Since $\frac{0}{0}$ is not a number, Bayes' rule is not well defined. This is because Bayes' rule conditions on a history that occurs with zero probability given the strategy profile. The definition of weak consistency only requires that beliefs obey Bayes' rule when the denominator is greater than zero. Accordingly, in our efforts to characterize an equilibrium with the message strategy $m(\theta) = a$, the complication of defining $b(\theta = a \mid m = b)$ is not insurmountable. Since weak consistency imposes no constraints on this belief, we are allowed to specify this posterior in whatever manner is required to help us construct equilibria. Let's say that $b(\theta = a \mid m = b) = \pi$. Given this specification of beliefs, the question of what receiver strategy is sequentially rational requires simply comparing expected utilities. Policy $a$ is more desirable if

$$\pi u_2(a, a) + (1 - \pi)u_2(a, b) \geq \pi u_2(b, a) + (1 - \pi)u_2(b, b)$$

Otherwise policy $b$ is more desirable. Accordingly, if

$$(8.1) \qquad \pi \geq \frac{u_2(b,b) - u_2(a,b)}{u_2(a,a) - u_2(a,b) - u_2(b,a) + u_2(b,b)}$$

then given these beliefs sequential rationality of $p$ is satisfied by $p(m) = b$. If the inequality in equation 8.1 is reversed then sequential rationality of $p$ is satisfied with $p(m) = a$. Finally, we must check that the message strategy is sequentially rational given the policy function. This step is trivial. Because the policy function is constant in $m$, any message function is a best response. Accordingly, we have found a pooling equilibrium of the signaling game.

Returning to the question of what $b(\theta = a \mid m = b)$ should be, we can see that the answer to this question is quite important. Assume that $\pi$ is sufficiently big so that the receiver's best response is $b$. Suppose, now that the $b(\theta = a \mid m = b) = 0$. In this case sequential rationality requires that $p(b) = a$. So changing the off-the path beliefs and the requirement of sequential rationality means that the off-the-path policy action must change. Can we construct an equilibrium with these new beliefs and a policy function that uses the pooling message function $m(\theta) = a$? Consider a sender observes $\theta = a$. In the conjectured equilibrium she is supposed to send message $a$. The receiver learns nothing from the message and selects $b$ because $a$ is more likely ($\pi$ is high). But now, if our sender deviates from this conjectured strategy and sends the message $b$, the receiver's response is to select policy $a$. If this conjectured strategy and belief profile is a PBE it must be the case that the sender never has an incentive deviate in this way. However, when $\theta = a$ the sender is better off deviating and eliciting $p = a$ instead of $p = b$. Thus, while the pooling message function was supportable in a PBE with some consistent beliefs, it is not supportable as a PBE with every set of consistent beliefs. The off-the-path beliefs matter as they create incentives for on-the-path behavior.

Another way of constructing pooling equilibria in this game is to appeal to mixed strategies, by the sender. Suppose now that regardless of her type the sender sends message $a$ with probability $\sigma \in (0,1)$. Given this Bayes' rule yields

$$b(\theta = a \mid m = a) = b(\theta = a \mid m = b) = \pi$$

As before characterizing sequentially rational strategies for the receiver hinges on $\pi$. One key difference between this pooling equilibria with mixed messages and the case of pure message strategies, is that with mixed strategies Bayes' rule pins down the beliefs at every observable information set.

In more general settings we can characterize equilibria along these lines, separating, pooling, and introduce another category called partially separating or partially pooling.[2]

DEFINITION 8.5. *In a general signaling game with multiple senders and one receiver in which each player has a type space $\Theta_i$, message space $M_i$ and the receiver has prior beliefs $\pi(\cdot)$ on $\Theta = \times_i \Theta_i$ we have the following: in a **separating** equilibrium on the equilibrium path the receivers posterior beliefs are concentrated at the true state. In a **pooling** equilibrium, on the equilibrium path the posteriors correspond to the priors, in a **partially separating** equilibrium neither of the above conditions are true.*

## 2. Application: Entry Deterrence in Elections

One of the most intriguing puzzles in the study of campaign finance is the question of why incumbent politicians exert so much effort to raise more campaign money than seems necessary to simply finance their campaigns. A standard explanation is that incumbents raise to these sums to deter entry by potential challengers. Thus, fundraising success is used to signal the incumbent's electoral strength. A formal model of this phenomenon has been developed by Epstein and Zemsky (1995). Suppose a challenger is deciding whether to run against an incumbent for office, but that the challenger only wishes to run when the incumbent is politically "weak." If the incumbent is "strong", the challenger would rather sit out the race. Conventional wisdom suggests that the incumbent may wish to signal to the challenger that she is strong by raising a large amount of campaign monies before the challenger decides to enter. If this "war chest" convinces the challenger that the incumbent is strong, he may be deterred from entering. To capture this intuition in a model, let $p_o$ be the prior probability that the challenger ($C$) places on the incumbent ($I$) being strong ($S$). Obviously, $1 - p_o$ is the probability that the incumbent is weak ($W$). Both the $C$ and $I$ get 1 unit of utility from office. Let the probability that $C$ wins against $W$ be $\pi_w$ while $C$ beats $S$ with $\pi_s$ where $\pi_w > \pi_s$. Let $k$ be $C$'s cost of running where to keep things interesting we assume that $\pi_w > k > \pi_s$. If $k > \pi_w$, $C$ would never enter and if $k < \pi_s$ he would always enter.

The key assumption of the model is that $S$ and $W$ have different costs of raising a war chest. To keep things simple, we model only

---

[2]Those who see the glass half-empty will probably prefer the term partially pooling. Others will find partially separating more in line with their philosophies of life.

whether or not the incumbent builds a war chest, but not his choice of size. Thus, the incumbent's strategy set is $S_I = \{WC, \tilde{\ }WC\}$ where $WC$ is the decision to build a war chest and $\tilde{\ }WC$ the decision to forgo one. By $s_I \in S_I$ we denote the strategy chosen by $I$. Thus, we assume that types $W$ and $S$ must pay $c_w$ and $c_s$ respectively to build a war chest where $c_s < c_w$. Note that the probability that the incumbent wins depends only on his type. There is no direct effect of the war chest except for its role in signaling the incumbent's strength to the challenger. After observing whether $I$ builds a warchest, $C$ decides whether to enter the race $E$ or sit it out $\tilde{\ }E$.

Figure 8.9 provides the extensive form of this game. We begin our analysis with the last stage of the game. Clearly $C$ will only enter if the expected utility of entering is greater than or equal $k$. Thus, entry requires that $\pi_s \Pr\{S|s_I\} + \pi_w \Pr\{W|s_I\} \geq k$.

### Insert Figure 8.9 Here

Note that the game in Figure 8.9 can also be in presented in a manner similar to Figure 8.8. We depict both approaches to familiarize readers with different graphical presentations.

**2.1. The First Period.** In the first period, $W$ and $S$ must choose strategies. There are three possible types of equilibria:

(1) *Separating*: $S$ and $W$ choose the different actions.
(2) *Pooling*: $S$ and $W$ choose the same strategies
(3) *Semi-pooling*: $S$ and $W$ choose different mixed strategies.

2.1.1. *The Separating Equilibria.* We first consider a separating equilibrium where the strong incumbent builds a war chest and the weak incumbent does not. Given these strategies, the challenger will learn the incumbent's type in equilibrium and will enter only if the incumbent does not build a war chest.

PROPOSITION 8.1. *If $c_s \leq \pi_s$ and $c_w \geq \pi_w$, the following strategies and beliefs constitute a perfect Bayesian equilibrium:$s(S) = WC$, $s(W) = \tilde{\ }WC$, $s(C) = E$ if $\tilde{\ }WC$ and $\tilde{\ }E$ otherwise, $\Pr\{S|WC\} = 1$, and $\Pr\{S|\tilde{\ }WC\} = 0$.*

Given these strategies, the application of Bayes' rule suggests that $C$'s equilibrium beliefs are $\Pr\{S|WC\} = 1$ and $\Pr\{W|\tilde{\ }WC\} = 1$. Since $\pi_w > k > \pi_s$, $C$'s strategy to enter only in the absence of a war chest is a best response. Now, we need to check that the incumbent's strategies are best responses. Incumbent $S$'s utility from $WC$ is $1 - c_s$ which is greater than her utility of $\tilde{\ }WC$ which is $1 - \pi_s$ since $c_s \leq \pi_s$. Now, we check whether $s(W) = \tilde{\ }WC$ is a best response. The payoff

for building a war chest is $1 - c_w$ while the utility for not building a war chest is $1 - \pi_w$. So ˜$WC$ is a best response since we assume that $c_w \geq \pi_w$.

It is worth trying to generate some intuition as to why the conditions $c_s \leq \pi_s$ and $c_w \geq \pi_w$ are required to support a separating equilibrium. It is clear that for a war chest to credibly signal strength it must be substantially more costly for the weak incumbent to build it. If this were not the case, weak incumbents would also build war chests and it would no longer be sequentially rational for $C$ to be deterred. But if war chests did not deter challenges, neither type would seek to build them, destroying the PBE.

Next we show that there exist no equilibria where the signal is reversed so that weak incumbent build warchests and strong incumbents do not.

PROPOSITION 8.2. $s(S) = $ ˜$WC$ and $s(W) = WC$ cannot be a Bayesian Nash equilibrium.

To see that the claim is true, note that $C$'s best response to these strategies would be {enter if $WC$}. Then $S$ would get a utility of 1 for ˜$WC$ and $1 - \pi_s - c_s$ for $WC$. So $S$'s strategy would be a best response. However, consider $W$'s best response. She would get $1 - \pi_w - c_w$ for $WC$ and 1 for ˜$WC$. Clearly, $s(W) = WC$ is not a best response.

The reason that no such separating equilibria exists can perhaps be understood in terms of the concept of incentive compatibility. We say that a messages $m$ and $m'$ are incentive compatible for types $\theta$ and $\theta'$ if and only if $EU(m|\theta) \geq EU(m'|\theta)$ and $EU(m'|\theta') \geq EU(m|\theta')$. In other words, each type must weakly prefer its own message to that of the other type. Clearly, incentive compatibility is a necessary condition for a PBE.

Note that for the reversed signals to be an equilibrium, the incentive compatibility requirements are

$$EU(˜WC|S) \geq EU(WC|S)$$
$$EU(WC|W) \geq EU(˜WC|W)$$

We can combine these requirements by adding the inequalities and moving things around so that

$$EU(˜WC|S) - EU(˜WC|W) > EU(WC|S) - EU(WC|W)$$

which becomes

$$\Pr\{E|˜WC\}(\pi_w - \pi_s) > \Pr\{E|WC\}(\pi_w - \pi_s) + c_w - c_s$$

Note that since $c_w > c_s$ and $\pi_w > \pi_s$, the incentive compatibility constraints require that $\Pr\{E|\tilde{}WC\} > \Pr\{E|WC\}$. However, given $C$'s beliefs, this cannot be a best response. Therefore, reversing the signal is not a PBE.

We will see this incentive compatibility approach used much more extensively when we look at models of mechanism design in chapter 11.

2.1.2. *Pooling Equilibria.* Now we turn to the analysis of pooling equilibria where both types of incumbents choose the same strategy. Generally, the easiest way to construct such equilibria is to specify the most unfavorable beliefs possible in the event that one of the senders chooses the out-of-equilibrium message. In this context, this requires that $C$ believe that the incumbent is weak following an out-of-equilibrium action. If these "pessimistic" beliefs support a strategy profile, then that profile constitutes a PBE. However, a number of slightly less pessimistic beliefs will also support that profile as a PBE. So its a good practice to compute the set of beliefs that support each PBE. In the current model, this requires us to specify the largest posterior on $\Pr\{S\}$ following the defection that supports the PBE.

PROPOSITION 8.3. *Suppose that $p_o \geq \frac{\pi_w - k}{\pi_w - \pi_s}$. The following strategies and beliefs are a perfect Bayesian equilibrium: $s(W) = s(S) = \tilde{}WC$; $s(C) = \{E\ if\ WC\}$, $\Pr\{S|\tilde{}WC\} = p_o$ and $\Pr\{S|WC\} \leq \frac{\pi_w - k}{\pi_w - \pi_s}$.*

Since $W$ and $S$ play the same strategy, Bayes' rule implies that $\Pr\{S|\tilde{}WC\} = p_o$. So $C$'s utility from entering is

$$\pi_s p_o + \pi_w (1 - p_o) - k = \pi_w - (\pi_w - \pi_s) p_o - k$$

so $C$ will enter if $p_o \leq \frac{\pi_w - k}{\pi_w - \pi_s}$ after observing $\tilde{}WC$. What should $C$ believe and do if he observes $WC$? Nash equilibrium is silent about what to do off the equilibrium path. Note that we cannot apply Bayes' rule because the denominator would be zero. So we assign the arbitrary beliefs $\Pr\{S|WC\} \leq \frac{\pi_w - k}{\pi_w - \pi_s}$ after observing $WC$. Thus, $C$ should enter if he observes $WC$. Now consider the strategies of $S$ and $W$. They both get 1 for not building a war chest and $1 - \pi_s - c_s$ and $1 - \pi_w - c_w$ respectively for building one. Thus, their strategies are best responses.

PROPOSITION 8.4. *Suppose that $p_o < \frac{\pi_w - k}{\pi_w - \pi_s}$. The following strategies and beliefs are a perfect Bayesian equilibrium: $s(W) = s(S) = \tilde{}WC$; $s(C) = \{E\}$, $\Pr\{S|\tilde{}WC\} = p_o$ and $\Pr\{S|WC\} \leq \frac{\pi_w - k}{\pi_w - \pi_s}$.*

Clearly, $C$'s strategy of entering is a best response since regardless on what the incumbent does $C$ will think it unlikely that the incumbent is $S$. On the equilibrium path, Bayes' rule implies that $\Pr\{S|\tilde{}WC\} = p_o$. Off the equilibrium path, we are free to assign $\Pr\{S|WC\} \leq$

$\frac{\pi_w - k}{\pi_w - \pi_s}$. Since $C$ always enters, a war chest by either type of incumbent is a waste of resources. Thus, neither type of incumbent builds one.

PROPOSITION 8.5. *Suppose that $p_o \geq \frac{\pi_w - k}{\pi_w - \pi_s}$. The following strategies and beliefs are a perfect Bayesian equilibrium: $s(W) = s(S) = $ ~$WC$; $s(C) = \{$~$E\}$, $\Pr\{S|$~$WC\} = p_o$ and $\Pr\{S|WC\} \geq \frac{\pi_w - k}{\pi_w - \pi_s}$.*

$C$'s best response given the beliefs on the equilibrium path is to stay out of the race. Suppose instead that the incumbent defected and built a war chest. In this equilibrium, $C$ believes that it is relatively likely that the incumbent is strong at the out-of-equilibrium information set. So $C$ still chooses ~$E$. Since $C$ always stays out, it is optimal for neither incumbent to build a war chest.

PROPOSITION 8.6. *Suppose that $p_o \geq \frac{\pi_w - k}{\pi_w - \pi_s}$, $c_s \leq \pi_s$, and $c_w \leq \pi_w$ The following strategies and beliefs are a perfect Bayesian equilibrium: $s(W) = s(S) = WC$; $s(C) = \{E$ if ~$WC\}$, $\Pr\{S|WC\} = p_o$ and $\Pr\{S|$~$WC\} \leq \frac{\pi_w - k}{\pi_w - \pi_s}$.*

Clearly, given equilibrium beliefs, $C$'s strategy is a best response. If $C$ observes ~$WC$, we assign a probability of $S$ sufficiently low so that she chooses $E$ following a defection. So consider incumbent $S$ who gets $1 - c_s$ in equilibrium but gets $1 - \pi_s$ by defecting. So $S$ will choose $WC$ so long as $c_s \leq \pi_s$. Incumbent $W$ gets $1 - c_w$ from the equilibrium and $1 - \pi_w$ from defecting so $WC$ is a best response if $c_w \leq \pi_w$.

PROPOSITION 8.7. *There is no equilibrium where $s(W) = s(C) = WC$ and $C$ always enters.*

Both types would defect to ~$WC$ since it does not change $C$'s behavior.

Clearly, there are two types of pooling equilibria to this game: one where both incumbents build warchests and one where neither does. Note that both reveal the same amount of information to the challenger, namely none, but differ in the costs incurred by the incumbent. Clearly, every player prefers the PBE where both types play ~$WC$ to the one where they both play $WC$. Thus, we can refer to the former as the efficient PBE and the latter as the inefficient one. Importantly, there is nothing intrinsic to the concept of PBE to predict which of these equilibria are more likely to be played. However, as we will see, various authors have proposed criteria for refining the set of PBE. Often these refinements select the efficient PBE.

2.1.3. *Partial Pooling.* The remaining possibility is that at least one of the incumbent types chooses a mixed strategy that occasionally separates from the other type and occasionally pools. There are

a couple of reasons for exploring such possibilities.    First of all, if
such equilibria exist, a full characterization of the set of PBEs requires
analysis of partial pooling equilibria.[3]   Secondly, authors sometimes
would like to analyze the most informative equilibrium i.e. the one in
which the receiver's posteriors are closest to the true distribution of
types.  In many cases, semi-pooling equilibria will exist for parameter
values for which there are no separating equilibria.  In this cases, the
partial pooling equilibria will be the most informative.

We present only one of the partial pooling equilibria, and leave the
other possibility to the reader as an exercise.  In this equilibrium, the
$S$ incumbent always builds a warchest.   However, the $W$ incumbent
also builds one with probability $q$.

PROPOSITION 8.8. *If $\frac{c_w}{\pi_w} > \frac{c_s}{\pi_s}$, then $s(S) = WC$, $s(W) = \{WC$ with
prob$= q\}$, $C = \{enter$ if $\tilde{}WC$ and enter with probability $r$ if $WC\}$ is
a perfect Bayesian equilibrium.*

Clearly, Bayes' rule implies that $\Pr\{S|\tilde{}WC\} = 0$ and thus follow-
ing $\tilde{}WC$ entry is a best response.  When $WC$ is observed, Bayes' rule
is more complicated:

$$\Pr(S \mid WC) = \frac{\Pr(WC \mid S)\Pr(S)}{\Pr(WC \mid S)\Pr(S) + \Pr(WC \mid W)\Pr(W)} = \frac{p_0}{p_0 + q(1 - p_0)}$$

Since $C$ plays a mixed strategy, he must be indifferent between entering
or not entering so that

$$\pi_s \Pr\{S|WC\} + \pi_w \Pr\{W|WC\} = k$$

or

$$\frac{\pi_s p_0}{p_0 + q(1 - p_0)} + \frac{\pi_w q(1 - p_0)}{p_0 + q(1 - p_0)} = k$$

Solving we find that

$$q^* = \frac{p_0(k - \pi_s)}{(1 - p_0)(\pi_w - k)}.$$

When $C$ chooses $r$, it must make $W$ indifferent between building $WC$
and not doing so.  If $W$ chooses $\tilde{}WC$, she gets $1 - \pi_w$.  Choosing $WC$
gets $r(1 - \pi_w) + (1 - r) - c_w$.   So the indifference condition implies
that

$$r^* = \frac{\pi_w - c_w}{\pi_w}$$

---

[3]The authors have revelaed their type.

Now we need only check that $S$ prefers $WC$. If she plays $\tilde{}WC$ she gets $1 - \pi_s$ whereas in equilibrium she gets $r(1 - \pi_s) + (1 - r) - c_s$. Thus,

$$EU_S(WC) - EU_S(\tilde{}WC) = \pi_s(1 - r) - c_s = \frac{\pi_s c_w}{\pi_w} - c_s > 0$$

by assumption.

We know from above that the separating equilibrium only exists when $\frac{c_s}{\pi_s} \leq 1 \leq \frac{c_w}{\pi_w}$ while the semi-pooling equilibrium that we have just considered exists when $\frac{c_s}{\pi_s} < \frac{c_w}{\pi_w}$. Thus, whenever the separating equilibrium exists, the semi-pooling equilibrium must also exist. However, the partial pooling equilibrium exists in many circumstances (i.e. $\frac{c_w}{\pi_w} < 1$) where the separating equilibrium does not.

## 3. Application: Information and Legislative Organization

One of the longest standing debates about legislative politics is the nature of the role of standing committees. While some scholars have focused on their role in stabilizing majority rule (Shepsle 1978), maintaining distributive coalitions (Shepsle and Weingast 1987, Weingast and Marshall 1988) while others have focused on their role in promoting the interests of the majority party (Cox and McCubbins 1994). However, Gilligan and Krehbiel (1987) present an alternative model based on the idea that the role of committees is to specialize in the development of policy-specific expertise. This quite influential theory is based on games of incomplete information.

The basis of these models is the idea that policymakers do not always know the exactly link between their policy choices and the ultimate policy outcomes. For example, legislators may not know how a specific agricultural policy will affect farmer's incomes because of a number of unforeseeable intervening variables such as the weather and competition from foreign producers. Since legislators would presumably value such information, they should desire institutional arrangements that facilitated its gathering and transmission. Gilligan and Krehbiel argue that a committee system with some limited parliamentary rights will help solve these informational problems.

To capture the distinction between policy choices and outcomes, Gilligan and Krehbiel assume that the policy outcome $x$ is a additive function of policy $p$ and random term $\omega$ so that

$$x = p + \omega.$$

Their game is played between two players: the floor $F$ and the committee $C$. In our version of the model, we assume that $C$ knows $\omega$

with certainty but that $F$'s prior beliefs are that:

$$\omega = \begin{cases} \theta \text{ with prob } .5 \\ -\theta \text{ with prob } .5 \end{cases}$$

They assume that each has quadratic spatial preferences over a single dimension and that $F$ has an ideal point of 0 and $C$ has an ideal point $c > 0$. Thus, the payoffs are $u_F(x) = -x^2$ and $u_C(x) = -(c - x)^2$ respectively. However, since $F$ does not know $\omega$, his utility from a given policy $p$ is given by the following expected utility:

$$.5u_F(p + \theta) + .5u_F(p - \theta) = -.5[p^2 + 2p\theta + \theta^2] - .5[p^2 - 2p\theta + \theta^2]$$
$$= -p^2 - \theta^2$$

Thus, $F$'s utility has two components. The first is $-p^2$ which reflect the difference between his ideal point and the expected policy while the second is $-\theta^2$ which reflects the variance in the policy shock $\omega$. Thus, $F$ is risk adverse in that he would be willing to "pay" to obtain information about $\omega$. Our concern here is whether $F$ may be willing to grant extra parliamentary rights to $C$ to provide it with the incentive specialize and provide information. Formally, $F$ chooses the procedures or rules under which legislation proposed by $C$ can be considered. To simply things, we assume that $F$ chooses between the following types of rules:

(1) Open rule: The committee reports a bill and the floor may freely amend it.
(2) Closed Rule: The $F$ must vote up or down on $C$'s proposal against the status quo, $SQ = 0$.

The closed rule obviously represents more extensive parliamentary rights for $C$ since she can make a take-it-or-leave to $F$. To solve this as a game of incomplete information, we need to specify strategies of both "types" of committees $(-\theta, \theta)$, beliefs of the floor following any proposal, and the floor's best response given these beliefs. Following Gilligan and Krehbiel, we will look for the most "informative" equilibrium. In our context, we will mainly be interested in whether or not a separating equilibrium exists.[4]

3.0.4. *Open Rule.* Under the open rule, the committee must worry about whether the information it provides can be used to "roll" it on the floor. Suppose that the committee revealed all of its information

---

[4]The reader who is paying attention should respond "aha! but haven't you warned us that informative semi-pooling equilibria can sometime exist when separating equilibria do not!" The attentive reader should consult excercise 8.6 and verify for herself that this is not one of those cases.

by proposing its ideal policy for each outcome. This means that the committee chooses $p^c$ so that $x = c$ or $p^c = c - \theta$ when $\omega = \theta$ and $p^c = c + \theta$ when $\omega = -\theta$. Since there are distinct proposals for each state of the world, these proposal strategies fully reveal all information to the floor. Thus, under an open rule, the floor will amend until $x = 0$ or $p^f = -\theta$ when $\omega = \theta$ and $p^f = \theta$ and $\omega = -\theta$. Therefore, utilities for this equilibrium are $u_F = 0$ and $u_C = -c^2$. Since we are interested in determining whether or not a separating equilibria exists, we will specify the beliefs most unfavorable to $C$ following an out-of-equilibrium proposal. Thus, we assume that $F$ treats any other bill as if it came from the "high" type: $\omega = \theta$ and adopts $p^f = -\theta$. Given $F$'s best responses, it is easy to see that there are three only outcomes that $C$ can generate $-2\theta, 0, 2\theta$.

Now we can check to see if $C$ like to defect from this separating equilibrium. If she defects by proposing $c + \theta$ when $\omega = \theta$, $F$ will pass $p^f = \theta$. Since the resulting outcome is $2\theta$, C's utility from the defection is $-(c-2\theta)^2$. Alternatively, by proposing $c - \theta$ when $\omega = -\theta$, $p^f = -\theta$ leading to $x = -2\theta$ and $u_c = -(c + 2\theta)^2$. Finally, consider $p^c \notin \{c + \theta, c - \theta\}$. Since all of these proposals generate $p^f = -\theta$, $u_c(\theta) = -c^2$ and $u_c(-\theta) = -(c + 2\theta)^2$.

Since $c > 0$ and $\theta > 0$, the highest utility defections are given in Table 1.

| Table 8.1: Possible Defections | | |
|---|---|---|
| Defection | Conditions | Utility |
| $p^c = c + \theta$ | if $\omega = \theta$ and $c > \theta$ | $-(c - 2\theta)^2$ |
| $p^c \notin \{c + \theta, c - \theta\}$ | if $\omega = \theta$ and $c < \theta$ | $-c^2$ |
| $p^c \neq c + \theta$ | if $\omega = -\theta$ | $-(c + 2\theta)^2$ |

Clearly the "low" type won't defect, since $-c^2 > -(c + 2\theta)^2$. Similarly if $c < \theta$, the "high" type will not defect since her equilibrium utility is the same as that of her most profitable defection. However, when $c > \theta$ and $\omega = \theta$, the high type prefers the utility of $-(c - 2\theta)^2$ to her equilibrium payoff. Therefore, the separating equilibrium does not exist whenever $c > \theta$. Absent a separating equilibrium, the only equilibrium is an uninformative pooling equilibrium where both types use the same mixed strategy across the set of proposals.

Thus, when the committee is an "outlier" i.e. $c > \theta$, no information can be revealed. In this case, the high type wants to convince the floor that $\omega = -\theta$ since it prefers policy $\theta$ to 0.

3.0.5. *Closed Rule.* Under the open rule, a major reason why information cannot be revealed by the "outlier" committee is that the floor

will use this information to "roll" the committee and move policy to the floor median. However, under the closed rule, the committee cannot be "rolled". The floor must accept or reject the proposal in favor of the status quo. So now suppose there is a separating equilibrium where $F$ learns the value of $\omega$. The utility of the status quo is $-\omega^2$ and of any other proposal $-(p+\omega)^2$. This implies that $F$ will accept any proposal in the interval between 0 and $-2\omega$. Thus, when $\omega = \theta$, the largest policy $F$ will accept is 0. This results in an outcome of $\theta$. When $\omega = -\theta$, the largest policy $F$ will accept is $2\theta$. This also results in an outcome of $\theta$.

So the committee can guarantee an outcome as high as $\theta$ by proposing $p^c = 0$ when $\omega = \theta$ and $p^c = 2\theta$ otherwise. If $c < \theta$, an outcome of $c$ can be guaranteed by proposing $p^c = c - \theta$ when $\omega = \theta$ and $p^c = c + \theta$.

To complete the specification of the perfect Bayesian equilibrium, we need to specify $F$'s behavior and beliefs following any out-of-equilibrium proposals from $C$. To support all possible separating equilibria, it is sufficient to have $F$ accept any proposal such that $-2\theta \leq p^c \leq 0$ and vote down out-of-equilibrium proposal $p^c > 0$ and . This strategy is a best response to the beliefs $\Pr\{\omega = \theta | p^c\} = 1$ following an out-of-equilibrium proposal.

Now we check that $C$ prefers its equilibrium strategy to any defection. When $c < \theta$, $C$ gets its ideal point so she cannot do better by defecting. Thus, a separating equilibrium exists, just as it did for the open rule. So we need only focus on the case where $c > \theta$. First suppose that $\omega = -\theta$, then $C$ gets $-(c - \theta)^2$ for $p^c = 2\theta$ and $-(c + \theta)^2$ for any other bill. $C$ will not defect in this case. Now suppose that $\omega = \theta$, then $C$ receives $-(c - \theta)^2$ for $p^c = 0$, $-(c - 3\theta)^2$ for $p^c = 2\theta$, and $-(c - \theta)^2$ for any other bill. Thus, so long as $c < 2\theta$, $C$ weakly prefers her equilibrium proposal $p^c = 0$. However, if $c > 2\theta$, the committee will defect so that no separating equilibrium exists. In this case, the only equilibria are uninformative ones such as those where both types of committee use the same mixed strategy over all proposals and the floor rejects all proposals except $p^c = 0$.

Since the closed rule can sustain a separating equilibrium in cases where the open rule cannot i.e. $2\theta \geq c > \theta$, the model predicts that restrictive rules might be employed to encourage greater information transmission from the committee to the floor.

**3.1. Committee Specialization.** In the preceding section, we showed that our version of the Gilligan-Krehbiel model predicts that restrictive rules can encourage more information transmission from informed committees than do open rules. We now turn to an analysis of

whether a commitment to restrictive rules by the floor can induce the committee to specialize in the first place. So now the game takes the following form:

(1) $F$ chooses whether or not the committee will report the bill under open or closed rule.
(2) $C$ decides whether to specialize by paying a cost $k$ to learn $\omega$.
(3) $C$ proposes $p^c$.
(4) Under closed rule, $F$ votes $p^c$ up or down against $SQ$. Under open rule, $F$ chooses $p^f$.

If the committee does not specialize, $F$ will decide based on its prior beliefs. Therefore, under the closed rule, she will veto any proposal other than $p^c = 0$. Under the open rule, $F$ passes $p^f = 0$. Thus, non-specialization results in a policy of $p = 0$ for both rules.

Now consider the committees specialization decision under open rule. If the committee specializes when $c \leq \theta$, $C$ and $F$ play the open rule separating equilibrium resulting in an overall payoff of $-c^2 - k$ for the committee. If the committee does not specialize, then its expected utility from $p = 0$ is $-c^2 - \theta^2$. So the committee specializes so long $k \leq \theta^2$. If $c > \theta$ and the committee specializes, the committee and floor will play the pooling equilibrium which leads to $p = 0$ and an overall committee payoff of $-c^2 - \theta^2 - k$. Thus, the committee will obviously not specialize under these circumstances.

Now consider the closed rule. If $c \leq \theta$, the separating equilibrium generates an outcome of $c$ so that the committee's utility of specializing is $-k$. Similarly, if $c > \theta$, the utility of specializing is $-(c-\theta)^2 - k$. In both cases, non-specialization leads to a payoff of $-c^2 - \theta^2$. Comparing these utilities, we find that it pays for the committee to specialize when $k \leq 2\theta c$ and $c > \theta$ and when $k < c^2 + \theta^2$ and $c \leq \theta$. Since both of these critical values for $k$ are higher than $\theta^2$, the committee often specializes under closed rule when it would not have under open rule.

**3.2. Implications.** Gilligan and Krehbiel draws several implications about institutional design by computing the floors utility under the different rules. The floor's expected payoffs under open rule are

$$E[u_F] = 0 \quad \text{if } c \leq \theta \text{ and } k \leq \theta^2$$
$$E[u_F] = -\theta^2 \quad \text{if } c > \theta \text{ and } k > \theta^2$$

while the payoffs under closed rule are

$$E[u_F] = -c^2 \quad \text{if } c \leq \theta \text{ and } k \leq c^2 + \theta^2$$
$$E[u_F] = -\theta^2 \quad \text{if } c > \theta \text{ and } k > c^2 + \theta^2$$

There are a couple of implications worth noting. First, $F$ always is always better when $C$'s ideal point is close to the floor's ideal 0.

This is the basis of Krehbiel's argument that majoritarian legislatures should attempt to appoint committees that are representative of the preferences of the chamber. He contrasts this implication with that of the distributive theory of legislatures which predicts that committees will be composed of high demanders for the policies in the committees jurisdiction.

Secondly, the model makes predictions about when $F$ will prefer a closed rule. Specifically, $F$ will choose a closed rule when $c \leq \theta$ and $\theta^2 \leq k \leq c^2 + \theta^2$. Thus, the committees most likely to receive closed rules are those whose ideal point does not diverge from the floor and those with intermediate specialization costs.

## 4. Application: Informational Lobbying

The traditional literature on lobbying in legislatures has uncovered a striking empirical regularity: interest groups almost always lobby their friends. Since these friends are likely to vote with the interest group anyway, this observation has been interpreted to mean that lobbying activities are not very consequential.

However, Austen-Smith and Wright (1992, 1994) develop a model where groups do indeed lobby their friends, but that their efforts are important. The main premise of their model is that interest groups have private information about the consequences of a legislative decision that are unknown to the legislator. "Lobbying" consists of a group making a speech to the legislator. Since lobbying is assumed to be costly, groups will choose whether or not lobby the legislator. The main result of their paper is that groups will often lobby friendly legislators to counteract the lobbying efforts of opposing groups.

Before analyzing the full model, we begin with a model with only one interest group and a legislator. The legislator is required to choose between two policies $A$ and $B$. However, she is uncertain as to which policies she prefers. Assume that there are two states $s = A$ or $B$ such that the legislator prefers policy $A$ in state $A$ and $B$ in state $B$. To keep things simple, we assume that $u_L(A) = 1$ in state $A$ and $u_L(A) = -1$ in state $B$. We assume $u_L(B) = 0$ in both states. The legislator believes that $s = A$ with probability $p < \frac{1}{2}$. Thus, in the absence of any additional information provided by lobbying, the legislator will choose $B$.

Suppose that there is an interest group $G_A$ who prefers $A$ to $B$. We assume that $u_{G_A}(A) = 1$ and $u_{G_A}(B) = 0$ in both states. If $G_A$ decides to lobby, it pays cost $c > 0$ to learn the true state with certainty. Its ex ante beliefs about the state are the same as $L$'s. Once informed,

the group sends one of two messages $m = A$ or $B$ where the messages are to be interpreted literally.

After observing message $m$, $L$ may attempt to verify the group's information by "auditing" the message. In doing so, $L$ incurs cost $\kappa$. If the message is found to be incorrect i.e. $m \neq s$, the group is penalized by an amount $\delta$.

Now consider the group's possible lobbying strategies. First, suppose $G_A$ always reports $A$ independent of $s$. Then the message would be uninformative and $L$ would always vote $B$ if she does not audit and will choose the optimal outcome following an audit. Since the utility of an audit is $1 - \kappa$, it is easy to show that $L$'s best response to an "always $A$" strategy involves auditing if and only if $p > \kappa$.

Since an "always $A$" lobbying strategy does not affect the outcome, $G_A$ would prefer not to lobby at all. Therefore, successful lobbying requires that $G_A$ tell the truth at least some of the time. Suppose $G_A$ told the truth all of the time. Then $L$ would always follow such advice. However, this would give $G_A$ the incentive to report $A$ even when $s = B$. Since $G_A$ has an incentive to deviate, this cannot be an equilibrium.

Thus, any perfect Bayesian equilibrium to this game has to be semi-pooling. So consider an equilibrium where $G_A$ always reports $A$ when $s = A$ and reports $A$ with probability $\mu$ when $s = B$. Such a strategy leads $L$ to update (using Bayes' Rule) her belief that $s = A$ to

$$\widehat{p} = \frac{p}{p + \mu (1 - p)}.$$

In such an equilibrium, $L$ must be indifferent to voting for $A$ and auditing. The expected utility of voting $A$ is $2\widehat{p} - 1$ whereas the utility of auditing is $\widehat{p} - \kappa$. So the group will choose $\mu$ so that

$$\frac{p}{p + \mu (1 - p)} = 1 - \kappa$$

or

$$\mu = \frac{\kappa p}{(1 - p)(1 - \kappa)}.$$

To close the model, $L$ must choose a probability of auditing $\alpha$ so that $G_A$ is indifferent between lying and being truthful when $s = B$. If $G_A$ is truthful, it gets $-1$. If $G_A$ lies, it gets $-\alpha(\delta + 1) + (1 - \alpha)1$. Therefore, $L$ sets

$$\alpha = \frac{1}{\delta + 2}.$$

Given these lobbying and auditing strategies, we can compute $G_A$'s expected utility to determine whether or not it will choose to become

informed and lobby.    Since $L$ always chooses $B$ in the absence of lobbying, the group's utility from not lobbying is $-1$. If it does lobby, the group gets $A$ for sure when $s = A$ and an expected utility of $-1$ when $s = B$.    Thus, the ex ante expected payoff from lobbying is $2p - 1 - c$.  Thus, the group lobbies if and only if $p > \frac{c}{2}$.

That this equilibrium requires $0 \leq \mu \leq 1$ suggests that we need $1 - p > \kappa > 0$.  Thus, it cannot be too costly for the legislator to audit.  If it is too costly, the legislator will not audit which guarantees that the group has the incentive to always lie.  However, if the group always lies, lobbying will be ineffective and the group will choose not to become informed.

Now consider the model with two groups.  Suppose now that there is a group $G_B$ with preference opposed to $G_A$ so that $u_{G_B}(B) = 1$ and $u_{G_B}(A) = -1$.   It also may learn the true state by paying a cost $c$. We assume that both groups decide whether to become informed and choose their messages, $m_A$ and $m_B$, simultaneously.

First, note that since $L$ chooses $B$ in the absence of lobbying there can be no equilibrium where $G_B$ lobbies but $G_A$ does not.  If $G_B$ lobbied alone, it would incur cost $c$ without altering the outcome.    Austen-Smith and Wright interpret this result as implying that groups only lobby "friendly" legislators to counteract the lobbying of other groups.

Since the outcome of $G_A$ lobbying alone is outlined above, we need only focus on the outcome when both groups lobby.    Consider the following equilibrium.    Both groups send truthful messages.  When $m_A = m_B = s$, the legislator believes that the true state is $s$.  However, if $m_A \neq m_B$, $L$ audits the message and chooses her optimal policy.  For auditing to be a best response to the out-of-equilibrium messages, we specify that $\widehat{p} \geq \kappa$ if $m_A \neq m_B$.

Now we check that $G_A$ will not deviate and choose an untruthful message.  Clearly, it has no incentive to choose $m_A = B$ when $s = A$. So consider whether it will choose $m_A = A$ when $s = B$.  Given $G_B$'s strategy, this results in a posterior of $p$.  Since this message will lead to an audit, it results in a policy of $B$ and a penalty of $\delta$ for a total payoff of $-1 - \delta$.  Since telling the truth leads to a payoff of $-1$, $G_A$ has no incentive to deviate.

Turning to $G_B$'s decision, we need to check that it will not choose $m_B = B$ when $s = A$.   As before, such a message leads to an audit and a penalty for $G_B$.  Thus, it will not defect.

Now we must verify the conditions under which both groups would prefer to lobby.  Since lobbying by both groups leads to the full information outcome, the equilibrium utilities of $G_A$ and $G_B$ are $2p - 1 - c$ and $1 - 2p - c$ respectively.   If only $G_A$ lobbies, its utility $2p - 1 - c$

so that its utility from lobby is independent of $G_B$'s choice. However, note that $G_B$'s utility from having $G_A$ lobby alone is

$$-p + (1-p)\left[\mu\left(2\alpha - 1\right) + (1-\mu)\right] = 1 - 2p - \left[\frac{\kappa p}{(1-\kappa)}\frac{2\delta + 2}{\delta + 2}\right]$$

Therefore, $G_B$ will only decide to lobby if $p > \frac{1-\kappa}{\kappa}\frac{\delta+2}{2\delta+2}c$.

Thus, we can characterize the equilibrium lobbying decisions of both groups as follows. If $p < \frac{1}{2}c$, neither group lobbies. If $\frac{1-\kappa}{\kappa}\frac{\delta+2}{2\delta+2}c > p > \frac{1}{2}c$, only $G_A$ lobbies. If $p > \frac{1-\kappa}{\kappa}\frac{\delta+2}{2\delta+2}c$, both groups lobby.

Austen-Smith and Wright interpret this perfect Bayesian equilibrium as implying the following hypotheses.

- Ceterus paribus, when a legislator is lobbied by groups from just one side of an issue, the only groups that lobby are those opposed to the legislator's ex ante position.
- The decision of a group to lobby an "unfriendly" legislator is independent of the lobbying decisions of opposing groups.
- Conditional on a "friendly" legislator being lobbied by an opposing group, a group's decision to lobby that legislator is purely counteractive.

## 5. Refinements of Perfect Bayesian Equilibrium*

**5.1. Sequential Equilibria.** By now it should be clear that weak consistency lives up to its name by not being a strong enough constraint on off the path actions. In fact PBE need not even be subgame perfect, a fact demonstrated by the example illustrated in Figure 8.10.

### Insert Figure 8.10 here

In this extensive form game, $C$ is a potential candidate for office, and must decide whether to run or not. After this decision, a popular media source such as the local newspaper decides whether to endorse $C$ or not. Following this decision, but without knowing the media's decision, the incumbent $I$ must decide how much campaign effort to exert.[5]  The payoffs are chosen to reflect the fact that $C$ prefers to enter only if he is endorsed and the incumbent chooses low effort, the media prefers to endorse $C$ only if the incumbent chooses low effort, and the incumbent prefers to exert effort only when $C$ is not endorsed.

In Figure 8.10, the shaded branches denote a strategy profile and the numbers $p = 0$ and $p = 1$ denote beliefs. It is not difficult to see that the figure depicts a PBE. Given the belief that "not endorse" occurs with probability 1, the incumbent optimally selects high effort.

---

[5]We can also interpret the endorsement and effort decisions as occuring simultaneously.

Moreover, given the expectation that endorsement and high effort will follow a decision to run, the optimal decision for $C$ is to not run for office. Given this strategy profile, the incumbent's information set is not reached and thus weak consistency does not restrict beliefs. Despite the fact that this is a PBE, the specified strategy profile is not even subgame perfect. To see this, consider the circled subgame that starts with $m$'s decision. Given that $m$ is choosing endorse, high effort is not a best response. The pathology exhibited by this game is that we can write down strategy profiles for which their are non-trivial information sets that are several moves away from the equilibrium path.

In a reasonably defined equilibrium, we would like $I$'s beliefs about $m$'s decision (conditional on reaching this information set) to be somehow consistent with player $m$'s strategy. Accordingly, we might conjecture that (i) if $C$ runs then either $m$ endorses and $I$ exerts low effort or (ii) $m$ does not endorse and $I$ exerts high effort or (iii) both $m$ and $I$ randomize. As we shall see combining sequential rationality and a stronger notion of consistency can avoid this pathology.

In this section we limit ourselves to finite games and present several stricter equilibrium concepts. All of the concepts defined will involve sequential rationality as defined in Definition 8.2. Where these concepts differ is in the restrictions imposed on beliefs.

First we need a bit of notation. Given a finite game $\Gamma^{EI}$ a mixed strategy profile $\sigma(\cdot)$ is a mapping that determines a lottery over available actions at each information set. For a finite game such a strategy profile can be written as a vector $(\sigma(1,1),....,\sigma(a,I)...)$ with generic coordinate $\sigma(a,I)$ denoting the probability that action $a$ is played at information set $I$. A completely mixed strategy profile selects every action at every information set with positive probability. A pure strategy is then a vector containing only 0's and 1's. A sequence of mixed strategies $\{\sigma(\cdot)^n\}_{n=1}^{\infty}$ is said to converge to a mixed strategy $\sigma(\cdot)$ for each $I, a$ $\sigma(a,I)^n$ converges to $\sigma(a,I)$. The notion of sequential equilibrium replaces weak consistency with a stronger condition.

DEFINITION 8.6. *Given a finite extensive form game with imperfectly observed actions,* $\Gamma^{EI}$ *a **sequential equilibrium** (SE) is a pair* $(\sigma(\cdot), b(\cdot))$ *such that: (1) the strategy profile* $\sigma(\cdot)$ *is sequentially rational relative to the belief* $b(\cdot)$*, and (2) there exists some sequence of completely mixed strategies*$\{\sigma(\cdot)^n\}_{n=1}^{\infty}$ *and beliefs* $\{b(\cdot)^n\}_{n=1}^{\infty}$*that converge to* $\sigma(\cdot), b(\cdot)$ *respectively and for some* $n'$ *if* $n > n'$ $b(\cdot)^n$ *is weakly consistent relative to the strategy profile* $\sigma(\cdot)^n$.

To demonstrate how this concept refines our notion of PBE, we first return to the game in Figure 8.10. Consider any completely mixed

profile $\sigma^n(\cdot)$ where we will assume that every pure strategy must be played with probability $\varepsilon_n$ where $\varepsilon_n \to 0$ as $n \to \infty$. If this sequence of mixed strategies is going to converge to the PBE in Figure 8.10, we require $\sigma^n(endorse, run)$ to be a mixed best response to $\sigma^n(run) = \varepsilon_n$ and $\sigma^n(high\ effort, run) = 1 - \varepsilon_n$. However, $m$'s expected utility from $endorse$ is therefore $\varepsilon_n$ while its payoff to $not\ endorse$ is $1 - 2\varepsilon_n$. Therefore, since $\varepsilon_n$ converges to 0, there must be some $N$ such that for all $n > N$, $m$ prefers to choose $\sigma^n(endorse, run) = \varepsilon_n \to 0$. Thus, the equilibrium in Figure 8.10 is not the limit of these completely mixed equilibria and is not a sequential equilibrium.

We can also see the restrictions that sequential equilibria place on beliefs. Weak consistency requires that the beliefs about the media's decision correspond to the strategy, $b^n(endorse) = \sigma^n(endorse, run)$. Accordingly for a sequence of completely mixed strategies that converges to a profile playing $not\ run$ and $endorse\ if\ run$ with probability 1, $b(\cdot)^n$ must converge to beliefs that put probability 1 on $endorse$. Therefore, there are no sequences of completely mixed strategies and weakly consistent beliefs that converge to the strategies and beliefs depicted in Figure 8.10. In fact it can be shown that every sequential equilibrium involves strategies that are subgame perfect. The proof of this result is left as an exercise.

While sequential equilibrium is an improvement over PBE, in many applications the two concepts turn out to have equivalent equilibrium sets. For example in the classic signaling games considered earlier the concepts coincide. Fudenberg and Tirole (1991) have shown that

PROPOSITION 8.9. *(Fudenberg and Tirole 1991) If $\Gamma^{EI}$ has only 2 periods or every player has at most two types then the set of PBE and SE coincide.*

**5.2. The Intuitive Criterion.** While the concept of sequential equilibrium can sometimes eliminate implausible PBE, it often continues to permit implausible out-of-equilibrium beliefs. To address these problems, Kreps and Cho's propose the *intuitive criterion* to further restrict the equilibrium set in many signaling games.

Since we wish categorize certain off-the path beliefs as implausible, Kreps and Cho postulate beliefs should be concentrated on those types of senders who have the greatest incentive to defect. They illustrate their concept with what has become known as the Beer and Quiche game. We follow this tradition because this is the simplest possible game upon which the concept can be demonstrated, but offer some political embellishments rephrasing the interaction as a game between Saddam Hussein and George Bush on the eve of the American invasion

of Iraq.   Consider a two-player signaling game in which the sender, Hussein has two possible types.  He can have weapons of mass destruction, $\theta = w$ or not, $\theta = \tilde{w}$.  The receiver, Bush can either attack $s_b = a$ or not $s_b = \tilde{a}$.  Bush has no concerns about winning if he attacks, but needs to justify the attack with the claim that Hussein has weapons of mass destruction.   Accordingly, he prefers to attack if $w$ and he prefers not to attack if $\tilde{w}$. Prior to Bush's decision, Hussein decides whether to allow weapons inspections $s_h = y$ or $s_h = \tilde{y}$.  However, the result of a weapons inspection would not be realized before Bush's decision of whether to attack or not.   We assume that regardless of Bush's decision type $w$ suffers a cost of 1 from inspection as he is shown to be in violation of UN resolutions. Type $\tilde{w}$ receives a benefit of 1 from inspection as he is publicly vindicated.  Regardless of Hussein's type he would prefer not to be attacked, and in fact the cost of being attacked is sufficiently large that Hussein of either type is willing to make either decision regarding UN inspections to avoid attack.  Bush on the other hand does not care directly about the inspections, but he prefers attacking as long as the probability of $\theta = w$ is sufficiently high (greater than $\frac{1}{2}$).   Figure 8.11 depicts the game.

## Insert Figure 8.11

It is not difficult to see that there are pooling equilibria in which weapons inspectors are allowed $(s_h = y)$ and there are pooling equilibria in which they are not allowed $(s_h = \tilde{y})$.  In any pooling equilibria, Bush's posterior must correspond to the prior.  Thus if the prior probability $w$ is sufficiently low, Bush will not attack.  To support pooling at $s_h = y$, it is necessary that both Hussein types prefer $y$ to $\tilde{y}$.  This results in the incentive compatibility conditions

$$3 \geq EU_H(\tilde{y}, \tilde{w})$$
$$2 \geq EU_H(\tilde{y}, w)$$

Letting $pr(a \mid \tilde{y})$ denote the probability that Bush attacks after observing $\tilde{w}$, the above conditions require that we have

$$3 \geq 0pr(a \mid \tilde{y}) + 2(1 - pr(a \mid \tilde{y}))$$
$$2 \geq 1pr(a \mid \tilde{y}) + 3(1 - pr(a \mid \tilde{y})).$$

This is true as long as $pr(a \mid \tilde{y}) \geq \frac{1}{2}$.   In order for Bush to use a strategy in which $pr(a \mid \tilde{y}) \geq \frac{1}{2}$ his posterior belief about Hussein's type conditional on the off-the path action $\tilde{y}$ needs to satisfy $pr(w \mid \tilde{y}) \geq \frac{1}{2}$.  This posterior is not pinned down by Bayes' Rule.  Recall that the on-the-path $pr(w \mid y) = \frac{1}{4}$ corresponds to the prior and is pinned down in a pooling equilibrium.  So we have shown that there is

a PBE (and by Proposition 8.9 as SE) in which both types of Husseins pool at $y$ and Bush does not attack.

To support the other pooling equilibrium where both types of Hussein select $\tilde{y}$ and Bush again does not attack, we need only specify off-the-path beliefs $pr(w \mid y) \geq \frac{1}{2}$. Following these beliefs, Bush attacks with probability at least $\frac{1}{2}$ and so both Hussein-types prefer to select $\tilde{y}$ and avoid attack.

While they are both PBE and SE, Kreps and Cho argue that only one of these pooling equilibria is reasonable. Consider the equilibria in which both Hussein types select $\tilde{y}$. This equilibria requires that the off-the-path belief satisfies $pr(\omega \mid y) \geq \frac{1}{2}$. Is it reasonable for Bush to believe that Hussein is more likely to have weapons if he allows inspections than if he doesn't? Recall that in this equilibrium $\theta = w$ is getting his maximal possible payoff. No inspections and no attack result in a payoff 3. However, the defection to $y$ and not attack results in a payoff of 2 for $w$. Such a deviation is not very desirable under the assumption that the deviation will not trigger attack. If an attack were triggered by $y$ then the defection is even less attractive. On the other hand consider the potential incentive for a deviation by a type $\tilde{w}$. In equilibrium he gets payoff 2. However if his defection did not result in an attack he would get utility of 3 (and thus improve his situation). Thus, it seems *intuitive* that if a defection were observed that it would most likely be committed by $\tilde{w}$. Kreps and Cho argue that Bush would be foolish to interpret $y$ as suggestive of $w$. Instead they imagine that the only justification for such an off the path deviation is that a type $\tilde{w}$ Hussein might deviate to $y$ and send the following justification.

> Dear W:
>
> Sorry for past squabbles with your old man. About this recent disagreement, I know that you are expecting me to choose $\tilde{y}$ and this doesn't tell you anything–its a pooling equilibrium after all (you remember pool tables from Yale don't you). But I am not going to do this, because I actually don't have any weapons and I want to show this to the world, so I am going to make myself even better off. Beside avoiding the tanks, special ops, and media barrage when you decide not to attack, I'm also going to allow weapons inspections in to show that I have been well-behaved. You should trust that this action indicates that I really have no weapons because if I did have weapons and I expected you not to attack if I didn't allow the inspections (which is the equilibrium we

are playing) then I would only hurt myself by letting in
the inspectors.

Sincerely

SH


Of course this type of communication is not modeled in standard
signalling games. The point is that given the Bush strategy, one type
can possibly gain from the off-the-path deviation, while the other type
can only lose. In such a setting, the off-the-path beliefs should be
concentrated on the type that stands to gain. Note that the $y$ pooling
equilibria is immune from this criticism. The only type that stands to
possibly gain from choosing $\tilde{y}$ is $w$. So the beliefs justifying Bush's
attack following $\tilde{y}$ are justified.

We now present the intuitive criterion in slightly a more rigorous
manner and so we require a bit more notation. Let $\Gamma^s$ denote a simple
signaling game with two periods  and two players. Player 1 has a
type space $\Theta$ and a message space $M$. Player 2 observes player 1's
message $m$ and selects an action from $A$. For simplicity assume that
all of these sets are finite. For the more complicated case of non-finite
sets, the following conditions can be modified but some technical issues
may be encountered. While player 1 knows her type, player 2 only
knows that 1's type is drawn from some probability mass function $f(\cdot)$
on $\Theta$ and player payoffs are given by utility functions $u_s(m, a, \theta)$ and
$u_r(m, a, \theta)$. Accordingly, a mixed strategy profile is a message function
$\sigma_s(\theta)$ that selects a lottery on $M$ for every $\theta$ and an an action function
$\sigma_r(m)$ that selects an action for each possible message. By $\sigma_s(m, \theta)$ and
$\sigma_s(a, m)$, we denote the probability that $m$ is played by a sender with
type $\theta$ and the probability that $a$ is played by an $r$ that has observed $m$
respectively. An equilibrium (PBE or SE) also involves a belief $\mu(\theta \mid
m)$. Given a signaling game and a sequential equilibrium to the game, let
$U_s^*(\theta)$ denote expected utility to player 1 of type $\theta$ from the equilibrium
profile. Finally let $\Delta$ denote the set of probability distributions on
$\Theta$ and let $BR_r(m) = \cup_{p(\theta) \in \Delta} \{\arg \max_{a \in A} \sum u_r(m, a, \theta) p(\theta)\}$ denote
the set of actions by $r$ that maximize the receiver's expected utility for
some beliefs about $\theta$. We say an action $r$ is rationalizable if it is an
element of $BR_r(m)$.


DEFINITION 8.7. *An SE $(\sigma_s(\cdot), \sigma_s(\cdot), F(\cdot \mid \cdot))$ satisfies the intuitive
criterion if for any message $m$ such that $\sum \sigma_r(m, \theta) f(\theta) = 0$, the pos-
terior belief $\mu(\theta \mid m) > 0$ only if $U_s^*(\theta) < \max_{a \in BR_r(m)} u_s(m, a, \theta)$.*

In words, an intuitive equilibrium requires that out-of-equilibrium beliefs put zero probability on types that could not gain from the observed deviation under some expectation that the receiver would respond to the deviation by playing a strategy from her set of best responses.

To further demonstrate the concept, we modify the entry deterrence game considered above. Now instead of restricting the message space to be $\{WC, \tilde{W}C\}$ we allow the incumbent to select a level $s_I \in \mathbb{R}^1_+$ at a cost $cs_I$ where $c$ is either $c_w$ or $c_s$ depending on the incumbent's type. Let the value of office be 1 for the incumbent, so if in equilibrium she accumulates $s'$ and wins with probability $\pi$ (because of randomness in the challenger's decision and the randomness associated with election in a contested race) her payoff is $\pi - cs'$. As you will see there are pooling, partially pooling and separating equilibria to this game. It can be shown however that there is exactly one intuitive equilibrium. This equilibrium is a separating equilibrium. We leave the analysis of this game as an exercise, and provide a solution in the back of the book. Students are strongly encouraged to devote the time to work through these two problems before looking at our solution.

The literature on refinements is quite large and refinements to the intuitive criterion have appeared in applications. Commonly, models with types spaces with more than two element require stronger refinements. This is because several types might stand to gain from a defection for different best responses by the receiver. In such situations, stronger refinements such as *universal divinity* (Banks and Sobel 1992).[6] Since universal divinity has been used in numerous political applications, we present a definition and an example.

DEFINITION 8.8. *An SE $(\sigma_s(\cdot), \sigma_s(\cdot), F(\cdot \mid \cdot))$ satisfies universal divinity if for any message $m$ such that $\sum \sigma_r(m, \theta) f(\theta) = 0$, the posterior belief $\mu(\theta \mid m) > 0$ only if there exists an action $a \in BR_r(m)$ such that $U^*_s(\theta) < u_s(m, a, \theta)$ and for every $\theta' \neq \theta$ $U^*_s(\theta') \geq u_s(m, a, \theta')$*

In comparing the two refinements, universal divinity is more stringent in the set of types that it allows the posteriors to place positive probability. In the case of universal divinity, the posteriors can only put weight on a type if there is rationalizable action that makes the deviation desirable for this type and not desirable for any other type. Informally, universal divinity requires that off the path beliefs put weight only on the types "most likely" to deviate.

---

[6]We suggest students to seek out the original Cho and Krep and Banks and Sobel pieces. In addition, Banks (1991) is an exemplary presentation of signaling games and refinements in political science.

To contrast the intuitive criterion and universal divinity, we consider a version Michael Spence's model of job signaling. We consider an application to the question of political reform as a signal to gain foreign investments or loan guarantees. Suppose a developing country has type $\theta \in \{1, 2, 3\}$ with $\theta$ measuring the nations potential to successfully repay debts (higher numbers are better). Let $\pi_1$ and $\pi_2$ denote the probability of types 1 and 2 (with type 3 occurring with probability $1 - \pi_1 - \pi_2$). The country must select a level of observable political reform $r \in \mathbb{R}^1_+$. The pain associated with reform is dependent on the nations type. After observing $r$ the IMF determines a financial package $f \in \mathbb{R}^1_+$ for the nation. We assume that the receiver's goal is to match $f$ with the nation's $\theta r$. Payoffs are as follows, given type $\theta$, reform $r$ and package $f$, the nation's utility is $f - \frac{1}{\theta} r^2$.

We first consider the case of $\pi_1 + \pi_2 = 1$ (so there are just two types of developing countries). In this case, there are pooling, partially pooling and separating sequential equilibria. However, the intuitive criterion selects a unique equilibrium. We sketch out the argument here. Consider Figure 8.12 which depicts indifference contours over pairs of messages and responses for senders of types 1 and 2. Since all senders prefer more funds and fewer reforms, movements to the northwest quadrant are desirable from their perspective.

## Insert Figure 8.12

Consider a pooling (or partially pooling) equilibrium in which both sender types are selecting the same level $r^p$ with positive probability. After observing $r^p$ the receiver knows that the posterior probability of $\theta < 2$ is greater than 0 and thus in any sequential equilibrium the package that corresponds with $r^p$, $f(r^p)$, must be less than $2r^p$. We show that there exists a message $r' > r^p$ such that if $f(r') = 2r'$ type 2 nations would prefer the deviation and type 1 nations would not. In this case the intuitive criterion implies that beliefs must place probability 1 on $\theta = 2$ if $r'$ is chosen. Given these beliefs following $r'$ the package $f(r') = 2r'$ is the unique sequentially rational package for the receiver. In the notation of our definitions above we have

$$U_r^*(2) = f(r^p) - \frac{r^p}{2}$$
$$u_r(r', 2r', 2) = 2r' - \frac{r'}{2}$$
$$U_r^*(1) = f(r^p) - r^p$$
$$u_r(r', 2r', 1) = 2r' - r'$$

Accordingly, $U_r^*(1) > u_r(r', 2r', 1)$ requires only $f(r^p) - r^p > r'$ while $U_r^*(2) < u_r(r', 2r', 2)$ requires only $f(r^p) - \frac{r^p}{2} < 2r' - \frac{r'}{2}$ or $\frac{r'-r^p}{2} < 2r' - f(r^p)$. Both of these inequalities can be simultaneously satisfied. See Figure 8.12 for the region of such values of $r'$. To recap, in a pooling or partially pooling equilibrium in which both types play $r^p$ with positive probability, the intuitive criterion requires that following the (possibly off the path) message of $r'$ beliefs assign probability 1 to type 2. Thus if $r'$ is played the financial package will be $2r'$. The value $r'$ was chosen so that type 2's strictly prefer message $r'$ to $r^p$ meaning that type 2's cannot put positive probability on the message $r^p$ contradicting the assumption that there is an intuitive equilibrium in which both types player $r^p$ with positive probability. Having ruled out all but separating equilibria, we claim that the intuitive criterion selects a unique separating equilibrium. We leave this as an exercise below.

## Insert Figure 8.13

We now assume that $\theta = 3$ occurs with probability $1 - \pi_1 - \pi_2 > 0$. The first question to address is whether the intuitive criterion still eliminates all pooling or partially pooling equilibria. The answer is no. To see this suppose that types 1 and 2 are both playing a message $r^p$ with positive probability, and $\theta = 3$ plays a pure strategy $r^3 > r^p$. It can be shown that this happens in some sequential equilibrium. Our argument before was that the intuitive criterion required that following some reform $r'$ that is higher than $r^p$ posterior beliefs are concentrated at the higher type. However, with the third type present this turns out not to be the case. We can now satisfy the intuitive criterion with posteriors putting weight $\theta = 3$ following a reform $r' > r^p$. With such beliefs, sequentially rational choices of $f$ might lead $\theta = 1$ to prefer the deviation to the equilibrium payoff. More specifically, the requirement that $f(r') \leq 2r'$ need no longer hold. Now it is just the case that $f(r') \leq 3r'$ needs to hold. With $f(r')$ this big, type $\theta = 1$ might be willing to deviate with the expectation that a deviation will result in $f(r')$. Accordingly in order for type $\theta = 2$ to signal that it is not type 1, it needs to send a message at least as high as $r^{\min}$ as depicted in Figure 8.13. However, for every level of $r > r^{\min}$, there are possible best responses $f(r)$ that make type 2 worse than under equilibrium. Notice that there is space between type 2's indifference curve and the line $f = 2r$. Since the intuitive criterion only requires that type 2 expect a response of $f > 2r$ for an $r$ greater than $r^{\min}$, it cannot be sure that the deviation is desirable.

However, if we test whether this partial pooling can happen in a universally divine equilibria the answer is different. Again suppose that

types 1 and 2 are both playing a message $r^p$ with positive probability, and $\theta = 3$ plays a pure strategy $r^3 > r^p$. Under universal divinity, a message of $r'$ that is slightly larger than $r^p$ must result in a posterior concentrated at $\theta = 2$. To see this, note that for values of $r$ to the right of $r^p$ type 1's indifference curve is above type 2's. This means that for pairs $(r', f(r'))$ that lie between the two indifference curves, we have $U_s^*(2) < u_s(r', f(r'), 2)$ and $U_s^*(1) > u_s(r', f(r'), 1)$. Moreover, since type 3 is getting $f = 3r$, his utility is higher in equilibrium than under the deviation. Accordingly, universal divinity requires that a message of $r'$ result in beliefs concentrated at $\theta = 2$ and thus sequential rationality requires that $f(r') = 2r'$, and so type 2 would gain from the deviation. It is left as an exercise to show that there is exactly one universally divine equilibrium in the game with 3 types.

## 6. Exercises

EXERCISE 8.1. *Consider the game of Figure 8.7, with the payoff to the path $B, NR$ being (5,0) instead of (4,0). Characterize all of the PBE (mixed and pure strategy) to the game.*

EXERCISE 8.2. *Consider the game of Figure 8.7, with the payoff to the path $ND$ being $(w, 5)$ instead of $(0, 5)$. Here $w$ is an exogenous parameter known to the agents that is ranging from $[-2, 5]$. For what regions of this range are there PBE in which $ND$ occurs with positive probability. In other words for what subset of $[-2, 5]$ are there PBE in which $ND$ is played.*

EXERCISE 8.3. *In the game depicted in Figure 8.8, show that there is not a PBE with $m(\theta) = \begin{cases} b \text{ if } \theta = a \\ a \text{ if } \theta = b \end{cases}$.*

EXERCISE 8.4. *Find all of the PBE of the game depicted in Figure 8.14.*

### Insert Figure 8.14 here

EXERCISE 8.5. *A Model of Political Repression*

Suppose that in each of two periods, society must decide whether to protest the policies of the state. The state may either acquiesce or repress. Society gets 1 if the state acquiesces, $-1$ the state represses, and 0 if it does not protest.

Suppose there are two types of states: Moderate and Hardline. The moderate state $(M)$ gets 0 if the protest does not take place, $-2$ if it acquiesces, and $-3$ if it represses. The hardline $(H)$ state gets 0 for no

protest, $-2$ for repression, and $-3$ for acquiescing. Let $p_0$ be the prior probability that the state is $M$.

a. In the second period, what is the critical value $p^*$ such that $S$ protests if $p_1 \geq p^*$ (where $p_1$ is $S$'s updated belief that the state is $M$)?

b. Is there a separating equilibrium with these strategies?

$M : \{acquiesce, acquiesce\}$

$H : \{repress, repress\}$

$S : \{protest, stay\ home\ if\ repressed\ in\ period1\}$

If so, what values of $p_0$ does it hold?

c. Is there a pooling equilibrium in the first period with these strategies?

$M : \{repress, acquiesce\}$

$H : \{repress, repress\}$

$S : \{protest, stay\ home\ if\ repressed\ in\ period1\}$

If so, what values of $p_0$ does it hold? Is it consistent with the intuitive criterion?

d. Is there a pooling equilibrium in the first period with these strategies?

$M : \{repress, acquiesce\}$

$H : \{repress, repress\}$

$S : \{stayhome, stay\ home\ if\ repressed\ in\ period1\}$

If so, what values of $p_0$ does it hold? Is it consistent with the intuitive criterion?

e. Compute a semi-pooling equilibrium where $M$ represses in the first period with probability $q$. For what values of $p_0$ does $S$ protest?

EXERCISE 8.6. *Show that there are no partial pooling equilibria in open rule version of the Gilligan-Krehbiel model.*

EXERCISE 8.7. *Compute the other partial pooling equilibria for the Warchest Game.*

EXERCISE 8.8. *Consider the open rule version of the Gilligan Krehbiel model described above, but suppose that their are two committee members with $c > \theta$ that observe the state and make simultaneous messages to the floor. Does a separating equilibrium exist. Now suppose*

*that their are three such committee members does a separating equilibrium exist?*

EXERCISE 8.9. *Show that given a finite extensive form game if $\sigma(\cdot), b(\cdot)$ is constitutes a sequential equilibrium then $\sigma(\cdot)$ is subgame perfect.*

EXERCISE 8.10. *Show that in a one sender, one receiver signaling game if the senders type space has two elements the set of PBE and SE coincide.*

EXERCISE 8.11. *Prove Proposition 8.9.*

EXERCISE 8.12. *Show that in the Hussein-Bush game the pooling equilibria with y on the path does not violate the intuitive criterion.*

EXERCISE 8.13. *Can you modify the probability of w in Figure 8.11 to support the observed path of play ( $\tilde{y}$ and a) as a PBE?*

EXERCISE 8.14. *Characterize the levels of $s_I$ and entry lotteries that are supportable in a PBE to the modified entry game with message space $\mathbb{R}_+^1$ ((note this question should be answered in the book))*

EXERCISE 8.15. *Characterize the unique intuitive equilibrium to the modified entry game with message space$\mathbb{R}_+^1$ ((note this question should be answered in the book)).*

EXERCISE 8.16. *In the reform for loan guarantees game with 2 types show that the unique intuitive equilibrium involves $r(1) = \arg\max\{r - r^2\} = \frac{1}{2}$ and $r(2) = \arg\max\{r - \frac{1}{2}r^2\} = 1$.*

EXERCISE 8.17. *In the reform for loan guarantees game with 3 types show that the unique universally divine equilibrium involves $r(1) = \arg\max\{r - r^2\} = \frac{1}{2}$ and $r(2) = \arg\max\{r - \frac{1}{2}r^2\} = 1$ and $r(3) = \arg\max\{r - \frac{1}{3}r^2\} = \frac{3}{2}$.*

CHAPTER 9

# Repeated Games

Many important models in political game theory are based situations consisting of agents playing the same game repeatedly over time. In many of these cases, the authors are interested in how certain social practices like conventions, norms, cooperation, and trust are sustained when actors may appear to have short run incentives to deviate from the expected behaviors. Another significant application of "repeated" is to understand how social dilemmas such as the Prisoner's dilemma can be solved without recourse to centralized authority (Taylor 1976).

The most interesting conceptual issue in such games is the extent to which repetition creates the opportunity to sustain more behavior as Nash equilibria than is possible in single-shot games. As we will see, in general, the set of Nash equilibria is much larger in repeated games than the corresponding static versions. This is because expectations about the future can lead to playing strategies that would not be optimal in a static context.

To generate some intuition as to how expectation about the future can influence behavior, consider the following abstract normal form game.

| Table 9.1 | | | |
|---|---|---|---|
| 1\2 | $L$ | $M$ | $R$ |
| $T$ | $8, 8$ | $0, 0$ | $1, 9$ |
| $M$ | $0, 0$ | $5, 5^*$ | $0, 0$ |
| $B$ | $9, 1$ | $0, 0$ | $3, 3^*$ |

Note that if this game is played only once there is only two Nash equilibria: $(M, M)$ and $(B, R)$. Even though the strategy profile $(T, L)$ provides the highest aggregate payoffs, it is not an Nash equilibrium since player 1 would defect to $B$ and player 2 would would defect to $R$. Now consider what happens if this game is played twice. Suppose player 1 chooses the strategy, "play $T$ in period 1 and play $M$ in period if player 2 plays $L$ in period 1. Otherwise play $B$ in period 2." Furthermore, suppose that player 2 chooses the strategy "play $L$ in period 1 and play $M$ in period 2 if player 1 plays $T$ in period 1. Otherwise play

$R$ in period 2." Note that this pair of strategies is a Nash equilibrium where the "good" outcome $(T, L)$ is played in the first period. To see this note if either player defects, her payoff 12 which is less than the equilibrium utility of 13. In fact, these strategies not only constitute a Nash equilibrium but the equilibrium is also subgame perfect. Since $(B, R)$ is an equilibrium of the one-shot game, it is an equilibrium to all of the subgames that follow actions other than $(T, L)$. Note that $(T, L)$ cannot be a Nash equilibrium to the second subgame.

By repeating the game the players can use a "norm" that if there is cooperation in the first period, the good equilibrium will be played in the second period. Otherwise, the bad equilibrium will be played. A second important point is that repeating the game can only improve the player's average utility. Since playing $(B, R)$ in both periods is also a subgame perfect Nash equilibrium, repetition cannot lower average utility below $(3, 3)$. As we will see almost any set of payoffs can be a Nash equilibrium if the game is long enough and the players care enough about the future.

## 1. The Repeated Prisoner's Dilemma

It is often argued that trade policy among nations is an example of the Prisoner's Dilemma played repeatedly over time. It is generally thought that the world economy does better when all nations agree to free trade, but that individual countries might do better by protecting their domestic economy. Given this tension, the question arises as to how to sustain free trade regimes. One answer is that free trade can be supported as an equilibrium in a repeated game where trade wars begin whenever a major country defects from the trade agreement. To illustrate this argument, consider the following representation of a trade policy dilemma between the US and the EU.

| Table 9.2: Free Trade Game | | |
|---|---|---|
| US\EU | *Free Trade* | *Protect* |
| *Free Trade* | $10, 10$ | $1, 12$ |
| *Protect* | $12, 1$ | $4, 4$ |

Obviously, if the game is just played once, the unique Nash equilibrium is the strategy profile (*Protect*, *Protect*). Suppose that it is played twice (ignoring the discounting of future payoffs), then the strategy sets for each player are (where the period 2 strategies depend on the strategy of the other player)

$$FT_1FT_2FT_2, FT_1FT_2P_2, FT_1P_2FT_2, FT_1P_2P_2, P_1FT_2FT_2, P_1FT_2P_2,$$
$$P_1P_2FT_2, P_1P_2P_2$$

where $FT_1FT_2P_2$ means "play $FT$ in period 1 and play $FT$ in period 2 if the other country plays $FT$ in period 1 otherwise play $P$."

| US\EU | $FT_1FT_2FT_2$ | $FT_1FT_2P_2$ | $FT_1P_2FT_2$ | $FT_1P_2P_2$ | $P_1FT_2FT_2$ | $P_1FT_2P_2$ | $P_1P_2FT_2$ | $P_1P_2P_2$ |
|---|---|---|---|---|---|---|---|---|
| $FT_1FT_2FT_2$ | 20,20 | 20,20 | 11,22 | 11,22 | 11,22 | 11,22 | 2,24 | 2,24 |
| $FT_1FT_2P_2$ | 20,20 | 20,20 | 11,22 | 11,22 | 13,13 | 13,13 | 5,16 | 5,16 |
| $FT_1P_2FT_2$ | 22,11 | 22,11 | 14,14 | 14,14 | 11,22 | 11,22 | 2,24 | 2,24 |
| $FT_1P_2P_2$ | 22,11 | 22,11 | 14,14 | 14,14 | 13,13 | 13,13 | 5,16 | 5,16 |
| $P_1FT_2FT_2$ | 22,11 | 13,13 | 22,11 | 13,13 | 14,14 | 5,16 | 14,14 | 5,16 |
| $P_1FT_2P_2$ | 22,11 | 13,13 | 22,11 | 13,13 | 16,5 | 8,8 | 16,5 | 8,8 |
| $P_1P_2FT_2$ | 24,2 | 16,5 | 24,2 | 16,5 | 14,14 | 5,16 | 14,14 | 5,16 |
| $P_1P_2P_2$ | 24,2 | 16,5 | 24,2 | 16,5 | 16,5 | 8,8 | 16,5 | 8,8[*] |

**Table 9.3: Two-Period Free Trade Game**

Unlike our first example, repeated the game only once does not effect behavior as $(P_1P_2P_2, P_1P_2P_3)$ is the only Nash equilibrium. This result can be generalized to any finite number of periods. In the last period, each country will like to protect. Since this is known in the penultimate period, each country will have an incentive to protect in this period as well. This process unravels until each country is protecting in every period.

The reason that we were able to induce some cooperation in our first example was that the first period behavior helped coordinate between multiple equilibrium in the second period. The good equilibrium was used as a reward while the bad equilibrium was used as a punishment. However, since the Prisoner's Dilemma has but one Nash equilibrium, it is impossible to encourage cooperation with the promise coordinating on a good equilibrium in the future.

However, if the game lasts an infinite number of periods, this ceases to be an issue. Now suppose that the "good equilibrium" is to free trade in every period while the "bad equilibrium" is to protect in every period. Since there is no last period, the good equilibrium does not unravel as it did in the finite case. Thus, in every single period, cooperation is sustained by the reward of the good equilibrium and the sanction of the bad.

## 2. The Grim Trigger Equilibrium

To see how infinite repetition eliminates the "last period" problem, consider the following strategy in the infinite period trade game: "Play free trade in every period until the other country protects, then protect forever." This is known as the *grim trigger strategy*, because any failure to cooperate leads to the non-cooperative equilibrium in all future periods. If each country plays this strategy, both receive 10 in every period. Assuming that both countries discount the future at a common rate $\delta$, the long-term utility of this strategy is $\frac{10}{1-\delta}$.[1] To show that this strategy is a Nash equilibrium, we must show that neither player is willing to defect. It is necessary and sufficient to show that no player is willing to defect for one period. Because of each stage game is identical, either a player will want to defect in every period or in no period. Defection from this strategy gives the defector 12 in the period of the defection. Since the other player will protect forever following the defection, the defector's best response is to protect forever following the defection. Therefore, the defector gets 4 in every period following the defection. The utility from defecting is therefore $12 + \frac{\delta 4}{1-\delta}$. Thus, the grim trigger strategies are a Nash equilibrium to the repeated prisoners dilemma if and only if $\frac{10}{1-\delta} \geq 12 + \frac{\delta 4}{1-\delta}$ or $\delta \geq \frac{1}{4}$. So as long as the players are sufficiently patient ($\delta$ large), the grim trigger strategy is a Nash equilibrium. We can also show that the grim trigger equilibrium is subgame perfect. A proper subgame to this game is also an infinitely repeated prisoner's dilemma. Since the grim trigger equilibrium is a Nash equilibrium for the whole game, it must be an equilibrium in each of the subgames.

Now consider a generalized Prisoner's Dilemma

| Table 9.4: Generalized Prisoner's Dilemma | | |
|---|---|---|
| 1\2 | Cooperate | Don't cooperate |
| Cooperate | $a, a$ | $d, c$ |
| Don't cooperate | $c, d$ | $b, b$ |

where $c > a > b > d$. Using exactly the same steps as above, the grim trigger strategy is a SPNE if and only if $\frac{a}{1-\delta} \geq c + \frac{\delta b}{1-\delta}$ or

$$\delta \geq \frac{c - a}{c - b}$$

---

[1]See chapter 3 for a discussion of time discounting and the calculation infinite sums of dicounted utilities.

Thus, we find that cooperation is harder to sustain (requires a higher discount factor) when:

(1) $c$ is large relative to $a$ and $b$.

(2) $a$ and $b$ are roughly equal.

## 3. Tit-for-Tat Strategies

The grim trigger strategy is not the only SPNE to the infinitely repeated prisoner's dilemma which sustains the cooperative outcome. Many authors find the grim trigger equilibrium unrealistic because it predicts that cooperation disappears forever following a single defection. Therefore, it is not very robust to mistakes by the players. It assumes that the player's cannot renegotiate a return to the cooperative phase which they would clearly have a incentive to do. However, if we assume that players can engage in such renegotiation, cooperation vanishes because the uncooperative equilibrium is no longer a deterrent.

An alternative SPNE is based on "tit for tat" strategies of the following form "cooperate until your opponent cheats. Then cheat until your opponent cooperates, then cooperate." Note that there are two possible subgames:

(1) A sub-game where both players are expected to cooperate in the next iteration. We call this the cooperation phase.

(2) A sub-game where the defector is supposed to cooperate and the other player is supposed to punish the defector by not cooperating in the next iteration. We call this the punishment phase.

Consider the first subgame. Again the equilibrium utility is $\frac{a}{1-\delta}$, but the utility from a defection is a bit more complicated. Defecting gives a one -period utility of $c$. In the next period, the defector gets $d$ from cooperating and being punished by the other player. In third period, the game returns to cooperation. Note that since we are intent upon establishing the existence of a Nash equilibrium, we need only check for one time deviations. Therefore, we assume that cooperation lasts forever following the punishment phase. Thus, the total utility is $c + \delta d + \delta^2 a + \delta^3 a + ... = c + \delta d + \frac{\delta^2 a}{1-\delta}$. Thus, a player will not defect during this type of subgame if $\frac{a}{1-\delta} \geq c + \delta d + \frac{\delta^2 a}{1-\delta}$ or

$$\delta \geq \frac{c - a}{a - d}$$

Now consider the second subgame. The equilibrium utility for the previous defector $d + \frac{\delta a}{1-\delta}$ which cooperating in every period while the

other player punishes him in the first period. If a player defects from this subgame, they get $b$ in the current period, $d$ in the next period, and $a$ after that. Thus, a player will not defect if $d + \frac{\delta a}{1-\delta} \geq b + \delta d + \frac{\delta^2 a}{1-\delta}$ or

$$\delta \geq \frac{b - d}{a - d}$$

Given these two results, we have established that if $\delta \geq \max \left\{ \frac{c-a}{a-d}, \frac{b-d}{a-d} \right\}$.the "tit for tat" strategies constitute a SPNE. Intuitively, we might expect tit-for-tat strategies to sustain more cooperation than the grim-trigger if the payoff of $d$ were a sufficiently large deterrent against defecting in the cooperation phase. However, this effect is counteracted by the fact that if $d$ is sufficiently bad, a defector has a stronger incentive to defect from the punishment phase. In fact we can show that $\max \left\{ \frac{c-a}{a-d}, \frac{b-d}{a-d} \right\} \geq \frac{c-a}{c-b}$ so that cooperation is always easier to sustain under the grim trigger strategy than under tit-for-tat.[2] However, note the following:

(1) The grim trigger strategy is not optimal when players may make mistakes. Tit-for-tat may be better because the effect of mistakes is not permanent.

(2) The grim trigger strategy is not *renegotiation proof*. Both players could do better by renegotiating to leave the punishment phase and return to the original equilibrium. However, if both players foresee this possibility the punishment phase will not be an effective deterrent. Tit-for-tat is renegotiation-proof. The punisher gets *higher* utility in the punishment phase and will not wish to renegotiate.

## 4. Intermediate Punishment Strategies

The grim trigger and tit-for-tat strategies represent just two of the possible supergame strategies that may sustain the cooperative outcome. We can generalize this class of strategies to include strategies that involve punishment phases of intermediate length. Consider the following strategies:

(1) Cooperate until your opponent defects. Then do not cooperate for $k$ periods. Return to cooperating after the punishment phase ends.

(2) Cooperate until your opponent defects. Then do not cooperate for $k$ periods if your opponent cooperates. If your opponent does not cooperate at any point during this punishment phase, begin a new punishment phase of $k$ periods.

---

[2]See Axelrod (1984) for evidence that real-life players typically choose strategies resembling tit-for-tat.

Strategy 1 is similar to the grim trigger strategy in that a punishment consists of a reversion to the strategy pair (don't cooperate, don't cooperate). However, now the punishment phase is finite. The second strategy is similar to tit-for-tat in that any defector is punished by cooperating while the other player does not.

We consider strategy 1 first. Clearly there is no incentive to defect during a punishment phase since mutual non-cooperation is a Nash equilibrium. We need only consider defections from a cooperation phase. Now the payoff from a single defection during a cooperative phase consists of the one period gain from defecting, the discounted utility of $b$ for $k$ periods, and the utility of getting $a$ every period of the end of the punishment phase.[3] Thus, using the rules for sums of discount factors we encountered in chapter (choice theory), we can write this utility as $c + \frac{\delta - \delta^{k+1}}{1-\delta} b + \frac{\delta^{k+1}}{1-\delta} a$ Thus, sustaining cooperation requires that $a\left(1 - \delta^{k+1}\right) \geq (1 - \delta) c + \left(\delta - \delta^{k+1}\right) b$. While we cannot get a closed form for the critical value of $\delta$, note that we can re-write this expression as

$$\delta > \frac{c - a}{c - b} + \delta^{k+1} \frac{a - b}{c - b}$$

Note that the first term on the right side of the inequality is the critical value for the grim trigger strategy while the second term is positive for any finite $k$. Thus, not surprisingly, it is harder to sustain cooperation with a finite punishment phase. However, in a model where players may make mistakes, this equilibrium may be preferred to the grim trigger strategy.[4]

Now we consider strategy 2. First lets consider a defection from the cooperation phase. The payoffs from a defection consist of a one period benefit $c$, a punishment payoff of $d$ for $k$ periods, and a return top cooperative payoffs $a$ at the end of the punishment. Summing all of these up generates $c + \frac{\delta - \delta^{k+1}}{1-\delta} d + \frac{\delta^{k+1}}{1-\delta} a$. Simple algebra reveals that this payoff is lower than the payoff from defection in the tit-for-tat equilibrium by $\frac{\delta^2 - \delta^{k+1}}{1-\delta}(a - d)$. Thus, increasing the length of the punishment phase decreases the incentive to defect from the cooperative phase.

However, increasing $k$ may not make such an equilibrium easier to sustain as it reduces the incentive to comply in the punishment phase. To see this, consider the payoffs from defecting from the punishment

---

[3]Again the logic of Nash equilibrium suggests we can ignore the possibility of future defections.

[4]Like the grim trigger SPNE, this one is not renegotiation proof.

phase. These payoffs consist of getting $b$ for one period, $d$ for $k$ periods, and then returning to $a$ or $b + \frac{\delta - \delta^{k+1}}{1-\delta}d + \frac{\delta^{k+1}}{1-\delta}a$. The payoffs from adhering to the equilibrium in the punishment phase depends on which period of the punishment phase the game is in. Since we have to verify compliance in each period, we need to ensure compliance in the period where the payoff to compliance is lowest, the first period of the punishment phase. Thus, the utility for complying with the punishment in this period is $\frac{1-\delta^k}{1-\delta}d + \frac{\delta^k}{1-\delta}a$. Thus, compliance with the punishment requires

$$\delta > \left(\frac{b-d}{a-d}\right)^{\frac{1}{k}}$$

This critical value is clearly diminishing in $k$. A SPNE in these strategies requires that both conditions on $\delta$ be satisfied.

## 5. The Folk Theorem*

A common theme of our examples is that so long as the agents are sufficiently patient outcomes that are not Nash equilibria in static games can be supported as SPNE of infinitely repeated games. This result generalizes significantly. In fact, any individually rational payoff to an infinitely repeated game can be sustained as a SPNE if agents are sufficiently patient. This important result has been well known for so long that no one knows who derived it first. It has therefore been afforded the status of a *Folk* theorem. In this section, we formally prove a version of this result.

The primitives of a repeated game are a normal form stage game $\Gamma = \langle N, S, u \rangle$ and vector of agent discount rates $\boldsymbol{\delta} = (\delta_1, ..., \delta_n)$. In each period $t \in \{1, 2, 3, ...\}$, the agents play the normal form game $\Gamma$. The game $\Gamma$ is often called the stage game to distinguish it from the repeated game. Before agent $i$ selects $s_i^t \in S_i$, her strategy in period $t$ she observes the strategy profile $s^{t-1}$ played in period $t-1$. Moreover, we maintain the assumption of perfect recall, meaning that $s_i^t$ can be conditioned on the history $h^{t-1} = (s^1, ....., s^{t-1}) \in S^{t-1} := \prod_{j=1}^{t-1} S$. The null history is $h^0 = \emptyset$. A pure strategy for player $i$ is then a sequence of mappings $\{s_i^t(h^{t-1}) : S^{t-1} \to S_i\}_{t=1}^{\infty}$. A mixed strategy is a sequence of mappings $\{\sigma_i^t(h^{t-1}) : S^{t-1} \to \Delta(S_i)\}_{t=1}^{\infty}$. Given a sequence of lotteries over stage game profiles $\{\sigma^t\}_{t=1}^{\infty}$ agent $i$'s expected utility is given by $\mathbb{E}U_i(\{\sigma^t\}_{t=1}^{\infty}) = (1 - \delta_i) \sum_{t=1}^{\infty} \delta_i^{t-1} \mathbb{E}_{\sigma^t} u_i(s^t)$ where $\mathbb{E}_{\sigma^t} u_i(s^t)$ takes the expectation of $u_i(s^t)$ over the mixture $\sigma^t$. The multiplier $(1 - \delta_i)$ is included so that for a constant sequence $\sigma^t$, $\mathbb{E}U_i(\{\sigma^t\}_{t=1}^{\infty}) = u_i(\sigma^t)$. We denote the repeated game induced by a stage game, by

$\Gamma^\infty = \langle N, S, u, \delta \rangle$. Of course a repeated game is also an extensive form game and our notions of NE and SPNE are well defined in the repeated game.

We now focus on repeated games generated by finite normal form stage games. Given Nash's theorem, we know that every such stage game has at least one mixed strategy NE. It is not surprising then that every such repeated game has as a mixed strategy SPNE the infinite repetition of the stage game mixed strategy NE.

PROPOSITION 9.1. *If $\sigma^*$ is a SPNE of the stage game then the repeated game profile $\sigma_i^t(h^{t-1}) = \sigma_i^*$ for every $(h^{t-1})$ for every $t$ for every $i$ is a SPNE of the repeated game.*

An interesting feature of repeated games is that the set of SPNE is usually very large. The class of results termed "Folk theorems" serve to quantify the set of equilibrium payoffs that are supportable in an equilibrium. We prove a particularly useful and simple Folk theorem. We first need several definitions.

DEFINITION 9.1. *The payoff vector $v \in \mathbb{R}^n$ is **individually rational** if*

$$v_i \geq \min_{s_{-i} \in S_{-i}} \left\{ \max_{s_i \in S_i} u_i(s_i, s_{-i}) \right\}.$$

The value $\min_{s_{-i} \in S_{-i}} \{\max_{s_i \in S_i} u_i(s_i, s_{-i})\}$ is the minimum stage game utility that player $i$ can attain when she plays a best response. This value is identified by letting the players $-i$ select $s_{-i}$ so as to minimize the utility to $i$ of playing a best response to $s_{-i}$.

DEFINITION 9.2. *The payoff vector $v \in \mathbb{R}^n$ is **feasible** if there is some sequence of pure strategy stage game profiles $\{s^t\}_{t=1}^\infty$ such that for each $i \in N$, $\mathbb{E}U_i(\{s^t\}_{t=1}^\infty) = v_i$.*

PROPOSITION 9.2. *For every feasible and individually rational payoff vector $v \in \mathbb{R}^n$ there is an $n$-tuple of discount rates $\delta$ s.t. the payoff vector $v$ occurs in a NE of the repeated game with the discount rates $\boldsymbol{\delta}$.*

PROOF. Assume that $v$ is feasible and individually rational. Let $\{s^{vt}\}$ be a strategy profile that calls for playing the strategy that attains the payoff vector $v$ as long as no one has previously deviated from this strategy or more than two players have deviated, and plays the strategy $\{s^{pt} = \arg\min_{s_{-i} \in S_{-i}} \{\max_{s_i \in S_i} u_i(s_i, s_{-i})\}\}$ which punishes the unique player that deviated in all subsequent periods. At any period $t$ the payoff to $i$ of playing $\{s^{vt}\}$ is $v_i$ and the payoff to deviating is bounded

by

$$(1 - \delta_i^t)v_i + \delta_i^t(1 - \delta_i) \max_{s \in S} u_i(s) + \delta_i^{t+1} \min_{s_{-i} \in S_{-i}} \left\{ \max_{s_i \in S_i} u_i(s_i, s_{-i}) \right\}.$$

This value is less than $v_i$ if

$$\delta_i \geq \frac{\max_{s \in S} u_i(s) - v_i}{\max_{s \in S} u_i(s) - \min_{s_{-i} \in S_{-i}} \left\{ \max_{s_i \in S_i} u_i(s_i, s_{-i}) \right\}}.$$

Since $\max_{s \in S} u_i(s) \geq v_i \geq \min_{s_{-i} \in S_{-i}} \left\{ \max_{s_i \in S_i} u_i(s_i, s_{-i}) \right\}$ the right hand side is strictly less than 1. Thus as long as this condition is satisfied for each $i \in N$ the conjectured strategy profile is a NE to the repeated game.■                                                                    □

The equilibria used in the proof need not be SPNE as the punishment might be very costly to impose. We can quantify a set of payoff vectors supportable in SPNE to the repeated game using reversion to stage game NE strategies as the punishment.

PROPOSITION 9.3. *If $v \in \mathbb{R}^n$ is a feasible payoff vector for which there is some mixed strategy stage game NE which yields the payoff vector $v'$ s.t. $v'_i < v_i$ for every $i \in N$ then there is a SPNE in the repeated game which yields the payoff vector $v$.*

## 6. Application: Interethnic Cooperation

Fearon and Laitin (1996) use infinitely repeated games to understand how inter-ethnic cooperation might be sustained. Consider two groups $A$ and $B$ both with $n$ (even) members. In each period $t$, players are randomly matched to play the following Prisoner's dilemma.

| Table 5: Inter-Ethic Cooperation Game | | |
|---|---|---|
| 1/2 | $Cooperate$ | $Defect$ |
| $Cooperate$ | 1,1 | $-\beta, a$ |
| $Defect$ | $\alpha, -\beta$ | 0,0 |

where $\alpha > 1$, $\beta > 0$, and $\frac{\alpha - \beta}{2} < 1$. Further suppose that each of the members has a common discount factor $\delta \in (0, 1)$. In each period $m$ members of each group are selected to be paired with members of the other group while the remaining $n - m$ are matched with members of their own group. This random matching process suggests that each player will have a $p = \frac{m}{n}$ probability of being matched in an "out-group" member.

To capture the dynamics of intergroup and intragroup interaction, Laitin and Fearon assume that within groups the entire history of play

is observed by all members of the group. However, the history of play for members of the other group are not observed. Thus, in their model inter-ethnic cooperation is hard to sustain because those who defect in inter-group interactions cannot be individually singled out for punishment by members of the other group. Nevertheless, Fearon and Laitin argue that cooperation can be sustained even in the absence of these direct sanctions. They consider two such equilibria to this game. The first is what they call the *spiral* equilibrium. In this equilibrium, cooperation is supported within groups by $k^{in}$ period punishments against individual defectors. However, intergroup cooperation is sustained by the threat of group specific punishment phases of $k^{out}$ periods. During these punishment phases, all members of a given group are punished by the other group if any has defected in an inter-group interaction. The second equilibrium is the *in group policing* equilibrium in which there is no cross-group punishments but each group punishes its own for defections against the other group. Below we analyze the in-group policing equilibrium and refer the reader to the original article for the discussion of the spiral equilibrium.

**6.1. The In-group Policing Equilibrium.** The strategy for the in-group policing equilibrium follows.

> Play $C$ in all out-group pairings. For in-group pairings, always play $C$ with any partner not in a punishment phase, and $D$ in a punishment phase. A player enters or restarts a punishment phase for $k^{gp}$ periods by defecting against the out-group member or against an in-group member

We will focus on group $A$ as the proof extends identically to the strategies of group $B$. Let $s_t = (k_1, k_2, ...k_n)$ be the *state* of the system where $k_i$ is the number of periods remaining in the punishment period for player $i$ at the beginning of period $t$. If $k_i = 0$, we say that player $i$ is a cooperator and that player $i$ is a defector if $k_i > 0$. For a given state $s_t$ and any integer $l > 0$, let $n_{t+l}$ be the number of members of group $A$ who will be cooperators in period $t + l$, assuming that each plays the equilibrium strategy from period $t$ to period $t+l$. Therefore, $q_{t+l} = \frac{n_{t+l}}{n-1}$ is the probability of facing a cooperator in an in-group interaction.

To demonstrate that these strategies constitute a subgame perfect Nash equilibrium, we need to verify the following conditions.

(1) A cooperator $i$ has no incentive to
   (a) to defect against any out-group player

      (b) to defect against an in-group
      (c) to cooperate with any in-group defector
    2. a defector $i$ has no incentive to
      (a) to defect against an out-group player
      (b) to defect against an in-group cooperator
      (c) to cooperate against an in-group defector

Clearly, conditions 1(c) and 2(c) will be always satisfied since those deviations lower utility and the current period without affecting strategies of any other player (these deviations do not trigger punishments). Further note that 1(a) and 1(b) reflect the same trade-offs since both deviations generate a payoff of $\alpha$ in period $t$ followed by $k^{gp}$ periods of punishment. Thus, we need only establish that there will be no incentive to deviate in cases 1(a), 2(a), and 2(b).

We can write the utility for cooperation against an in-group cooperator or out-group member as

$$1 + \sum_{l=1}^{\infty} \delta^i (p + (1 - p)(q_{t+l} + (1 - q_{t+l})\alpha)$$

while the utility from deviations 1(a) and 1(b) is
(9.1)
$$\alpha + \sum_{l=1}^{k^{gp}} \delta^l(p + (1-p)(-q_{t+l}\beta + (1-q_{t+l})0) + \sum_{l=k^{gp}+1}^{\infty} \delta^i(p + (1-p)(q_{t+l} + (1-q_{t+l})\alpha)$$

The net utility of cooperating is therefore

(9.2) $$1 - \alpha + \sum_{l=1}^{k^{gp}} \delta^l((1 - p)(q_{t+l}(1 + \beta) + (1 - q_{t+l})\alpha)$$

To show that deviations 1(a) and 1(b) will not occur, we need equation 9.2 to be positive for all states and the resulting sequences of $q_{t+l}$. If $\alpha > 1 + \beta$, equation 9.2 is minimized by $q_{t+l} = 1$ for $l = 1, k^{gp}$. This is the path following $s_t = (0, 0, ..., 0)$. The net utility is positive following this state if and only if

(9.3) $$\delta^{k^{gp}} \leq 1 - \frac{(1 - \delta)(\alpha - 1)}{\delta(1 - p)(1 + \beta)}$$

Now consider the case where $1 + \beta > \alpha$. The the net utility would be minimized at $q_{t+l} = 0$ for $l = 1, k^{gp}$. However, given the definition of $q$ this is an infeasible sequence since all players are assumed to cooperate in their punishment phases and terminate their punishments after $k^{gp}$ periods. Thus, the minimizing sequence is one where all players defect in time $t - 1$ and return to cooperation status in period $t + k^{gp} - 1$.

The sequence of $q$ is therefore $q_{t+l} = 0$ for $l = 1, k^{gp} - 1$ and $q_{t+k^{gp}} = 1$. Thus, we now require (after some algebra) that

$$(9.4) \qquad \frac{(\delta - \delta^{k^{gp}})}{1 - \delta} \frac{\alpha}{(1 + \beta)} + \delta^{k^{gp}} \geq \frac{\alpha - 1}{(1 - p)(1 + \beta)}$$

Now we need to check to see whether a defector at time $t$ will wish to make deviation 2(a). Suppose that a defector with $k_i$ is paired against an out-group player. The utility of cooperating is

$$1 + \sum_{l=1}^{k_i-1} \delta^l(p+(1-p)(-q_{t+l}\beta+(1-q_{t+l})0) + \sum_{l=k_i}^{\infty} \delta^i(p+(1-p)(q_{t+l}+(1-q_{t+l})\alpha)$$

while the utility of the deviation is given by equation 9.1. Thus, the defector will cooperate with an out-group member so long as

$$\sum_{l=k_i}^{k^{gp}} \delta^i((1-p)(q_{t+l}(1+\beta)+(1-q_{t+l})\alpha) \geq \alpha - 1$$

The right side of this inequality is minimized when $k_i = k^{gp}$ so that we require

$$\delta^{k^{gp}}((1-p)(q_{t+k^{gp}}(1+\beta)+(1-q_{t+k^{gp}+l})\alpha) \geq \alpha - 1$$

Since all players play according to the equilibrium strategy, $q_{t+k^{gp}} = 1$ for all $s_t$. Therefore, we require

$$(9.5) \qquad \delta^{k^{gp}} \geq \frac{\alpha - 1}{(1 - p)(1 + \beta)}$$

Finally, we need to check deviation 2(b). So assume a defector with $k_i = 1$ is paired against as cooperator. The utility of cooperating is

$$-\beta + \sum_{l=1}^{k_i-1} \delta^l(p+(1-p)(-q_{t+l}\beta+(1-q_{t+l})0) + \sum_{l=k_i}^{\infty} \delta^i(p+(1-p)(q_{t+l}+(1-q_{t+l})\alpha)$$

Again the utility for defecting is given by

$$\sum_{l=1}^{k^{gp}} \delta^l(p+(1-p)(-q_{t+l}\beta+(1-q_{t+l})0) + \sum_{l=k^{gp}+1}^{\infty} \delta^i(p+(1-p)(q_{t+l}+(1-q_{t+l})\alpha)$$

so that the net utility of cooperating is $\sum_{l=k_i}^{k^{gp}} \delta^i((1-p)(q_{t+l}(1+\beta)+(1- q_{t+l})\alpha) - \beta$. Using the same argument as we did on 2(a), our SPNE requires

$$(9.6) \qquad \delta^{k^{gp}} \geq \frac{\beta}{(1 - p)(1 + \beta)}$$

Now we have a full set of equilibrium conditions. First, consider the case of $\alpha > 1 + \beta$, we require 9.3,9.5,and 9.6. First, note that equation 9.6 holds whenever 9.5 does. Less obviously, we can show that equation 9.5 implies equation 9.3. We can re-write 9.3 as

$$(9.7) \qquad \delta \left( \delta^{k^{gp}} - \frac{\alpha - 1}{(1-p)(1+\beta)} \right) \leq \delta - \frac{(\alpha - 1)}{(1-p)(1+\beta)}$$

Note that if equation 9.5 holds both sides of equation 9.7 are positive and the right side must be larger since $1 > \delta > \delta^{k^{gp}}$. Thus, 9.5 is necessary and sufficient for the in-group policing strategies to be a SPNE if $\alpha > 1 + \beta$.

Now consider the case $\alpha < \beta + 1$ where we require equations 9.4, 9.5,and 9.6 to hold, but now 9.6 implies 9.5. Also, if equation 9.6 holds, note that

$$\delta^{k^{gp}} \geq \frac{\beta}{(1-p)(1+\beta)} > \frac{\alpha - 1}{(1-p)(1+\beta)} > \frac{\alpha - 1}{(1-p)(1+\beta)} - \frac{\left( \delta - \delta^{k^{gp}} \right)}{1 - \delta} \frac{\alpha}{(1+\beta)}$$

Therefore, equation 9.6 implies equation 9.4 so that 9.6 is necessary and sufficient for the in-group punishments to constitute a SPNE. We have established the following proposition.

PROPOSITION 9.4. *The in-group punishment strategy with $k^{gp}$ period punishments is a SPNE if and only if $\delta^{k^{gp}} \geq \min \left\{ \frac{\alpha-1}{(1-p)(1+\beta)}, \frac{\beta}{(1-p)(1+\beta)} \right\}$*

Note some important features of the SPNE. First, if $\delta^{k^{gp}} \geq \min \left\{ \frac{\alpha-1}{(1-p)(1+\beta)}, \frac{\beta}{(1-p)(1+\beta)} \right\}$ holds for $k^{gp} > 1$, it must hold for $k^{gp} = 1$. Thus, no more than a single period of punishment is required to sustain the equilibrium. In fact, longer punishments are counterproductive since they lower the incentives of defectors to cooperate in order to end the punishments.

A second important point about the SPNE is that they can only be sustained if $p$, the probability of out-group interactions is low enough. Since the SPNE requires $\min \left\{ \frac{\alpha-1}{(1-p)(1+\beta)}, \frac{\beta}{(1-p)(1+\beta)} \right\} \leq 1$, it can never exist if $p > \min \left\{ \frac{1}{1+\beta}, 1 - \frac{\alpha-1}{1+\beta} \right\}$. When the probability of interaction with the out-group is large, the probability of punishment for any deviation is low since punishments are meted out only from in-group players. This has the somewhat counterintuitive implication that inter-ethnic cooperation is impeded by too much inter-ethnic interaction. Fearon and Laitin argue that this result provides an endogenous rationale for groups wanting to preserve ethnic boundaries.

## 7. Application: Trade Wars

Consider a generalization of the free trade game as presented in Table 9.6.

| Table 9.6: Generalized Free Trade Game | | |
|---|---|---|
| $1\backslash 2$ | Free Trade | Protection |
| Free Trade | $\Theta, \Theta$ | $0, \Theta + \rho$ |
| Protection | $\Theta + \rho, 0$ | $\rho, \rho$ |

We now interpret $\Theta$ as the value to each country of the other countries open markets and $\rho$ as each countries gain from protecting its own markets. From before, we know that Free Trade can be supported by the grim trigger strategies if an only if $\frac{\Theta}{1-\delta} \geq \Theta + \rho + \frac{\delta\rho}{1-\delta}$ or

$$\delta \geq \frac{\rho}{\Theta}$$

Supporting this equilibrium depends crucially on each country perfectly observing the policies of the other countries. This may not be a very realistic assumption since countries may use invisible trade barriers. Also since trade flows will vary with a number of market conditions unrelated to trade policy, each country will not know for certain whether the fall in trade is due to malfeasance by the other side.

To model these issues, we assume that each country cannot directly observe the policies of the other country, but observes only the value of its trade $\Theta$ which is random variable. To keep things as simple as possible, we let $\Theta_i = \theta > 0$ with probability $\pi$ when country $j$ engages in free trade and 0 with probability $1 - \pi$. When country $j$ protects its markets, $\Theta_i = 0$ with probability 1. Consequently, country $i$ knows for sure if $j$ chooses free trade if $\Theta_i = \theta$ but if uncertain of $j$'s policies when $\Theta_i = 0$. We will assume that $\pi\theta > \rho$ so that each county prefers the free trade outcome to mutual protectionism in expectation.[5]

Clearly, just as before, protection is a SPNE to this game, but we would like to see if there are equilibria which will sustain some level of free trade. Obviously, such an equilibrium will require some form of punishment when $\Theta_i = 0$ is observed even though it cannot be verified with certainty that country $j$ actually defected.

First, consider a grim trigger strategy in which country $i$ protects forever whenever it observes $\Theta_i = 0$ . Thus, the payoffs to free trade in the first period are $\pi\theta + (1 - \pi)0$. Free trade continues to the next

---

[5]This model is based loosely on Green and Porter's (1984) model of imperfect collusion and price wars in economic cartels.

period so long as $\Theta_1 = \Theta_2 = \theta$ which occurs with probability $\pi^2$. Thus, country $i$'s payoffs in the second period are $\pi^2(\pi\theta + (1 - \pi)0) + (1 - \pi^2)\rho = \pi^3\theta + (1 - \pi^2)\rho$. Continuing the same logic to period 3, we get $\pi^5\theta + (1 - \pi^4)\rho$. Thus, the infinite discounted sum of utilities from free trade are

$$V^{FT} = \pi\theta + \delta\left(\pi^3\theta + (1 - \pi^2)\rho\right) + \delta^2\left(\pi^5\theta + (1 - \pi^4)\rho\right) + ...$$

$$V^{FT} = \pi\theta(1 + \delta\pi^2 + \delta^2\pi^4 + ...) + \delta(1 - \pi^2)\rho + \delta^2(1 - \pi^4)\rho + ...$$

$$V^{FT} = \frac{\pi\theta - \delta\pi^2\rho}{1 - \delta\pi^2} + \frac{\delta\rho}{1 - \delta}$$

The utility from defecting to protection is more straightforward. The one period payoff is $\pi\theta + \rho$ while the future payoff is $\frac{\delta\rho}{1-\delta}$ so that $V^P = \pi\theta + \frac{\rho}{1-\delta}$ Thus, country $i$ will choose free trade if and only if $V^{FT} \geq V^P$ or

$$\delta > \frac{\rho}{\pi^3\theta}$$

For comparison, note that if policies were observable, a SPNE in grim trigger strategies would exist so long as $\delta > \frac{\rho}{\pi\theta}$. Thus, the grim trigger strategy is significantly more difficult to sustain when policies are unobservable. In fact, if $\rho > \pi^3\theta$, grim trigger strategies would not constitute a SPNE for any value of $\delta$. The grim trigger strategy is also very costly in the sense that infinite punishments can be generated by variation in $\Theta$, independent of policy.

So now we will follow Green and Porter (1984) and consider finite trigger strategies. Now if either country observes $\Theta = 0$, a trade war begins in which both countries protect their markets for $k \geq 1$ periods. We know check conditions under which free trade is the optimal policy if there is no trade war going on. [6] Let $V_i^k$ be the value of the payoffs for country $i$ for a $k$ period trade war. It is easy to see that

$$V_i^k = \frac{(1 - \delta^k)\rho}{1 - \delta}$$

Let $V^{FT}$ be the payoff beginning a period in which the countries are in a free trade phase. Therefore,

$$V^{FT} = \pi\left(\theta + \pi\delta V^{FT}\right) + \left(1 - \pi^2\right)\delta\left(V_i^k + \delta^k V^{FT}\right)$$

Thus, the equilibrium payoff to free trading is

$$V^{FT} = \frac{\pi\theta + \left(1 - \pi^2\right)\delta V_i^k}{1 - \pi^2\delta - (1 - \pi^2)\delta^{k+1}}$$

---

[6]Since mutual protection is a Nash equilibrium, we do not need to check the optimality of protecting during a trade war.

We can compute the value of a deviation as $V^P$

$$V^P = \pi\theta + \rho + \delta \left( V_i^k + \delta^k V^{FT} \right)$$

An equilibrium requires that $V^{FT} \geq V^P$ or

(9.8)
$$\frac{\delta - \delta^{k+1}}{1 - \delta^{k+1}} > \frac{\rho}{\pi^3\theta}$$

The left side of this expression is increasing in $k$ so let $k^{\min}(\delta)$ be the smaller integer such that the inequality holds for a given $\delta$. Thus, we have established the following proposition.

PROPOSITION 9.5. *If $\delta > \frac{\rho}{\pi^3\theta}$ and $k \geq k^{\min}(\delta)$, the following strategies is a SPNE.*

   (1) *Begin the game, free trading.*
   (2) *Free trade until $\Theta_i = 0$ for either country.*
   (3) *Following a period in which $\Theta_i = 0$, protect for $k$ periods.*
   (4) *After $k$ periods, return to free trade.*

Note that a SPNE can be supported with trade wars of any length greater that $k^{\min}(\delta)$. However, if we assume that the countries can coordinate on the optimal duration of trade wars, the model provides a theory of their duration. Intuitively, since trade wars are costly, the countries should coordinate on the minimal length war sustaining cooperation, $k^{\min}(\delta)$. This intuition can be verify by checking that $V^{FT}$ is strictly decreasing in $k$ so long as $\pi\theta > \rho$. Thus, we can derive empirical predictions by examining equation 9.8. Recall that the left side is increasing in $k$, this implies that $k^{\min}(\delta)$ is increasing in the value of protectionism and decreasing in the value of free trade. Clearly this makes sense. Trade wars should be longer when the incentive problems are more severe. We can also see that the duration of trade conflict is decreasing $\pi$. This is a very counter-intuitive result. Suppose that we interpreted $1 - \pi$ as the volatility of trade flows (the probability of low trade during a free trade regime). This interpretation suggests that trade volatility increases the duration of trade wars. This is necessary to keep countries from enacting barriers and blaming the results on natural volatility. However, in equilibrium, the countries never protect outside trade wars so that they know with certainty that $\Theta_i = 0$ was caused by natural volatility. Yet they must engage in costly, length trade wars to ensure that barriers are not erected.

## 8. Exercises

EXERCISE 9.1. *Assume that there are three groups with the following preferences over three policies.*

|   $A$   |   $B$   |   $C$   |
| :---: | :---: | :---: |
|   $x$   |   $z$   |   $y$   |
|   $y$   |   $x$   |   $z$   |
|   $z$   |   $y$   |   $x$   |

We will be analyzing this as a repeated game where in each period groups may make a counter proposal to the status quo. For example, suppose that the status quo is $x$ then both $B$ and $C$ wish to propose $z$ which passes an becomes the new status quo. In the next period $A$ and $C$ wish to propose $y$ which passes. Assume that each group discounts the future by $\delta$.

    a. What is each group's utility from the infinite cycle (starting at $x$) of policies that results?

    b. Now suppose that $A$ and $B$ decide to form a political party to implement policy $x$ forever. Find a critical value $\delta^*$ such that if $\delta > \delta^*$ there is a subgame perfect Nash equilibrium where $A$ and $B$ vote for $x$ in every period, $C$ proposes $z$, and if $A$ or $B$ ever defects a policy cycle starts.

EXERCISE 9.2. *Prove Proposition 1.*

EXERCISE 9.3. *Prove Proposition 3.*

EXERCISE 9.4. *Find conditions for the existence of the Spiral SPNE to the Fearon and Laitin's model. This equilibrium is based on the following strategies:*

> In in-group pairings, always play $C$ with cooperator, and always play $D$ against a defector, regardless of ones status. A player enters or restarts the in-group punishment phase for $k^{in}$ periods by defecting against a cooperator. In out-group pairings, play $C$ if neither group is in an out-group punishments phase. Otherwise play, $D$. A group enters the out-group punishment phase for $k^{out}$ periods if any member defects in a cross-group pairing when neither group is in the out-group punishment phase.

EXERCISE 9.5. *Consider the model of trade wars. Construct the following "probabilistic grim trigger SPNE." Instead of reverting to protectionism forever the first time $\Theta_i = 0$ is observed, assume that country i plays a mixed strategy and protects forever with probability $\mu$.*

CHAPTER 10

# Bargaining Theory

If political science is the study of "who gets what, what, when and how" then bargaining theory lies at its foundation.[1]  Legislators and executives bargain over new legislation.  States bargain to reach new international agreements and to settle crises.  Political parties bargain over coalition governments. And so on.

Not surprisingly given its importance, the application of game theoretic models of bargaining to study political processes has been a very active area of research.  These models have focused on two sets of issues.  The first are the questions of distribution – "who wins" and "who loses." Does the president get his preferred legislation?  Which country gets to control the disputed region?  Which parties received government portfolios?  The second important question concerns the efficiency of political bargaining.  Does the bargaining process itself consume resources or fail to reach outcomes that make everyone better off?  Does legislative bargaining end in gridlock or a veto even though there are policy compromises that all prefer?  Do international disputes end in costly militarized conflicts and wars?  Why does it take so long to form new coalition governments?

In this chapter, we review some of the most important bargaining models and their application to political science.

## 1. The Nash Bargaining Solution

One of the earliest attempts to model bargaining was the framework developed by John Nash.  His approach was axiomatic in that he stipulated a number of features that should characterize the outcome of any bargaining situation.  Before discussing his axiomatic requirements, we describe his "solution" to the bargaining problem   Our discussion closely mirrors that of Muthoo (1999).

Suppose that two players $A$ and $B$ are negotiating over the allocation of $X$ units of some resource.  We assume that $X$ is infinitely divisible so that the feasible allocations are all $x_A$ and $x_B$ such that $x_A + x_B \leq X$.  Each player receives utility based on their allocations,

---

[1]See Lasswell (1936).

$U_A(x_A)$ and $U_B(x_B)$. We assume that $U_i$ is strictly increasing and concave for both players $i = A, B$. In the event that no agreement is reached, each player receives a default utility, *disagreement value* or *outside option* of $\underline{u}_i > u_i(0)$. Finally to ensure that the bargaining problem is non-trivial, we assume that there exists at least one allocation $(x_A, x_B)$ such that $U_i(x_i) > \underline{u}_i$ and $x_A + x_B \leq X$. This ensures that there is feasible allocation that both players prefer to their disagreement values.

In analyzing Nash's solution to this problem, it is useful convert it into one of allocations of utilities $(u_A, u_B)$ rather than one of allocations of $X$. Therefore, we define the feasible utility allocations as the set $\Omega = \{(u_A, u_B) : u_A(x_A) = u_A, u_A(x_B) = u_B, \text{ and } x_A + x_B \leq X\}$. Given our assumptions about the utility functions, the boundary of this feasible set can be represented as a locus of points such as the one in Figure 10.1. We define this locus as the function $g(u_A) = U_B(X - U_A^{-1}(u_A))$. Muthoo (1999) provides a proof that $g$ is both decreasing and concave in $u_A$. To simplify our exposition, we assume that it is twice-differentiable.

Now we can state Nash's solution to the bargaining problem. Based on the axioms we discuss below, his solution is the utility allocation $(u_A, u_B) \in \Omega$ that maximizes

$$(u_A - \underline{u}_A)(u_B - \underline{u}_B)$$
$$\textit{subject to } u_A \geq \underline{u}_A \text{ and } u_B \geq \underline{u}_B$$

The requirement that $u_A \geq \underline{u}_A$ and $u_B \geq \underline{u}_B$ is illustrated by the dotted lines in Figure 10.1. Thus, the constraint set is $g(u_A)$ on the range $[\underline{u}_A, g^{-1}(\underline{u}_B)]$. Since $g$ is concave and decreasing, the feasible set is convex. It is easy to see that the *Nash product* $(u_A - \underline{u}_A)(u_B - \underline{u}_B)$ is quasi-concave in both $u_A$ and $u_B$. Thus, we can represent its level curves in the region that the product in positive as the heavy dotted lines in Figure 10.1.

### Insert Figure 10.1 Here

Thus, there is a unique Nash bargaining solution at the tangency of $g$ and the iso-product curves. Mathematically, the solution to the constrained optimization problem is given by

$$-g'(u_A) = \frac{u_B - \underline{u}_B}{u_A - \underline{u}_A}$$
$$u_B = g(u_A)$$

Before moving to a general results about the Nash bargaining solution, it is useful to consider some special cases. First, assume that $X = 1$

and $u_i(x_i) = x_i$. The Nash Bargaining solution for this model is

$$u_A = x_A = \frac{1 + \underline{u}_A - \underline{u}_B}{2} \text{ and } u_B = x_B = \frac{1 - \underline{u}_A + \underline{u}_B}{2}$$

It is easy to see two important features of the solution. First, each player does better when disagreement provides it with a higher utility and worse when their opponent has a better outside option. Second, if each player has an equally valuable outside option, the resources are split evenly. Another way to interpret Nash's solution is to note that the bargainers insist upon their disagreement values and equally split the surplus $1 - \underline{u}_A - \underline{u}_B$ which gives each a utility of $\underline{u}_i + \frac{1 - \underline{u}_A - \underline{u}_B}{2}$.

Now we turn to the general case. First, we can state the Nash bargaining solution in terms of shares.

PROPOSITION 10.1. *The Nash bargaining shares are given by the solution to*

$$\frac{U_A(x_A) - \underline{u}_A}{U_A'(x_A)} = \frac{U_B(X - x_A) - \underline{u}_B}{U_B'(X - x_A)}$$

PROOF. Direct application of previous result using the fact that $g'(u_A) = -U_A^{-1'}(X - U_A^{-1}(u_A)) \cdot U_B'(X - U_A^{-1}(u_A)) = -\frac{U_B'(X - U_A^{-1}(u_A))}{U_A'(X - U_A^{-1}(u_A))}$ and $U_A^{-1}(u_A) = x_A$. $\square$

A direct implication of this result is that if the disagreement values and utility functions are the same for both players, the Nash bargaining shares are $x_A = x_B = \frac{1}{2}X$. Finally, we show that given our assumptions about $g$ payoffs increase in one's own disagreement value and decline in the opponent's.

PROPOSITION 10.2. *Assume that $g$ is twice-differentiable, then let* $\frac{\partial u_i}{\partial \underline{u}_i} > 0$ *and* $\frac{\partial u_j}{\partial \underline{u}_i} < 0$ *for* $i \neq j$.

PROOF. Since $g$ is twice differentiable, we can use implicit differentiation to the solution $\frac{g(u_A) - \underline{u}_B}{u_A - \underline{u}_A} + g'(u_A) = 0$. Since the second order condition is satisfied i.e. $\frac{g'(u_A)(u_A - \underline{u}_B) - (u_A - \underline{u}_B)}{(u_A - \underline{u}_A)^2} + g''(u_A) < 0$, the result follows from $\frac{-1}{u_A - \underline{u}_A} < 0$, $\frac{g(u_A) - \underline{u}_B}{(u_A - \underline{u}_A)^2} > 0$, and $g'(u_A) < 0$. $\square$

**1.1. Application: Risk Aversion and the Nash Bargaining Solution.** Intuitively, risk is an important component of bargaining. Bargainers always have to contend with the possibility that an agreement will not be reached and they will be left with their outside options. Also we should expect that if a player makes a more aggressive demand, she increases the probability that the negotiations will collapse. Consequently, it seems natural to think that bargainers who

are more willing to tolerate risk should do better because they are willing to make tougher demands and more aggressive reject offers. While the Nash bargaining model, "black boxes" the negotiation process, the solution is consistent with this intuition.

To see this, assume that each player has a utility function given by $U_i(x_i) = x_i^{\alpha_i}$ where $0 < \alpha_i < 1$, disagreement values $\underline{u}_i = 0$, and $X = 1$. The different values of $\alpha$ capture the players risk aversion, the lower $\alpha$ the greater the risk aversion.[2] It is easiest to compute the equilibrium shares using the formula from Proposition 10.1. The solution is

$$x_A = \frac{\alpha_A}{\alpha_A + \alpha_B} \text{ and } x_B = \frac{\alpha_B}{\alpha_A + \alpha_B}$$

These results imply that each bargainer's share is decreasing in their own risk aversion and increasing in the risk aversion of their opponent. This effect is consistent with our intuition that bargainers who are risk-acceptant enough to take tough positions (i.e. increase the likelihood of disagreement) should receive larger allocations.

**1.2. Nash's Axioms.** In this section we outline the axioms that underlie Nash's bargaining solution. Informally, the axioms are intended to encapsulate the following principals.

(1) The bargainers are expected utility maximizers.
(2) Bargaining should be efficient. The players should fully allocate all of the available resources and no player should do worse that their disagreement value.
(3) The allocation should depend only on the player's preferences and disagreement values.
(4) The bargaining solution should no be effected by eliminating from consideration allocations other than the solution.

To formalize this axioms, recall that $\Omega$ is the set of feasible utility applications $(u_A, u_B)$ that can be reached through some allocation of $X$. We now define the set of Pareto optimal allocations as $\Omega^e = \{\omega \in \Omega : u_A \geq \underline{u}_A \text{ and } g(u_A) \geq \underline{u}_B\}$. We can define a generic bargaining situation as a pair $(\Omega, \underline{u})$ where $\underline{u}$ is the vector of disagreement values. We denote the set of all bargaining games as $\Sigma$ and the bargaining solution as a correspondence $F : \Sigma \rightrightarrows R^2$. We let $F_i$ denote the utility allocated to agent $i$.

The following axioms form the basis of Nash's solution.

---

[2]A standard measure of risk aversion is $-\frac{u''}{u'}$. For these utility functions, $-\frac{u''}{u'} = -\frac{\alpha(\alpha-1)x^{a-2}}{\alpha x^{\alpha-1}} = \frac{(1-\alpha)}{x}$.

AXIOM 10.1. *Invariance to equivalent utility representations: Let $U_i' = \alpha_i U_i + \beta_i$ and $\underline{u}_i' = \alpha_i \underline{u}_i + \beta_i$ for $\alpha_i > 0$ and define $\Omega'$ accordingly. Then $F_i(\Omega', \underline{u}') = \alpha_i F_i(\Omega, \underline{u}) + \beta_i$ for $i = A, B$.*

Affine transformations of utility functions and disagreement utilities should not alter the bargaining outcomes. Since the utility allocations are adjusted by the same transformations as the utility functions, it is easy to show that the resource allocations $x_i = F_i^{-1}(\Omega, \underline{u})$ and $x_i' = F_i^{-1}(\Omega', \underline{u}')$ for $i = A, B$. As we know saw in chapter 3, this axiom implies that the players are expected utility maximizers.

AXIOM 10.2. *Pareto efficiency: If $F(\Sigma) = (u_A, u_B)$, then there are no other allocations $(u_A', u_B') \in \Omega$ such that $u_i' > u_i$ for some $i$, $u_j' \geq u_j$ for $j \neq i$, and $u_i' \geq \underline{u}_i$ for all $i$.*

The Pareto axiom holds that the bargainers should not be able to improve upon the bargaining solution by choosing an allocation that makes one of the bargainers better off without reducing the utility of the other.

AXIOM 10.3. *Symmetry: Let $\underline{u}_A = \underline{u}_B$ and assume that $(u_1, u_2) \in \Omega$ if and only if $(u_2, u_1) \in \Omega$. Then $F_A(\Omega, \underline{u}) = F_B(\Omega, \underline{u})$.*

The basic idea of this axiom is that if neither player is advantaged by having a better disagreement outcome or a utility level unreachable by her opponent, then the bargainers should get the same utility allocations.

AXIOM 10.4. *Independence of Irrelevant Alternatives. Consider two bargaining situations $(\Omega, \underline{u})$ and $(\Omega', \underline{u})$ such that $\Omega' \subset \Omega$ and $F(\Omega, \underline{u}) \subset \Omega'$. Then $F(\Omega, \underline{u}) = F(\Omega', \underline{u})$*

The intuition behind the IIA axiom is that, holding the disagreement points constant, a smaller feasible set of allocations should only change the bargaining solution if it makes the original allocation infeasible.

From our analysis of the Nash bargaining solution in the previous section, it is clear that it satisfies all of these axioms. However, the next proposition establishes that it is the only solution which satisfies all four axioms.

PROPOSITION 10.3. *A bargaining solution $F : \Sigma \rightrightarrows R^2$ satisfies axioms 1-4 if and only if it is the Nash bargaining solution.*

PROOF. see Muthoo (1999). □

## 2. Non-cooperative Bargaining

While it does make a number of reasonable empirical predictions, the Nash bargaining solution is best interpreted as a normative argument about what bargaining outcomes should look like rather than a positive theory about how actual bargaining will take place. In this section, we turn to non-cooperative game theoretic models which deduce behavior of bargainers under different extensive forms.

The starting point for the application of non-cooperative game theory to bargaining is the model of Rubinstein (1982). Suppose that two players are trying to decide how to divide \$1. The players will take turns making offers so that player 1 proposes in periods $0, 2, 4$, etc. and player 2 makes proposals in the other periods. The game continues (possibly infinitely) until a proposal is accepted by the other player.

In each period that she is the proposer, player 1 can make an offer $(x_1, x_2)$ where $x_1$ is player 1's share and $x_2$ is player 2's share where $x_1 + x_2 \leq 1$. If player 2 accepts, the game ends and the dollar is divided accordingly. If player 2 rejects, then she gets to make an offer $(x_1, x_2)$ and the game continues if player 1 rejects. To simply matters, we assume that both players have linear utility functions $u_1 = x_1$ and $u_2 = x_2$. Each player has a discount factor $\delta_i$ so that players value a proposal of $(x_1, x_2)$ $t$ periods in the future as $(\delta_1^t x_1, \delta_2^t x_2)$.

Just as in the bargaining game we encountered in chapter 7, there are lots of Nash equilibria to this game. For example, consider the strategies "Player 1 demands $x_1 = 1$ and refuses all other offers, while player 2 always offers $x_1 = 1$ and accepts any offer". However, this equilibrium is not subgame perfect. If player 2 rejected player 1's first offer, and offered $x_1 > \delta_1$ player 1 should accept because the best it can get is the whole dollar next period. So we will focus on subgame perfect Nash equilibrium.

**2.1. Subgame Perfect Equilibria.** Rubinstein shows that there is a unique SPNE to this game based on playing the following strategies in every period:

Player 1 proposes $\left( \frac{1-\delta_2}{1-\delta_1\delta_2}, \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2} \right)$ and accept player 2's offer if and only if $x_1 \geq \frac{\delta_1(1-\delta_2)}{1-\delta_1\delta_2}$.

Player 2 proposes $\left( \frac{\delta_1(1-\delta_2)}{1-\delta_1\delta_2}, \frac{1-\delta_1}{1-\delta_1\delta_2} \right)$ and accept player 1's offer if and only if $x_2 \geq \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2}$.

We begin by simply verifying that these strategies are in fact a SPNE. First we check whether player 1 has an incentive to defect in any subgame. Consider a subgame beginning with a proposal by player

1(i.e. an even period). In the equilibrium, player 1 proposes the split $(\frac{1-\delta_2}{1-\delta_1\delta_2}, \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2})$ which is accepted by player 2. Clearly, player 1 cannot gain by lowering $x_1$ as it will be accepted but she gets a lower share. If player 1 raises $x_1$, then she must lower $x_2$ so that the proposal is feasible. However, any $x_2 < \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2}$ will be rejected. Following such a rejection, player 2 proposes $x_1 = \frac{\delta_1(1-\delta_2)}{1-\delta_1\delta_2}$ which player 1 accepts. Thus, player 1's utility of this defection is $\frac{\delta_1^2(1-\delta_2)}{1-\delta_1\delta_2}$ which is less than her equilibrium utility of $\frac{1-\delta_2}{1-\delta_1\delta_2}$ since $\delta_1 < 1$.

Now consider whether player 1 will defect when player 2 is the proposer (i.e. an odd period). Player 2 proposes $(\frac{\delta_1(1-\delta_2)}{1-\delta_1\delta_2}, \frac{1-\delta_1}{1-\delta_1\delta_2})$. Note that accepting and rejecting the offer lead to the same utility as the best that player one can do is have $x = \frac{1-\delta_2}{1-\delta_1\delta_2}$ accepted one period later.

Showing that player 2 will not defect is entirely similar.

**2.2. Computing the Equilibrium.** The problem with the preceding proof is that it does not give much of a sense of how the result is derived. Now we consider a more constructive proof. Let $v_1$ and $v_2$ be the utilities of player 1 and 2 for subgames in which they are the proposer. For example, if player makes a proposal $x_1$ that is accepted $v_1 = x_1$. If player 1's proposal is rejected, $v_1$ is the discounted values of the maximum of what player 2 offers and what it gets by rejecting and proposing in his next turn. Given that the postulated strategies are the same in every period, these values are independent of $t$. We will call these *continuation values* since they also reflect the utility of rejecting a proposal and moving to the next subgame. Consider a subgame where player 1 is the proposer. She must offer player 2 at least $\delta_2 v_2$. Thus, $x_1 = 1 - \delta_2 v_2$. Since this offer is accepted $v_1 = x_1$ so that $v_1 = 1 - \delta_2 v_2$. Consider a subgame where player 2 is the proposer. She must offer at least $\delta_1 v_1$ so that $v_2 = 1 - \delta_1 v_1$.

Solving these two equations leads to

$$v_1 = \frac{1-\delta_2}{1-\delta_1\delta_2}$$

$$v_2 = \frac{1-\delta_1}{1-\delta_1\delta_2}$$

These continuation values are consistent with the strategies presented in the last section.

**2.3. Uniqueness.** While we have shown that Rubinstein's equilibrium is a subgame perfect Nash equilibrium, we have not ruled out

the possibility that there are others. We now show that this equilibrium is the unique SPNE by proving that $v_1$ and $v_2$ above are the only continuation values consistent with a SPNE. Suppose there are more than one SPNE. Let $\overline{v}_i$ and $\underline{v}_i$ be player $i$'s highest and lowest SPNE continuation values for any subgame where player $i$ is the proposer. Let $\overline{w}_i$ and $\underline{w}_i$ be player $i$'s highest and lowest SPNE continuation values for any subgame where player $i$ is not the proposer.

When player 1 makes a proposal, she never need to offer more than $\delta_2 \overline{v}_2$ since player 2 cannot expect more than $\overline{v}_2$ by rejecting and making her own proposal in the next round. Thus, her lowest possible continuation value must satisfy $\underline{v}_1 \geq 1 - \delta_2 \overline{v}_2$. By the symmetric argument, $\underline{v}_2 \geq 1 - \delta_1 \overline{v}_1$. Since we now know that the other player will never offer anything greater than $\delta_i \overline{v}_i$ then we know that $\overline{w}_i \leq \delta_i \overline{v}_i$.

Now consider player 1's strategy. When she proposes the best she can do is either to pay $\delta_2 \underline{v}_2$ or trigger a rejection to get $\delta_1 \overline{w}_1$. Thus, we know that her continuation value satisfies $\overline{v}_1 \leq \max\left\{1 - \delta_2 \underline{v}_2, \delta_1 \overline{w}_1\right\} \leq \max\left\{1 - \delta_2 \underline{v}_2, \delta_1^2 \overline{v}_1\right\} = 1 - \delta_2 \underline{v}_2$. Similarly, $\overline{v}_2 \leq 1 - \delta_1 \underline{v}_1$. Thus, we have the following four inequalities

$$\underline{v}_1 \geq 1 - \delta_2 \overline{v}_2$$
$$\underline{v}_2 \geq 1 - \delta_1 \overline{v}_1$$
$$\overline{v}_2 \leq 1 - \delta_1 \underline{v}_1$$
$$\overline{v}_1 \leq 1 - \delta_2 \underline{v}_2$$

Using the first and third inequalities, we see that $\underline{v}_1 \geq 1 - \delta_2(1 - \delta_1 \underline{v}_1)$ which implies that $\underline{v}_1 \geq \frac{1-\delta_2}{1-\delta_1\delta_2}$. Similarly, using the second and fourth, we get $\overline{v}_1 \leq 1 - \delta_2(1 - \delta_1 \overline{v}_1)$ or $\overline{v}_1 \leq \frac{1-\delta_2}{1-\delta_1\delta_2}$. This implies that $\overline{v}_1 = \underline{v}_1 = \frac{1-\delta_2}{1-\delta_1\delta_2}$. Similarly, we can derive that $\overline{v}_2 = \underline{v}_2 = \frac{1-\delta_1}{1-\delta_1\delta_2}$. Thus, there is a single continuation value for each player. Thus, the postulated strategies are the only SPNE.

**2.4. Implications.** The model suggests a very simple path of play. In period zero, player 1 proposes $(\frac{1-\delta_2}{1-\delta_1\delta_2}, \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2})$, player 2 accepts, and the game ends. Since the whole dollar is allocated and there is no delay, the subgame perfect Nash equilibrium is efficient. It is easy to see that the SPNE has the following implications.

(1) If both players have the same discount factor, there is a first mover advantage since $\frac{1-\delta}{1-\delta^2} > \frac{\delta(1-\delta)}{1-\delta^2}$. Intuitively, since player 2 discounts the future, player 1 only need offer her a fraction of what she would get for being the proposer next period. Since both players are identical, this means that player 2 is getting only a fraction of what player 1 gets.

(2) Both players shares are increasing in their discount factors and declining in their opponent's. It pays to be patient. When player 2's discount factor is high, player one has to offer her more to secure immediate agreement. Conversely, when player 1's discount factor is high, player 2 will have to offer him more to reach agreement in the event that player 2 gets to make an offer. Thus, rejecting player 1s offer is less valuable for player 2 suggesting that player 1 gets to keep more in the first period.

(3) Let $\delta_1 = \delta_2 = \delta$, the both players shares converge to $\frac{1}{2}$ as $\delta$ converges to 1. As both players become perfectly patient, they are less willing to accept offers that are less than what they can get as the proposer next period. In the limit, they demand exactly what they expect to get next period which is satisfied by the proposal $\left(\frac{1}{2}, \frac{1}{2}\right)$. One way to think about the discount rates converging to one is to consider a situation in which offers and counter- offers can be made very quickly so that rejecting an offer creates only infinitessimal delay. In such a case, the equilibrium is equal division and corresponds exactly to the Nash bargaining solution for this problem.

**2.5. Asymmetric Disagreement Values.** In the canonical Rubinstein game, the players get 0 in any period for which there is no agreement. We now modify the game in two ways. First, we assume that the players receive an allocation of $(d_1, d_2)$ in each period prior to an agreement where $d_1 + d_2 < 1$. After an agreement, $(x_1^*, x_2^*)$ is reached, the bargainers receive this allocation in every period over an infinite horizon. This contrasts with the model of the last section where the allocation is "consumed immediately."[3] To keep things simple, we assume that $\delta_1 = \delta_2 = \delta$. Thus, the utilities of reaching agreement $(x_1^*, x_2^*)$ in period $t$ are $\left(\frac{\left(1-\delta^{t-1}\right)d_1+\delta^t x_1^*}{1-\delta}, \frac{\left(1-\delta^{t-1}\right)d_2+\delta^t x_2^*}{1-\delta}\right)$.[4]

Let $v_i$ be the continuation values for periods in which $i$ proposes. If an agreement $(x_1^*, x_2^*)$ is reached in such a period, $v_i = \frac{x_i^*}{1-\delta}$. Consider player 2's decision to accept or reject an offer of $x_2$. If she accepts, she

---

[3]This modification rules out strategies where the bargainers delay infinitely in the hopes that the discounted sum of $d_i$ exceeds the one-period agreement. We can easily adjust the original model to correspond to the assumption that the agreement is over a flow of utilities rather than one-shot consumption. We would simply use the original model and assume that the players were allocating $\frac{1}{1-\delta}$.

[4]We assume that any agreement results in the same allocation in each period. However, since the players are risk neutral, there might be agreements to random allocations that generate the same payoffs.

gets a value of $\frac{x_2}{1-\delta}$ whereas if she rejects she gets $d_2$ in the current period and a continuation value $v_2$ in the next. Thus, she will accept so long as $x_2 > (1 - \delta)(d_2 + \delta v_2)$. Now consider player 1's choice. If he makes the minimal acceptable offer $x_2 = (1 - \delta)(d_2 + \delta v_2)$, his continuation value is $v_1 = \frac{1-(1-\delta)(d_2+\delta v_2)}{1-\delta} = \frac{1}{1-\delta} - d_2 - \delta v_2$. Similarly, assuming that player 2 wished to secure an agreement with her proposals we require $v_2 = \frac{1}{1-\delta} - d_1 - \delta v_1$. The solution to these two equations is given by

$$v_1 = \frac{1 - d_2 + \delta d_1}{1 - \delta^2} = \frac{d_1}{1 - \delta} + \frac{1 - d_1 - d_2}{1 - \delta^2}$$

$$v_2 = \frac{1 - d_1 + \delta d_2}{1 - \delta^2} = \frac{d_2}{1 - \delta} + \frac{1 - d_1 - d_2}{1 - \delta^2}$$

To show that these are in fact equilibrium continuation values, we must show that each player prefers to make their equilibrium proposal rather than defect and get their disagreement value for an additional period. Thus, we require $v_1 > d_1(1 + \delta) + \delta^2 v_1$ or $v_1 > \frac{d_1}{1-\delta}$ which is easily verified. Similarly, our equilibrium requires that $v_2 > \frac{d_2}{1-\delta}$ which is also satisfied. The techniques of Section 2.3 can easily be generalized to show that this is the unique SPNE.

This equilibrium has a number of qualitative similarities to the Nash bargaining solution. Note that each player's continuation value increases in her disagreement value and decreases in their opponent's. This equilibrium also has a surplus-splitting interpretation. Note that each player's continuation value has two components. The first is $\frac{d_i}{1-\delta}$ which is the utility each player can guarantee herself in the absence of any agreement. The second component $\frac{1-d_1-d_2}{1-\delta^2}$ corresponds to the equilibrium continuation value in a game to split $1 - d_1 - d_2$ when the players have outside option values of 0. Thus, a useful interpretation of this equilibrium is that both players take what they are entitled and bargain over the rest.

## 3. Majority Rule Bargaining Under Closed Rule

A key feature of the Rubinstein model is that unanimous consent is required to reach an agreement on the allocation. This rules out a number of important political settings where only a simple or super-majority is required for agreement. Baron and Ferejohn (1989) have extended Rubinstein's model to simple majority rule.

Suppose that there are $N$ (odd) players bargaining and any proposal requires $n = (N + 1)/2$ votes. Instead of assuming alternating offers, Baron and Ferejohn consider a bargaining protocol with a *Random Recognition Rule*. According to this protocol, in each period, every

player has an equal probability $(1/N)$ of being chosen to be the proposer. In this section, we focus on bargaining under *closed rule* where the proposer makes a take-it-leave-it offer for the current legislative session. The proposer in each period makes an offer $(x_1, x_2, \ldots, x_N)$ such that $x_i$ is the share for player $i$ and we require $\sum x_i \leq 1$. If this proposal is rejected, the session ends, discounting occurs, and a new proposer is chosen at the beginning of the next session. In a later section, we consider open rule bargaining where proposals can be amended within the current session. To keep things simple, we assume that each player has the same discount factor $\delta$.

This game has lots of subgame perfect equilibria. In fact for large $N$ and $\delta$, there is a SPNE that can support any division of the dollar. This is due to the fact that if the players are patient enough, they can design punishment strategies to guarantee \$0 to any defector. However, these strategies require that each player know the whole history (possibly infinite) of the game so as to know which actions are consistent with the prescribed punishment. Thus, following Baron and Ferejohn, we will analyze only *stationary* equilibria. A stationary equilibrium to this game is one in which:

(1) A proposer proposes the same division every time she is recognized regardless of the history of the game.
(2) Voters vote only on the basis of the current proposal and expectations about future proposals. Because of assumption 1, future proposals will have the same distribution of outcomes in each period.

These two assumptions imply that the game essentially starts over in every period. Therefore, the continuation value of each player is exactly the expected utility of the game. Let $v_i$ be the continuation value for player $i$. We will focus on symmetric equilibria so that $v_i = v$ for all $i$. Finally, we will focus only on equilibria in which voters do not choose weakly dominated strategies in the voting stage. Therefore, a voter will accept any proposal that provides at least as much as the discounted continuation value. Therefore, any voter who gets $x_i \geq \delta v$ will vote in favor of the proposal while any voter who receives less than $\delta v$ will vote against.

Give these voting strategies, a proposer knows that she must propose $\delta v$ to $n - 1$ other players and 0 to the rest. Let $z$ be the amount that the proposer keeps so that

$$z = 1 - (n - 1)\delta v$$

We will assume (although we can show that it must be true), that the proposer chooses her coalition partners randomly. Now we can compute $v$. Since the continuation value is just the expected value of the game starting next period, it is simply $z$ times the probability of being chosen as proposer $\frac{1}{N}$, $\delta v$ times the probability of being included in the winning coalition $\frac{n-1}{N}$, and 0 times the remaining probability. Thus,

$$v = \frac{z}{N} + \frac{n-1}{N}\delta v.$$

Substituting for $z$ we obtain

$$v = \frac{1}{N}.$$

Thus, the continuation value is just a proportional share of the dollar. Since $v$ is also the expected utility of the game, this result implies that bargaining is efficient in the sense that the sum of player utilities is maximized. As the reader will discover in the exercises, this efficiency result may not hold if voters are risk-averse.

Finally, given our solution for $v$, we can compute the proposer's share:

$$z = 1 - \delta\frac{n-1}{N} = 1 - \delta\frac{N-1}{2N}.$$

To ensure that the proposer will prefer to make an acceptable proposal, we must check that $z > \delta v$, otherwise a proposer would prefer punt and wait for the next period. This condition is easily verified.

Among the important implications of the model is its predictions about proposal power. First, note that proposal power is increasing in $N$. When $N$ increases, the proposer has more potential coalition partners to play off of one another. This increases the competition for inclusion in the winning coalition and drives down what the proposer must pay. Secondly, proposal power is decreasing in $\delta$. When $\delta$ is higher, the voters are more willing to vote down proposals and wait for a chance to propose themselves. Thus, the proposer must be relatively more generous to secure agreement.

**3.1. Supermajority Rule.** We can also easily extend the model to capture situations where more that a simple majority is required for passage of the bill. Now assume that $k > n$ votes are required. If is easy to see that the proposer's share is now

$$z = 1 - (k-1)\delta v$$

and the continuation values are now given by

$$v = \frac{z}{N} + \frac{k-1}{N}\delta v$$

Simple algebra reveals that once again $v = \frac{1}{N}$. This is not terribly surprising given that the supermajority rule preserves the same symmetry we found in the majority rule game. However, the proposers equilibrium share is now lowered to $z = 1 - \delta\frac{k-1}{N}$. Thus, the primary consequence of supermajority rules is to mitigate the proposer's advantage.

**3.2. Asymmetric Proposal Power.** A limitation of the preceding model is that it assumes that all legislators have the same probability of being recognized to make the proposal. This assumption would ignore real world legislative institutions such as committees and parties which may affect the probability that an individual legislator gets to make a proposal.

To show how the model generalizes, suppose that the members are divided into two parties $A$ and $B$. Party $A$ has $N - m > n$ members so that it is the majority. Each member of $A$ has a proposal power $p > 1/N$. Alternatively, there are $m$ members of $B$ who have proposal power $q < 1/N$. For consistency, we require that $(N - m)p + mq = 1$.

Again we will assume symmetry so that every legislator with the same recognition probability has plays the same strategies and therefore has the same continuation value. The members of the two parties have continuation values $v_A$ and $v_B$, respectively. We conjecture for now (and prove later) that $v_A > v_B$. Given these continuation values, a member of party $A$ will vote for any proposal that provides her at least $\delta v_A$ and a member of party $B$ will vote for a proposal giving at least $\delta v_B$. Given these strategies and the assumption that $v_A > v_B$, a proposer from party $A$ will give $\delta v_B$ to the $m$ members of party $B$ and $\delta v_A$ to $n - m - 1$ members of party $A$. Thus,

$$z_A = 1 - (n - m - 1)\delta v_A - m\delta v_B$$

A member of $B$ will give positive allocations to $m - 1$ members of $B$ and $n - m$ members of $A$ so that

$$z_B = 1 - (n - m)\delta v_A - (m - 1)\delta v_B$$

Note that $z_A > z_B$. We can now compute $v_A$ and $v_B$

$$v_A = pz_A + p(n - m - 1)\delta v_A + qm(n - m)\delta v_A/(N - m)$$
$$v_B = qz_B + (1 - q)\delta v_B$$

Thus, we have 4 equations with 4 unknowns. It is straightforward to solve, but its messy. So we will consider a simple example. Let $N = 3, m = 1$. Note that $q = 1 - 2p < 1/3$. Therefore, the equilibrium conditions are the following:

$$z_A = 1 - \delta v_B$$
$$z_B = 1 - \delta v_A$$
$$v_A = p z_A + q \delta v_A / 2$$
$$v_B = q z_B + (1 - q) \delta v_B.$$

After some tedious algebra, we find that

$$v_A = \frac{(1 - q)(1 - \delta)}{2 + q\delta - 2\delta}$$
$$v_B = \frac{q(2 - \delta)}{2 + q\delta - 2\delta}.$$

We still need to check our assumption that $v_A \geq v_B$. This occurs when

$$q < \frac{1 - \delta}{3 - 2\delta} \leq \frac{1}{3}.$$

Since this upper bound is always less than $1/3$ when $\delta > 0$, the asymmetry in proposal power must be substantial to give an advantage to party $A$. The reason is that its greater proposal power makes members of $A$ unattractive coalition partners. Thus, the likelihood of being the proposer must be large enough to offset this effect. However, it is easy to show that that $v_A$ is decreasing and $v_B$ is increasing in $q$.

To complete our analysis, we need to consider what happens when $\frac{1-\delta}{3-2\delta} < q \leq \frac{1}{3}$. We can rule out $v_B > v_A$ as this would imply that the member of $B$ is never in a coalition with the proposer. Thus,

$$v_B = q z_B = q(1 - \delta v_A)$$
$$v_A = p z_A + (1 - p)\delta v_A = p(1 - \delta v_A) + (1 - p)\delta v_A.$$

This leads to $v_A = \frac{(1-q)}{2(1-\delta q)}$ and $v_B = \frac{q(2-\delta-\delta q)}{2(1-\delta q)}$. Note that if $v_B > v_A$ only if $q \geq \frac{1+2\delta}{3-2\delta} \geq \frac{1}{3}$ which violates our original assumption about $q$. Thus, the only possible outcome for $\frac{1-\delta}{3-2\delta} < q \leq \frac{1}{3}$ is $v_A = v_B$. To support this equilibrium, proposers from $A$ must choose a mixed strategy that randomizes between formed a coalition with the remaining member of $A$ and the member of $B$. We leave computation of the equilibrium mixed strategy to an exercise.

**3.3. Asymmetric Veto Powers.** Another institutional variation in legislative institutions is that certain players are privileged with the ability to block legislation such as the president, an upper chamber, or a court. In this section, we provide a simple example of how to incorporate vetoes into the Baron-Ferejohn model.[5] Now suppose that one member of our three person legislature has an absolute veto power in that she must approve every proposal. Let party $B$ have the veto player. To keep things simple, we return to the case of equal proposal powers.

Since $B$ has an absolute veto, any proposer must include $B$ in her coalition so that

$$z_A = 1 - \delta v_B$$
$$z_B = 1 - \delta v_A$$

Computing the continuation values:

$$v_A = \frac{1}{3}z_A + \delta\frac{1}{3}v_A$$
$$v_B = \frac{1}{3}z_B + \delta\frac{2}{3}v_B$$

Thus, we can solve for $v_A = \frac{3(1-\delta)}{\delta^2-9\delta+9}$ and $v_B = \frac{3-2\delta}{\delta^2-9\delta+9}$. Note that $v_A < v_B$ so long as $\delta > 0$.

**3.4. The Baron-Ferejohn Model under Open Rule.** In the preceding sections, we have focused exclusively on models where proposals cannot be amended within the current legislative session. We now show how the model can be extended to capture the possibility that proposals can be amended before a final passage vote. We now assume that following each proposal a member is selected at random from the remaining $N-1$ legislators. The selected legislator has two choices. First, she may *call the question* and bring about a final passage vote on the previous proposal. Alternatively, she may make a new offer or *amendment*. The amendment is paired against the current offer. The winner of this vote is the proposal on the floor at the beginning of the next session. In the next session, a new legislator is chosen to either amend or call the question.

Now a legislative proposer has two considerations. First, just as before, a simple majority must receive their discounted continuation values in order to support the proposal on final passage. Secondly, the proposer must craft a proposal for which deters others from amending.

---

[5]This variation of the Baron-Ferejohn game was developed in McCarty (2000a,2000b).

This can be accomplished by allocating sufficient resources such that
the next proposer prefers to move the initial proposal rather than have
her own proposal on the floor at the beginning of the next session.

To keep things simple, we will focus again on $N = 3$. First, we
consider a scenario where the proposer keeps $z$, provides $\frac{1-z}{2}$ to both
other legislators, and each legislator moves the question. To solve for
the optimal $z$, define $v_i^2(z)$ as the continuation value of beginning a
session with a proposal giving $z$ to player $i$ and $\frac{1-z}{2}$ to the other two
legislators. Since we are focusing on symmetric equilibria, we can
suppress the subscript $i$. Thus, $v^2(z)$ is expected utility of this strategy
of the first proposer and that of any proposer who successfully amends
a proposal.

Given this definition, a proposer must give each legislator at least
$\delta v^2(z)$ to induce them to call the question   Otherwise, she would
defect to a proposal giving herself $z$ for the next period. Therefore,
the equilibrium requires that $\frac{1-z}{2} \geq \delta v^2(z)$. So long as this condition
holds, the proposer gets $z$ with probability 1 so that $v^2(z) = z$. Thus,
the proposer will choose to maximize $z$ such that $\frac{1-z}{2} \geq \delta v^2(z)$. This
leads to a solution of

$$v^2(z) = z = \frac{1}{1 + 2\delta}$$

While the proposer can secure $z = \frac{1}{1+2\delta}$ with certainty, it may prefer
to secure the support of only one legislator and risk the defeat of their
proposal if the other is selected to make an amendment. So now assume
that the proposer keeps $z$, gives $1 - z$ to some other legislator, and $0$
to the third legislator. The legislator who receives $1 - z$ moves the
question if selected. The legislator who receives $0$ offers an amendment
giving $z$ to herself, $0$ to the original proposer, and $1 - z$ to the other
legislator. Such an amendment carries with the votes of the legislators
receiving positive allocation in the amended proposal.

To compute the optimal $z$, we need to consider two values. Let
$v_i^1(z)$ be the value to legislator $i$ of beginning the period with a proposal
giving $z$ to $i$ and $1 - z$ and $0$ to the others. Similarly, let $v_i^1(0)$ be the
value to $i$ of the game starting from a proposal that gives her $0$ and
$z$ and $1 - z$ to the others. Again due to symmetry we can drop the
subscripts.

First, we compute $v(z)$. With probability $\frac{1}{2}$, the proposal is moved
and approved giving the proposer $z$. However, with probability $\frac{1}{2}$,
the proposal is amended so that the original proposer gets $0$ in the
proposal in play at the beginning of the next session. Therefore, $v^1(z) =
\frac{1}{2}z + \frac{1}{2}\delta v^1(0)$. Now consider the value of starting the period with $0$.

With probability $\frac{1}{2}$, the proposal is moved and passed leading to a payoff of 0. However with probability $\frac{1}{2}$, the member is selected and can amend the proposal so that she gets $z$ in the standing proposal at the beginning of the next session. Therefore, $v^1(0) = \frac{1}{2}\delta v^1(z)$. Putting these two values together, we get $v^1(z) = \frac{1}{2}z + \frac{1}{4}\delta^2 v^1(z)$ or

$$v^1(z) = \frac{2z}{4 - \delta^2}$$

Finally, we have to insure that the legislator receiving $1 - z$ prefers to move the question rather than amend. This requires that $1 - z \geq \delta v^1(z)$ or $z \leq 1 - \delta v^1(z)$. Therefore, the proposer will choose $z$ to maximize $v^1(z)$ subject to this constraint. The solution is $z = \frac{4-\delta^2}{4+2\delta-\delta^2}$ which leads to a continuation value of

$$v^1(z) = \frac{2}{4 + 2\delta - \delta^2}$$

To determine which strategy the proposer will choose, we simply need to compare $v^1(z)$ and $v^2(z)$. Straightforward algebra shows that $v^1(z) > v^2(z)$ when $\delta > \delta^* \equiv \sqrt{3} - 1$. Intuitively, when players are patient and value the future, it is very expensive to inhibit amendments from both legislators. Therefore, the proposer prefers to buy off only one member and take its chances with an amendment.

There are two interesting features of the open rule model. First, it is possible that coalition are greater than minimum winning. This occurs when $\delta < \delta^*$ so that the proposer spreads resources sufficiently to deter all amendments. Secondly, there can be equilibrium delay in agreement. This occurs when $\delta > \delta^*$ and the proposer gives 0 to one member. If that member is then selected, they make a successful amendment which precludes agreement in the first session.

It is useful to compare the equilibrium allocations from the open rule with those for the closed rule. The literature has paid particular attention to the proposer's share.[6] Recall that for $N = 3$, the proposer keeps $\frac{3-\delta}{3}$. This share is always greater than $v^2(z)$ and greater than $v^1(z)$ when $\delta > \delta^*$. Thus, the open rule lowers the proposer's advantage. However, proposal power can also be mitigated by the use of supermajority rules. Consider the case where $k = N = 3$. The proposer's share is $\frac{3-2\delta}{3}$ which is always lower than $v^1(z)$. Thus, when $\delta > \delta^*$, the unanimity rule lowers proposal power below that of the open majority rule without incurring costly delay.

---

[6]This is not only an important equity consideration. In models of allocating the benefits of costly projects, institutions which limit the proposer's share reduce the incentive to pass inefficient projects (Baron 1992, McCarty 2000b, Primo 2004).

## 4. Bargaining with Incomplete Information

In all of the bargaining models that we have seen so far all of the agents know the disagreement or continuation values of their opponents. Consequently, they know with certainty which offers will be accepted and rejected.[7] This assumption is obviously problematic in many political contexts. A legislature doesn't know whether an executive will sign a particular bill, states do not know whether the peace terms will be accepted or if the opponent will prefer to continue fighting, and so on. In this section, we provide a bare bones model of bargaining with incomplete information. We then elaborate the model with examples from executive-legislative bargaining and crisis bargaining in international relations.

**4.1. A Basic Model.** Consider the setup we used for the Nash bargaining problem where two players are negotiating over the division of $X$. However, now there is uncertainty about player $A$'s disagreement value. With probability $\pi$, $\underline{u}_A = 0$ and with probability $1 - \pi$ $\underline{u}_A = d > 0$. We will refer to $\underline{u}_A = 0$ as the "weak" type and $\underline{u}_A = d$ as the "strong" type. To keep things as simple as possible, we assume that player $B$ makes a take-it-or- leave it offer to $A$ of $(u_A, u_B) \in \Omega$. If $A$ accepts, the offer is implemented. However, if $A$ rejects, the payoffs are $(\underline{u}_A, \underline{u}_B)$.

Clearly, $A$ will only accept an offer that gives her $u_A \geq \underline{u}_A$. However, since $B$ does not know the value of $\underline{u}_A$ she does not know how much utility to transfer to $A$ to secure her agreement. If she offers less than $A$'s disagreement value, $A$ will reject leading to $(\underline{u}_A, \underline{u}_B)$. Let $P(u_A | \underline{u}_A)$ be the probability that $A$ will accept $u_A$ when her disagreement value is $\underline{u}_A$. Therefore, we can write $B$'s expected utility function as a function of her offer

$$EU_B(u_A) = [\pi P(u_A|0) + (1 - \pi) P(u_A|d)] \, g(u_A) + [1 - \pi P(u_A|0) - (1 - \pi) P(u_A|d)] \, \underline{u}_B$$

Since $P(u_A | \underline{u}_A) = 1$ if $u_A \geq \underline{u}_A$ and $0$ otherwise, we can re-write $B$'s utility as

$$EU_B(u_A) = \begin{cases} \pi g(u_A) + (1 - \pi)\underline{u}_B \text{ if } d > u_A \geq 0 \\ g(u_A) \text{ if } u_A \geq d \end{cases}$$

Given that $g(u_A)$ is decreasing in $u_A$, the only possible solutions are $u_A = d$ or $u_A = 0$. We will call $u_A = 0$ the aggressive offer and $u_A = d$

---

[7]In the open rule Baron-Ferejohn game, the uncertainty is about which player will be selected to offer an amed,nment, but not whether a particular player prefers an amendment to the proposal.

the accommodating offer.  $B$ will choose the aggressive offer if and only if $\pi g(0) + (1 - \pi)\underline{u}_B > g(d)$

$$\pi > \frac{g(d) - \underline{u}_B}{g(0) - \underline{u}_B}$$

While very simple, the model makes a number of sensible predictions. First, $B$ is more likely to make the "aggressive offer" of $u_A = 0$ when

- the probability that $A$ is the weak type is high.
- her disagreement value $\underline{u}_B$ is good
- the utility difference between the aggressive and accommodating offers, $g(0) - g(d)$, is large.

Given the assumptions of one-sided incomplete information and a single take-it-or-leave-it offer, this model is very limited.  However, rather than generalize this abstract model, we will review a number of political science applications which relax these assumptions.

## 5. Application: Veto Bargaining

In Chapter 7, we studied the application of the Romer-Rosenthal agenda setting model to the presidential veto.  While this complete information model of the presidential veto provides an excellent tool for studying veto power, it cannot provide a basis for studying vetoes, for the obvious reason that it predicts vetoes will not occur.  We now turn to a simple model for studying vetoes, rather than veto power.  In this model, vetoes do occur.  This simple incomplete information model in turn provides the foundation for building more complex models of veto bargaining that incorporate reputation, learning, and dynamics.[8]

If one wants to explain the fact that vetoes occur, one must dispense with at least one of the assumptions underlying the basic model. While the model presented in Chapter 7 has a number of very restrictive assumptions, few of them are actually consequential in the prediction of no vetoes.  One important exception is the assumption that $C$ has complete information about the preferences of $P$ and $O$. When there is such uncertainty, vetoes may occur because the legislature overestimates its ability to extract concessions from the president or the override pivot.

Relaxing the assumption of complete information has been the starting point for most of the recent work on veto bargaining (Matthews 1989, McCarty 1997, and Cameron 2000).  To present the basic flavor of these models, we consider a model without an override possibility so

---

[8]This section draws heavily on Cameron and McCarty (2004).

that $q$ remains the policy in the event of a veto. To capture the uncertainty that the proposer $C$ faces about the receiver $P$'s preferences, we assume she believes $P$ is one of two preference "types," a moderate with ideal point $m$ or an extremist with ideal point $e$. We assume all agents have linear preferences given by $-|x - i|$ for policy $x \in \mathbb{R}$ and ideal point $i \in \{c, e, m\}$ where $e < m < c$. Let $\pi$ be the probability that $P$ is the extreme type.

The main implication of the uncertainty about preferences is that $C$ no longer knows for sure which bills the president will accept and which he will veto. To see this, consider Figure 10.2 where we assume that $q < e$. Here the set of bills the extremist type of receiver is willing to accept over the status quo is only a subset of those the moderate type is willing to accept. Thus, $C$ can force a more attractive bill (from her perspective) on the moderate receiver than she can on the extremist one. $C$'s dilemma is whether to propose a bill she finds relatively less attractive but that both types will accept – a bill like $b_e$ – or be more aggressive and propose a bill – like $b_m$ – she finds more attractive but only the moderate receiver will accept. Clearly, the attractiveness of the gamble depends on $C$'s beliefs about $P$'s type. If $\pi$ is high (so $C$ believes $P$ is probably an extremist), $C$ will likely be deterred from making the aggressive proposal. On the other hand, if $\pi$ is low (so $C$ believes $P$ is probably a moderate), $C$ may well find an attractive gamble. If she offers it, on occasion it will prove a poor choice: the receiver will turn out to be the extreme type and will veto it.

## Insert Figure 10.2 Here

Now we will compute the necessary conditions for an equilibrium veto to occur. First, assume the preference configuration of Figure 10.2 holds, i.e. $q < e < m < c$. Let $B_t(q)$ be the sets of bills that each type $t \in \{e, m\}$ is willing to accept over the status quo. Similar to our analysis of the complete information version of the model, these sets are $[q, 2t - q]$ if $t > q$ and $[2t - q, q]$ otherwise. Notice that for any $q$, president $m$ is willing to accept a higher bill that is $e$. Since $e > q$, so that $b_e(q) = [q, 2e - q] \subset b_m(q) = [q, 2m - q]$– any bill that $e$ accepts $m$ will accept, but the converse is not true. Therefore, $C$ faces a tradeoff. It can propose $2e - q$ which both types accept, or can propose $2m - q$ which $e$ will veto. Given $C$'s beliefs the latter strategy results in a veto with probability $\pi$.

**Case 1:** $c > 2m - q$. Given $C$'s linear preferences, her utility from $b = 2e - q$ is $2e - q - c$ while her expected utility from $b = 2m - q$ is $\pi q + (1 - \pi)(2m - q) - c$. Thus, if $\pi \leq \frac{m-e}{m-q}$, she prefers $b = 2m - q$ and a veto occurs with probability $\pi$.

**Case 2:** $2e - q < c < 2m - q$. $C$'s payoff from $b = 2e - q$ remains $2e - q - c$, but now $m$ will accept $b = c$. Thus, proposing her ideal point leads to an expected utility of $\pi(q - c)$. Thus, $C$ will propose $b = c$ if $\pi \leq \frac{c+q-2e}{c-q}$. Note that the critical value of $\pi$ is lower than in case 1 making a veto less likely for this preference configuration with $c$ closer to $m$.

**Case 3:** $c < 2e - q$. Now both types will accept $b = c$. So $C$'s proposes its ideal point for all values of $\pi$ and no vetoes occur.

The punchline of this simple model is that vetoes are much less likely to occur when $C$'s preferences are closer to $m$ and $e$. Empirically, Cameron (1999) has found that vetoes are less likely to occur during periods of unified party control of Congress and presidency which he interprets as evidence for this prediction.

**5.1. Models with Reputation, Learning, and Dynamics.** An interesting feature of the incomplete information model is that a moderate receiver $P$ does better if the proposer $C$ believes $P$ is the extreme type. This raises the possibility that $P$ might attempt to manipulate $C$'s beliefs about his type, his reputation. In this section, we examine three models in which the actors try to manipulate $P$'s reputation. All are signaling models, because an informed player takes an action that conveys information about $P$'s type. In the first two models, the veto threat and sequential veto bargaining (SVB) models, the informed player is the receiver $P$ himself. In the third model, the blame game veto model, both $C$ and $P$ take actions to convey information to uninformed voters.

5.1.1. *Veto Threats.* Ranging from the dramatic "read my lips" variety to the much more mundane "statements of administration policy" routinely produced by the Office of Management and Budget, the veto threat is an important feature of legislative politics in the U.S. However, none of the models reviewed thus far provide any leverage on understanding this phenomena. Matthews (1986) however provides an influential model of veto threats where the president may use a costless signal or "cheap talk" to reveal information about preferences and veto intentions.

To illustrate this model, it is helpful to increase the number of presidential types from two to four. Therefore, in addition to $m$ and $e$, we add the two following types: $r$ the "recalcitrant" type and $a$ the "accommodating" type. We assume that each type and $C$ have linear preferences and that $r < q < e < m < c < a$ as in Figure 10.3. President $r$ is called recalcitrant because he will veto any bill that C prefers to the status quo while $a$ is accommodating because he prefers

$c$ to the status quo. We will also assume that the probability of these types are $\pi_r, \pi_e, \pi_m$, and $\pi_a$. In this game, the president first makes a "speech" which is simply a costless signal to the legislature. Each of these messages has no literal meaning, just a contextual one derived from the equilibrium that is being played. Following the speech, $C$ updates her beliefs about the president's preferences and then makes a proposal which the president can either accept or reject.

### Insert Figure10.3

As a baseline, first consider an equilibrium where the president's speeches contain no information because each type makes the same speech. In this babbling equilibrium, $C$ will simply choose the bill from $b_r = q$, $b_e = 2e - q$, $b_m = 2m - c$, or $b_a = c$ to maximize her utility. For any such choice, those with lower types will veto. For example, if $C$ chooses $b_m$, types $e$ and $r$ will veto so that the veto probability will be $\pi_e + \pi_r$. Rather than present the formulae for the conditions for each proposal, they are illustrated graphically in Figure 10.4. The first figure shows which proposal will be made in the babbling equilibrium for different values of $\pi_m$ and $\pi_e$ for given values of $\pi_a$ and $\pi_r$. Note that the proposal $b_r = q$ is never made since $C$ does at least as well with a vetoed proposal. Note that this equilibrium is somewhat bad from the president's perspective. If the president is type $a$, there is a utility loss associated with the fact that $C$ may propose the less desirable policies $b_m$ and $b_e$. For president $m$, there are losses associated with the fact that $C$ might propose $c$ (which he then vetoes) rather than his preferred $b_e$. Since $r$ and $e$, only get their status quo utility from all proposals, they are not affected. $C$ is also affected by the lack of information as it may force her either to accommodate more than necessary or to risk a veto.

### Insert Figure 10.4 Here

So given the bad outcomes from the babbling equilibrium, it is reasonable to ask whether there are other equilibria where more information is transmitted. Matthews shows that some information can be revealed in presidential speeches, but not all of it. First, consider why a separating equilibrium where every presidential type gives a distinct speech cannot be an equilibrium. If $C$ could learn the president's type from the speech, she would optimally propose $b_r$ to $r$, $b_e$ to $e$, etc. However, since $m$ prefers $b_e$ to $b_m$, $m$ would prefer to defect and give $e$'s speech. Thus, a separating equilibrium cannot exist. Matthew's shows that the most informative equilibrium is one where type $a$ reveals his type with an "accommodating" speech and the other types all make the same "threatening" speech. Following an accommodating speech,

$C$ correctly infers that the president will accept her ideal point and thus proposes $c$. Type $a$ is willing to make the accommodating speech since she clearly prefers $c$ to $b_m$ or $b_e$..Following the threatening speech, $C$ learns that the president is not $a$ and updates her beliefs accordingly. Given these beliefs, $C$ chooses between $b_m$ and $b_e$. The second panel of Figure 10.4 illustrates the optimal proposal as a function of $\pi_m$ and $\pi_e$ for given values of $\pi_a$ and $\pi_r$. There are two important things to note. First, it is more likely that $C$ proposes $b_e$ because the knowledge that the president is not type $a$ makes the probability that $b_m$ will be vetoed much higher. This leads to the prediction that $C$ makes a larger concession to the president's preferences after a threatening speech than after an accommodating speech.

It is important to note that an informative equilibrium is not guaranteed to exist. Suppose type $a$ preferred $b_m$ to $c$ to $b_e$, an informative equilibrium would exists only if $C$'s best response to the threatening message was $b_e$. Otherwise, $a$ would defect to the threatening speech. Similarly, if $a$ prefers $b_e$ to $c$, no informative equilibrium can exist.

It is possible for some configurations of preferences that the veto threat is simply a bluff. Consider what would happen if $m$ were moved in Figure 10.3 sufficiently to the right that he preferred $c$ to $q$ (thus became an accommodator) but still preferred $b_e$ to $c$. In the informative equilibrium, $m$ would still give the threatening speech, but it is a bluff in the sense that he would have signed $C$'s ideal point.

The informative equilibrium makes $C$ better off (if it didn't she could just turn off the TV and ignore the speech). However, it is possible that some presidential types will be worse off. Suppose that $a$ were repositioned so that his preference ordering were such that he preferred $b_m$ to $c$ to $b_e$. Further, suppose that the babbling equilibrium produced $b_m$ while a threat in the more informative equilibrium produced $b_e$. Then $a$ would clearly prefer the outcome of the babbling equilibrium to the $c$ she gets from making her accommodating speech in the informative equilibrium.

5.1.2. *Sequential Veto Bargaining with Incomplete Information.* Often, the proposer can make multiple offers, learning about the receiver as she does so. For example, if the receiver rejects a tough offer early, the proposer may believe the receiver is genuinely tough. If so, the proposer's next offer is apt to be more accommodating. This "haggling" dynamic is very common in many types of bargaining, and one might well expect to see it in veto bargaining as well. But a complicating factor is misdirection: the proposer will often have an incentive to reject early offers in order to build a reputation that leads to better later offers. But knowing this, why should the proposer actually make

the compromises? The sequential veto bargaining model (SVB) model explores these questions about learning and credibility.

A simple example conveys many of the basic ideas. First, consider a situation in which $q = 0$, $e = \frac{1}{4}$, $m = .6$, and $c = 1$. By now it should be clear that in a one-shot game (without a veto threat) $b_e = .5$ and $b_m = c$. Using the results of the one-shot incomplete information model, it is easy to see that $C$ will offer $b_m = c$ if $\pi < \frac{1}{2}$ and $b_e = .5$ otherwise. But suppose this is not a one-shot game, so that $C$ may make a second offer if the first is rejected. More specifically, suppose bargaining breaks down with probability $\rho$, but otherwise a second offer can be made (The probability of a bargaining breakdown reflects the inherent uncertainty of the legislative and other political processes.) Is a haggling equilibrium possible, that is, one in which $C$ first makes a tough offer then, following a veto and no bargaining break down, makes a more accommodating offer?

In such a haggling equilibrium, the moderate president must accept the tough offer in the first round (if both types rejected the tough offer, then $C$ should make the accommodating offer lest a break down saddle her with the unappealing status quo). Therefore, the following "incentive compatibility constraint" must hold:

$$(m - c) \geq (1 - \rho)(m - 2e + q) + \rho(q - m)$$

or

$$\rho \geq \frac{c - q - 2(m - e)}{2(e - q)}.$$

The incentive compatibility constraint indicates that accepting the tough offer in the first round is better for the moderate type than rejecting the offer and holding out for the more proximate accommodating offer, taking into account the probability of a bargaining breakdown. In the example, the critical value for the break down probability is .6. Let $\mu(e)$ be $C$'s belief that the president is the extreme type, following a veto. Note the following: in a haggling equilibrium, it must be case that $\mu(e) \geq \frac{1}{2}$, otherwise, following a veto, $C$ will make the tough offer again in the final period (this was proven above). But if the probability of a breakdown is greater than .6, then the moderate type accepts the initial offer, so that by Bayes' Rule $\mu(e) = 1$, following a veto, and $C$ will indeed make the accommodating offer in the second round.

There remains an additional incentive compatibility constraint to examine, however. Congress must find it more appealing to make a tough offer followed by an accommodating offer (conditional on a veto and no break down), rather than make an initial accommodating offer

that would be surely accepted. This requires that

$$\pi\left[(1-\rho)(2e-q-c)+\rho(q-c)\right]\geq\rho(2e-q-c)$$

or

$$\pi\leq\frac{c-2e-q}{c-2e(1-\rho)+q(1-2\rho)}$$

In the example, this condition becomes $\pi\leq\frac{1}{1+\rho}$. We can now indicate a haggling equilibrium in the two period, two type sequential veto bargaining model with the ideal points indicated earlier. If $.6\leq\rho\leq1$ and $\pi\leq\frac{1}{1+\rho}$,then $C$ offers $b_1=b_m=c$ and $b_2=b_e=.5$. Presidential type $m$ accepts both $b_e$ and $b_m$ in both periods, while type $e$ accepts offer $b_e$ and vetoes $b_m$ in both periods. Finally, $C$'s belief that the president is an extreme type is $\mu(e)=1$ following a veto.

5.1.3. *Bargaining over Multiple Bills.* While the last section shows that incomplete information can effect the dynamics of bargaining on a single issue, McCarty (1997) considers how informational and reputational incentives alter the bargaining across multiple issues over time. He considers a model of veto bargaining with incomplete information where $P$ and $C$ bargain over a series of policies with status quo points $q_1$ and $q_2$. In each of the two periods, $C$ proposes $b_t$ and the president decides whether to accept or reject it. . Thus, bargaining over each policy is modeled as one-shot such that if $P$ vetoes $b_t$ the status quo $q_t$ is the policy outcome. Since the president's ideal point is assumed to be constant across policies, the outcome on policy 1 may provide information to C prior to her making an offer on policy 2. Since in the last period, the game is identical to the one-shot incomplete information game described above, type $m$ does better on the second policy by having $C$ believe that he is the extreme type if preferences are such as those given in panel a or b. Thus, given those preference configurations, $m$ may be willing to use his first period veto to build a reputation as the extreme type in order get a better outcome on policy 2. This involves rejecting bills that he, but not type $e$, prefers to $q_1$. Thus, reputational incentives increase the likelihood of a veto on policy 1. Given that $C$ understands these incentives, she may be willing to be sufficiently accommodating on the first policy to discourage type $m$ from vetoing on reputational grounds. Thus, McCarty's model predicts a "honeymoon" pattern of accommodating policies early in the president's term followed by less accommodating policies toward the end when reputational incentives are diminished. However, he notes that since the existence of reputational incentives depends on preference configurations such as those in panel a and b, this honeymoon

effect is unlikely when the expected difference $P$ and $C$ is small such as during unified governments.

**5.2. Blame Game Vetoes.** A recent model argues that vetoes are less a product of legislative uncertainty than of electoral politics. Groseclose and McCarty (2000) examines a model in which the legislative agenda setter uses its proposal power to signal that the president has policy views that are out of step with the voters. In this "blame game" model, vetoes occur when the agenda setter receives a larger payoff from signaling that the president has extreme preferences than she does from enacting new policy. Thus, in this model, the electorate's uncertainty about the president is critical, not the uncertainty of legislators.

To illustrate a simple version of this model, consider a new actor $V$, the voter. We assume $V$ also has linear preferences and an ideal point $v$. Following the notation of the last section, $V$ believes the president is type $e$ with probability $\pi$ and type $m$ otherwise. We focus on the case where $e < m < v$. We assume the voter evaluates the president based on the expected distance between the president's ideal point and her own ideal point. Therefore, the voters evaluation is just

$$w(e, m, \pi; v) = -\pi|v - e| - (1 - \pi)|v - m| = \pi e + (1 - \pi) m - v$$

An important feature of this model is that $P$ and $C$ care how much expected utility $V$ gets from the president's position. The most interesting case is one of conflict, in which the president gets greater utility when the voter believes he is a moderate and Congress gets greater utility when the voter believes the president is an extremist. Such a case would plausibly arise when Congress and the presidency are controlled by different political parties or factions, especially when those parties are highly polarized ideologically, and voters are generally more moderate. In such a case, $C$ and $P$ trade gains from enacting policy with gains from political posturing. More specifically, the president would like to take actions that lead the public to lower $\pi$ while the legislature would like to take actions that lead the public to increase $\pi$. We allow $C$ and $P$ to value these trade-offs differently by letting $\lambda_c$ and $\lambda_p$ be the respective weights each place on policy. Therefore, the utility functions for $C$ and $P$ become:

$$-\lambda_c|x - c| + (1 - \lambda_c)(\pi e + (1 - \pi) m - v)$$

and

$$-\lambda_p|x - p| - (1 - \lambda_p)(\pi e + (1 - \pi) m - v)$$

An important assumption of this model is that while $V$ is relatively uninformed about $P$'s preferences, $C$ is fully informed. Therefore, $C$ may be able to credibly communicate its information about $\pi$ through its choice of bill. Similarly, the president's decision whether to veto particular proposals may also provide information to voters about his preferences.

A particularly interesting equilibrium is one in which $C$ proposes an acceptable bill when $P$ is moderate and submits a bill that will be rejected when the president is extreme. McCarty (2002) shows that such an equilibrium is the only one in which vetoes occur and it exists if and only if the following two conditions hold:

$$(10.1) \qquad \frac{\lambda_p - \lambda_c}{\lambda_p \lambda_c} (1 - \pi)(m - e) \geq 2(e - q)$$

and

$$(10.2) \qquad 2 \geq \frac{\lambda_p - \lambda_c}{\lambda_p \lambda_c} \pi$$

These conditions produce a number of predictions about the occurrence of vetoes.[9] First, note that the first condition cannot be satisfied if $m = e$ or $\pi = 1$. Thus, voter uncertainty about the president's preferences is crucial. Without this uncertainty, orchestrating a veto has no signaling value to $C$ so she might as well make acceptable proposals to both types. Next, note that both conditions are easier to satisfy when $\pi$ is lower. Since the *ex ante* evaluation of the president is decreasing in $\pi$ (the probability he is extreme), the model suggests that vetoes will occur more likely when the public believes the president is moderate (that is, believes the president is ideologically proximate). Intuitively, Congress finds the blame game most attractive when it has negative information about the president's policy preferences that is inconsistent with the voter's beliefs.

The next three prediction are based on $C$ and $P$'s willingness to trade policy gains for political gains. Figure 10.5 illustrates how each of the conditions are affected by the policy weights $\lambda_p$ and $\lambda_c$. The area under the higher solid line represents the combinations of $\lambda_p$ and $\lambda_c$ that satisfy the first condition. Alternatively, the area above the lower dashed line are those satisfying the second condition. The blame game equilibrium described above exists in the intersection of these regions. First, note that the first condition can be met only when $\lambda_p > \lambda_c$,

---

[9]These are conditions are necessary for the case of $c > 2m - q + \frac{1 - \lambda_p}{\lambda_p}(m - e)$. Different positions of $c$ result in slightly modified but qualitatively similar conditions.

suggesting that the president must put relatively more weight on the policy outcome than does Congress. If this were not the case, $C$ would prefer to achieve policy gains by passing mutually attractive bills rather than seek purely electoral advantage by passing bills the president will reject. However, the second condition puts an upper bound on the difference in policy weights. If $\lambda_p$ is much greater than $\lambda_c$, $C$ loses the ability to signal credibly through its proposals. One final prediction emerges from the fact that only extreme types ever veto in the blame game model. Since only type $e$ vetoes, every veto is followed by a reduction of voter support.

### Insert Figure 10.5 Here

## 6. Application: Crisis Bargaining

One of the limitation of bargaining theory is that solutions are generally highly dependent on the bargaining protocol and are therefore not robust to changes in the extensive form of the game. In the context of veto bargaining, this is not such a large problem because its protocol is often codified in constitutional provisions and well-established legislative procedures. However, in the case of crisis bargaining among sovereign states, it is clearly less desirable to have bargaining solutions depend heavily on particular extensive forms since the relevant protocols are generally more informal, non-codified, and unobservable due to secrecy concerns.

Recognizing this problem, Banks (1990) considers what equilibria of a large class of crisis bargaining games have to have in common. Consider the following crisis bargaining scenario. Two states 1 and 2 bargain over 1 unit of territory. Let $x$ be the share that country 1's. Following Banks, we assume that both countries are risk neutral so that we can define country 1's payoffs from a settlement as $x$ and country 2's as $1 - x$. Failure to reach an agreement on the division of the territory leads to a war.[10] We assume that country 1's expected utility of a war is $u$ and country 2's is $v$. These expected utilities encapsulate expectations about the probability of winning the war, the benefits of winning and losing, and the allocation of territory that the winner can secure.

Banks assumes that country 1 has an informational advantage vis a vis country 2 about the values of $(u, v)$. Following the usual practice, he models this asymmetric information by assuming the country 1's type is some $t \in T$ defined such that country 1's expected benefits of

---

[10]We include in the set of possible agreement that the the status quo ex ante remains intact.

war, $u(t)$, are increasing in its type. While country 1 learns $t$ prior to negotiations, country 2 has only a common knowledge prior $f(t)$ over $T$.

Given this framework, standard game theoretic models would specify a set of decisions available to the countries and a (probabilistic) outcome function specifying the probability of a war and the distribution of settlements as a function of these decisions. From this model, equilibrium strategies $(\sigma_1(t), \sigma_2)$ from which the equilibrium probability of war $p(t)$ and expected settlement $x(t)$ can be derived. In such an equilibrium, country 1's payoffs are

$$U(t; x, p) = p(t)u(t) + (1 - p(t))u(t)$$

Clearly, not every $p(t)$ and $x(t)$ can arise from a Bayesian equilibrium. In particular, Bayesian equilibrium places two requirements. The first is incentive compatibility. Type $t$ cannot prefer the outcomes $p(t')$ and $x(t')$ to $p(t)$ and $x(t)$. Otherwise it would defect from its equilibrium strategy $\sigma_1(t)$. Thus, incentive compatibility requires

$$(10.3) \qquad p(t)u(t) + (1 - p(t))x(t) \geq p(t')u(t) + (1 - p(t'))x(t')$$

$$(10.4) \qquad p(t')u(t') + (1 - p(t'))x(t') \geq p(t)u(t') + (1 - p(t))x(t)$$

The second condition imposed by Bayesian equilibrium is individual rationality.[11] It cannot be the case that $u(t)$ is greater than $x(t)$ unless $p(t) = 1$. Otherwise, $t$ would withdraw from the agreement and start a war with probability 1. The individual rationality constraint is

$$(10.5) \qquad p(t)u(t) + (1 - p(t))x(t) \geq u(t)$$

While incentive compatibility and individual rationality are minimal requirements, we will see that they impose quite a bit of structure on bargaining outcomes. The most important feature of the Bayesian equilibrium concern the monontinicity of $p$, $x$, and $U$ in $t$.

LEMMA 10.1. *If $p$ and $x$ are incentive compatible and individually rational, then $p(t)$ is weakly increasing on $T$.*

PROOF. Let $t' > t$. Note that we can subtract the right side of equation 10.3 from the left side of equation 10.4 and the left side of equation 10.3 from the right side of 10.4 to produce

$$p(t')[u(t') - u(t)] \geq p(t)[u(t') - u(t)]$$

Since $u(t)$ is strictly increasing, $u(t') - u(t) > 0$ so that it must be the case that $p(t') \geq p(t)$. $\qquad\qquad\qquad\square$

---

[11]The similarity of Banks' approach to mechanism design should be obvious to the attentive reader.

This lemma shows that in any Bayesian equilibrium the probability of war cannot decrease as country 1's expected utility of war increases. Not only is this a feature of strategic models, it is consistent with the assumptions of a number of decision-theoretic models of war.

For the next result, we need to define the set of types who resolve the dispute through bargaining with a positive probability for some outcome $p, x$. Let $T_b(x, p) = \{t \in T : p(t) < 1\}$. Note that individual rationality requires that $x(t) \geq u(t)$ for any $t \in T_b$.

LEMMA 10.2. *If $p$ and $x$ are incentive compatible and individually rational, then $x(t)$ is weakly increasing on $T_b$.*

PROOF. Let $t, t' \in T_b$ and $t' > t$ that that Lemma 1 implies $1 > p(t') \geq p(t)$. Since $x(t) \geq u(t)$ for all $t \in T_b$, we know that

$$(10.6) \qquad p(t)u(t') + (1 - p(t))x(t') \geq p(t')u(t') + (1 - p(t'))x(t')$$

We can combine this result with equation 10.4 to produce

$$p(t)u(t') + (1 - p(t))x(t') \geq p(t)u(t') + (1 - p(t))x(t)$$

which reduces to $x(t') \geq x(t)$ after dividing by $(1 - p(t))$ which we know is positive since $t \in T_b$.                    $\square$

Not surprisingly, country 1 must do at least as well in the bargaining outcome when its war utility improves. Thus, taken together Lemmas 2 and 3 suggest that in any Bayesian equilibrium that higher types may get better bargaining outcomes, but that this comes at a greater risk for war.[12]     However, incentive compatibility requires that these trade-off benefit higher types (or else they would mimic lower types). Let $T_w = \{t \in T : p(t) > 0\}$ so that $T_w$ is the set of types that go to war with some probability.

LEMMA 10.3. *If $p$ and $x$ are incentive compatible and individually rational, then $U(t; x, p)$ is weakly increasing on $T$ and strictly increasing on $T_w$.*

PROOF. Let $t' > t$, and $t', t \in T_w$.   Suppose (contra the lemma) that $U(t) \geq U(t')$ so that

$$p(t)u(t) + (1 - p(t))x(t) \geq p(t')u(t') + (1 - p(t'))x(t')$$

Since $u(t') > u(t)$, this implies that

$$(10.7) \qquad p(t)u(t') + (1 - p(t))x(t) \geq p(t')u(t') + (1 - p(t'))x(t')$$

_____

[12]Banks also shows that if $x$ and $p$ are incentive compatible and individually rational then $x(t') > x(t)$ if and only if $p(t') > p(t)$ for $t' > t$ and $t, t' \in T_b$.

Since $t', t \in T_w$, $p(t)$ and $p(t')$ are greater than zero so that equation 10.7 violates equation 10.4. Thus, $U(t') > U(t)$. However, this strict equality does not hold on $T/T_w$. If $t, t' \in T/T_w$, then $p(t) = p(t') = 0$ so that equations 10.3 and 10.4 clearly imply that $x(t) = x(t')$ and $U(t) = U(t')$.                                                                $\square$

Thus, having a better expected utility of war cannot make country 1 worse off. In fact, for types that go to war with a non-zero probability, higher war payoffs leads to strictly higher equilibrium payoffs.[13].

While the incentive compatibility approach can go a long way towards telling us what predictions are generic to crisis bargaining models, there are a number questions that it cannot resolve. For example, we do not learn how country 2's perceptions of country 1's war utilities, as measured by the priors $f(t)$ effect the likelihood of war or the bargaining settlement. For that we will have to turn to more explicit models of crisis bargaining.

**6.1. Models of Crisis Bargaining.** Fearon's (1995) seminal article explores several models in the class covered by Bank's results. He gives a specific form for $(u, v)$. In particular, Fearon assumes that each country has a cost of war $c_i > 0$, country 1 wins any war with probability $\pi \in (0, 1)$, and the winner of the war can impose its most preferred settlement ($x = 1$ for country one and $x = 0$ for country 2). Therefore, $u = \pi - c_1$ and $v = 1 - \pi - c_2$. Let $x_0$ be the status quo allocation of the territory.

Given this framework, we can define the set of agreements that each side will accept in lieu of going to war. For country 1, we require that $x > \pi - c_1$ and for country 2 we require $1 - x > 1 - \pi - c_2$. Therefore, any allocation $x \in [\pi - c_1, \pi + c_2]$ avoids conflict. Since the costs of war are positive, the set of peaceful agreements is non-empty. Therefore, under perfect information, we should expect that one of these agreements will be reached and war will be avoided.[14]

Fearon considers a simple model with incomplete information about country 1's costs. While he assumes a continuous distribution of $c_1$, it suffices for us to consider a cost where $c_1$ can take on only two values $\overline{c} > \underline{c}$. The common knowledge prior is that $c_A = \underline{c}$ with probability $\lambda$. Fearon first considers a model where country 2 makes a single take-it-or-leave-it offer to country 1. If country 1 rejects it, war ensues.

---

[13]Banks also shows that $U(t; x, p)$ is continuous in $t$, but we refer the reader to his article.

[14]This statement of course assumed that the territory is infinitely divisable so that an agreement in thus region is feasible. This may not be the case in some disputes especially when the terrotory involves religious or ideational significance.

In analyzing this model, note that if country 2 offers $x \geq \pi - \underline{c}$ both country 1 types will accept and war will be avoided. Clearly, country 2 has no incentive to pay higher than $\pi - \underline{c}$, so let $\underline{x} = \pi - \underline{c}$. If country 2 offers $x \in (\pi - \overline{c}, \pi - \underline{c}]$, only the low cost type will accept it so that war starts with probability $\lambda$. Of these offers, country 2 prefers $\overline{x} = \pi - \overline{c}$. Finally, if country 2 offers $x < \pi - \overline{c}$ both types reject and a war starts with certainty. This generates a payoff to country 2 of $v = 1 - \pi - c_2$. Thus, country 2's choice boils down to a choice of three utilities $1 - \underline{x}$, $\lambda v + (1 - \lambda)(1 - \overline{x})$, or $v$. We can easily dismiss the third option. Since the interval $[\pi - \underline{c}, \pi + c_2]$ is non-empty, we know that $1 - \underline{x} > v$. Thus, country 2 will never choose to sabotage the negotiations to generate a war with probability 1. Now we can determine country 2's preferences over the remaining offers. Clearly, $\lambda v + (1 - \lambda)(1 - \overline{x}) \geq 1 - \underline{x}$whenever

$$\lambda \leq \frac{\overline{c} - \underline{c}}{c_2 + \overline{c}}$$

Thus, when the probability that country 1 has low costs is sufficiently low, country 2 will take an aggressive bargaining stance that risks war in the event that country 1 actually does have low costs. Note that the critical threshold is decreasing in country 2's costs since it is less willing to take such a risk when its military capabilities are low.

We can easily check that Bank's results hold trivially for this model. When $\lambda > \frac{\overline{c} - \underline{c}}{c_2 + \overline{c}}$, both types get the same allocation and have zero probabilities of going to war. When $\lambda \leq \frac{\overline{c} - \underline{c}}{c_2 + \overline{c}}$, both types are still offered the same allocation but $\underline{c}$ goes to war with probability one.

As Fearon notes there are reasons skeptical of informational explanations for war. Perhaps opportunities for communication should resolve such informational asymmetries and avoid war. However, given that the countries have diametrically opposed preferences over the allocations, it is easy to show that cheap talk will not influence bargaining or the probability of war. Now let country one announce $H$ or $L$ as a signal of its costs $\overline{c}$ and $\underline{c}$, .respectively [15] Following the message, let $\lambda^*$ be country 2's updated beliefs about 1's costs. Clearly, based on the these updated beliefs, country 2 will use the same cutpoint rule as before. First we consider whether there are separating equilibria where type $\overline{c}$ reports $H$ and $\underline{c}$ reports $L$. In such a case, $\lambda^*(H) = 0$, $\lambda^*(L) = 1$, $x(H) = \overline{x}$, and $x(L) = \underline{x}$. The probability of war would be zero. Therefore, separating messages requires the incentive compatibility conditions $x(H) \geq x(L)$ and $x(L) \geq x(H)$. These conditions

---

[15]The restriction to two messages or enbuing them with literal meaning is not consequential.

clearly fail since $\underline{x} > \overline{x}$. We leave it to the reader to verify that there are no partially information semi-pooling equilibria.

**6.2. A Model of Escalation\*.** This section is based Fearon (1994) who develops a version of the war of attrition to explore how "audience costs" imposed on states who back down in international disputes effect the dynamics of crisis escalation.

Assume that two states 1 and 2 are in a dispute over a prize worth $v > 0$. The game is played in continuous time beginning at $t = 0$. At every instant each state can choose among three strategies: *attack, quit,* or *escalate.* The game continues until one or both of the states quits or attacks. If both states choose *escalate,* the game continues. If either state *attacks* before the other *quits* they both receive their expected payoffs from war $w_i < 0$. Fearon interprets these payoffs as resolve since the state with the higher $w_i$ is relatively more willing to engage in military conflict to settle the dispute.

Fearon's interest is in understanding how the sanctions imposed on leaders who back down during disputes effect crisis behavior. Therefore, he assumes that if state $i$ quits before state $j$ at time $t$ if suffers audience costs $a_i(t)$ which are strictly increasing in $t$. The dependence on $t$ is designed to reflect the intuition that it is more costly to back down during a protracted dispute than in a short one.

A pure strategy in this game specifies for any subgame beginning at time $t'$ a rule specifying a finite time $t \geq t'$ at which to *attack* or *quit.*[16] We can write these strategies as $\{t, attack\}$ meaning "escalate until $t$ and then attack" or $\{t, quit\}$ to represent "escalate until time $t$ and then quit."

Before turning to the more general model where each side is uncertain of the other's resolve, it is instructive to consider the case of complete information. For each state, we can compute the time $t$ for which it strictly prefers to attack rather than back down. Clearly, this occurs when $w_i \geq -a_i(t)$. Let

$$\overline{t}_i = -a_i^{-1}(w_i)$$

Suppose, either because state 1's resolve or audience costs are higher, that $\overline{t}_1 < \overline{t}_2$. Thus, at $\overline{t}_1$, state 2 prefers still prefers to quit than be attacked, thus it will quit. Let $Q_i(t; t')$ be the probability that state $i$ quits before time $t$ conditional on it not quitting before $t'$ and $Q_i(t)$ be the unconditional probability of quitting by time $t$.

---

[16]Not allowing "never quit" to be part of the stragy set eliminates the uninteresting equilibrium where both sides choose this strategy and escalate forever.

Thus, at every subgame $0 \leq t' < \bar{t}_1$, state 2 receives $-a_2(t')$ for quitting immediately and $Q_1(t; t')v - (1 - Q_1(t; t'))a_2(t)$ for $\{t, \ quit\}$. However, consider state 1's strategy. From subgame $0 \leq t' < \bar{t}_1$, strategy $\{\bar{t}_1, \ attack\}$ has a payoff of $v$ while $\{t, \ quit\}$ has a payoff of $Q_2(t; t')v - (1 - Q_2(t; t'))a_2(t) < v$. Therefore, $Q_1(t; t') = 0$ for all $t < \bar{t}_1$. Now we can see that state 2's payoff from $\{t, \ quit\}$ is $-a_2(t) < -a_2(t')$. Thus, state 2 will quit immediately at every subgame $0 \leq t' < \bar{t}_1$ including 0. The equilibrium with complete information therefore involves state 2 stopping immediately and state 1 claiming the prize.

The complete information equilibrium has the property that both high resolve and audience costs lead to better crisis bargaining outcomes. However, it has the unsettling prediction that no crises ever occur, because the weaker side capitulates immediately. Therefore, Fearon also considers an incomplete information version of the game where each side is uncertain of the others resolve. We assume that the resolve of state $i$ is distributed according to the cumulative distribution function $F_i$ on the interval $[\underline{w}_i, 0]$.

Just as in the complete information game, the equilibrium depends on defining a time point after which neither state will wish to quit. Fearon refers to such a time point as the *horizon* of the crisis game. Formally, we define this horizon point as $t_h$, the earliest time point at which that $Q_i(t)$ is not increasing for $t > t_h$ for $i = 1, 2$.

Fearon observes that the following must be true in any PBE[17]:

(1) Both states quit simultaneously with probability zero. Suppose states 1 quits with positive probability mass at $t'$. Then clearly state 2 has an incentive to wait at until at least $t' + \varepsilon$ before quitting as this increasing its probability of winning $v$ substantially with only an infinitesimal increase in its audience costs. Similarly, there are no PBE where a state quits contemporaneous with an attack from the other state. Again, if state 2 expects that state 1 will quit with probability mass at point $t'$, it should hold off its attack until $t' + \varepsilon$. This implies that both states cannot plan to quit at $t_h$.

(2) State $i$ will not attack during at time $t'$ if $Q_j(t)$ is increasing at $t'$. Since attacks have negative expected utility, it pays to wait longer in the hopes that the opponent will drop out prior to the attack.

(3) Both states will quit with positive probability in time intervals arbitrarily close to $t_h$. By the definition of $t_h$, at least one state

---

[17]For formal statememnts of these observations and their proofs, see Fearon 1994.

must quit with positive probability in the arbitrarily small in-
terval before $t_h$. Suppose that this is true for state $i$. Now
suppose to the contrary that state $j$ quits with zero proba-
bility after time $t' < t_h$. Clearly from observation 2, state
2 will not attack between $t'$ and $t_h$. Therefore, quitting with
positive probability after $t'$, state 1 unnecessarily increases its
audience costs – it knows it will quit before state 2 yet it keeps
escalating.

(4) Both states attack with probability zero for $t < t_h$. This
follows directly from observations 2 and 3. Therefore, the
utility of strategy $\{t > t_h, \text{attack}\}$ for state $i$ is

$$(10.8) \qquad U_i^a(t, w_i) = Q_j(t_h)v + (1 - Q_j(t_h))w_i$$

while the utility of $\{t < t_h, \text{quit}\}$ is

$$(10.9) \qquad U_i^q(t) = Q_j(t)v - (1 - Q_j(t_h))a_i(t)$$

Since $U_i^q(t)$ does not depend on $w_i$, $\{t < t_h, \text{quit}\}$ can only
be a best response if it is constant for all $t$, say $k_i$. Suppose
this were not true and let $U_i^q(t') > U_i^q(t)$ for some $t' < t_h$
and all $t < t_h$. Then all types such that $U_i^q(t') > U_i^a(t_h, w_i)$
would quit exactly at $t'$. State $j$'s best response would then
be $\{t > t', \text{quit}\}$ making $U_i^q(t') = -a_i(t')$ which contradicts
$U_i^q(t') > U_i^q(t)$ for all $t < t_h$.

Based on these observations, the following lemmas help to charac-
terize the perfect Bayesian equilibrium for this game.

LEMMA 10.4. *In any equilibrium in which both states choose to
escalate with positive probability, there must exist a finite horizon $t_h$.*

The logic of this lemma straightforward. Suppose to the contrary
that there were a PBE where $Q_i(t)$ were increasing for all $t$. By
observation 2, state $j$ will not attack with probability 1 which turn (by
observation 4) suggests that $U_i^q(t)$ is constant for all $t$. This implies
that

$$Q_j(t) = \frac{k_i + a_i(t)}{v + a_i(t)}$$

However, since $j$ never attacks, it must be the case that $\lim_{t \to \infty} Q_j(t) = 1$.
This can only be true if $k_i = v$ or $\lim_{t \to \infty} a_i(t) = 1$. If $k_i = v$, $Q_j(0) = 1$
implying that $j$ does not escalate with probability 1. If $\lim_{t \to \infty} a_i(t) = 1$,
$i$ will not wish to choose $\{t, \text{quit}\}$ for arbitrarily large $t$.

LEMMA 10.5. *In any equilibrium with $t_h$ as the horizon and in which
escalation may occur, (1) if state $i$ chooses $\{t, \text{attack}\}$ it must be the*

*case that $t \geq t_h$; and (2) state $i$ will choose $\{t, attack\}$ where $t \geq t_h$ if $w_i > -a_i(t_h)$ and only if $w_i \geq -a_i(t_h)$.*

Part 1 of this lemma follows directly from observations 2 and 3. Ignoring several technical complications, part 2 follows from the fact that $U_i^a(t, w_i) \geq U_i^q(t)$ if and only if $w_i \geq -a_i(t_h)$.[18]

From lemma (2) we know that the ex ante probability that state $j$ will attack at $t_h$ is $(1 - F_j(-a_j(t_h)))$. Thus, state $i$'s ex ante utility of escalating up to $t_h$ and then backing down is

$$u_i(t_h) = F_j(-a_j(t_h))v - (1 - F_j(-a_j(t_h)))\, a_i(t_h)$$

We can then define $t_i^*$ such that $u_i(t_i^*) = 0$. Thus, $t_i^*$ has the property that state $i$ is indifferent between escalating to time $t_i^*$ and conceding immediately.[19]

PROPOSITION 10.4. *Let $t_i^*$ be the unique solution $u_i(t_i^*) = 0$ and $t^* = \min\{t_1^*, t_2^*\}$. For any equilibrium in which escalation occurs with positive probability, the horizon must be $t^*$.*

If $t_h$ were greater than $t^*$, the state with the lower $t_i^*$ would have an incentive to quit with probability 1 before $t_d$. This would contradict observation 3. If $t_d$ were greater than $t^*$, then both states would have an incentive to bluff a little longer at $t_d$ before quitting. This contradicts the definition of $t_d$.

These lemmas lead directly to the main result.

PROPOSITION 10.5. *Label the players so that $t^* = t_2^* < t_1^*$. Let $k_1 = u_1(t^*) > 0$. The following describes equilibrium strategies for state $i = 1, 2$ as a function of type $w_i$:*

*For $w_i \geq -a_i t^*$, state $i$ plays $\{t, attack\}$ for any $t > t^*$.*

*For $w_i < -a_i t^*$, state $i$ plays $\{t, quit\}$ with any pure strategies that yield the following cumulative distributions*

$$\Psi_1(t) = \frac{1}{F_1(-a_1(t^*))}\frac{a_2(t)}{v + a_2(t)}$$

$$\Psi_2(t) = \frac{1}{F_2(-a_2(t^*))}\frac{k_1 + a_1(t)}{v + a_1(t)}$$

*For $t \leq t^*$, state $i$ believes that the probability that $j$ will not back down is given by*

$$\Pr(w_j \geq -a_j t^* | t) = \frac{v + a_i(t)}{v + a_i(t^*)}$$

---

[18]The technical complications involve ruling situations were $Q_i(t)$ has mass points.

[19]We implicitly assume that the range of $a_i(t)$ is sufficiently large that there is a unique solution to $u_i(t_i^*) = 0$.

*For $t > t^*$, state $i$'s beliefs follow Bayes' Rule in accord with the opponents's strategy for attacking. For any $t > t^*$ off the equilibrium path, let $i$ believe that $w_j > -a_j(t^*)$ and is distributed according to $F_j$ truncated at $-a_j(t^*)$.*

PROOF. Let $Q_i(t)$ be the unconditional probability that state $i$ quits by time $t$ From lemma (2) and Proposition (1), $Q_i(t^*) = F_i(-a_i(t^*))$. The utility to state $i$ of $\{t, \text{quit}\}$ is therefore

$$Q_j(t)v - (1 - Q_j(t))a_i t$$

To ensure that $i$ is indifferent between quitting and continuing for any $t < t^*$, we require that $Q_j(t)v - (1 - Q_j(t))a_i(t) = u_i(t^*)$ or

$$Q_j(t) = \frac{u_i(t^*) + a_i(t)}{v + a_i(t)}$$

Since only types $w_j < -a_j(t^*)$ ever quit, these types must be quit at rates $\frac{1}{F_j(-a_j(t^*))}Q_j(t)$ or $\Psi_j(t) = \frac{1}{F_j(-a_j(t^*))}\frac{u_i(t^*)+a_i(t)}{v+a_i(t)}$.

We know that by time $t$, $\Psi_j(t)$ of the types in the interval $[\underline{w}_j, a_j t^*)$ will have dropped out so that

$$\Pr(w_j \geq -a_j(t^*)|t) = \frac{1 - F_j(-a_j(t))}{1 - F_j(-a_j(t)) + F_j(-a_j(t))(1 - \Psi_j(t))}$$

$$\Pr(w_j \geq -a_j(t^*)|t) = \frac{1 - F_j(-a_j(t))}{1 - Q_j(t)} = \frac{(1 - F_j(-a_j(t)))(v + a_i(t))}{v - u_i(t^*)}$$

$$\Pr(w_j \geq -a_j(t^*)|t) = \frac{v + a_i(t)}{v + a_i(t^*)}$$

$\square$

In this PBE, types with low resolve from state $i$ drop out at a rate designed to keep the low resolve types state $j$ indifferent between dropping out and escalating through time $t^*$. At time $t^*$, both states attack because they know that all of the low resolve types for the other state would have dropped out by then and escalating would simply lead to larger audience costs.

It is instructive to explore why low resolve types of state $j$ quit at a rate to make low resolve types of state $i$ indifferent between dropping out and escalating. If they dropped out at a faster rate, all low-resolve types in state $i$ would conclude at each $t$ that state $j$ is more likely to be a strong type. This would lead low-resolve types of state $i$ to quit more quickly. Then state $j$ would then begin to infer that the pool of state $i$ types is stronger, and begin to drop more quickly. In the limit, all low-resolve types would quit at $t = 0$, but this cannot be an equilibrium (recall Observation 1).

Now consider what happens if low resolve types from state $i$ drop out at a slower rate than the equilibrium. Then at each $t$, state $j$ would infer that the pool of remaining types is weaker than the corresponding equilibrium pool. This would lead state $j$ types to also drop out at a slower rate, which in turn induces state $i$ to escalate more, and so on. Such a dynamic would lead to all types to prefer escalating until $t^*$. However, this cannot be an equilibrium because the low resolve types would have clear preference for dropping out before $t^*$ than attacking at $t^*$.

An important substantive feature of the model is that it turns out to be beneficial to be able to incur large audience costs. Such a state is better able to convince its opponent that it is "locked-in" to the conflict since a larger set of types are willing to escalate to the horizon and then attack. This runs counter to a simple intuitive prediction that those with the most to lose by backing down will back down earlier. However, conditional starting a crisis, the signaling value outweighs this effect. By modifying his model slightly to include an explicit initiation phase, Fearon argues that this framework provides a justification for why democracies (highly sensitive to audience costs) might be less likely to initiate conflict, but will be more likely to prevail.

## 7. Exercises

EXERCISE 10.1. *Let $U_i(x_i) = \ln x_i$ for $i \in \{A, B\}$. Solve for the Nash bargaining solution as a function of the disagreement values.*

EXERCISE 10.2. *In the Rubinstein bargaining model with $\delta_1 = \delta_2$ and $d_1 = d_2 = 0$, assume that $u_i(x_i) = x_i^\alpha$ where $0 < \alpha < 1$. Compute the SPNE shares. What is the effect of risk-aversion (lower $\alpha$)?*

EXERCISE 10.3. *Consider the Baron-Ferejohn model where $u_i(x_i) = x_i^\alpha$. Show that the initial proposer's share is decreasing in $\alpha$.*

EXERCISE 10.4. *In the model described in Section 3.2 with $N = 3$ and $m = 1$, suppose that $\frac{1-\delta}{3-2\delta} < q \leq \frac{1}{3}$. Compute a mixed strategy equilibrium where $v_A = v_B$.*

EXERCISE 10.5. *In the model described in Section 3.2, compute $v_A$ and $v_B$ for generic values of $N$ and $n$.*

EXERCISE 10.6. *Consider an extension of the model described in Section 4.1. First assume that there is two rounds of bargaining so that A makes a counter offer if it rejects B's initial offer. Assume that the payoffs are discounted by a factor $\delta$ if agreement is reached in the second round. What is the PBE to this game? Now assume that there*

*is incomplete information about B's disagreement value where $\underline{u}_B = 0$
with probability $\pi$ and $\underline{u}_B = d$ withe probability $1 - \pi$? Construct a
PBE equilibrium to this game. Is it unique?*

EXERCISE 10.7. *Consider an extension of the model considered in
Section 3.3. Assume that there are two groups of bargainers, 1 and
2. Let $m_i$ be the number of members of group $i$ so that $m_1 + m_2 = N$.
Suppose that all members of group $i$ have a qualified veto power in that
if they object to the proposal $k_i$ votes are required to override their veto
where $0 < k_i < N$. Assume that $k_1 > k_2$. Compute continuation
values for members of each group.*

EXERCISE 10.8. *In the model described in Section 6.1, assume that
$c_1$ is distributed uniformly on the interval $[\underline{c}, \overline{c}]$. Compute the perfect
Bayesian equilibrium. Now consider the extension with pre-bargaining
cheap talk. Show that if $\overline{c}$, there is a PBE where the high cost types
reveal information about their cost to country 2.*

# CHAPTER 11

# Mechanism Design and Agency Theory

So far we have discussed techniques for analyzing how strategic agents behave in specific games. In social settings where the "rules" are fairly clear, this approach to game is a powerful source of empirical predictions about the outcomes of strategic interactions.

An alternative approach is to ask a slightly different question: given a desired outcome, what game should be designed so that strategic agents will produce it? The field of game theory that asks such questions is called mechanism design. Here a designer or principle selects a Bayesian game, or *mechanism*, for an agent or group of agents to play. Examples include the selection of tax codes to induce agents to reveal their willingness to fund public projects, the design of auctions to maximize revenue, and selection of reelection functions by voters to create incentives for government officials to "behave" while in office.

Typically, we model the choice of mechanisms as a maximization problem given the designer's preferences over agent types and outcomes. If the designer's preferences can be interpreted as an objective notion of social preferences then the problem of selecting a mechanism can be viewed as a normative exercise. A classic example is the selection of rules to determine the provision of public goods so as to maximize the sum of individual utilities. Given this normative interpretation, mechanism design is closely related to social choice theory. A version of mechanism design known as implementation theory seeks to uncover classes of choice rules (mappings from agent types to collective decisions) for which there exists a mechanism which will achieve the choice function. Choice functions of this type are said to be implementable. The Gibbard-Satherwaite theorem was an example of this type of work.

While the applications of mechanism design are often normative or prescriptive, we also may use it to make positive predictions. For example, mechanism design is often used to investigate whether a poorly informed principal (e.g. legislature or executive) can create incentives so that well informed agents (e.g. committees or bureaucrats) take actions to achieve ends which the principal desires.

In most applications in economics, it is assumed that the designer has a very rich set of games to choose from. She is usually free to commit to contracts with very elaborate reward and punishment schemes. Such assumptions seem entirely reasonable in economic settings where third-parties such as courts may be counted on to enforce complex agreements and large monetary rewards and sanctions are considered legitimate. However in applications to politics, it is often unreasonable to believe that principals can pre-commit to reward scheme because third-party enforcement is often unavailable. Also, the monetary incentives may be legally or socially proscribed. Accordingly, after presenting some basic concepts and results of implementation theory and mechanism design, we will focus much of our attention on the design of incentive mechanism when the principal is more constrained.

To get a feel for mechanism design, let's consider one of the classic examples – a political science department of $n$ members and a chair that is deciding whether to purchase a nice Saeco espresso maker. The fancy coffee maker has a cost $c$ and the chair wants to see if the department members value the machine enough to justify the expense. Assume that each member's valuation of the coffee maker is $\theta_i \in \mathbb{R}^1_+$. The chair, a Benthamite and non-coffee drinker[1], wishes to purchase the machine if and only if $\sum_{i=1}^{n} \theta_i \geq c$. Unfortunately, the chair does not know the valuations of individual department members. Instead she believes that each members type is drawn from the probability distribution $F(\cdot)$. What should the chair do? One solution might be to privately ask each member for his valuation and purchase the espresso maker with department funds if the total announced valuations exceed $c$. The problem is that some colleagues might find it advantageous to inflate their valuations in order to increase the likelihood that the machine is purchased. Thus, this scheme induces a game where each faculty member's best response is to report a valuation that exceeds his true preference. So this solution may not do a very good job in determining whether the espresso maker should be purchased. The chair may decide that a better solution is to ask each member to contribute her valuation and then if the total contributions exceed $c$, purchase the maker, and keep the surplus to pay for coffee beans. If the contributions do not reach $c$, the chair would return them. This mechanism is also flawed. Now the strategic scholars may decide to understate their valuations

---

[1]Jeremy Bentham (1748-1832) was a prominent British philosopher who advocated an ethical system based on pursuing the "greatest good for the greatest number." Historians record that he was also a coffee drinker.

hoping to free-ride off the contributions of their colleagues. This rule creates the classic collective action problem.

While neither of the preceding schemes works well, we can use the theory of mechanism design to uncover a class of particularly simple mechanisms that can be used by the chair to learn the faculty's preferences. Groves (1973) and Clarke (1971) show that the following mechanism has good properties.

- Ask each faculty member to e-mail her valuation $m_i$ to the chair.
- If $\sum_{i=1}^{n} m_i \geq c$ purchase the coffee maker, otherwise do not.
- If the coffee maker is purchased collect from faculty member $i$ the amount $t_i(m_i.m_{-i}) = c - \sum_{j \neq i} m_j$.
- If the coffee maker is not purchased, collect no money.

We can demonstrate that given this mechanism each faculty member has an incentive to reveal her true valuation regardless of the other members' valuations. A key property of this mechanism is that member $i$'s message only indirectly effects her contribution through its effect on the ultimate decision about whether to purchase the espresso maker. The amount that each member pays depends on the messages of all the other members.

To see that all members will offer truthful messages, consider the decision of member $i$ with type $\theta_i$. Suppose that she were to lie with the announcement $m_i' < \theta_i$. This understatement affects the outcome only if affects the likelihood that the machine is purchased, or if $\sum_{j \neq i} m_j + m_i < c < \sum_{j \neq i} m_j + \theta_i$. Under these circumstances, $c - \sum_{j \neq i} m_j < \theta_i$ so that the contribution required from truthful announcement, $c - \sum_{j \neq i} m_j$ is less than the member's value of the new espresso maker. This deviation from a truthful response can only make the department member worse off. Now consider whether a member has an incentive to overstate her demand for espresso with a message $m_i' > \theta_i$. This fabrication only effects $i$'s utility if the inflated message results in purchasing the machine where the truthful message would not have. Such a scenario requires that $\sum_{j \neq i} m_j + m_i > c > \sum_{j \neq i} m_j + \theta_i$

so that $c - \sum_{j \neq i} m_j > \theta_i$. Thus, member $i$'s contribution is more than her value the coffee maker. So lying doesn't pay.

Beyond promoting honesty in departmental affairs, the Groves-Clarke mechanism has the desirable property that the espresso maker is purchased if and only if the aggregate valuation of the department exceeds its cost. However, it has a less desirable property that it is not "budget balancing." When the machine is purchased, the chair collects $\sum_{i=1}^{n} t_i(m_i.m_{-i}) = nc - \sum_{i=1}^{n} \sum_{j \neq i} m_j$ which is greater than or equal to $c$. Of course, this isn't much of a problem as far as the chair is concerned – a little compensation for having to send and read all those e-mail messages.

## 1. The Mechanism Design Problem

Now we consider the mechanism design problem in a more abstract setting. Consider a set $N$ of $n$ agents and a mechanism designer (denoted agent 0). The designer ultimately selects a policy $x \in X$. Each agent has a private type $\theta_i \in \Theta$, with the joint type vector $\theta$ drawn from the joint distribution function $F(\theta)$. Agents also have Bernoulli utility functions $u_i(x, \theta) : X \times \Theta^n \to R^1$ that depend on the chosen policy and the agents type. The mechanism designer has a Bernoulli utility function $u_0(x, \theta)$. In many application agents care only about their own type, but we allow agent's payoffs to be a function of the entire profile. The primitives of a mechanism design problem are therefore $\langle \Theta, F(\theta), X, u \rangle$.

In a typical application, the mechanism designer will elicit a public vector of signals from the agents. The designer chooses a message space for each agent $M_i$ and a policy function $p(m) : \prod_{i \in N} M_i \to X$, that selects a policy for every possible profile of messages $m = (m_1, m_2, ...., m_n) \in M = \prod_{i \in N} M_i$. Accordingly a mechanism is a pair $\langle M, p(\cdot) \rangle$.

For a given choice of message spaces and policy function, the $n$ agents play the Bayesian normal form game with the strategy sets $S_i = M_i$ and payoffs given by the composition of $u$ and $p$, $u_i(p(m), \theta)$.

It is straightforward to see how the espresso mechanism maps into this general framework.. Clearly, the chair is the designer who selects the message space $M = \Theta$ and implements the policy "buy if $\sum_{i=1}^{n} m_i \geq c$

and charge $c - \sum_{j \neq i} m_j$ to agent $i$." The faculty members then play a Bayesian normal form game with strategies and payoffs determined by the chair's choices.

Given a mechanism, determining how agents will behave requires specifying a form of rationality. One possibility is to make predictions only if the agents have dominant strategies in the induced game. These are known as *dominance solvable* mechanisms. Alternatively, we could assume that agents select a Bayesian Nash equilibrium. Clearly, the question of how well the mechanism performs rests on assumptions about how the agents play the induced game.

In this chapter we will focus on Bayesian Nash equilibria. A large literature exists in economics using other solution concepts. For example, the Groves-Clarke mechanism originated in literature on implementation in dominant strategies.. Given the focus on Bayesian Nash equilibria, we wish highlight the types of choice functions $g : \Theta \to X$ that satisfy the following condition: there exists a mechanism $\langle M, p(\cdot) \rangle$ such that if agents play a Bayesian Nash equilibria to the mechanism then the final outcome corresponds to the policy that would be selected by the choice function, $p(m(\theta)) = g(\theta)$.

From the mechanism designer's perspective the mechanism is instrumental to achieving a particular choice function. If the designer wishes to implement the function $g(\theta)$ in Bayesian Nash strategies then she must select a mechanism $\langle M, p(\cdot) \rangle$ such that in the corresponding game has a Bayesian Nash equilibrium in which the agents use strategies $m_i^*(\theta_i)$ so that $p(m_1^*(\theta_1), m_i^*(\theta_i), ..., m_n^*(\theta_n)) = g(\theta)$. Thus the choice of a mechanism is informed by knowledge of the incentives the mechanism creates, and the designer anticipates how these incentives will shape behavior by anticipating that agents will play equilibrium strategies to the mechanism.

DEFINITION 11.1. *Given a mechanism design problem $\langle \Theta, F(\theta), X, u \rangle$ we say the choice function $g(\theta)$ is implementable in Bayesian Nash strategies if there exists a mechanism $\langle M, p(\cdot) \rangle$ which has a Bayesian Nash equilibrium $m_i^*(\cdot)$ in which $p(m_1^*(\theta_1), m_i^*(\theta_i), ..., m_n^*(\theta_n)) = g(\theta)$ for every $\theta \in \Theta$.*

If a game is dominance solvable, the surviving strategy profile will be a Bayesian Nash equilibrium. This means that given a mechanism design problem the set of choice functions which are implementable in dominant strategies is a subset of the set of choice functions which are implementable in Bayesian Nash strategies.

Our first result, the so-called *revelation principal,* dramatically simplifies the search for choice functions which are implementable by allowing us to focus on a smaller set of possible mechanisms. *Direct mechanisms* are mechanisms in which the agents are asked to report their types directly. Thus, direct mechanisms have $M_i = \Theta_i$. The revelation principal says that if there exists a mechanisms to implement choice function $g(\theta)$ then there must exist direct mechanism that implements $g(\theta)$. This powerful result tells us that we need not consider all possible mechanisms just the direct ones.

While revelation principal is quite general, its proof is very straightforward. It can also be extended to a very large class of equilibrium concepts.

PROPOSITION 11.1. *(Revelation Principal in Bayesian Nash strategies). Given a mechanism design problem $\langle \Theta, F(\theta), u \rangle$ , if the choice function $g(\theta)$ is implementable in Bayesian Nash strategies then there exists a direct mechanism $\langle \Theta, p(\cdot) \rangle$ that implements $g(\theta)$ in Bayesian Nash strategies.*

PROOF. We begin by assuming that there is a non-direct mechanism $\langle M, p'(\cdot) \rangle$ that implements $g(\theta)$ in Bayesian Nash strategies. We use this mechanism to construct a direct mechanism that also implements the choice function $g(\theta)$. Let $s_i(\theta_i)$ denote the strategy that player $i$ deploys in one of the Bayesian Nash equilibria to the game induced by $\langle M, p'(\cdot) \rangle$ . Consider the direct mechanism in which agents are asked to announce messages $m_i \in \Theta_i$ and then the policy is chosen by the function $p(\theta) = p'(s_1(\theta_1), ..., s_i(\theta_i), ..., s_n(\theta_n))$. We need only verify that under the direct mechanism, truthful announcements of $m_i(\theta_i) = \theta_i$ form a Bayesian Nash equilibria. First, suppose that all agents $N \backslash i$ are playing truthful strategies. If agent $i$ also uses a truthful strategy then the final outcome will be $g(\theta)$. Now suppose that there is a desirable deviation $m_i' \neq \theta_i'$ in the direct mechanism for agent $i$ with type is $\theta_i' \in \Theta_i$. If agent $i$ can select $m_i'$ in the direct mechanism, it must be the case that $m_i' \in \Theta_i$. However, since $s_i(\cdot) : \Theta_i \rightarrow M_i$ in the game induced by $\langle M, p'(\cdot) \rangle$,it must be the case that $s_i(m') \in M_i$ exists. This implies that

$$\int_{\theta_{-i} \in \Theta_{-i}} u_i(p'(s_1(\theta_1), ..., s_i(m'), ..., s_n(\theta_n)) dF(\theta_{-i} \mid \theta_i) >$$

$$\int_{\theta_{-i} \in \Theta_{-i}} u_i(p'(s_1(\theta_1), ..., s_i(\theta_i'), ..., s_n(\theta_n)) dF(\theta_{-i} \mid \theta_i)$$

This expression contradicts the fact that $s_i(\cdot)$ is a best response in the Bayesian game induced by $\langle M, p'(\cdot)\rangle$. Thus we have established the result.                                                                       □

One cautionary note is in order. We focus only on the existence of a mechanism which implements a choice function as the outcome of Bayesian Nash Equilibrium. Clearly, there may be other equilibria to the Bayesian game mechanism that result in different collective choices.

Given the revelation principal, the question of whether a particular choice function is implementable can be answered by focusing on truthful direct mechanisms. If a choice function cannot be implemented by such a mechanism, it cannot be implemented by any mechanism.

We now consider a few examples before returning to the development of the theory.

## 2. Applications

**2.1. Polling.** Suppose that there are $N = \{1, 2, .n\}$ ($n$ odd –not surprisingly) voters with symmetric single-peaked preferences on $\mathbb{R}^1$. The mechanism designer does not know the agent's ideal points but she may ask each voter a question and then select a policy $x \in \mathbb{R}^1$. Each ideal point $\theta_i$ is drawn from a distribution $F(\cdot)$ on $\mathbb{R}^1$. The natural question to ask is whether there are any mechanisms that induce the agents to reveal their ideal points in dominant strategies. The answer is yes. Consider the mechanism that asks each agent to announce their ideal point $m_i \in \mathbb{R}^1$ and then chooses policy equal to the median announcement, $x(m) = median(m)$. To see that truthful response is a best response, consider a respondent with ideal point $y_i$. Let $x(m_i, m_{-i})$ denote the median of the profile of messages that includes $m_i$ and the responses of the other $n-1$ respondents. Further, let $\underline{x}(m_{-i})$ and $\overline{x}(m_{-i})$ be the lower and upper of $median(m_{-i})$.[2]

Suppose that agent $i$ reports $y_i$. If $y_i < \underline{x}(m_{-i})$, then $\underline{x}(m_{-i})$ becomes the median report so that $x(y_i, m_{-i}) = \underline{x}(m_{-i})$. Similarly, if $y_i > \overline{x}(m_{-i})$, $x(y_i, m_{-i}) = \overline{x}(m_{-i})$. Finally, if $y_i \in [\underline{x}(m_{-i}), \overline{x}(m_{-i})]$, $y_i$ is the median report so that $x(y_i, m_{-i}) = y_i$. Thus, the deviation to $y_i$ can result in only three types of outcomes $\underline{x}(m_{-i}), \overline{x}(m_{-i})$, and $y_i \in [\underline{x}(m_{-i}), \overline{x}(m_{-i})]$.

Clearly, the defection cannot pay in the state of the world where if $\theta_i \in [\underline{x}(m_{-i}), \overline{x}(m_{-i})]$ since she obtains her ideal point by reporting $m_i = \theta_i$. So suppose that $\theta_i < \underline{x}(m_{-i})$. Now agent $i$'s best feasible

---

[2]Recall that if a set of real numbers has an even number of distinct elements, it has two medians.

outcome is $\underline{x}(m_{-i})$ which can be obtained by any message less that $\underline{x}(m_{-i})$ which includes $\theta_i$. Similarly, if $\theta_i > \overline{x}(m_{-i})$, announcing $m_i = \theta_i$ weakly maximizes her utility. We have thus seen that regardless of the responses of the other players, agent $i$'s best strategy is to announce $m_i = y_i$.

Relating this example back to the Gibbard-Satherwaite Theorem of chapter 4, Herve Moulin (1980) obtained a mechanism design analogue to results about the existence of preference aggregation rules that do not violate the Arrow's conditions when preference are single-peaked. Moulin focused on mechanisms in which respondents are asked to announce a number (interpreted as their ideal point) and set up a small number of criterion, anonymity and efficiency. The first condition anonymity requires, simply, that the mechanism treat individuals identically.

DEFINITION 11.2. *A mechanism* $g : \mathbb{R}^n \to X$ *is anonymous if for any permutation* $\pi : N \to N$, $g(m_1, .., m_i, .., m_n) = g(m_{\pi(1)}, ..., m_{\pi(i)}, ..., m_{\pi(n)})$.

The efficiency condition is essentially identical to Arrow's Pareto condition.

DEFINITION 11.3. *A mechanism* $g : \mathbb{R}^n \to X$ *is efficient, if for every profile of types* $y = (y_1, ..., y_n)$ $g(y)$ *is Pareto efficient, that is there is not some other policy* $z \neq g(y)$ *such that every agent weakly prefers* $y$ *to* $g(y)$ *and some agent strictly prefers* $y$ *to* $g(y)$.

Moulin showed the following result.

PROPOSITION 11.2. *If preferences are single-peaked then every efficient, and anonymous, strategy proof mechanism is of the following form: Take the announcements* $m_1, m_2, ....m_n$ *and add* $k$ *fixed numbers* $a_1, ..., a_k$, *and select the median of this longer list,* $x(m, a)$.

We have already sketched out the argument for why such a mechanism is strategy-proof. The proof of Moulin's result is left as an exercise.

**2.2. Auctions.** Because of the role of auctions in allocating everything from broadcast spectra to chatchas on E-Bay, economists have developed large body of theory both on how optimally structure auctions to generate maximal revenue and to achieve allocative efficiency. While political scientists are generally not concerned with those issues, we present some of the basics of auction theory for a couple of reasons. First, several important aspects of politics can be thought of as particular types of auctions. Second, auction theory demonstrates how mechanism design can be used to study the choice of institutions.

However, it is not that auction theory has been ignored in political science. A classic application of auction theory are models of favor-buying in which competing interest groups offer bribes to secure public contracts or policy concessions from politicians. Recently one of us has argued for modeling electoral competition as a form of an auction (Meirowitz 2004).

The standard auction problem involves a seller (agent 0) and a population of potential bidders, $N = \{1, .., n\}$. Each bidder places a valuation of $\theta_i \in \mathbb{R}^1_+$ on the item to be sold. These valuations are the private information of each buyer. We assume for simplicity that the common prior is that bidder valuations are independent and identical draws from a twice differentiable distribution function $F(\cdot)$. The utility to a bidder $i$ that has valuation $\theta_i$ and wins the item by paying $b_i$ is just $\theta_i - b_i$. The payoff to losing bidders is 0. We begin with a few commonly studied auction mechanisms.

2.2.1. *Second Price and Ascending Price Auctions.* Consider two auction designs. In a second price auction, each participant submits a sealed bid $b_i$ and the one who bids the highest wins the object. However, the winner the amount equal to the second highest bid. In an ascending price auction, all participants begin with their placards up, and the auctioneer announces an ascending sequence of prices, $10, $11, $12,..... When a price is announced causing the second to last placard to fall, the item is sold to that bidder with the upright placard at the announced price.

Both of these auctions are commonly used and the experience of bidding under these schemes may seem quite different. Nonetheless, from a game-theoretic perspective the mechanisms are identical. Note that the auctions induce different games (one is a Bayesian game with simultaneous moves and the other is a Bayesian game with sequential moves). Nonetheless, the incentives are identical. To see this imagine that each bidder entered their valuation $\theta_i$ into a computer and the computer did two things: (1) it submitted the valuations into a second price auction, i.e. it uses strategy $b_i = \theta_i$. (2) it has a robot play the following strategy in an ascending price auction -hold the placard up until the announced price exceeds the valuation, $\theta_i$. In both cases the bidder with the highest valuation would win the item and pay the second highest price.

It is easy to see that in the second price auction bidding $b_i = \theta_i$ is a dominant strategy. Consider any other bid $b_i \neq \theta_i$. Let $b^{\max}_{-i} = \max_{j \neq i} b_j$ denote the highest bid submitted by all the other agents. There are three possibilities, either $b^{\max}_{-i} = \theta_i$, $b^{\max}_{-i} < \theta_i$ or $b^{\max}_{-i} > \theta_i$. In the first case use of the strategy $b_i = \theta_i$ will result in a ties (which

has expected utility equal to 0). The strategy $b_i > \theta_i$ will result in victory but at the price $b_{-i}^{\max} = \theta_i$ which results in utility of 0. Finally under the strategy $b_i < \theta_i$ you do not win the item and receive utility of 0. In the second case, under the strategy of $b_i = \theta_i$ you win and pay $b_{-i}^{\max}$ receiving payoff $\theta_i - b_{-i}^{\max} > 0$. Under the strategy $b_i > \theta_i$ you win and also pay $b_{-i}^{\max}$ again receiving payoff $\theta_i - b_{-i}^{\max}$. Under the strategy of $b_i < \theta_i$ you do not win and receive a payoff of 0. In the last case, under the strategy of $b_i = \theta_i$ you do not win and receive payoff of 0. Under the strategy $b_i > \theta_i$ you either win and pay $b_{-i}^{\max}$ for an item that you only value at $\theta_i < b_{-i}^{\max}$ and thus you receive a negative payoff, or you do not win and receive a payoff of 0. Thus the strategy of $b_i = \theta_i$ does at least as well as any other strategy and avoids two types of regret: not winning an item at a price you would be willing to pay and winning an item at a price that you do not want to pay.

Similarly, lowering the placard once the price exceeds $\theta_i$ is a dominate strategy in the ascending price auction. Is there a reason to lower the placard before the price reaches $\theta_i$? Doing so insures that you do not win the item, thus this deviation results in a payoff of 0. Thus, early resignation cannot do better than the conjectured strategy of lowering the placard at a price equal to $\theta_i$. Moreover, if it is the case that all other agents will lower the placard by a price $p' < \theta_i$ than the conjectured strategy would result in victory at a price of $p'$ and utility $\theta_i - p_i$. Now consider the potential consequences of deviating to a strategy of keeping the placard up after the price exceeds $\theta_i$. In this case either you are not the winner and you receive a utility of 0 or you are the winner and you pay a price $p'' > \theta_i$ to win an object worth $\theta_i$ resulting in a negative payoff.

We now consider different types of auctions.

2.2.2. *First Price and Descending Price Auctions\*.* In a first price auction, agents simultaneously submit bids and the highest bidder wins the object and pays her bid. In a descending price auction each bidder has a placard and the auctioneer begins with a high price and lowers it as the auction progresses $100, $99, $98,.... The bidder that raises her placard first wins the item and pays the most recent price. As before the simultaneous move and extensive form versions of the game are similar. We focus only on the analysis of the first price auction, leaving as an exercise the construction of a subgame perfect equilibrium to the descending price auction. We follow Krishna (2002) and present an informal derivation of the equilibrium strategies. We begin with a conjecture about equilibrium strategies and verify that the conjecture turns out to be consistent with a PBE.

Suppose that bidders $j \neq i$/have strategies that are represented by the differentiable and strictly increasing function $b(\theta)$. Since a bidders expected utility is 0 if she pays her valuation or loses given the random variable $b_{-i}^{\max}$ the expected utility to a bidder with type $\theta_i$ from using strategy $b_i$ is

$$pr(b_{-i}^{\max} < b_i)\left[\theta_i - b_i\right].$$

Since we have assumed that the other bidders use the strategy $b(\theta)$, the term $pr(b_{-i}^{\max} < b_i)$ is just the probability that $n - 1$ draws of $\theta_j$ from $F(\cdot)$ each have values lower than $b^{-1}(b_i)$. Given the independence of types the expression for this probability is just $F(b^{-1}(b_i))^{n-1}$. Accordingly, the expected utility is

$$F(b^{-1}(b_i))^{n-1}\left[\theta_i - b_i\right].$$

An optimal $b_i$ must solve the first order necessary condition

$$(n-1)F(b^{-1}(b_i))^{n-2}f(b^{-1}(b_i))\frac{db^{-1}(b_i)}{db_i}\left[\theta_i - b_i\right] - F(b^{-1}(b_i))^{n-1} = 0$$

In the conjectured equilibrium $b(\theta_i) = b_i$ and thus we have

$$(n-1)F(\theta)^{n-2}f(\theta)\theta = \frac{db(\theta)}{d\theta}F(\theta)^{n-1} + (n-1)F(\theta)^{n-2}f(\theta)b(\theta)$$

Note that the right hand side is $\frac{d}{d\theta}\left(F(\theta)^{n-1}b(\theta)\right)$ allows us to re-express the above as

$$\frac{d}{d\theta}\left(F(\theta)^{n-1}b(\theta)\right) = (n-1)F(\theta)^{n-2}f(\theta)\theta$$

The first theorem of calculus allows us to re-express this as

$$(11.1) \qquad b(\theta) = \frac{\int_0^\theta \theta'(n-1)F(\theta')^{n-2}f(\theta')d\theta'}{F(\theta)^{n-1}}$$

which is the conditional expectation of $b_{-i}^{\max}$ given that $b_{-i}^{\max} < \theta$. Verifying that this is in fact an equilibrium amounts to determining whether the local extrema characterized by equation 11.1 is in fact a global maximum. We leave this as an exercise.

2.2.3. *The Revenue Equivalence Principle\**. We see that the first price and second price auctions result in different strategies. Given this, a natural question is which auction the seller would prefer to use. This question is answered by a fundamental result in auction theory, the *revelation equivalence principal* (Riley and Samuelson 1981; Myerson 1981). This principal guarantees that when the types are independent the expected revenue of an auction depends only on the probability that a seller wins as a function of her type $\theta_i$ and the utility realized by the seller from the lowest possible type. An immediate consequence of

this result is that, assuming independent evaluations, all of the above auctions yield the same expected revenue to the seller. In this section we develop this logic.

Returning to the notation from the beginning of this chapter, let $M = \Theta$ denote the space of possible messages. We assume that each of the $n$ players have a valuation which is an independent draw from the differentiable distribution $F(\cdot)$ on $\Theta$ with density $f(\cdot)$. For convenience we assume that $\Theta$ is an interval of the form $[0, k]$. By $\Delta(N)$ we denote the set of lotteries over the bidders $N$. Let $w(m_1, ...., m_n) :$ $M^n \to \Delta(N)$ and $t(m_1, ...., m_n) : M^n \to \mathbb{R}^n_+$ denote a mechanism which specifies for each profile $m$ of messages a lottery over the identity of the winner and a profile of transfers. Thus, the first price auction has winner

$$w(m) = \begin{cases} 1 \text{ if } i = \arg\max\{m_j\} \\ \quad 0 \text{ otherwise} \end{cases}$$

and transfers

$$t_i(m) = \begin{cases} m_i \text{ if } i = \arg\max\{m_j\} \\ \quad 0 \text{ otherwise.} \end{cases}$$

An all-pay, first-price auction (which we analyze below) has $t_i(m) = m_i$. The subject of the revenue equivalence result is that the expected revenue generated by Bayesian Nash play in the auction. Standard auctions have the "winner takes it with certainty" mapping $g(m)$ defined above. Given an increasing bid function $b(\theta_i)$ the expected revenue is $ER = \sum_{i=1}^{n} \int_0^k t_i(b(\theta_i)) dF(\theta_i)$. We can now state the result.

THEOREM 11.1. *Assume valuations are independently and identically distributed. In every standard auction, every symmetric and increasing equilibrium in which the expected payment of a bidder with value 0 is 0 yields the same value of $ER$.*

PROOF. We consider a fixed standard auction, and symmetric and increasing equilibrium characterized by the function $b(\cdot)$. Let $t^+(\theta_i)$ denote the expected payment of a bidder with value $\theta_i$. Suppose that bidder $i$'s valuation is $\theta'_i$ and consider the bid $z = b(\theta''_i)$. When all other bidders use the strategy $b(\cdot)$ the expected utility of bid $z$ to bidder $i$ is

$$\theta'_i F(\theta''_i)^{n-1} - t^+(\theta''_i)$$

Differentiating this expected utility with respect to the type $\theta''_i$ yields the first order condition

$$(n-1)\theta'_i f(\theta''_i)^{n-2} = \frac{\partial t^+(z)}{\partial z} \Big|_{z = \theta''_i}$$

Since $b(\cdot)$ is an equilibrium it must be the case that $z = b(\theta_i'') = b(\theta_i')$ solves this condition, yielding

$$(n-1)\theta_i' f(\theta_i')^{n-2} = \frac{\partial t^+(z)}{\partial z}\Big|_{z=\theta_i'}\ .$$

But this differential equation yields the solution

$$t^+(\theta_i') = t^+(0) + \int_0^{\theta_i'} (n-1)\theta_i' f(\theta_i')^{n-2} d\theta_i'.$$

Thus under the assumption that $t^+(0) = 0$ we have $t^+(\theta_i') = \int_0^{\theta_i'}(n-1)\theta_i' f(\theta_i')^{n-2}d\theta_i'$ which does not depend on the details of the auction. Since

$$ER = \sum_{i=1}^{n}\int_0^k t_i(b(\theta_i))dF(\theta_i) = n\int t^+(\theta_i'))f(\theta_i')d\theta_i'.$$

this completes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

2.2.4. *Contests and All-Pay Auctions\**. An all-pay auction requires that contestants pay their bids regardless of whether they win. Many contests take this form. Two examples of interest in political science are models of interest group influence in which special interests pay bribes and the group that makes the largest payment receives a policy and electoral campaigns in which candidates compete for office by mounting costly campaigns. We develop this second application in some detail. We use the revenue equivalence principle to reach some interesting conclusions.

Consider a pool of candidates, $C = \{1, 2, ...., c\}$. Each candidate has a non-negative real-valued type $\theta_p \in \mathbb{R}_+^1$ that corresponds to their overall efficiency at governing. The social goal is to select the most efficient candidate $p^* = \arg\max_{p\in C}\theta_p$. We assume that the quality types $\boldsymbol{\theta} = (\theta_1, ...., \theta_c)$ are private information – only candidate $p$ knows $\theta_p$. For simplicity we assume that types are independent and identical draws from the distribution $F(\cdot)$ with density function $f(\cdot)$. Let $M_c = \max\{\theta_j\}_{j\neq 1}$ denote the maximum type from $c-1$ draws of $\theta$. Given $\theta_p$ the probability that $M_c \le \theta_p$ is $F_c(\theta_p) \equiv F(\theta_p)^{c-1}$. Differentiating yields the density of $M_c$, $f_c(\theta_p) = (c-1)F(\theta_p)^{c-2}f(\theta_p)$.

Our baseline model of an election involves two periods. In the first period each candidate simultaneously selects a level of campaign effort $a_p$. In the second period the candidate with the highest level of effort wins office. The key assumption is that a candidate's cost of campaign effort is decreasing in her governing efficiency. That is, we assume that if candidate 1 has a lower cost of campaign effort than

candidate 2, candidate 1 is likely to be a more effective leader than is
candidate 2. For simplicity, we consider the case where campaigning
costs are inversely proportional to efficiency, thus $\beta_p = \frac{1}{\theta_p}$. Candidate
payoffs are then

$$(11.2) \qquad Eu(a_p, \theta_p) = \begin{cases} 1 - \frac{a_p}{\theta_p} \text{ if } a_p > \max_{j \neq p}\{a_j\} \\ -\frac{a_p}{\theta_p} \text{ if } a_p < \max_{j \neq p}\{a_j\} \\ \frac{1}{\#\{j:a_j=a_p\}} - \frac{a_p}{\theta_p} \text{ if } a_p = \max_{j \neq p}\{a_j\} \end{cases}$$

where the last line of the expression follows from assuming that ran-
domization is used in case of ties. Note that multiplying each candi-
dates' utility function by $\theta_p$ translates this payoff function into that of
a standard all-pay auction where the prize has value $\theta_p$ and bids have
unitary cost.[3] We now characterize a symmetric pure strategy equi-
librium in which a candidate's effort level is a strictly increasing and
differentiable function of her efficiency type $\theta_i$. If each player is using
a strictly increasing and differentiable strategy $\alpha(\theta)$ then candidate $p$
with type $\theta_p$ that selects effort $a_p$ has expected utility

$$(11.3) \qquad \pi(a_p, \theta_p) = F(\alpha^{-1}(a_p))^{c-1} - \frac{a_p}{\theta_p}$$

In order for $a_p$ to be optimal it must solve the first order condition

$$(11.4) \qquad (c-1)F(\alpha^{-1}(a_p))^{c-2}f(\alpha^{-1}(a_p))\frac{d\alpha^{-1}}{da_p} = \frac{1}{\theta_p}.$$

In a symmetric equilibrium it must be the case that $\alpha(\theta_p) = a_p$, so that
$\alpha^{-1}(a_p) = \theta_p$. Thus the condition reduces to

$$(11.5) \qquad \frac{d\alpha^{-1}}{da_p} = \frac{1}{\theta_p(c-1)F(\theta_p)^{c-2}f(\theta_p)}.$$

This means that

$$(11.6) \qquad \frac{d\alpha}{d\theta} = \theta_p(c-1)F(\theta_p)^{c-2}f(\theta_p).$$

Integration yields the required solution,

$$(11.7) \qquad \alpha(\theta) = \int_0^\theta x(c-1)F(x)^{c-2}f(x)dx$$

$$= \int_0^\theta xf_c(x)dx.$$

This function is strictly increasing in $\theta$ so it remains only to verify that
this solution satisfies the sufficient second order condition. This result

---

[3]See for example Milgrom and Weber (1985) for results on standard auction
models.

follows form Theorem 2 of Krishna and Morgan (1997) so we do not reproduce the proof.

PROPOSITION 11.3. *In the basic model a symmetric equilibrium in which accumulations $\alpha(\theta)$ are strictly increasing in efficiency exists. In this equilibrium the best candidate $p^* = \arg\max_{p \in C}(\theta_p)$ is chosen with probability one.*

A few comments are in order. First $a(\theta_i) < \theta_i$.

CONCLUSION 1. *The winning candidate achieves a strictly positive payoff.*

Second, as a function of $c$ the effort strategies have derivative

$$\frac{\partial \alpha(\theta)}{\partial c} = \int_0^\theta x(c-1)\ln(F(x))F(x)^{c-2}f(x)dx$$

which is negative as $F(x) \leq 1$. Thus for $c < c'$ we have $\alpha_c(\theta) > \alpha_{c+1}(\theta)$ where $\alpha_n(\theta)$ denotes the equilibrium effort function when $c = n$.

CONCLUSION 2. *Candidate efforts are decreasing in the number of candidates.*

Third, since the model is equivalent to a first-price all-pay auction with independent values and $\alpha(0) = 0$, the revenue equivalence principle proven above implies that he expected total amount of effort

$$A_c = c \int_0^\infty \alpha(x)f(x)dx$$

must be the same as the expected payment of the winner in a second price auction in which player values are drawn from $f(\cdot)$. This value is just the expected value of of the second highest of $c$ draws from $f(\cdot)$ which is

$$\int_0^\infty xc(c-1)(1-F(x))F(x)^{c-2}f(x)dx.$$

which is increasing in $c$.

CONCLUSION 3. *The total effort is increasing in the number of candidates*

It is reasonable to think that the social objective is maximization of $\theta_{p_c^*}$. If in addition the actual effort of campaigning is viewed as wasteful then it is natural to think that the social objective involves a trade-off between increasing $\mathbb{E}\theta_{p_c^*}$ and decreasing $\mathbb{E}\sum_{p=1}^c \frac{\alpha(\theta_p)}{\theta_p}$ where $\mathbb{E}$ is the expectations operator. Letting $v(\cdot)$ and $u(\cdot)$ denote twice

differentiable increasing functions with $v'' < 0$ and $u'' > 0$, it is natural to consider a social welfare function of the form

$$V_c = \mathbb{E}v(\theta_{p_c^*}) - \mathbb{E}u\left(\sum_{p=1}^{c} \frac{\alpha(\theta_p)}{\theta_p}\right).$$

While the all-pay auction aspect of campaigns does an excellent job selecting the best quality candidate, unlike a typical auction where the bidding is viewed as revenue to the seller, in the campaign context the costs $\sum_{p=1}^{c} \frac{\alpha(\theta_p)}{\theta_p}$ are not desirable. This perspective motivates a natural question. How does one select good candidates while not inducing excessive wasteful campaigning? In the next section we demonstrate that one answer to this question involves imperfect voting–a phenomena that many believe is empirically reasonable.

The baseline model assumes that elections are perfect screening devices selecting the candidate that selects the highest level of $a$. In reality voting is imperfect, and candidates likely recognize this. In this section we consider elections with a random decision. For simplicity assume that the candidate with the highest level $a$ wins with probability $q + \frac{1-q}{c}$ and that the remaining candidates win office with probability $\frac{1-q}{c}$. To be sure a large number of other models of probabilistic selection can be incorporated. In this setting the candidate payoffs are then

$$(11.8) \quad Eu(a_p, \theta_p) = \begin{cases} q + \frac{1-q}{c} - \frac{a_p}{\theta_p} & \text{if } a_p > \max_{j \neq p}\{a_j\} \\ \frac{1-q}{c} - \frac{a_p}{\theta_p} & \text{if } a_p < \max_{j \neq p}\{a_j\} \\ \frac{q}{\#\{j:a_j=a_p\}} + \frac{1-q}{c} - \frac{a_p}{\theta_p} & \text{if } a_p = \max_{j \neq p}\{a_j\} \end{cases}$$

If each player is using a strictly increasing and differentiable strategy $\alpha(\theta)$ then candidate $p$ with type $\theta_p$ that selects accumulation $a_p$ has expected utility

$$(11.9) \qquad \pi(a_p, \theta_p, q) = qF(\alpha^{-1}(a_p))^{c-1} - \frac{a_p}{\theta_p} + \frac{1-q}{c}$$

In order for $a_p$ to be optimal it must solve the first order condition

$$(11.10) \qquad q(c-1)F(\alpha^{-1}(a_p))^{c-2}f(\alpha^{-1}(a_p))\frac{d\alpha^{-1}}{da_p} = \frac{1}{\theta_p}.$$

In a symmetric equilibrium it must be the case that $\alpha(\theta_p) = a_p$, so that $\alpha^{-1}(a_p) = \theta_p$. Thus the condition reduces to

(11.11) $$\frac{d\alpha^{-1}}{da_p} = \frac{1}{\theta_p q(c-1)F(\theta_p)^{c-2}f(\theta_p)}.$$

This means that

(11.12) $$\frac{d\alpha}{d\theta} = \theta_p q(c-1)F(\theta_p)^{c-2}f(\theta_p).$$

Integration yields the required solution,

(11.13) $$\alpha(\theta; q) = \int_0^\theta xq(c-1)F(x)^{c-2}f(x)dx$$

$$= q \int_0^\theta x f_c(x)dx$$

(11.14) $$= q\alpha(\theta; 1).$$

Accordingly, when the election is imperfect, each candidate shades down her effort and the efforts are proportional to the imperfection parameter $q$.

For the remainder of this section we will focus on two candidate contests and consider the optimal level $q \in [0, 1]$. In this case the strategies are

$$\alpha(\theta; q) = q \int_0^\theta x f(x)dx.$$

The social welfare associated with $q$ is

$$V(q) = q \int_0^\infty v(x)2F(x)f(x)dx + (1-q) \int_0^\infty v(x)2(1-F(x))f(x)dx$$

$$- \int_0^\infty u\left(\frac{2q}{\theta} \int_0^\theta x f(x)dx\right) df(\theta)$$

The optimal value $q$ satisfies the first order condition

$$\int_0^\infty v(x)F(x)f(x)dx - \int_0^\infty v(x)(1-F(x))f(x)dx =$$

$$\int_0^\infty \left(\left[\frac{1}{\theta}\int_0^\theta x f(x)dx\right]u'(\frac{2q}{\theta}\int_0^\theta x f(x)dx)\right) f(\theta)d\theta$$

The optimal value $q^*$ is less than 1 if

$$\int\limits_0^\infty v(x)F(x)f(x)dx - \int\limits_0^\infty v(x)(1 - F(x))f(x)dx <$$

$$\int_0^\infty \left( \left[ \frac{1}{\theta} \int_0^\theta xf(x)dx \right] u'(\frac{2}{\theta} \int_0^\theta xf(x)dx) \right) f(\theta)d\theta$$

and greater than 0 if

$$\int\limits_0^\infty v(x)F(x)f(x)dx - \int\limits_0^\infty v(x)(1 - F(x))f(x)dx >$$

$$\int_0^\infty \left( \left[ \frac{2}{\theta} \int_0^\theta xf(x)dx \right] u'(0) \right) f(\theta)d\theta$$

Accordingly, if $v(\cdot)$ is relatively flat and $u(\cdot)$ is relatively steep then then imperfect elections are efficiency enhancing. When $v$ and $u$ are linear the first order condition does not depend on $q$ and the optimal $q$ is either 1 or 0, as it is either efficient to maximize the probability that the best candidate serves or to minimize the campaigning costs depending on which effect has a larger marginal impact on social welfare.

## 3. Incentive Compatibility and Individual Rationality

In the contribution and polling examples, we considered mechanisms for which truthfully reporting one's type was a best response. While the revelation principal tells us that focusing on direct mechanisms will not limit the choice functions that we can implement, it is silent about what types of choice functions are implementable. In the next two sections we consider this question. From the revelation principal, we know that we need only focus on direct mechanisms. Consequently, the question becomes "which types of mechanisms will induce agents to be truthful?" In a direct mechanism agents are asked to report their private information type, and then this information is used to make a decision. Incentive compatibility is the requirement that given the mechanism and the belief that all other agents are being truthful, agent $i$ prefers being truthful to lying.

We begin with the general problem. Consider a problem with agents $N$, choice space $X$, type space $\Theta$, prior joint density function over types $f(\theta)$, and state-contingent utility functions $u_i(x, \theta)$ for each $i \in N$. A direct mechanism is a mapping $p(\theta) : \Theta \to X$. In order for there to be a Bayesian Nash equilibrium with truthful strategies to the

game that the mechanism induces the following Incentive compatibility condition must be true:

DEFINITION 11.4. *(Incentive Compatibility condition) For every $i \in N$ and every $\theta_i \in \Theta_i$*
(IC)
$$\int_{\Theta_{-i}} u_i(p(\theta_i, \theta_{-i}), \theta_i) f_{-i}(\theta_{-i}) d\theta_{-i} \geq \int_{\Theta_{-i}} u_i(p(\theta'_i, \theta_{-i}), \theta_i) f_{-i}(\theta_{-i}) d\theta_{-i}$$
*for every $\theta'_i \in \Theta_i$.*

Informally, the condition requires that truthful messages are a best response if everyone else is using truthful messages. We now focus on gaining additional leverage on the types of mechanisms $p(\cdot)$ that can satisfy this condition. The easiest way to ensure that the IC is satisfied is to make $\int_{\Theta_{-i}} u_i(p(m_i, \theta_{-i}), \theta_i) f_{-i}(\theta_{-i}) d\theta_{-i}$ constant in $m_i$.

Consider the following example. Let $\omega$ is a random state variable that effects the players payoffs from various policies. Assume that $n \geq 3$ agents observe $\omega$ so that $\theta_i = \omega$ with probability. The mechanism designer's job can implement a choice function $x(\omega)$ which maps a profile of messages into policies $x$ in the following simple type of mechanism
$$p(\theta) = \begin{cases} x(\omega') \text{ if } \#\{j \in N : \theta_j = \omega'\} \geq n-1 \\ x(w^*) \text{ otherwise} \end{cases}$$
where $\#\{j \in N : \theta_j = \omega'\}$ denotes the number of individuals announcing $\theta_j = \omega'$ and $w^*$ is an arbitrary value of $\omega$. Since a single defection does not alter the policy choice, this mechanism satisfies incentive compatibility.

The mechanism in the above example, however, does not lead to a strong incentive to be truthful. In the auction example incentive compatibility is satisfied by a second price auction as (ignoring ties),
$$u_i(p(m_i, \theta_{-i}), \theta_i) = \begin{cases} \theta_i - \max_j \theta_j \text{ if } m_i > \max_j \theta_j \\ 0 \text{ otherwise} \end{cases}$$
is constant in $m_i$ if $m_i > \max_j \theta_j$ and constant in $m_i$ for $m_i < \max_j \theta_j$. The non-constant part of the function jumps from 0 to $\theta_i - \max_j \theta_j$ which is positive only if $m_i > \max_j \theta_j$. In contrast the first price auction is not incentive compatible.

To generate some intuition for incentive compatibility conditions we return to the coffee machine problem. Consider an arbitrary transfer schedule $t_i(m_i.m_{-i})$ which maps a message profile into an amount charged to member. Let $p(m)$ be a policy function which maps the message profile into the probability that coffee maker is purchased..

Given this transfer schedule, the expected utility to department member $i$ from announcement $m_i$ is

$$Eu(m_i, \theta_{-i}) \equiv \int \left[ \theta_i p(m_i, \theta_{-i}) - t_i(m_i, \theta_{-i}) \right] f_{-i}(\theta_{-i}) d\theta_{-i}$$

If we assume that $p$ and $t$ are differentiable functions, incentive compatibility requires that $m_i = \theta_i$ maximize $Eu(m_i, \theta_{-i})$. Therefore, we require that a *local incentive compatibility condition* based on the first-order condition for maximization:

$$\frac{\partial Eu(m_i, \theta_i)}{\partial m_i} \Big|_{m_i = \theta_i} = 0.$$

Interchanging the order of integration and differentiation in our example leads to the condition

$$\theta_i \int \frac{\partial p(m_i . \theta_{-i})}{\partial m_i} f_{-i}(\theta_{-i}) d\theta_{-i} \Big|_{m_i = \theta_i} = \int \frac{\partial t(m_i . \theta_{-i})}{\partial m_i} f_{-i}(\theta_{-i}) d\theta_{-i} \Big|_{m_i = \theta_i}$$

Intuitively, this means that an incentive compatible mechanism requires that the expected decrease in transfers associated with a slight under reporting of $\theta_i$ is exactly offset by the reduction in the expected likelihood of coffee maker purchase–weighted by the value of the purchase $\theta_i$.

While incentive compatibility requires that players be willing to reveal their private information, we also require that participants be willing to play the game. Analysis of these constraints is somewhat simpler and *ad hoc*. The key intuition is players will rationally participate only if she gets a higher payoff in equilibrium that her payoff from not participating. In some settings, these constraints can be trivial as it is reasonable to assume that agents have no choice but to participate. In settings where the constraints are taken seriously, it is necessary to be explicit about the value to players of not participating. Formally,

DEFINITION 11.5. *(Individual Rationality constraint) if the value to a player from not participating in the mechanism is give by $v_i(\theta, \theta_{-i})$ and the value to being truthful in a direct mechanism is $u_i(\theta_i, \theta_{-i})$ then*

$$\int u_i(\theta_i, \theta_{-i}) dF(\theta_{-i}) \geq \int v_i(\theta_i, \theta_{-i}) dF(\theta_{-i})$$

*for all $\theta_i \in \Theta_i$ for each $i \in N$.*

## 4. Constrained Mechanism Design

Classical mechanism design allows the planner to commit to one of a large number of mechanisms. Incentive compatibility, and potentially individual rationality are the only constraints. In many settings of

political science the ability to commit is considered unreasonable. As such the choice of mechanisms is somewhat limited. Restrictions include an inability to provide transfers, and limits on the potential commitments that can be made. The remainder of this section develops some examples that illustrate how political scientists use constrained mechanism design to address problems of institutional design. With few exceptions, the models in political science are within principal-agent paradigm. This means nothing more than the recognition that there is a principal or boss (or several) and an agent or subordinate (or several). The principal would like to have "good" agents do "good" jobs and the agents tend to have a desire to keep the principal in the dark about whether they are of the "good" type or doing a "good" job. The principal is generally assumed to have a limited number of possible instruments which he can control. In the language of mechanism design, the principal is the planner and doing a "good" job or revealing whether one is a "good" type corresponds to selecting appropriate messages in the context of a direct mechanism. The limited number of levers corresponds to constraints on the mechanism that can be enacted. Before turning to some interesting applications, we demonstrate these concepts in a simple delegation problem.

Suppose there is a principal and two agents. Agents have one of two possible types $\theta_i \in \{good, bad\}$. Each agent is a good type with probability $\pi$. In addition if an agent is chosen to perform a task she can devote one of two levels of effort $a_i \in \{high, low\}$. We suppose that $x_i = high$ imposes a cost of $c$ on the agent while $a_i = low$ is costless. The principal must select one of the two agents to perform a task and the chosen agent must decide which effort level to choose. The fundamental question of **principal agent models** (or agency models) can then be phrased as follows. How can the principal design institutions to select a *good* type agent (if one exists) and induce the chosen agent to select *high*? The former aspect, designing institutions to select *good* types is often termed **adverse selection**. The latter aspect, designing institutions to create incentives for *high* effort is termed **moral hazard**.

If the principal could observe a label on the agent's shirt that indicated his type, or if the choice of effort were readily observable then we would say there is no monitoring or observability problem. In this case, things are not too challenging for our principal. She would simply hire only good types and fire them when they give low effort. More interesting are situations where agent types are private information and effort is imperfectly observed. Ferejohn (1986) models accountability in repeated elections as a principal agent problem. To clarify these

concepts, we sketch and work through a simplified version of Ferejohn's model.

Ferejohn focused only on the moral hazard problem, so for now we will ignore the different types. There are two identical parties that can serve in office. Suppose that the in government party selects an effort level, but the voters observe a noisy policy outcome $x \in \{high, low\}$ that has the following form. If $a = high$ then $x$ is $high$ with probability $q > \frac{1}{2}$ and $low$ with probability $1 - q$; if $a = low$ then $x$ is low with probability $q$ and high with probability $1 - q$. Finally, suppose the government party receives a rent $r$ from being in office each period and discounts the future with discount rate $\delta$. The opposition party gets a payoff of zero for each period it is out of power.

The voter gets to select a reelection rule specifying which values of $x$ will result in the decision to reelect the incumbent and which values of $x$ will result in removal. Can she select such a rule that creates incentives for the government to always select $a = high$? One simple rule would be to retain if $x = high$ and replace otherwise. If the voter uses this rule in every period, then the government party faces a simple decision problem. If it selects $a = high$ in the current period then her expected utility is

(11.15) $$r - c + q\delta V_I + (1 - q)\delta V_O$$

where $V_I$ is the value of the game starting next period if the party is in office at the beginning of next period, while $V_O$ is value of the game starting next period if the party is out of office at the beginning of next period.[4] If the government selects $a = low$ then her expected utility is

(11.16) $$r + (1 - q)\delta V_I + q\delta V_O$$

If the voter's rule works so that is the incumbent finds it best to select $a = high$ in every period, then we also have

$$V_I = r - c + q\delta V_I + (1 - q)\delta V_O$$

and

$$V_O = (1 - q)\delta V_I + q\delta V_O$$

which imply that

(11.17) $$V_I = \frac{r - c + cq\delta - qr\delta}{2q\delta^2 - \delta^2 - 2q\delta + 1}.$$

and

(11.18) $$V_O = \frac{r\delta - c\delta + cq\delta - qr\delta}{2q\delta^2 - \delta^2 - 2q\delta + 1}.$$

---

[4]See chapter 3 for a discussion of Bellman equations and value functions.

A version of incentive compatibility (for all governments to select $a = high$,) requires that equation 11.15 is no less than equation 11.16. We can write this incentive constraint as

(11.19) $$V_I - V_o \geq \frac{c}{(2q-1)\delta}.$$

Substituting equations 11.17 and 11.18 leads to the condition

(11.20) $$r - c - r\delta + c\delta \geq \frac{c(2q\delta^2 - \delta^2 - 2q\delta + 1)}{(2q-1)\delta}$$

which is necessary in order for the voters rule to induce $a_i = high$ in every period. In other words if the exogenous parameters satisfy this inequality the moral hazard problem is solvable with the rule that retains only if $x = high$. In this setting the voters do not really need to commit to the rule because the rule represents an equilibrium strategy given the decision rule that the government officers are using as long as the voter prefers $x = high$ to $x = low$. This is true because all governments will select $a_i = high$ and all governments have the same type (i.e. there is no adverse selection problem). Accordingly, if the above condition is satisfied then we have characterized a Nash equilibrium to the game between the two parties and the voter: If in office select $a_i = high$ and when deciding how to vote, retain only if $x = high$.

Now we consider an adverse selection version of the problem. Suppose the government does not select an effort level but her type is private information. For simplicity, we assume that each of the parties has a type that is $high$ with probability $\pi > \frac{1}{2}$. We assume these draws are independent. In each period the voter only observes $x$ which takes the value $high$ with probability $z > \frac{1}{2}$ if the government is a high type and the value of $low$ with probability $1 - z$ if the government is a high type. If the government is a low type then $x = high$ with probability $1 - z$ and $low$ with probability $z$. Once again the voter might try the simple rule: retain if $x = high$ and replace otherwise. However, this rule runs the risk of throwing out quality governments who are temporarily unlucky. A better rule uses information from previous periods. For any finite number of periods $k_i$ in which party $i$ was in office let $h_i$ denote the number of these periods in which $x$ was $high$. It follows that in $k_i - h_i$ of those periods we had $x = low$. The posterior probability that party $i$ is of type $high$ after $h_i$ realizations of $high$ from $k_i$ trials is given by Bayes' rule as

$$\Pr(\theta_i = high \mid k_i, h_i) = \frac{\pi z^{h_i}(1-z)^{k_i - h_i}}{\pi z^{h_i}(1-z)^{k_i-h_i} + (1-\pi)(1-z)^{h_i}z^{k_i-h_i}}$$

An optimal rule is to keep the original incumbent until $\Pr(\theta_i = high \mid k_i, h_i)$ falls below $\pi$, the expected quality of party $j$ before it serves. Then the voter should dump party $j$ as soon as $\Pr(\theta_i = high \mid k_i, h_i) > \Pr(\theta_j = high \mid k_j, h_j)$. Recalling the law of large numbers we see that this rule will eventually select the optimal party (they may both be optimal).

This discussion of agency has been meant only as an introduction to the concepts. We now consider several applications of agency theory to the study of delegation. The basic concepts of trying to create incentives to limit moral hazard and trying to select the best agents are recurrent.

**4.1. A Model of Delegation to Bureaucrats.** One of the key questions in the study of the modern administrative state is the trade-off between political control of an agency and the autonomy that an agency requires to apply its expertise to policy problems. Principal agent theory has been a natural approach to this question. Originally applied in politics to study when and how the U.S. congress delegated rule-making authority to regulatory agencies, principal agent applications have spread to many other political systems and to many other types of political bodies such as political parties and international organization. In this section, we consider a version of the model that Epstein and O'Halloran (1994) applied to the study of statutory delegation in the United States.

Suppose that a legislature $L$ is considering how much authority to delegate to a bureaucratic agency $A$. We assume that the policy space $X$ is single dimensional and a subset of $\mathbb{R}$. Each of the $n$ legislators have quadratic policy preferences with ideal points $l_1 < .... < l_n$. The agency is treated as a unitary actor with quadratic policy preferences and an ideal point $a$.

To motivate why $L$ would delegate policy making authority to $A$ rather than set policy itself, we assume that the members of $L$ are uninformed about the consequences of various policy choices $p \in X$ whereas $A$ is fully informed. We model this uncertainty just as Gilligan and Krehbiel do in their study of legislative committees. We assume that the policy outcome $x$ is a function of the policy choice $p$ and an error term $\varepsilon$. To keep things as simple as possible, let $x = p - \varepsilon$. Each legislator has a common knowledge prior that $\varepsilon$ is mean zero and is distributed according to a distribution function $F(\varepsilon)$ with density $f(\varepsilon)$. The agency knows $\varepsilon$ with certainty. Thus, this informational structure captures the idea that bureaucrats are policy experts on whom legislators rely on to improve policy formulation.

Rather than solve for the optimal mechanism from the $L$'s perspective, we limit the analysis to a very simple mechanism where $L$ chooses a set of allowable policies $P \subset X$. We assume that the sanction against an agency who chooses $p \notin P$ are so large that it never chooses to do so. Therefore, given a set of allowable policies $P$ and the state of the world $\varepsilon$, $A$ will choose $p$ to maximize

$$-(a - p + \varepsilon)^2 \text{ subject to } p \in P$$

Thus, we can see that whenever $a + \varepsilon \in P$, the agency can get her ideal outcome by choosing $p = a + \varepsilon$. When $a + \varepsilon \notin P$, $A$ will choose the point in $P$ closest to $a + \varepsilon$.

Given $A$'s best response, we turn to the legislature's choice of $P$. While in principle we can allow $P$ to be any subset of $X$, the following result establishes that $P$ will always be a closed interval $\left[\underline{p}, \overline{p}\right] \subset X$.

PROPOSITION 11.4. *$P^*$ will always be a closed interval* $\left[\underline{p}, \overline{p}\right] \subset X$

.

We leave the details of the proof to the reader, but here's a hint. Suppose that $P$ has a "hole" in it e.g. $P = \left[\underline{p}, p'\right] \cup [p'', \overline{p}]$ where $p' < p''$. Then whenever $p' < a + \varepsilon < \frac{p'+p''}{2}$, $A$ chooses $p'$ and when $\frac{p'+p''}{2} < a + \varepsilon < p''$, she chooses $p''$. Thus, the policy outcome as a function of $\varepsilon$ will appear as the solid line of Figure 11.1. It is easy to see the variance in the policy outcome can be lowered by moving $p'$ and $p''$ closer together. Since all legislators are risk-adverse, they would want to reduce the variance so long as the expected policy outcome does not change. The hint is over. It is now up to the reader to show that it is always possible to move $p'$ and $p''$ closer together without changing the mean policy outcome.[5]

### Insert Figure 11.1 Here

Given the proposition, it is clear that the legislature's collective choice problem is to choose $\underline{p}$ and $\overline{p}$.and the agency's best response function is

$$p^* = \begin{cases} \underline{p} \text{ if } a + \varepsilon < \underline{p} \\ \overline{p} \text{ if } a + \varepsilon > \overline{p} \\ a + \varepsilon \text{ otherwise} \end{cases}$$

The solid line of Figure 11.1 gives this best response.

A complication arises in modelling $L$'s collective choice in that it must decide over two dimensions, the lower bound and the upper bound. Fortunately, the reader can check that the Plott conditions

---

[5]Of course, eliminating holes does not prove that the interval must be closed. Closedness is required to make $A$'s best response well-defined.

will hold so that the majority rule decision will be the $P^*$ preferred by the legislator with the median ideal point $l_m$. To get started, note that given $A$'s best response, $l_i$ prefers the combination of $\underline{p}$ and $\overline{p}$ that maximizes

$$-\int_{\overline{p}-a}^{\infty} (l_i - \overline{p} + \varepsilon)^2 f(\varepsilon)d\varepsilon - \int_{\underline{p}-a}^{\overline{p}-a} (l_i - a)^2 f(\varepsilon)d\varepsilon - \int_{-\infty}^{\underline{p}-a} (l_i - \underline{p} + \varepsilon)^2 f(\varepsilon)d\varepsilon$$

The first order conditions are

$$\frac{\partial}{\partial \overline{p}} = 2\int_{\overline{p}-a}^{\infty} (l_i - \overline{p} + \varepsilon)f(\varepsilon)d\varepsilon = 0$$

$$\frac{\partial}{\partial \underline{p}} = 2\int_{-\infty}^{\underline{p}-a} (l_i - \underline{p} + \varepsilon)f(\varepsilon)d\varepsilon = 0$$

Note that whenever $l_i < a$, $\frac{\partial}{\partial \underline{p}} < 0$ for any finite $\underline{p}$. Thus, the optimal choice is $\underline{p}^* = -\infty$. The intuition is that when $l_i < a$ the agency always wants a higher policy than legislator $i$ would want if she were informed. Thus, legislator $i$ never finds it in her interest to constrain $A$ from choosing low policies. Similarly, if $l_i > a$, $\frac{\partial}{\partial \overline{p}} > 0$ and $\overline{p}^* = \infty$.

PROPOSITION 11.5. *The majority rule outcome for $P^*$ is the closed interval $[\underline{p}, \overline{p}]$ preferred by the legislator with ideal point $l_m$.*

Since the majority rule outcome is the median's ideal statute, the delegation game becomes one between the agency and legislative median. Thus, the allowable policies for the agency are given by the solutions:

$$\underline{p}^* = -\infty$$

$$\int_{\overline{p}^*-a}^{\infty} (l_m - \overline{p}^* + \varepsilon)f(\varepsilon)d\varepsilon = 0$$

if $l_m < a$ and

$$\overline{p}^* = -\infty$$

$$\int_{-\infty}^{\underline{p}^*-a} (l_m - \underline{p}^* + \varepsilon)f(\varepsilon)d\varepsilon = 0$$

if $l_m > a$.

Having already discussed why the median legislator will not want to constrain the agency on one end of the spectrum, consider the intuition for the expressions for the other constraint. Consider the case where $l_m < a$. We can re-write the expression for $\overline{p}^*$ as

$$l_m = \int_{\overline{p}^* - a}^{\infty} (\overline{p}^* - \varepsilon) \frac{f(\varepsilon)}{1 - F(\overline{p}^* - a)} d\varepsilon$$

This condition implies that the expected outcome conditional on $A$ being constrained has to be at the median legislator's ideal point. If this expected outcome were greater than $l_m$, the median could do better in expectation by further constraining $A$'s choice to generate lower policies. Similarly, we can write the condition for $\underline{p}^*$ when $l_m > a$ as

$$l_m = \int_{-\infty}^{\underline{p}^* - a} (\underline{p}^* - \varepsilon) \frac{f(\varepsilon)}{F(\underline{p}^* - a)} d\varepsilon$$

To generate some more specific results, assume that $\varepsilon$ is distributed uniformly on $[-E, E]$ so that $F(\varepsilon) = \frac{\varepsilon + E}{2E}$ and $f(\varepsilon) = \frac{1}{2E}$. Then if $l_m < a$,

$$l_m = \overline{p}^* - \frac{\overline{p}^* - a + E}{2}$$

$$\text{or}$$

$$\overline{p}^* = 2l_m - a + E.$$

We can see that if the agency and the median get closer together (e.g. by raising $l_m$ or decreasing $a$) then the legislature will pass a more permissive statute with a greater upper bound. Intuitively, the legislature will delegate more authority to an agency that shares its preferences. Also, when there is more policy uncertainty (e.g. $E$ is larger), the agency is granted more authority. Thus, when information asymmetries are greater, the legislature will be more dependent on the informed agency to formulate policy.[6]

**4.2. Bureaucratic Capacity.** One of the important assumptions of the Epstein and O'Halloran and most other models of delegation is that the agency can implement its policy choice perfectly without error. This may be a reasonable assumption for advanced democracies with

---

[6]When $l_m > a$, the solution is $\underline{p}^* = 2l_m - a - E$. The reader can verify that this lower bound is less restrictive when $l_m$ and $a$ are close together and when $E$ is large.

cadres of professional, highly trained bureaucrats, but it is far less applicable in many developing states and in earlier historical eras.

To address this issue, Huber and McCarty (2004) develop a model in which bureaucracies vary in their capacity to implement policies. In that model, if $A$ attempts to implement policy $p$, the resulting policy is $\widetilde{p} = p + \omega$ where $\omega$ is an implementation error which is assumed to have mean 0 and variance $\sigma_\omega^2$. Bureaucracies with high capacity are assumed to be better able to implement policies and therefore have lower values $\sigma_\omega^2$. Conversely, low capacity bureaucracies implement policy with imprecision so that $\sigma_\omega^2$ is high. Let $G(\omega)$ be the distribution function for $\omega$ and $g(\omega)$ be the associated density.

Huber and McCarty embed this model of capacity into a delegation model very similar to that of Epstein and O'Halloran. The legislature would like to delegate to the agency because the agency is better informed about the consequences of various policy choices. As above, the agency knows $\varepsilon$ but the legislature knows only that it is distributed uniformly on $[-E, E]$. They also assume that the members of $L$ and $A$ have quadratic preferences over policy space $X$.

The legislature moves first and creates a statute specifying the set of admissible policies $P = \left[\underline{p}, \overline{p}\right]$.[7] The agency then moves and attempts to implement policy $p$ which results in policy $\widetilde{p} = p - \omega$ and outcome $x = \widetilde{p} - \varepsilon$. Huber and McCarty assume that compliance with the statute requires that $\widetilde{p} \in P$. Thus, even if the agency attempted to comply i.e. $p \in P$, its implementation errors might lead to non-compliance. If $\widetilde{p} \notin P$, the agency incurs a cost $\delta$ as a sanction for non-compliance.[8] Unlike the Epstein and O'Halloran model, the agency may actually choose to be non-compliant. Alternatively, non-compliance could be purely a result of implementation errors. Since it is assumed that $L$ cannot observe $p$, the sanction for non-compliance must be the same regardless of the ultimate cause. Thus, $A$ will be sanctioned when $p - \omega > \overline{p}$ or when $p - \omega < \underline{p}$. Given a choice of $p$, the probability of sanction is $G(p - \overline{p}) + 1 - G(p - \underline{p})$.

To facilitate the exposition, note that a number of features of the Epstein-O'Halloran model generalize to this model. First, it can be shown that the majority rule choice of $P$ maximizes the utility of the legislator with ideal point $l_m$. We will assume throughout this section

---

[7]A note to the industrious reader who refers to the original publication is in order. In Huber and McCarty, the political principal is a generic politician rather than a legislature. We have adjusted the nomenclature and notation to paralell that of our discusion of Epstein and O'Halloran.

[8]Actually, Huber and McCarty assume that non-compliance is detected probabilitisically. However, our simplification doesn't alter the results.

that $a > l_m$. Secondly, we can show in the current model that $\underline{p}^* = -\infty$ if $a > l_m$.

Given this setup and the assumption that $a > l_m$, we can write $A$'s utility function as[9]

$$- \int_{-\infty}^{\infty} (a - p + \omega + \varepsilon)^2 g(\omega) d\omega - \delta \left[ G(p - \overline{p}) \right]$$

$$= - (a - p + \varepsilon)^2 - \sigma_\omega^2 - \delta \left[ G(p - \overline{p}) \right]$$

Her first order conditions for a maximum are

$$2(a - p + \varepsilon) - \delta g(p - \overline{p}) = 0$$

The first term in this expression represents the marginal benefit of moving the expected policy closer to the agency's ideal point $a$. The second term represents the net marginal cost of increasing the intended policy in terms of the probability of sanction. Increasing $p$ increases the probability of $p + \omega > \overline{p}$ by $g(p - \overline{p})$. Clearly, $A$ will choose $p$ to equate the marginal policy benefits with marginal sanction costs. To get an explicit solution for $p^*$, Huber and McCarty assume that $g(\omega) = \frac{\Omega - |\omega|}{\Omega^2}$. This density is "tent-shaped" on the interval $[-\Omega, \Omega]$ and implies that $\sigma_\omega^2 = \frac{\Omega^2}{6}$. Thus, $\Omega$ represents a measure of bureaucratic incapacity.

Given these assumptions about functional form, Figure 11.2 plots the marginal policy benefit curve and the marginal compliance costs for three values of $\overline{p}$ as a function of $p$. The marginal benefit line is the bold downward sloping line. The marginal policy benefit is independent of the location of $\overline{p}$, and declines as the Bureaucrat's action approaches the Bureaucrat's most-preferred action, $a + \varepsilon$. The marginal cost curves depend on the location of $\overline{p}$, and are depicted in Figure 11.2 by the three triangles centered at $\overline{p}_1, \overline{p}_2$, and $\overline{p}_3$. These triangles represent the function $\delta g(p - \overline{p})$. If $\overline{p}$ were too high or too low, the marginal costs would be zero at the Bureaucrat's ideal intended action, $a + \varepsilon$. This would lead the Bureaucrat to choose its ideal action. For non-extreme statute's, the Bureaucrat's best response lies at the intersection of the marginal benefit curve with the relevant marginal cost curve.[10] Given

---

[9]The second line follows from a useful fact expected utilities of quadratic functions: $\int_{-\infty}^{\infty} (x + \phi)^2 f(\phi) d\phi = (x - E(\phi))^2 + var(\phi)$

[10]McCarty and Huber assume that $\Omega^2 > \delta$ which guaratees that there is a unique intersection of the marginal benefit and cost curves and it represents a global maximum. Given their interest in systems where bureacratic capcity and the ability to sanction non-compliance is low, this assumption seems reasonable.

$\overline{p}_1$, for example, the optimal action is $p_1^*$. For any $p > p_1^*$, the marginal policy benefits of increasing the policy action (toward the Bureaucrat's most preferred) exceed the marginal compliance costs, and for any $p < p_1^*$, the reduction in compliance costs of moving the action away from the Bureaucrat's most-preferred action exceed the policy loses.

## Insert Figure 11.2 about here

We can see from Figure 11.2 that the effect of changes in $\overline{p}$ on the Bureaucrat's best-response depends on whether the apex of the "compliance cost" triangle is to the left or right of the policy benefit line (i.e., to the left or right of $\overline{p}_2$ in Figure 11.2). For $\overline{p}_1 < \overline{p} < \overline{p}_2$, $p^* > \overline{p}$. In this range, increases in $\overline{p}$ increase the marginal compliance costs of any $p > \overline{p}$, inducing the Bureaucrat to move toward the Politician's ideal point. At $\overline{p}_2$, however, this effect reverses. For $\overline{p} > \overline{p}_2$, $p^* < \overline{p}$, and increases in $\overline{p}$ decrease the marginal compliance cost of any $p^* < \overline{p}$, inducing the Bureaucrat to adjust his action closer to his ideal point. Consequently, the minimum action that the politicians can induce is given by $p^*(\overline{p}_2)$.

Huber and McCarty show that formal solution to the agent's maximization problem is

$$
p^* = \begin{cases}
a + \varepsilon & \text{if } \overline{p} - \varepsilon \leq a - \Omega \\
\frac{\Omega^2(a+\varepsilon) - \delta(\overline{p}+\Omega)}{\Omega^2 - \delta} & \text{if } a - \Omega \leq \overline{p} - \varepsilon \leq a - \frac{\delta}{\Omega} \\
\frac{\Omega^2(a+\varepsilon) - \delta(\overline{p}-\Omega)}{\Omega^2 + \delta} & \text{if } \overline{p} - \varepsilon \leq a + \Omega \leq \overline{p} - \varepsilon \leq a + \Omega \\
a + \varepsilon & \text{if } \overline{p} - \varepsilon \geq a + \Omega
\end{cases}
$$

A few features of $A$'s best response are worth noting. Notice that for extreme statutes ($\overline{p} \leq a - \Omega + \varepsilon$ or $\overline{p} \geq a + \Omega + \varepsilon$), $A$'s best response is to attempt to implement her ideal point. This is because if the statute is too lax or too constraining, the marginal compliance cost at $A$'s ideal policy is zero. Secondly, note that $l_m$ can only induce policies in the the interval $\left[a - \frac{\delta}{\Omega} + \varepsilon, a + \varepsilon\right]$. He cannot induce a lower policy because $A$ would comply less often under a more restrictive statute. Note that this minimum policy is increasing in $\Omega$. Thus, $l_m$'s ability to control $A$ decreases and capacity gets lower. The intuition is that a low capacity bureaucracy is non-compliant a large part of the time regardless of the policies it chooses and this probability of non-compliance is not very responsive to the agencies choices. Therefore, the agency will choose to implement policies closer to its ideal point since there is little additional penalty for doing so. Therefore, the model identifies one important effect of low bureaucratic capacity: bureaucrats are harder to control through statutes.

Given their interest in the bureaucratic politics of low-capacity system, Huber and McCarty focus on the optimal statute when bureaucratic capacity is sufficiently low.[11]  They find that under these conditions the optimal statute is

$$\overline{p}^* = a - \frac{\delta}{\Omega} + \frac{\delta E}{\Omega^2}$$

Just as the Epstein-O'Halloran model, the Huber-McCarty model predicts that $L$ will delegate more authority when $E$ is larger.  However, its prediction about preference divergence is the exact opposition.  In the Huber-McCarty model, the statute is more permissive when $a$ and $l_m = 0$ are further apart.  This is a consequence of the fact that low capacity bureaucrats are more likely to defect to their ideal point in response to restrictive statutes.  This defection is extremely costly to $L$ when $a$ is far from $l_m$.  Thus, $L$ is willing to grant more latitude to extreme bureaucrats to give them stronger incentives to comply with the statute.  Finally, the Huber-McCarty model makes another prediction at odds with standard models of bureaucratic delegation.  Other models have shown that when *ex post* sanctions are high, the principal is willing to delegate more.  In the Huber-McCarty model, high $\delta$ is associated with a more restrictive statute.  The rationale is straightforward.  High sanctions can induce even low capacity bureaucrats to comply.  Thus, $L$ no longer needs to grant more discretion solely to induce compliance.

**4.3. Some Generalizations.**  Most game theoretic treatments of delegation (included the ones outlined above) maintain a number of stylized assumptions.  First, principals and agents are assumed to be risk averse.  In fact in most applications, they are assumed to have quadratic preferences.  Secondly, the policy space and random shock are assumed to be single dimensional.  Finally, most models assume that policy outcomes are additive functions of policy choices and shocks.  Bendor and Meirowitz (2004) note that while these simplifying assumptions allow us to specify parsimonious models of the issues at hand, they often come at a price in terms of generalizability.  In particular, the assumptions have limited our ability to determine specifically which features are most important in the decision to delegate, the selection of and agent, and the choice of monitoring and control mechanisms.  In this section, we sketch the key points of Bendor and Meirowitz's argument.

---

[11]"Sufficiently low capacity" means $\Omega > \min\left\{E, \frac{\delta}{a}, \sqrt{\delta}\right\}$.

Assume that there is a single principal and $n$ subordinates. All agents have ideal points in $\mathbb{R}^d$ and without loss of generality the principals's ideal point is assumed to be the 0 vector. Preferences over outcomes are assumed to be represented by a utility function of the form

$$u_i(x) = h(-\|x - y_i\|)$$

where $h(\cdot)$ is a strictly increasing continuous function, $\|z\|$ is the Euclidean norm and $y_i$ is $i$'s ideal vector in $\mathbb{R}^d$. Euclidean preference over a single dimensional policy is clearly a special case of this assumption. As in Epstein-O'Halloran and Huber-McCarty, Bendor and Meirowitz assume that the principal is less informed than the subordinates, but they allow for (1) arbitrary functional forms, and (2) heterogeneity in the uncertainty associated with different policy selections. Formally this is captured by assuming that for any policy, outcome $x(p)$ is a random variable given by the conditional distribution $F(x \mid p)$. The principal knows only the family of conditional distributions and informed subordinates know deterministic mappings from $p$ into $x$. They assume that "perfect shock absorption" is possible in that an informed agent can implement any policy outcome $x$ by choosing the appropriate $p$. Epstein and O'Halloran's assumption that $x = p + \varepsilon$ is a special case of this assumption. However, because of the implementation shocks in McCarty and Huber, the shock absorption assumption holds only in expectation.

Unlike Huber-McCarty model where bureaucratic capacity refers to variation in the ability of agents to implement their intended policies, Bendor and Meirowitz model capacity as variation in the agent's expertise. They assume that with probability $q_i$ agent $i$ learns the random shock and can select $p$ to attain any $x$. However, with probability $1 - q_i$ subordinate $i$ is uniformed, knowing no more than the principal.

In their basic delegation model, the principal decides whether to delegate or not. If she does not delegate, she selects policy based on her priors about the policy shock. If she delegates, she selects an agent gives that agent complete discretion over the policy choice. The selected agent chooses policy $p$, and the game ends. We now highlight a few of the key findings. For now assume that all agents have high capacity, $q_i = 1$ for all $i \in N$.

PROPOSITION 11.6. *There exists a ball $B(\varepsilon, 0)$ centered at 0 with radius $\varepsilon$ containing ideal points of the agents to whom the principal is willing to delegate. $B(\varepsilon, 0)$ is known as the delegation set.*

The construction of $B(\varepsilon, 0)$ is straightforward. If agent $i$ is given control of policy, she would like to select a policy to enact the outcome

$x = y_i$. Accordingly, the principal is willing to delegate to $i$ only if the outcome $y_i$ is preferred to the best lottery that the principal can attain by selecting policy herself based on her priors. As long as the principal faces uncertainty, her utility of implementing policy herself, $u_0'$, is less than the utility associated with reaching $x = 0$ with probability 1, $h(0)$. This follows from the fact that $x = 0$ is the principal's ideal point so that any lottery on $\mathbb{R}^d$ has lower expected utility than getting her ideal policy $x = 0$ with certainty. So clearly, the principal will delegate to an agent with ideal point $y_i = 0$. The set of ideal points that are sufficiently close to 0 to merit delegation solve the inequality

$$h(-\|y_i\|) \leq u_0'$$

and given that $h(\cdot)$ is strictly increasing and continuous, for any value of $u_0'$ the set $\{y : h(-\|y\|) \leq u_0'\}$ is a closed ball.

Since so little structure has been placed on $h(\cdot)$, we see that the informational rationale for delegation in spatial settings hinges only on the desire to avoid bad outcomes: risk preferences, dimensionality, or the nature of uncertainty are secondary issues. It is not difficult to see that relaxing the assumption that all agents are perfectly competent $q_i = 1$ does not have a qualitative effect. However, as competency decreases the delegation set shrinks–if you are going to give authority to someone else that is not likely to know anything more than you, they had better have preferences very close to yours.

However, if several agents have ideal points in the delegation set, the choice of who to delegate to can be subtle. The traditional literature has often stressed the **ally principal** which says that if the principal delegates, she will select agent whose preferences most closely match hers. With homogenous competence, $q_i$ the same for all $i$ and the stylistic assumption that $x = p - \varepsilon$ (a multidimensional version of the assumption in Epstein-O'Halloran) the ally principal holds. The proof is left as an exercise.

However, with heterogeneity in $q_i$ or more general policy outcome function, the ally principal can fail. As an example of the first problem, consider the case of $x = p - \varepsilon$ and two agents with $\|y_1\| < \|y_2\|$. Would the principal ever choose to delegate to agent 2 instead of the more proximate agent 1? Bendor and Meirowitz show that if $q_2 > q_1$ then possibly yes. We can conclude only that the principal will never select an agent that is dominated by another agent in the sense that agent $i$ dominates agent $j$ if $\|y_i\| \leq \|y_j\|$ and $q_i \geq q_j$ with one of the inequalities strict. A more subtle finding is that once we relax that assumption that $x = p - \varepsilon$ even if $q_1 = q_2$ agent 2 may still be chosen over agent 1. The key to this possibility is that the uncertainty associated with

different policies need not be the same. In the general model it is possible that an uninformed agent 1's most preferred policy, $p_1$, results in more uncertainty than an uniformed agent 2's most preferred policy, $p_2$. This can be the case if attempting to enact certain types of outcomes (say ones far from the status quo) is harder and thus subject to larger possible errors than attempts to enact other types of outcomes (say ones close to the status quo). As an example suppose the policy and outcome spaces are $\mathbb{R}^1$. As an example of this point, consider two agents with ideal points $y_1 = -1 + \delta$ and $y_2 = 1$. Thus, agent 1 is closer to the principal. Suppose that $F(x \mid p)$ is as follows. If $p > 0$ then $x = p + .1$ or $x = p - .1$ with equal probability and if $p < 0$ then $x = p + .2$ or $x = p - .2$ with equal probability. In this case, either agent selects $p = y_i$ if they do not learn the shock. This means that agent 1's uniformed policy choice entails more outcome risk. Accordingly if the principal is risk averse meaning $h(\cdot)$ is strictly concave) and $\delta$ is sufficiently small then the principal would prefer to delegate to agent 2 as opposed to 1 even though her ideal point is closer. An example of this argument is left as an exercise.

Another potential violation of the ally principal arises from free-riding among agents when information acquisition is costly. Suppose now that the shock can be learned at a cost $c$, by any agent or principal. If the cost is incurred, the shock is observed with probability one. Further suppose that the principal can select an agent giving her authority and then observe whether she invests the cost $c$ to learn the shock. If the agent chooses not to learn then the principal can retake control and decide whether to invest $c$ herself to learn the shock and select policy. Alternatively the principal can select policy in ignorance. In this setting it turns out that the delegation set is a multidimensional doughnut. If the principal delegates to an agent with an ideal point very close to her own, the agent has the choice of paying $c$ to get the outcome utility $h(0)$ or not investing, with the knowledge that the principal will then retake control and invest $c$ herself. Thus, learning implies utility $h(0) - c$ while free-riding yields utility $h(-\|y_i\|)$. Thus for agents with ideal points closer than $h^{-1}(h(0) - c)$, free-riding on the principal is preferred. Accordingly the principal will not delegate to agents that are very close. Of course agents that are very distant will select undesirable outcomes.

PROPOSITION 11.7. *If information acquisition has cost c, then the delegation set consists of agents in the original delegation set with ideal points farther from $0$ than $d = h^{-1}(h(0) - c)$*

Bendor and Meirowitz also consider the effect of competition by agents in settings where there are many agents and one principal. Suppose that the agents are perfectly competent and simultaneously announce outcomes in $X$. The principal then selects an agent and the agent selects policy (knowing the shock) to enact the outcome she announced. If the outcome space is one dimensional and there are agents on either side of the principal this game looks like Downsian competition and in every equilibrium at least two agents announce that they will enact the principal's ideal outcome. It turns out that this conclusion holds regardless of the dimensionality of the policy space.

DEFINITION 11.6. *We say that preferences satisfy diversity if either (1) there does not exist a vector $s \in X$ such that for all $i \in N$ $x_i = \lambda_i s$ for some $\lambda_i \in R^1$ or (ii) if such a vector does exist then there must be two agents $i$ and $j$ with $\lambda_i > 0$ and $\lambda_j < 0$.*

Diversity requires that either preferences are not colinear or if they are that there are agents on either side of the principal.

PROPOSITION 11.8. *If preferences satisfy diversity then in every equilibrium at least two agents commit to enacting $x = 0$ and the principal accepts one of these offers.*

It is clear that if one agent promises to enact 0 the commitments of the other agents are payoff irrelevant. So any strategy profile in which at least 2 agents make this commitment is an equilibrium. Now suppose that no agents are making this commitment. With at least two agents at least one of the agents can move the final outcome closer to her ideal point by committing to an outcome that is closer to the principals than the closest commitment of the remaining agents. There cannot be an equilibrium in which the principal does not get her ideal outcome.

## 5. Mechanism Design and Signaling Games

So far our presentation of mechanism design proceeded quite independently of our analysis of signaling. This isn't surprising given that this is the way the literature has generally developed. However, we think much can be gained by exploring the connections between the two normally disjoint topics.

We begin with a generic setup for a signaling model. Suppose that the sender's (player $s$) type is $\theta \in \Theta = [0,1]$ and the message space is $M = [0,1]$. It is common knowledge that $\theta$ is drawn from a distribution $F(\cdot)$ on $\Theta$. Following the $s$'s message $m \in M$, the receiver (player $r$) selects a policy $p \in X = [0,1]$. Even without

specifying payoffs, we can use the concept of incentive compatibility to specify necessary conditions the existence of an equilibrium in which the senders message is fully-revealing such that $m^{-1}(m(\theta)) = \theta$. In an equilibrium in which $m(\theta)$ is one-to-one, consistent beliefs must be concentrated at the correct $\theta$ i.e. beliefs can be represented by the probability distribution

$$B(\theta \mid m') = \begin{cases} 1 \text{ if } \theta \geq m^{-1}(m') \\ \quad 0 \text{ otherwise} \end{cases}$$

Given a fully revealing message, sequential rationality by $r$ requires that $p(m) \in P(\theta) \equiv \arg\max u_r(p, \theta)$. Sequential rationality by $s$ requires that she not have an incentive to mislead $r$ by behaving as if her type were $\theta'$ when it is $\theta''$. Given the mapping $p(m)$, this incentive compatibility condition is $u_s(p(\theta''), \theta'') \geq u_s(p(\theta'), \theta'')$ for all $\theta', \theta'' \in \Theta$. Alternatively, the receiver's best response requires $u_r(p(\theta''), \theta'') \geq u_r(p(\theta'), \theta'')$ for all $\theta', \theta'' \in \Theta$. Thus, a requirement of a separating equilibrium is that $p(\theta)$ maximize both $u_s(p, \theta)$ and $u_r(p, \theta)$.

PROPOSITION 11.9. *A separating PBE exists only if the preferences of the players are similar–specifically for every $\theta \in \Theta$ it must be the case that $\{\arg\max_{p \in X} u_r(p, \theta)\} \cap \{\arg\max_{p \in X} u_s(p, \theta)\}$ is non empty.*

Given this result, it is clear that truthful revelation in cheap talk signaling with one sender requires strong similarity between sender and receiver payoffs. As an example, recall our version of the open rule Gilligan-Krehbiel model from chapter 8. There we showed that the best response by $F$ to informative signals is $p(-\theta) = \theta$ and $p(-\theta) = \theta$. The proposition specifies that a separating equilibrium exists if and only if

$$u_F(\theta| - \theta) \geq u_F(-\theta| - \theta)$$
$$u_F(-\theta|\theta) \geq u_F(\theta|\theta)$$
$$u_C(\theta| - \theta) \geq u_C(-\theta| - \theta)$$
$$u_C(-\theta|\theta) \geq u_C(\theta|\theta)$$

We know the first two inequalities holds since $p(-\theta) = \theta$ and $p(-\theta) = \theta$ so the crucial conditions are $-c^2 \geq -(2\theta + c)^2$ and $-c^2 > -(2\theta - c)^2$. Note that these both hold if $c < \theta$. This is exactly the condition of preference divergence that we derived before.

Now we extend the general model so that there are two senders 1 and 2 who each observe $\theta$ but have possibly different preferences. The question of whether there is a PBE in which the receiver learns $\theta$ can be viewed as a type of mechanism design problem where the receiver's

choice of a mapping $p(m_1, m_2) : M^2 \to X$ is analogous to selecting a mechanism. In contrast to mechanism design, a PBE of a signaling game requires that the receivers decision be sequentially rational given consistent beliefs. Thus, we are not free to choose just any mechanism that satisfies the sender's incentive compatibility conditions. However, Baron and Meirowitz (2004) show that in many cases the constraint that the receiver's actions must be sequentially rational is not very limiting.

Returning for the moment to the mechanism design problem, suppose the receiver wishes induce truthfulness by punishing the senders if their messages do not coincide. Suppose that there exists a bad policy $p^b$ which both senders like less that the receiver's best response to any truthful pair of messages. Formally, $p^b$ is defined so that for every $\theta$

$$(11.21) \qquad u_1(\arg\max_{p \in X} u_r(p, \theta), \theta) \geq u_1(p^b, \theta)$$

$$u_2(\arg\max_{p \in X} u_r(p, \theta), \theta) \geq u_2(p^b, \theta).$$

Given this definition of the $p^b$, the following mechanism satisfies the incentive compatibility conditions for both senders to be truthful:

$$p(m_1, m_2) = \begin{cases} \arg\max_{p \in X} u_r(p, m) \text{ if } m_1 = m_2 \\ \qquad p^b \text{ otherwise.} \end{cases}$$

Given this policy function, condition 11.21 implies that sender 2's best response to a truthful announcement by sender 1 is a truthful announcement to avoid the dreaded $p^b$. A similar argument applies for sender 1's incentive to be truthful. Within the mechanism design framework, the mere existence of $p^b$ is enough to induce truthfulness regardless of the receiver's utility from $p^b$. However, in signaling models, we need to worry about whether choosing $p^b$ is sequentially rational for the receiver. Accordingly it must be the case that $p^b$ is an optimal policy for the receiver given the beliefs she forms at all information sets where $m_1 \neq m_2$. This suggests that our mechanism design trick will not be very compelling to those committed to the signaling tradition. But recall that weak consistency only constrains beliefs at information sets that occur with positive probability. In an equilibrium in which the senders are truthful $m_1 \neq m_2$ does not occur. Accordingly, satisfying sequential rationality and credibly committing to enact $p^b$ if $m_1 \neq m_2$ is not that challenging. All that is required is that there exists some distribution $b^b(\cdot)$ on $\Theta$ such that $p^b \in \arg\max_p \int u_r(p, \theta) d\, b^b(\theta)$. Adding a state $\theta^b$ in which $p^b \in \arg\max_p u_r(p, \theta^b)$ would suffice. From this argument we reach the following conclusion

THEOREM 11.2. *(Baron and Meirowitz)* *With at least two senders that observe* $\theta$ *as long as there is a policy* $p^b \in X$ *satisfying condition 11.21 and a distribution* $b^b(\cdot)$ *on* $\Theta$ *such that* $p^b \in \arg\max_p \int u_r(p, \theta)d\, b^b(\theta)$ *a truthful PBE exists.*

In response to Gilligan and Krehbiel's (1989) work on heterogeneous legislative committees, Krishna and Morgan (2001) demonstrated that with two senders there are PBE in which the receiver learns the state $\theta$. Their equilibrium did not hinge on a punishment policy $p^b$ but rather used out-of-equilibrium beliefs to rationalize a policy that punished any agent that would have an incentive to lie. Krehbiel 2001 criticized this approach on the grounds off the path responses to some messages were highly discontinuous and tended to move in the wrong direction. While in equilibrium, high policies are best reponses to low messages, the Krishna and Morgan's PBE call for low policies in response to high out-of-equilibrium messages.

However, Battaglini (2002) showed that if in fact the policy and state spaces are multidimensional ($X = \Theta = [0,1]^2$) then truthful equilibria can be constructed that do not depend on beliefs in such a peculiar manner. As an example of this model, consider a receiver and two senders also with Euclidean preferences over two dimensions. The receiver has an ideal point of $(0,0)$, sender 1 has ideal point $(1,0)$ and sender 2 has ideal point $(0,1)$. In this model, each sender observes the shock $\theta \in \Theta$ perfectly, and the outcomes $x$ are a random function of policy $p$ where $x = p + \theta$. We denote a message by $s \in \{1,2\}$ by the vector $m_s = (m_s^1, m_s^2)$. Similarly outcomes and policies are vectors ($x = (x^1, x^2)$, $p = (p^1, p^2)$). A particularly simple direct mechanism can be constructed once we realize that while neither sender has the same preferences as the receiver (and thus by Proposition 11.9 there is not a truthful equilibrium in the 1 sender game), each sender has the same preferences as the receiver over a particular dimension of the problem. Sender 1 and the receiver both want the second dimension of the outcome as close to 0 as possible and sender 2 and the receiver both want the first coordinate of the outcome as close to 0 as possible. Suppose that sender 1 with ideal point $(1,0)$ is given complete control of the second coordinate, so that $p^1 = -m_1^1$ and that sender 2 with ideal point $(0,1)$ is given complete control of the first coordinate, so that $p^2 = -m_2^2$. Given this mechanism announcement of $m_s(\theta) = \theta$ is a best response, and thus the mechanism is incentive compatible. Can this receiver strategy be supported in a PBE? If so we need to find weakly consistent beliefs for which this policy function is sequentially

rational. Since the mapping $p(m)$ described above is a direct mechanism at every information set that is reached Bayes' rule results in concentrated beliefs $\theta = m_1 = m_2$. It remains to specify beliefs for information sets in which $m_1 \neq m_2$ that make policy function $p(m)$ sequentially rational for the receiver and make the truthful messages sequentially rational for the senders. One easy way to do this is to let the beliefs ignore $m_2^1$ and $m_1^2$. In other words the beliefs are concentrated at $\theta = (m_1^1, m_2^2)$. Battaglini shows that in the Euclidean preferences setting with at least 2 dimensions and 2 senders separating PBE exist as long as the ideal points of the 3 players are not on a line.

So we have seen that under particular types of preferences with 2 senders separating PBE exist. With three perfectly informed senders $\{1, 2, 3\}$ we do not need to make any assumptions about preferences as a particularly simple separating PBE can be characterized. Let $p^*(\theta)$ denote a selection from the correspondence $\arg\max u_r(p, \theta)$ and suppose that $p^+ \in p^*(\theta)$ for some $\theta$. Similar to above, suppose the receiver could commit to the following mechanism that depends on the messages $m_1, m_2, m_3$

$$p(m) = \begin{cases} p^*(\theta) \text{ if } \theta = m_i = m_j \text{ for some } i, j \in \{1, 2, 3\} \\ p^+ \text{ otherwise} \end{cases} .$$

Then truthful messages are a best response because if $i$ and $j$ are truthful then $k$'s message is outcome inconsequential. Supporting this policy function as sequentially rational is easy. Any belief mapping that is concentrated on $\theta$ if $m_1 = m_2 = m_3 = \theta$ is weakly consistent, and any beliefs that are concentrated at $\theta$ if $\theta = m_i = m_j$ for some $i, j \in \{1, 2, 3\}$ make the policy mapping a best response for the receiver.

The conclusion of Baron and Meirowitz is that any direct mechanism that (1) selects an optimal policy for the receiver following truthful messages and (2) following non-truthful messages selects a policy which is optimal given some belief on $\Theta$ can be supported in a PBE to the signaling game.

## 6. Exercises

EXERCISE 11.1. *Suppose that instead wishing to maximize the welfare of her department, the department chair wanted to maximize her surplus (contributions minus expenditures on the espresso maker). Would she still want to implement the Groves-Clarke mechanism?*

EXERCISE 11.2. *Suppose that a polling mechanism were used such that $x(m) = \frac{1}{n} \sum m_i$ so that average ideal point is the implemented policy. Show that this mechanism is not strategy-proof.*

EXERCISE 11.3. *Prove Moulin's result (Proposition 11.2).*

EXERCISE 11.4. *Prove that the strategies described by equation 11.1 form an equilibrium to the first price auction.*

EXERCISE 11.5. *Demonstrate that use of the strategy, "Hold up placard until the price exceeds $b(\theta)$" constitutes an equilibrium in the descending price auction.*

EXERCISE 11.6. *Find the expected revenue of a standard auction if $F(\cdot)$ is the uniform distribution on $[0, 1]$.*

EXERCISE 11.7. *Consider Ferejohn's model with moral hazard. Assume that equation 11.20 does not hold. Construct a mixed strategy equilibrium where the government sometimes chooses $a = low$ and the voter always removes the government when $x = low$ and occasionally removes it when $x = high$.*

EXERCISE 11.8. *Prove that in the Epstein-O'Halloran model $P^*$ will always be a closed interval $\left[\underline{p}, \overline{p}\right] \subset X$ .*

EXERCISE 11.9. *In the Epstein-O'Halloran model, show that the majority rule outcome for $P^*$ is that preferred by the legislator with ideal point $l_m$.*

EXERCISE 11.10. *In the Epstein-O'Halloran model, assume that $\varepsilon$ is distributed $N(0, \sigma^2)$. Compute the optimal statute $P$. How does $P$ depend on $l_m$, $a$, and $\sigma^2$  Hint: $E(\varepsilon | \varepsilon < m) = -\sigma \frac{\phi\left(\frac{m}{\sigma}\right)}{\Phi\left(\frac{m}{\sigma}\right)}$.*

EXERCISE 11.11. *Augment the Esptein-O'Halloran model by assuming that governor $G$ with ideal point $g > l_m$ appoints $A$ (i.e. selects $a$) prior to $L$ choosing $P$. Assume that $A$ learns $\varepsilon$ but that $G$ and $L$ believe that $\varepsilon$ is distributed uniformly on $[-E, E]$. Show that the governor's optimal appointment is $a^* \in (l_m, g)$. How does $a^*$ depend on $l_m$ and $E$?. What if the governor appoints $A$ after $P$ is selected?*

EXERCISE 11.12. *In the context of the Bendor-Meirowitz model, prove that $\{y : h(\|y\|) \leq u'_0\}$ is a closed ball.*

EXERCISE 11.13. *In the Bendor-Meirowitz model, show that if $q_i$ is the same for all $i$ and $x = p - \varepsilon$ (with $\varepsilon$ having distribution $F(\cdot)$ on $\mathbb{R}^d$), the ally principal holds.*

EXERCISE 11.14. *In the Bendor-Meirowitz model, suppose that $q_1 = q_2 = \frac{3}{4}$, and that $u(x) = -x^2$. Find the minimum value of $\delta$ such that the principal prefers delegation to agent 2.*

EXERCISE 11.15. *Construct a fully separating PBE in a version of the Battaglini model in which sender 1 has ideal point $(1,1)$ and sender 2 has ideal point $(0,1)$.*

CHAPTER 12

# Mathematical Appendix

Mathematics can be thought of as a language constructed to facilitate the derivation of logical propositions from basic axioms. Mathematical propositions are "true" if and only if the underlying axioms and assumptions are "true". Euclidean geometry is the set of all statements that are true if parallel lines never intersect. Other geometries are based on different axioms and generate other theorems.

Mathematical arguments have two main advantages over forms of non-formal argumentation:

- The basic assumptions of the proponent are more clearly laid out. This makes it clear which of the propositions are logically derived and which are assumed to be true.

In many non-formal theories of politics, it is not clear what types of assumptions about individual and institutional behavior are being made to produce the empirical predictions. The result is often predictions incompatible with any coherent theory of behavior. In mathematical models, the "micro- foundations" are well specified. The foundations are no more or less appropriate when they are formalized. However, understanding of the relationships between foundations and findings is made easier by the use of this type of reasoning.

- It is often easier to derive logical propositions within the framework of mathematics than outside it.

While mathematics was designed and has evolved to economize on the production of logical arguments, English and other spoken languages serve so many other purposes that their ability to produce logical argument is compromised.

**0.1. Mathematical Statements and Proofs.** Just as any other language, the fundamental unit of mathematics is the statement, which for now we denote $P$. Here we review the types of mathematical statements that readers are likely to encounter.

: *Universal Statement:* $P$ is always true within a given mathematical system.

Consider the following example of a universal statement.

: Let $x$ be a real number. $\forall x$, $x \leq |x|$ where the symbol ∎$\forall$ means "for all" or "for every".

To prove such a statement, we require it to be proven for a generic $x$ where we can only use the properties common to every value of $x$.

: Existential Statement: There are conditions under which $P$ is true.

The following is such an example

: $\exists x$ such that $x = |x|$ where $\exists$ means "there exist(s)."

To prove an existential statement, we must only find a value of $x$ in the given system for which $P$ is true. Of course in this example, $x \geq 0$ is the needed condition.

Mathematics also has very well defined procedures for verifying that a given statement is true. Now we consider the types of *proofs* that will be encountered in the text.

Deduction. Proofs by deduction are those in which a statement is true because it is logically connected to a statement known or presupposed to be true. Suppose we know $P$ to be true, then we can establish that $Q$ is true is we can prove "if $P$, then $Q$" or $P \Rightarrow Q$. Obviously, this can take place via a number of steps like:

$$P \Rightarrow R$$
$$R \Rightarrow S$$
$$S \Rightarrow Q$$

Sometimes when we work it out in our mind, showing $S \Rightarrow Q$ may be the first step. However, when communicating the logic of the proof, it should be written according in the order given above. Deduction can be used to prove both existential and universal statements.

The Contrapositive. Sometimes it is easier to establish, $P \Rightarrow Q$ by formulating it in terms of the negative statements $\tilde{}Q$ and $\tilde{}P$ where $\tilde{}$ means "not". It is logically true that $(\tilde{}Q \Rightarrow \tilde{}P) \Rightarrow (P \Rightarrow Q)$.

EXAMPLE 12.1. *If 7m is an odd number, then m is an odd number.*

Thus, $P = \{7m \text{ is odd}\}$, $Q = \{m \text{ is odd}\}$, $\tilde{}P = \{7m \text{ is even}\}$, and $\tilde{}Q = \{m \text{ is even}\}$. We wish to show $P \Rightarrow Q$ by showing $\tilde{}Q \Rightarrow \tilde{}P$.

$\tilde{}Q \Rightarrow m = 2k$ for some integer $k$
  $\Rightarrow 7m = 7(2k)$
  $\Rightarrow 7m = 2(7k)$
  $\Rightarrow 7m = 2n$ for some integer $n$
  $\Rightarrow \tilde{}P$

Contradiction. One way is to prove that the statement $P$ is true is to demonstrate that $\tilde{}P$ is false. Proving a statement is false is quite easy

we only need to provided a counterexample. One counterexample will show $\tilde{}P$ to be false and that $P$ is true. Note that this procedure works best when we are arguing against a universal statement or in favor of an existential statement.

EXAMPLE 12.2. *Let $n$ be any integer and let $P = \{there\ exists\ n > 0\ such\ n^2 + n + 17\ is\ not\ a\ prime\ number\}$ or $P = \{\exists n > 0 \ni n^2 + n + 17\ is\ not\ a\ prime\ number\}$  where $\exists$ means "there exists" and $\ni$ means "such that."*

We can now construct $\tilde{}P = \{\forall\ n > 0,\ n^2 + n + 17\ \text{is a prime number}\}$. But $\tilde{}P$ is false since $n = 17$ implies that $n^2 + n + 17 = 19 \cdot 17$ and is thus not prime. Thus, $P$ is true.

We may also establish $\tilde{}$P is false by deriving a series of implications from $\tilde{}P$ that lead to a false statement.

   $\tilde{}P \Rightarrow n + 1 + 17/n$ is not an integer for all $n$

      $\Rightarrow 17/n$ is not an integer for all $n$

      $\Rightarrow 1$ is not an integer

The final statement is obviously false.

## 1. Sets and Functions

**1.1. Sets.** A set is a collection of distinct objects be they numeric values or anything else. The objects of sets are called elements. We may denote sets both by enumerating them or by describing them. Suppose $S$ is the set of all positive integers less than 10. We may write $S$ as $S = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ or $S = \{n | n$ is an integer and $0 < n < 10\}$. The second statement is read as "$S$ is the set of all numbers that are integers and are greater than 0 and less than 10." Sets can either have a finite number of elements as above or an infinite such as $I = \{n\ | n$ is integer greater than 7$\}$ or $J = \{x \mid 0 < x < 1\}$. Infinite sets are either countable or uncountable. $I$ is countable since a integer can be associated with each element. $J$ on the other hand is uncountable since it is impossible to associate an integer with each element. We can prove that $J$ is uncountable as follows. Associate elements in $J$ with integers such that $x_1$ is the first element, $x_2$ is the second element and so forth. As long as $x_1 < x_2$ there exists an $x$ that is an element of $J$ such that $x_1 < x < x_2$ which is not associated with an integer. This is a contradiction.

EXAMPLE 12.3. *A rational number is one that can be written as a fraction $p/q$ where $p$ and $q$ are integers. Prove that $Q = \{x \mid x$ is rational$\}$ is infinite but countable.*

We may denote that some element belongs to a particular set with the symbol "$\in$". Therefore, the following are truthful statements: $3 \in S$, $9 \in I$, $.345678 \in J$, and $.8 \in Q$. We can also designate which elements are not in a particular set with "$\notin$" so that $10 \notin S$, $7 \notin I$, $0 \notin J$, and $\pi \notin Q$.

**1.2. Set Relations and Operations.** The following are some useful relationships between different sets.

Equality: $S_1 = S_2$ implies that $x \in S_1$ if and only if $x \in S_2$. In other words, if $S_1 = S_2$, $S_1$ and $S_2$ contain exactly the same elements.

Subset: $S_1$ is said to be a subset of $S_2$ or $S_1 \subseteq S_2$ if for all $x \in S_1$, $x \in S_2$. Thus, $S_1$ is a subset of $S_2$ if all of the elements of $S_1$ are also in $S_2$. Note that if $S_1 = S_2$, then $S_1 \subseteq S_2$ and $S_2 \subseteq S_1$. We can also define a "proper" subset which rules out the possibility of equality. $S_1 \subset S_2$ implies that for all $x \in S_1$, $x \in S_2$ and there exist $y \in S_2$ such that $y \notin S_1$. In other words, all elements of $S_1$ are in $S_2$ but $S_2$ has additional elements not found in $S_1$. We may use the symbols $\supset$ and $\supseteq$ to write such statements in the opposite order. Finally, the symbol $\nsubseteq$ means "not a subset of."

Disjoint: Two sets are said to be disjoint is they have no elements in common. Formally, $S_1$ and $S_2$ are disjoint if $x \in S_1$ implies that $x \notin S_2$.

There is a special set $\varnothing$ known as the "null set" defined by the following properties: $\varnothing$ contains no elements and for all $S$, $\varnothing \in S$. Clearly, $\varnothing$ is the smallest possible set. There is also a largest set $U$ such that for all $S$, $S \subseteq U$. We will call this set the universal set.

A number of mathematical operations on sets will be useful.

Unions: The union of two or more sets is the total set of elements contained by all of the sets. Formally, we write the union of $S_1$ and $S_2$ as $S_1 \cup S_2 = \{x \mid x \in S_1 \text{ or } x \in S_2\}$. We can write the union of a large number of sets indexed by $i$ as $\bigcup_{i=1}^{n} S_i = S_1 \cup S_2 \cup ... \cup S_n$.

Intersections: The intersection of two or more sets is the set of elements common to all of the sets. Formally, we write the intersection of $S_1$ and $S_2$ as $S_1 \cap S_2 = \{x \mid x \in S_1 \text{ and } x \in S_2\}$. We can write the intersections of a large number of sets indexed by $i$ as $\bigcap_{i=1}^{n} S_i = S_1 \cap S_2 \cap ... \cap S_n$. If $S_1$ and $S_2$ are disjoint, $S_1 \cap S_2 = \varnothing$. Since they have no elements in common, the only subset in the intersection must be the null set.

Complements: If we have a universal set $U$ and a subset $S$ we can define the complement of $S$ as $S^c = U/S = \{x \mid x \in U \text{ and } x \notin S\}$.

Let $A$, $B$, $C$, and $D$ be sets. Then the operations on these sets must satisfy the following properties..

: Communitive: $A \cup B = B \cup A$ and $A \cap B = B \cap A$
: Associative: $(A \cup B) \cup C = A \cup (B \cup C)$ and $(A \cap B) \cap C = A \cap (B \cap C)$
: Distributive: $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ and $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$

**1.3. Correspondences and Functions.** Another way to relate two sets to one another is to specify which elements of each set "correspond" or "go with each other". In general a correspondence is a rule, $f$ , that links the elements of $S_1$ to $S_2$. Formally, we may write $f : S_1 \rightarrow\rightarrow S_2$ where the set $S_1$ is called the domain (or pre-image) set and $S_2$ is the range (or image) set. As an example, consider Figure 12.1 where the correspondence $f$ relates $x \in S_1$ to elements $y, z \in S_2$.

**Insert Figure 12.1 Here**

A function is a special type of correspondence which relates each element of the domain to a unique point of the range. So if the correspondence $f : S_1 \rightarrow\rightarrow S_2$ is actually a function, for every $x \in S_1$ it is the case that $f(x)$ is a single element of $S_2$. For functions we drop one of the arrows and write $f : S_1 \rightarrow S_2$. With a function, multiple points from the domain may map into the same point in the range. In Figure 12.2, the function relates $x, w, v \in S_1$ to elements $y, z \in S_2$..Since each element of the domain maps into a single element in range, we may without ambiguity write a function as $y = f(x)$ for $x \in S_1$ and $y \in S_2$. We may also represent the function by a set of ordered pair such as $\{(y, x) \mid y = f(x)$ for all $x \in S_1$ and $y \in S_2\}$.

**Insert Figure 12.2 Here**

Consider a function $f : A \rightarrow B$. Two important properties are:

: Injectivity: For all $a_1$ and $a_2 \in A$, $f(a_1) = f(a_2)$ if and only if $a_1 = a_2$.. Each point in the range is associated with a single point in the domain. This property is also known as "one-to-one"
: Surjectivity: $\forall b \in B$ there exist $a \in A$ such that $f(a) = b$ This property is also known as "on to"

If a function has these two properties it is known as *bijection* and there exists and inverse function mapping points in $B$ to points in $A$. We can write this inverse function as $f^{-1}:B \rightarrow A$ or $f^{-1}(b) = a$ for $b \in B$ and $a \in A$ .

EXAMPLE 12.4. $y = 2x$ *is a bijection. Since every $y$ maps to a single $x$, we can write the inverse function $f^{-1}(x) = y/2$.*

EXAMPLE 12.5. $y = x^2$ *is not injective since $x$ and $-x$ produce the same $y$. However, if we restrict the domain to positive numbers, we can write $f^{-1}(x) = \sqrt{y}$*

## 2. The Real Number System

The real number system $\mathbb{R}$ consists of all the integers as well as the rational numbers (ratios of integers) and the irrational numbers (numbers that are not the ratio of integers). The real number system is actually a set of numbers and some additional structure. The system includes two operators, $+$ and $\times$ which map from $\mathbb{R} \times \mathbb{R}$ into $\mathbb{R}$ and a weak ordering $\geq$ which is a subset of $\mathbb{R} \times \mathbb{R}$. These are the familiar operators of addition and multiplication and the ordering is our old fried "greater than or equal to." Axiomatically, the system is characterized by 14 axioms. For our purposes it is sufficient to highlight only a subset of these conditions. The Real number system is a field, which means that the operations $+$ and $\times$ behave the way we expect them to: the order of addition (or multiplication) does not matter, multiplication is distributive (meaning that $\forall x, y, z \in \mathbb{R}$ $x(y + z) = xy + xz$), multiplication and addition by 0 and 1 have the expected consequences and every number has a multiplicative inverse (so that $x \times \frac{1}{x} = 1$). In addition the real number system satisfies order axioms which insure that $\geq$ behaves the way we expect it to. What is probably unfamiliar territory is one particular axiom of the real number system which we must emphasize.

DEFINITION 12.1. **Completeness axiom:***For every non-empty subset $S \subset \mathbb{R}$ if there exists an upper bound $b$ of $S$ (meaning $x \in s \implies X \leq b$) then there exists a least upper bound $c$ (meaning $c$ is an upper bound of $S$ and if $z$ is an upper bound of $S$ then $c \leq z$). In other words every set with an upper bound has a least upper bound.*

An example of a space that is not complete is $\mathbb{R} \backslash \mathbb{Z}$ where $\mathbb{Z}$ is the set of integers. In this space the set $(0, 1)$ has an upper bound (example $\frac{3}{2}$) but it does not have a least upper bound.

**2.1. Limits of Real numbers.** A sequence of real numbers $\{x_n\}_{n=1}^{\infty}$ sometimes just denoted $\{x_n\}$ we mean an infinite list of real numbers. More precisely a sequence is a function which maps the counting numbers $(1, 2, 3, ...)$ into the real numbers. In this sense $x_n$ is the value of this function evaluated at integer $n$.

DEFINITION 12.2. *The number $l \in \mathbb{R}$ is a limit of the sequence $\{x_n\}$ if for every $\varepsilon > 0$ there is an $N$ such that for all $n > N$ we have $|x_n - l| < \varepsilon$. If $l$ is a limit of the sequence $\{x_n\}$ we write $l = \lim x_n$.*

PROPOSITION 12.1. *A sequence has at most one limit.*

Proof: Suppose otherwise, then $a = \lim x_n = b$ and $a \neq b$. Since $a \neq b$ there exists some $\varepsilon > 0$ such that (1) $|a - b| > 2\varepsilon$. Since $a = \lim x_n = b$ it must be the case that for some $N$ if $n > N$ (2) $|x_n - a| < \varepsilon$ and (3) $|x_n - b| < \varepsilon$. Without loss of generality assume that $a < b$ if $x_n < a < b$ then we have contradicted 1 or 3, if $a < b < x_n$ then we have contradicted 1 and 2, if $a < x_n < b$ then 1 implies that either 2 or 3 are violated. One of these three cases must be true.∎

DEFINITION 12.3. *A sequence $\{x_n\}$ is a Cauchy sequence if for every $\varepsilon > 0$ there is an $N$ such that for all $n, m > N$ we have $|x_n - x_m| < \varepsilon$.*

PROPOSITION 12.2. *A sequence has a limit if an only if it is a Cauchy sequence.*

DEFINITION 12.4. *We say a sequence $\{x_n\}$ converges to infinity $\infty$ $(-\infty)$ if for any $b \in R$ there is some $N$ such that for all $n > N$ we have $x_n > (<)b$.*

DEFINITION 12.5. *The number $l \in \mathbb{R}$ is a cluster point of the sequence $\{x_n\}$ if for every $\varepsilon > 0$ and every $N$ there exists some $n > N$ such that. $|x_n - l| < \varepsilon$.*

The sequence $x_n = -1^n$ has two cluster points 1 and $-1$ but no limit.

PROPOSITION 12.3. *$l$ is a cluster point of $\{x_n\}$ if and only if it is the limit of a subsequence $\{x_{n'}\}$.*

DEFINITION 12.6. *The number $l \in \mathbb{R}$ is the limit superior (limsup) of the sequence $\{x_n\}$ if (1) for any $\varepsilon > 0$ there is an $N$ such that. for all $n > N$ we have $x_n < l + \varepsilon$, and (2) for any $\varepsilon > 0$ and any $N$ there exists some $n > N$ we have $x_n > l - \varepsilon$. We write $l = \limsup x_n$*

DEFINITION 12.7. *$l$ is the limit inferior (liminf) of the sequence $\{x_n\}$ if $l = -\limsup(-x_n)$.*

An alternative definition which may be more intuitive is,

DEFINITION 12.8. *The limsup is the greatest cluster point and the liminf is the least cluster point.*

PROPOSITION 12.4. *For any sequence $\{x_n\}$ $\limsup x_n \geq \liminf x_n$ and if equality holds then $\lim x_n$ exists and $\lim x_n = \limsup x_n = \liminf x_n$.*

## 3. Points and sets

We now move beyond the real numbers and consider arbitrary spaces that are endowed with particular structures. These spaces are called metric spaces. Typically, we think about a notion of distance in $\mathbb{R}$ and specify that the distance between two points $x$ and $y$ is $|x - y|$. More generally, we can think about arbitrary spaces sets endowed with a distance function.

DEFINITION 12.9. *A metric space $(X, d)$ is a set of points $X$ and a distance function $d(x, y) : X \times X \to \mathbb{R}$, satisfying the conditions:*
1. *$d(x, y) \geq 0$*
2. *$d(x, y) = 0$ if and only if $x = y$*
3. *$d(x, y) = d(y, x)$ for any $x, y \in X$*
3. *$d(x, z) \leq d(x, y) + d(y, z)$ for any $x, y, z \in X$.*

In addition it is convenient to think about balls around points. So for a scaler $\varepsilon > 0$ and a point $x \in X$ we say the $\varepsilon$-ball around $x$ is $B(x, \varepsilon) = \{y \in X : d(x, y) < \varepsilon\}$. There are two properties of sets that we are concerned with.

DEFINITION 12.10. *A set $A \subset X$ is open if for every $x \in A$ there is some $\varepsilon > 0$ such that. $B(x, \varepsilon) \subset A$. A set $A \subset X$ is closed if its complement $X \backslash A$ is open.*

Out of convention, we think of $X$ and $\varnothing$ as both open and closed. Several results about open and closed sets can be established.

PROPOSITION 12.5. *1. If $O_1$ and $O_2$ are open then $O_1 \cap O_2$ is open.*

THEOREM 12.1. *2. Given a collection of open sets $O_1, O_2, ....$ the set $\cup_i O_i$ is open.*
*3. If $C_1$ and $C_2$ are closed then $C_1 \cap C_2$ is closed.*
*4. Given a collection of closed sets $C_1, C_2, ....$ the set $\cap_i C_i$ is closed.*

PROOF. 1. Assume that $O_1$ and $O_2$ are open and pick an arbitrary point $x \in O_1 \cap O_2$. Since $O_1$ and $O_2$ are open there is some $\varepsilon_1, \varepsilon_2 > 0$ such that. $B(x, \varepsilon_1) \subset O_1$ and $B(x, \varepsilon_2) \subset O_1$. Letting $\varepsilon = \min\{\varepsilon_1, \varepsilon_2\}$ we have $B(x, \varepsilon) \subset O_1$ and $B(x, \varepsilon) \subset O_2$ implying that $B(x, \varepsilon) \subset O_1 \cap O_2$

2. Assume that $O_1, O_2, ....$ are open and pick an arbitrary point $x \in \cup_i O_i$. Since $x \in O_i$ for some $i$ there is some $\varepsilon > 0$ such that. $B(x, \varepsilon) \subset O_i$. But this means that $B(x, \varepsilon) \subset \cup_i O_i$.

3. and 4 follow from De Morgan's laws $X\backslash\{A \cup B\} = \{X\backslash A\} \cap \{X\backslash B\}$ and $X\backslash\{A \cap B\} = \{X\backslash A\} \cup \{X\backslash B\}$.∎ $\hspace{2cm}$ □

It is not the case that the infinite intersections of open sets is open. An example is the collection of open sets $(-\frac{1}{n}, \frac{1}{n})$. Each such set is open but the intersection is just the set $\{0\}$ which is not open.

Another property of sets that surfaces is

DEFINITION 12.11. *A set $A \subset X$ is bounded if there exists some finite scaler $k$ such that for every $x, y \in A$ we have $d(x, y) < k$.*

If $X$ is a subset of finite dimensional Euclidean space $\mathbb{R}^n = \{(x_1, x_2, ..., x_n) : x_i \in \mathbb{R}\}$ the we have the following definition.

DEFINITION 12.12. *A set $A \subset \mathbb{R}^n$ is compact if $A$ is closed and bounded.*

In arbitrary metric spaces a more general definition of compactness is needed. The more general (or Topological) definition of compactness deals with open covers.

DEFINITION 12.13. *Given a set $A$ an open covering of $A$ is a collection of sets $\{O_\theta\}_{\theta \in \Theta}$ where $\Theta$ is an arbitrary index set and $O_\theta$ is open for every $\theta \in \Theta$ such that. $A \subset \{\cup_{\theta \in \Theta} O_\theta\}$ (in other words if $x \in A$ then there is some $\theta \in \Theta$ such that. $x \in O_\theta$).*

A set is compact if every open covering has a finite sub covering.

DEFINITION 12.14. *A set $A$ is compact if $\{O_\theta\}_{\theta \in \Theta}$ an open covering of $A$ implies that for some finite set $B \subset \Theta$, $\{O_\theta\}_{\theta \in B}$ is a covering of $A$.*

Given our metric space is also a field (so that $+$ is defined) which is the case of $\mathbb{R}^n$ we have the additional useful definition:

DEFINITION 12.15. *A set $A$ is convex if for every $x, y \in A$ and every scaler $\lambda \in [0, 1]$ the point $\lambda x + (1 - \lambda)y$ is also in $A$.*

PROPOSITION 12.6. *If $A_1, A_2, ....$ is a collection of convex sets then $\cap_i A_i$ is a convex set.*

PROOF. Since $A_i$ is convex for any $x, y \in A_i$ and any weight we have $\lambda x + (1 - \lambda)y \in A_i$. Since any $x, y \in \cap_i A_i$ must be in every $A_i$ we know that if $x, y \in \cap_i A_i$ then $\lambda x + (1 - \lambda)y \in A_i$ for every $i$ and thus $\lambda x + (1 - \lambda)y \in \cap_i A_i$.∎ $\hspace{2cm}$ □

## 4. Continuity of Functions

We now consider two sets $X$ and $Y$ that are each subsets of metric spaces (not necessarily the same spaces). We refer to the metrics for these spaces as $d_X$ and $d_Y$. The motivating example is $X = Y = \mathbb{R}$, but the following applies to any functions that map one metric space into another. Of central importance in analysis is continuity.

DEFINITION 12.16. *A function $f : X \rightarrow Y$ is continuous at $x \in X$ if for any $\varepsilon > 0$ there is some $\delta > 0$ such that. for all $y \in X$ with $d_X(x, y) < \delta$ it is the case that $d_Y(f(x), f(y)) < \varepsilon$. A function is continuous if it is continuous at every point in its domain.*

Another way to say the same thing is offered by the following definition.

DEFINITION 12.17. *A function $f : X \rightarrow Y$ is continuous if for every open set $B \subset Y$ the inverse image $f^{-1}(B) := \{x \in X : f(x) \in B)\}$ is open*

To clarify in this definition whether or not a set $B \subset Y$ is open depends on the metric $d_Y$ and whether or not an inverse image set $f^{-1}(B)$ is open depends on $d_X$. Of particular interest are functions for which the range is a subset of the real line. These functions are called real-valued functions. For real valued functions another definition of continuity is often convenient.

DEFINITION 12.18. *Given a function $f : X \rightarrow \mathbb{R}$, the upper contour sets is a collection of sets of the form $U_\alpha = \{x \in X : f(x) \geq \alpha\}$ for every $\alpha \in R$. The lower contour sets are the sets $L_\alpha = \{x \in X : f(x) \leq \alpha\}$.*

Continuity can be restated in terms of the contour sets.

PROPOSITION 12.7. *The function $f : X \rightarrow \mathbb{R}$ is continuous if and only if all of the upper and lower contour sets are closed*

**4.1. Extrema, Solutions and Fixed Points*.** The following is a crucial result indicating sufficient conditions for optimization problems to have solutions.

THEOREM 12.2. *If $f : X \rightarrow Y$ is continuous and $X$ is compact and non-empty then there exists a point $x^* = \arg\max_{x \in X}\{f(x)\}$.*

Another crucial result follows.

THEOREM 12.3. *(Bolzano Intermediate Value theorem) If $F : [a, b] \rightarrow R$ is continuous with $f(a) < y < f(b)$ [or $f(b) < y < f(z)$] then there is a $c \in (a, b)$ with $f(c) = y$.'*

PROOF. Consider the lower contour set of $y$, $L_y = \{x \in [a, b] :$ $f(x) \leq y\}$. Now $L_y$ is non-empty as $a \in L_y$. This set is also bounded so by completeness it has a least upper bound. Call this point $c$. Either $c \in L_y$ (that is $f(c) \leq y$) or $c$ is a cluster point. If $c$ is a cluster point then there is some sequence $\{x_n\}$ of numbers in $L_y$ with $\lim x_n = c$. Since $f$ is continuous this implies that $\{f(x_n)\}$ converges to $f(c)$. Since $f(x_n) < y$ for every $n$ it is the case that $f(c) \leq y$. Thus we know that $c \in L_y$. Now consider the upper contour set of $y$ $U_y$. Now $U_y$ is non-empty as $b \in U_y$. This set is also bounded so by completeness it has a greatest lower bound. Call this point $c$. Either $c \in U_y$ (that is $f(c) \geq y$) or $c$ is a cluster point. If $c$ is a cluster point then there is some sequence $\{x_n\}$ of numbers in $U_y$ with $\lim x_n = c$. Since $f$ is continuous this implies that $\{f(x_n)\}$ converges to $f(c)$. Since $f(x_n) > y$ for every $n$ it is the case that $f(c) \geq y$. Thus we know that $c \in U_y$. Thus $c \in U_y \cap L_y$ implying that $f(c) = y$.■                                         □

Of central importance to the analysis of games are fixed points.

DEFINITION 12.19. *Given a function $f : X \to X$ a fixed point is a point $x \in X$ such that. $f(x) = x$.*

A key result is Brouwer's fixed point theorem.

THEOREM 12.4. *(Brouwer 1910) If $X \subset \mathbb{R}^n$ is compact, convex and non-empty and $f : X \to X$ is continuous then it has a fixed point.*

While the proof for $n > 1$ is beyond the scope of this review, one can prove the one-dimensional version with the intermediate value theorem.

PROPOSITION 12.8. *If $f : [a, b] \to [a, b]$ is continuous then it has a fixed point.*

PROOF. Define the function $g(x) = f(x) - x$. This is a continuous function from $[a, b]$ into $[a, b]$. If for some $a', b' \in [a, b]$ we have $g(a) > 0$ and $g(b) < 0$ or $g(a) < 0$ and $g(b) > 0$ then the intermediate value theorem implies that for some $c \in [a', b'] \subset [a.b]$ we have $g(c) = 0$ so that $f(c) = c$ and $c$ is a fixed point. The remaining cases are $g(x) > 0$ for all $x \in [a, b]$ or $g(x) < 0$ for all $x \in [a, b]$. These cases involve $f(x) > x$ for all $x$ or $g(x) < x$ for all $x$. But since $b = \sup\{x \in [a, b]\}$ $= \sup\{f(x) : x \in [a, b]$ and $a = \inf\{x \in [a, b]\} = \inf\{f(x) : x \in [a, b]\}$ this is not possible.■                                         □

When $X$ is a field the following condition is of relevant.

DEFINITION 12.20. *A function $f : X \to \mathbb{R}$ with $X$ a convex set is quasi-concave if the upper contour sets are convex. That is for every $t \in \mathbb{R}$ and $x, x' \in X$ and every $\lambda \in (0, 1)$ it is the case that $f(x) \geq t$*

and $f(x') \geq t$ implies $f(\lambda x + (1 - \lambda)x') \geq t$.   If the last inequality is always strict the function is strictly quasi concave.

A useful property of quasi-concave objective functions is easily obtained.

THEOREM 12.5. If $X$ is convex and $f : X \to \mathbb{R}$ is strictly quasi concave then $\arg\max_{x \in X}\{f(x)\}$ contains at most one point.

PROOF. By way of a contradiction assume otherwise, so that there are two distinct points $x, y \in \arg\max_{x \in X}\{f(x)\}$ and $f(\cdot)$ is strictly-quasi concave.   This means that for $\lambda \in (0, 1)$ it is the case that $f(\lambda x + (1 - \lambda)y) > f(x) = f(y)$ contradicting the fact that $x, y \in \arg\max_{x \in X}\{f(x)\}$.■                                                                 □

## 5.  Correspondences*

Some important concepts about correspondences should also be considered.

DEFINITION 12.21. A correspondence $f : X \to\to Y$ is convex-valued if for each $x \in X$ the set $f(x)$ is convex.

Notions of continuity may also be extended to correspondences. First we define the upper and lower images.

DEFINITION 12.22. The upper image of $E \subset Y$ under $f$ (denoted $f^+(E)$, is defined by $f^+(E) = \{x \in X : f(x) \subset E\}$.

The upper image of a set $E$ is the set of points in $X$ that map into subsets of $E$.

DEFINITION 12.23. The lower image of $E \subset Y$ under $f$ (denoted $f^-(E)$, is defined by $f^-(E) = \{x \in X : f(x) \cap E \neq \varnothing\}$.

The lower image of a set $E$ is the set of points in $X$ that map into sets that intersect $E$.   Just as continuity of functions pertains to properties of contour sets continuity of correspondences relates to properties of these image sets.

DEFINITION 12.24. A correspondence $f : X \to\to Y$ is upper hemi-continuous if for each $x \in X$, whenever $x \in f^+(E)$ for $E$ an open set in $Y$ there exists an open ball $B(x, \varepsilon)$ with $B(x, \varepsilon) \subset f^+(E)$.

DEFINITION 12.25. A correspondence $f : X \to\to Y$ is lower hemi-continuous if for each $x \in X$, whenever $x \in f^-(E)$ for $E$ an open set in $Y$ there exists an open ball $B(x, \varepsilon)$ with $B(x, \varepsilon) \subset f^+(E)$.

DEFINITION 12.26. *A correspondence $f : X \rightarrow\rightarrow Y$ is-continuous if it is both upper and lower hemi-continuous.*

For most problems we care about we can settle with an alternative condition which is more intuitive then upper-hemi continuity.

DEFINITION 12.27. *A correspondence $f : X \rightarrow\rightarrow Y$ is closed at $x \in X$ if $x_n \rightarrow x$, $y_n \in f(x_n)$ and $y_n \rightarrow y$ imply $y \in f(x)$. If a correspondence is closed at each point in its domain it is closed.*

PROPOSITION 12.9. *If $Y$ is compact then $f : X \rightarrow\rightarrow Y$ is upper hem-continuous if and only if it is closed.*

The following result is an early version of what is called the Theorem of the maximum. Alternative versions exist, but the basic point for formal theory is clear, the solutions of well-behaved optimization problems respond smoothly to changes in parameters.

THEOREM 12.6. *(Berge 1997) If $u : X \rightarrow \mathbb{R}$ is a continuous function and $\Gamma : Y \rightarrow\rightarrow X$ such that for each $y \in Y$, $\Gamma(y) \neq \varnothing$ then*
*(1) the function $v : Y \rightarrow \mathbb{R}$ defined by $v(y) = \max\{u(x)$ such that $x \in \Gamma(y)\}$ is continuous and*
*(2) the correspondence $a : Y \rightarrow\rightarrow X$ defined by $a(y) = \arg\max_{x \in \Gamma(y)}\{u(x)\}$ is upper hemi-continuous.*

Our final result is a generalization of Brouwer's fixed point theorem

THEOREM 12.7. *(Kakutani 1941) Let $A \subset R^n$ be compact and convex and let $f : A \rightarrow\rightarrow A$ be closed (or upper-hemi continuous) with non-empty and convex values then $f$ has a fixed point.*

## 6. Calculus

The above analysis results represent insights about problems that can be gained based only on knowledge of topological features (compactness, continuity) and convexity features. With more structure and the use of calculus finer conclusions can be reached through analysis. In this section we provide a quick review of basis concepts of calculus which will prove useful throughout the book. Readers are also referred to Gill (2004), Chiang (2004), and Simon and Blume (1994).

**6.1. Calculus in $\mathbb{R}^1$.** Many of the questions we ask in empirical political science involve what happens to variable $y$ when we change variable $x$. If the variables are related by a function so that $y = f(x)$, the *derivative* allows us to describe and quantify the effects the variables have on one another. Suppose that $y = f(x)$. What happens

if we increase $x$ to $x + h$? The change in $y$ per unit change in $x$ is then given by

$$\frac{\Delta y}{\Delta x} = \frac{f(x+h) - f(x)}{h}$$

which is just the slope of the line drawn from $f(x + h)$ to $f(x)$. The difficulty of this measure is that it depends on $h$ as as illustrated by the two heavy dotted lines corresponding to $h_1$ and $h_2$ in Figure 12.3. We would prefer a measure that does not depend on $h$ and describes the behavior of the function as close to $x$ as possible. Such a measure is

$$\frac{dy}{dx} = \lim_{h \longrightarrow 0} \frac{f(x+h) - f(x)}{h}$$

which is known as the *derivative* of $f$ with respect to $x$. This is the solid heavy line in Figure 12.3. We also may use the notation $f'(x)$. Note that while the numerator of this limit goes to zero, the denominator goes to infinity so it can converge to any value. There is no guarantee that such a limit exists. If it does exists, we say the function is *differentiable*. The limit cannot exist if $f$ is not continuous at $x$, however if may be continuous and still not be differentiable. For example consider the function : $f(x) = |x|$ which is continuous but not differentiable at $x = 0$. To see this, note that

$$\lim_{h \longrightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \longrightarrow 0} \frac{|h|}{h}$$

Such a limit does not exist because a sequence of $h < 0$ converges to $-1$ while sequences with $h > 0$ converge to 1.

## Insert Figure 12.3 Here

Since derivative is a measure of the rate of change in $y$ given a change in $x$, we can use it to determine whether or not a function is increasing or decreasing. If $f'(x) > 0$, the function is increasing while if $f'(x) < 0$ the function is decreasing.

6.1.1. *Some Special Derivatives.* Many ordinary functions have derivatives with well known forms. We now list those for reference.

(1) Constant: If $f(x) = c$ then $f'(x) = 0$.
(2) Linear: If $f(x) = a_0 + a_1 x$ then $f'(x) = a_1$.
(3) Polynomial: If $f(x) = ax^n$ then ; $f'(x) = nax^{n-1}$
(4) Exponential: If $f(x) = e^{ax}$ then $f'(x) = ae^{ax}$
(5) Natural logarithm: If $f(x) = a \ln(bx)$; $f'(x) = a/x$

6.1.2. *Derivatives of Composite Functions.* We can take derivatives of more complicated functions especially if we can break them down into composite functions say $f(x)$ and $g(x)$. The following rules help to compute such derivatives.

(1) The Addition and Subtraction Rule:  $\frac{d(f+g)}{dx} = f' + g'$ and
$\frac{d(f-g)}{dx} = f' - g'$

(2) The Product Rule:  $\frac{d(f \cdot g)}{dx} = f' \cdot g + f \cdot g'$

(3) The Quotient Rule:  $\frac{d\left(\frac{f}{g}\right)}{dx} = \frac{f' \cdot g - fg'}{g^2}$

(4) The Chain Rule:  Let $z = g(y)$ and $y = f(x)$ that $z = g(f(x))$, then $\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx} = f'(x)g'(y)$

6.1.3. *Higher Derivatives.* Since derivatives of $f(x)$ (when they exists) are themselves functions of $x$, we can take derivatives of derivatives to learn more about the properties of the function. We represent the derivative of $f'(x)$, or the second derivative of $f(x)$ as $\frac{d^2 f}{dx^2} = f''(x)$. As before, if $f'' > 0$, $f'$ is increasing and if $f'' < 0$, $f'$ is decreasing. The second derivative can also tell us about the behavior of the original function. If

- $f' > 0$, $f'' > 0$, then $f(x)$ is increasing at an increasing rate
- $f' > 0$, $f'' < 0$, then $f(x)$ is increasing at a decreasing rate
- $f' < 0$, $f'' > 0$, then $f(x)$ is decreasing at a decreasing rate
- $f' < 0$, $f'' < 0$, then $f(x)$ is decreasing at an increasing rate

Figure 12.4 plots a function that exhibits each of these properties at different ranges of $x$.

In principal, we can take $n^{th}$ order derivatives, provided that they exist. We denote these as $\frac{d^n f}{dx^n} = f^{(n)}(x)$.

### Insert Figure 12.4 Here

6.1.4. *Maxima and Minima of Functions.* Much of the mathematical analysis in political game theory involves maximizing or minimizing functions. Voters maximize utility functions and politicians maximize votes. States minimize the number of deaths in combats. The derivative is very handy in locating the local (as opposed to global) extrema of functions.

Intuitively, the local maximum (minimum) is the point where the function ceases to increase (decrease) and begins to decrease (increase). Therefore, the derivative must equal zero unless the local extremum is global and located on the boundary of range. Figure 12.5 illustrates the distinctions between global and local maxima and minima and the intuition as to why derivatives must be zero at local extrema. However, the derivative may be zero at a point that is not a maximum or a minimum as demonstrated by the function in Figure 12.6 which contains a "saddle point." Thus, a second order condition must be satisfied to guarantee that a point satisfying the first order condition is indeed an extremum. Approaching a local maximum, the derivative is

positive and becomes negative after reaching it. So the derivative must be decreasing which means second derivative cannot be positive. Conversely, at a local minimum the second derivative cannot be negative. At a saddle point, the second derivative is zero.

### Insert Figure 12.5 and 12.6 Here

Some Formal Definitions. The following definitions will be useful. Let $f : D \to \mathbb{R}$ then

- $f(x^*)$ is a global maximum if $f(x^*) \geq f(x)$ for all $x \in D$
- $f(x^*)$ is a global minimum if $f(x^*) \leq f(x)$ for all $x \in D$
- $f(x^*)$ is a local maximum for all $\epsilon > 0$, $f(x^*) \geq f(x)$ if $|x^* - x| < \epsilon$
- $f(x^*)$ is a local minimum for all $\epsilon > 0$, $f(x^*) \leq f(x)$ if $|x^* - x| < \epsilon$
- If $f(x^*)$ is a maximum, then $x^*$ is known as $\arg\max_{D} f(x)$
- If $f(x^*)$ is a minimum, then $x^*$ is known as a $\arg\min_{D} f(x)$
- If $f'(x^*) = 0$, then $x^*$ is a critical point of $f$.

6.1.5. *Application: Bureaucratic Resource Allocations.* A bureaucrat has a budget $B$ to spend on two activities that contribute to the output of the agency. The output of the agency is given by $O = \sqrt{x_1 x_2}$ where $x_1$ and $x_2$ are the expenditures on activities 1 and 2 respectively. Since $B = x_1 + x_2$, we can replace $x_2$ with $B - x_1$ so that $O = \sqrt{x_1 (B - x_1)}$. Now we wish to find the expenditure $x_1$ that maximizes the agencies output. First, we will compute the critical values $x_1^*$ to look for local maxima. The derivative of the output function is

$$O' = \frac{\left(\frac{1}{2}B - x_1^*\right)}{\sqrt{x_1^* (B - x_1^*)}}$$

Setting $O'$ to 0, reveals that the only critical value is $x_1^* = \frac{1}{2}B$. To determine whether this is indeed a maximum, we must compute the second derivative and evaluate it at $x_1^*$. The second derivative is $O'' = -\left((x_1^* (B - x_1^*))^{-\frac{1}{2}} - \left(\frac{1}{2}B - x_1^*\right)^2 (x_1^* (B - x_1^*))^{-\frac{3}{2}}\right)$. If we evaluate this second derivation at $x_1^* = \frac{1}{2}B$, it reduces to $O'' = -2 < 0$. Thus, $x_1^* = \frac{1}{2}B$ is a local maximum and produces an output of $\frac{1}{2}B$. It is easy to see that it is also a global maximum since $O(x_1^*) = \frac{1}{2}B$ is greater that $O(0) = O(B) = 0$.

6.1.6. *Concavity and Convexity of Functions.* Two important properties of functions are concavity and convexity. To illustrate these concepts, consider Figure 12.7. A function that curves downward like $f_1$ is known as concave. We can verify it is concave if for any points

like $x_1$ and $x_2$, the line between $f(x_1)$ and $f(x_2)$ lies below the function between those two points. Formally, $f : D \to \mathbb{R}$ is concave over the set $D$ if and only if $f(\lambda x_1 + (1 - \lambda)x_2) \geq \lambda f(x_1) + (1 - \lambda)f(x_2)$ for all $\lambda \in [0, 1]$ and $x_1, x_2 \in D$.

Alternatively, a function that curves upward like $f_2$ is convex. We can verify it is convex if for any points like $x_1$ and $x_2$, the line between $f(x_1)$ and $f(x_2)$ lies above the function between those two points. Formally, $f : D \to \mathbb{R}$ is convex over the set $D$ if and only if $f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$ for all $\lambda \in [0, 1]$ and $x_1, x_2 \in D$. We can extend the definition to the case of strict concavity and convexity by replacing the weak inequalities with strict ones.

Concave and convex functions are critical because of the following:

- If $f : D \to \mathbb{R}$ is concave and $f'(x^*) = 0$, $x^*$ is a global maximum.
- If $f : D \to \mathbb{R}$ is convex and $f'(x^*) = 0$, $x^*$ is a global minimum.

These statements are true because if $f'$ exists concavity implies $f'' < 0$ and convexity implies that $f'' > 0$.

## Insert Figure 12.7 Here

6.1.7. *Integral Calculus.* Let $F(x)$ be a function such that $F'(x) = f(x)$. Then we say that $F$ is the anti-derivative of $f(x)$. We typically write anti-derivatives in terms of the indefinite integral:

$$F(x) = \int f(x)\,dx$$

The laws of differentiation lead to the following results (where $C$ is an arbitrary constant). Check by differentiating the left side of each.

: $\int af(x)dx = a \int f(x)\,dx$
: $\int (f + g)\,dx = \int f dx + \int g dx$
: $\int x^n dx = \frac{x^{n+1}}{n+1} + C$
: $\int \frac{1}{x}dx = \ln x + C$
: $\int e^x dx = e^x + C$
: $\int e^{f(x)} f'(x)dx = e^{f(x)} + C$
: $\int (f(x))^n f'(x)dx = \frac{f(x)^{n+1}}{n+1} + C$
: $\int \frac{f'(x)}{f(x)}dx = \ln f(x) + C$

The most common use of the integral will be to measure the area under a function. If $F$ is the antiderivative of $f$, then the area underneath $f$ between points $a$ and $b$ is given by the definite integral

$$\int\limits_{a}^{b} f(x)\, dx = F(a) - F(b)$$

Differentiation of the Definite Integral. The rules for differentiating definite integrals are:

$\; : \; \frac{d}{dx} \int\limits_{a}^{b} f(x) dx = \int\limits_{a}^{b} f'(x) dx$

$\; : \; \frac{d}{db} \int\limits_{a}^{b} f(x) dx = f(b)$

$\; : \; \frac{d}{da} \int\limits_{a}^{b} f(x) dx = -f(a)$

$\; : \; \frac{d}{d\alpha} \int\limits_{a(\alpha)}^{b(\alpha)} f(x(\alpha)) dx = \int\limits_{a}^{b} f'(x(\alpha)) \frac{\partial x}{\partial \alpha} dx + f(b(\alpha)) \frac{\partial b}{\partial \alpha} - f(a(\alpha)) \frac{\partial a}{\partial \alpha}$

**6.2. Calculus of Several Variables.** Consider the function $y = f(\mathbf{x})$. It is often useful; to know how $y$ changes given a change in one of the elements of $\mathbf{x}$. Typically, we will look at the partial effects of $x_i$: that, is how a change in $x_i$ effects $y$ while assuming that the other elements of $\mathbf{x}$'s are held fixed. This is equivalent to examining the behavior of the function within a given "slice". Formally, the partial derivative is

$$\frac{\partial f}{\partial x_i} = \lim_{h \to 0} \frac{f(\mathbf{x} + \mathbf{h}_i) - f(\mathbf{x})}{h}$$

where $\mathbf{h}_i$ is a vector of zeros except for an $h$ in the $i^{th}$ position. Partial derivatives are as easy to take as regular derivatives by virtue of the fact that we can treat all of the other variables as constants.

EXAMPLE 12.6. *Let* $f(x_1, x_2) = \frac{x_1}{x_2}$. *Then* $\frac{\partial f}{\partial x_1} = \frac{1}{x_2}$, $\frac{\partial f}{\partial x_2} = -\frac{x_1}{x_2^2}$.

We often write the collection of partial derivatives as the *gradient vector*

$$D_{\mathbf{x}} f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, ..., \frac{\partial f}{\partial x_n} \right)'$$

The gradient vector evaluated at $\mathbf{x}$ describes the behavior of the function near $\mathbf{x}$.

6.2.1. *Higher Order and Cross Partial Derivatives.* Just as with functions of a single variable, we can take higher-order partial derivative to characterize the behavior of partial derivatives. The second partial

derivative with respect to $x_i$ is written as

$$\frac{\partial}{\partial x_i}\left(\frac{\partial f}{\partial x_i}\right) = \frac{\partial^2 f}{\partial x_i^2}$$

We can interpret exactly the same way as in the case of a single variable. However, the case of more than a single variable, we may want to know how derivatives change when other variables change. How does changing $x_j$ effect the partial derivative with respect to $x_i$? We can write the cross partial derivative as

$$\frac{\partial}{\partial x_j}\left(\frac{\partial f}{\partial x_i}\right) = \frac{\partial^2 f}{\partial x_j \partial x_i}$$

EXAMPLE 12.7. Let $f(x_1, x_2) = \frac{x_1}{x_2}$. Then $\frac{\partial^2 f}{\partial x_1^2} = 0$, $\frac{\partial^2 f}{\partial x_2^2} = -\frac{2x_1}{x_2^4}$, $\frac{\partial^2 f}{\partial x_1 \partial x_2} = -\frac{1}{x_2^2}$, and $\frac{\partial^2 f}{\partial x_2 \partial x_1} = -\frac{1}{x_2^2}$.

Note that $\frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial^2 f}{\partial x_2 \partial x_1}$. This is true generally as $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$. Thus, the order of partial differentiation does not matter.

Often we will denote the collection of second and cross derivatives in the form of the Hessian matrix:

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

6.2.2. *Implicit Function theorem\**. Many equilibrium characterizations involve finding a value of $\mathbf{x} \in \mathbb{R}^n$ that solves a system like

$$f(\mathbf{x}; \mathbf{y}) = 0$$

for a particular value of the parameters $\mathbf{y} \in \mathbb{R}^k$. When a closed form solution for a solution $\mathbf{x}^*$ exists we get an explicit relationship of the form

$$\mathbf{x}^* = g(\mathbf{y}).$$

If $g$ is a differentiable function than comparative statics analysis, (finding out how changes in $\mathbf{y}$ effect $\mathbf{x}$) is straightforward. Sometimes however, we can prove that a solution $\mathbf{x}^*$ exists for each $\mathbf{y}$ but we cannot directly solve for the function $g(\cdot)$. For example a fixed point theorem may tell us that a solution to the system

$$f(\mathbf{x}; \mathbf{y}) - \mathbf{x} = 0$$

exists, but we may not be able to analytically solve for the vector $\mathbf{x}$.

Under suitable conditions the implicit function theorem lets us implicitly characterize the derivative $D_{\mathbf{y}}\mathbf{x}^*$. We first present the result in the case of one endogenous and one exogenous variable. Let $x^*$ solve $f(x, y) = 0$.

PROPOSITION 12.10. *(Implicit Function Theorem) Let $x^*$ solve the system at $y^*$. If $f(\cdot, \cdot)$ is continuously differentiable and $\frac{\partial f(x^*, y^*)}{\partial x} \neq 0$ then for some open set $A$ containing $x^*$ and an open set $B$ containing $y^*$ there exists a continuously differentiable function $\phi : B \to A$ with $f(\phi(y), y) = 0$ and the derivative of this function at $y^*$ is given by*

$$\frac{\partial \phi(y^*)}{\partial y} = -\frac{\frac{\partial f(x^*, y^*)}{\partial y}}{\frac{\partial f(x^*, y^*)}{\partial x}}$$

To present the result in the more general case, we consider endogenous vectors of the form $\mathbf{x} = (x_1, ..., x_n) \in \mathbb{R}^n$ and exogenous vectors of the form $\mathbf{y} = (y_1, ..., y_k) \in \mathbb{R}^k$. Suppose the system $f(\mathbf{x}, \mathbf{y}) = 0$ is of the form

$$f_1(x_1, ..., x_n; y_1, ..., y_k) = 0$$

$$.$$

$$.$$

$$f_n(x_1, ..., x_n; y_1, ..., y_k) = 0.$$

The Jacobian matrix of this system with respect to the endogenous variables is then the $n$ by $n$ matrix that stacks the transpose of the Gradient vectors up

$$J = \begin{bmatrix} D_{\mathbf{x}} f_1' \\ . \\ . \\ D_{\mathbf{x}} f_n' \end{bmatrix}.$$

PROPOSITION 12.11. *(Implicit Function Theorem) Given a pair $(\mathbf{x}^*, \mathbf{y}^*)$ for which $\mathbf{x}^*$ is a solution to the system at $\mathbf{y}^*$, if it is the case that $f_1(\cdot)$ through $f_n(\cdot)$ are continuously differentiable in each coordinate of $\mathbf{x}$ and $\mathbf{y}$ and the Jacobian matrix of the system with respect to the endogenous variables, is non singular, (i.e. the determinant of $J$ is non-zero), then for some open set $A$ containing $\mathbf{x}^*$ and an open set $\mathbf{B}$ containing $y^*$ there exists a continuously differentiable function $\phi : B \to A$ with $f(\phi(\mathbf{y}), \mathbf{y}) = 0$ and the derivative of this function at $\mathbf{y}^*$ is given by,*

$$D_{\mathbf{y}}\phi(\mathbf{y}^*) = -\left[D_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y}^*)\right]^{-1} D_{\mathbf{q}}f(\mathbf{x}^*, \mathbf{y}^*)$$

6.2.3. *Optimization in $\mathbb{R}^n$.* Recall that if we want to maximize $f :$ $\mathbb{R} \to \mathbb{R}$, we need to look for values of $x$ for which $f'(x^*) = 0$.[1]  If this condition did not hold, some other $x$ in a neighborhood of $x^*$ would produce a larger value of $f(x)$.

The same logic holds for optimizing multivariate functions.  However, now we need the derivative with respect to each element of $\mathbf{x}$ to be zero.  Suppose this were not trues and that $\frac{\partial f}{\partial x_i} > 0$.  Then, the value of the function would increase for a small increase in $x_i$ and decrease for a small decrease.  Similarly, we cannot have $\frac{\partial f}{\partial x_i} < 0$ at an interior optimum.

Given this discussion, it is clear that a necessary condition for $\mathbf{x}^*$ to optimize $f : \mathbb{R}^n \to \mathbb{R}$ is that

$$Df(\mathbf{x}^*) = \mathbf{0}$$

The second order conditions for maxima and minima are based on the Hessian matrix defined in the last section and require some more advanced concepts in matrix algebra.[2]  However, we can state the conditions for $\mathbb{R}^2$.

DEFINITION 12.28. $\mathbf{x}^* \in \mathbb{R}^2$ *is a local maximum if $Df(\mathbf{x}^*) = \mathbf{0}$,* $\frac{\partial^2 f}{\partial x_1^2} < 0, \frac{\partial^2 f}{\partial x_2^2} < 0,$ *and* $\frac{\partial^2 f}{\partial x_1^2}\frac{\partial^2 f}{\partial x_2^2} > \left(\frac{\partial^2 f}{\partial x_1 \partial x_2}\right)^2$.

DEFINITION 12.29. $\mathbf{x}^* \in \mathbb{R}^2$ *is a local minimum if $Df(\mathbf{x}^*) = \mathbf{0}$,* $\frac{\partial^2 f}{\partial x_1^2} > 0, \frac{\partial^2 f}{\partial x_2^2} > 0,$ *and* $\frac{\partial^2 f}{\partial x_1^2}\frac{\partial^2 f}{\partial x_2^2} > \left(\frac{\partial^2 f}{\partial x_1 \partial x_2}\right)^2$.

For a discussion of the general case, we refer the readers to Chiang (2004), Simon and Blume (1994).

Example:  Party Resource Allocations. Suppose a political party wants to allocate its funds across two elections.  The party values each of these seats by $W_1$ and $W_2$ respectively (the party gets zero for each seat it loses).  Let $\frac{x_i}{1+x_i}$ be the probability that the party wins seat $i$ where $x_i$ is the amount of money it spends in election $i$.  The cost of spending $x_i$ is simply $x_i$.  Therefore, the party wishes choose $(x_1, x_2)$ to maximize:

$$\frac{x_1}{1+x_1}W_1 + \frac{x_2}{1+x_2}W_2 - x_1 - x_2$$

---

[1]Since the domain is $\mathbb{R}$, we need not worry about corner solutions.

[2]The sufficient condition for a maximum (minimum) is that $H$ is positive (negative) definite.  A $NxN$ matrix $M$ is positive definite if for all vectors $v \in \mathbb{R}^n, v'Mv > 0$.  It is negative definite if $v'Mv < 0$ for all $v \in \mathbb{R}^n$.

The first order conditions are

$$\frac{W_1}{(1+x_1)^2} - 1 = 0$$

$$\frac{W_2}{(1+x_2)^2} - 1 = 0$$

while the Hessian is

$$\begin{bmatrix} -\frac{W_1}{(1+x_1)^3} & 0 \\ 0 & -\frac{W_2}{(1+x_1)^3} \end{bmatrix}$$

From the first order conditions, we can see that there are four possible critical values: $\left(-\sqrt{W_1}-1, -\sqrt{W_2}-1\right)$, $\left(\sqrt{W_1}-1, -\sqrt{W_2}-1\right)$, $\left(-\sqrt{W_1}-1, \sqrt{W_2}-1\right)$, and $\left(\sqrt{W_1}-1, \sqrt{W_2}-1\right)$.

Note however that the second order conditions require $-\frac{W_i}{(1+x_i)^3} < 0$ or $x_i > -1$. Thus, since $W_i > 0$, the only critical value that satisfies the second order condition is $\left(\sqrt{W_1}-1, \sqrt{W_2}-1\right)$.

6.2.4. *Concave and Convex Functions.* The definitions of concavity and convexity that we have already encountered generalize easily to $\mathbb{R}^n$.

DEFINITION 12.30. *Let $U \subseteq \mathbb{R}^n$ and $f : U \to \mathbb{R}$. The function $f$ is concave if for all $\mathbf{x}$, $\mathbf{y} \in U$ and $\lambda \in [0,1]$, $f\left(\lambda\mathbf{x}+(1-\lambda)\mathbf{y}\right) \geq \lambda f\left(\mathbf{x}\right) + (1-\lambda) f\left(\mathbf{y}\right)$.*

DEFINITION 12.31. *Let $U \subseteq \mathbb{R}^n$ and $f : U \to \mathbb{R}$. The function $f$ is convex if for all $\mathbf{x}$, $\mathbf{y} \in U$ and $\lambda \in [0,1]$, $f\left(\lambda\mathbf{x}+(1-\lambda)\mathbf{y}\right) \leq \lambda f\left(\mathbf{x}\right) + (1-\lambda) f\left(\mathbf{y}\right)$.*

Just as before, concavity guarantees that the critical values generate global maxima while convexity guarantees global minima.

THEOREM 12.8. *Let $f : U \to \mathbb{R}$ be be a twice differentiable function where $U$ is an open and convex subset of $\mathbb{R}^n$. If $f$ is a concave function on $U$ and $Df(\mathbf{x}^*) = 0$ for $\mathbf{x}^* \in U$, then $\mathbf{x}^*$ is a global maximum of $f$ on $U$. If $f$ is a convex function on $U$ and $Df(\mathbf{x}^*) = 0$ for $\mathbf{x}^* \in U$, then $\mathbf{x}^*$ is a global minimum of $f$ on $U$.*

6.2.5. *Constrained Maximization.*

Equality Constraints. In a number of contexts in political game theory, it will be useful to solve constrained maximization problems. Such constraints may arise either because of feasibility constraints on agents choices or due to the behavior of other agents. Such problems

will be assumed to take the form of

$$\max f\left(\mathbf{x}\right) \text{ subject to } g_1\left(\mathbf{x}\right) = 0$$
$$g_1\left(\mathbf{x}\right) = 0$$
$$.$$
$$g_k\left(\mathbf{x}\right) = 0$$

where the function $f : R^n \rightarrow R^1$ and each of the functions $g_j : R^n \rightarrow R^n$ are assumed to be twice differentiable. Here we take a cook book approach demonstrating how to analyze problems of this type.

The solution to this constrained optimization problem can be found by setting up a somewhat different unconstrained optimization and solving this translated problem. The trick is to incorporate the constraints as part of the objective function.

The Langrangian

$$L\left(\mathbf{x}, \boldsymbol{\lambda}\right) = f\left(\mathbf{x}\right) - \sum_{j=1}^{k} \lambda_j g_j\left(\mathbf{x}\right)$$

represents this translated objective function. Note that it depends on both the choice variables $\mathbf{x}$ from our original problem and a new vector of $k$ variables. These new variables are called the constraint multipliers (as each constraint gets its own multiplier). Note that we started with a real valued objective function and $k$ constraints and translate the problem into an objective function that is the sum of $k+1$ real valued functions. The first order conditions for optimization of the Lagrangian are

$$\frac{\partial f\left(\mathbf{x}\right)}{\partial x_i} = \sum_{j=1}^{k} \lambda_j \frac{\partial g_j\left(\mathbf{x}\right)}{\partial x_i} \text{ for each } i = 1, ..., n$$
$$g_j\left(\mathbf{x}\right) = 0 \text{ for each } j = 1, ..., k$$

Analysis of the first $n$ conditions yields necessary conditions on $\mathbf{x}$ for a solution to the constrained problem. More formally,

PROPOSITION 12.12. *(Lagrangian Theorem) Assume that the gradient vectors of the $k$ constraint functions are linearly independent vectors. If $\mathbf{x}^*$ solves the constrained problem then there exists a vector of Lagrangian multipliers $\boldsymbol{\lambda} \in R^k$ for which $\left(\mathbf{x}^*, \boldsymbol{\lambda}\right)$ solve the above first order conditions.*

The motivation for translating the constrained problem to this unconstrained problem is best obtained by inspecting the first $n$ first order conditions of the Lagrangian. They require that any increase

in the value of $f$ obtained by changing $\mathbf{x}$ (from a solution to the first order conditions) results in a corresponding change in the value of at least one of the constraint functions $g$. In other words, if $\mathbf{x}$ solves the Lagrangian than any improvement in $f$ would come at the expense of violating the constraint. Note that the independence requirement for the procedure to work is that the Jacobian of the constraints with respect to the variables $x$ have rank $k$, (that is the gradient vectors of the $k$ constraints are independent). Without this *constraint qualification* condition it need not be the case that a change in $\sum_{j=1}^{k} \lambda_j \frac{\partial g_j(\mathbf{x})}{\partial x_i}$ corresponds to a violation of the constraint.

Application: Party Resource Allocations. Suppose now we assume that the party has a budget constraint that it must satisfy when allocating funds across districts. Now the party wishes to maximize:

$$\frac{x_1}{1+x_1} W_1 + \frac{x_2}{1+x_2} W_2$$

$$\text{subject to } B = x_1 + x_2$$

The Lagrangian is $\frac{x_1}{1+x_1} W_1 + \frac{x_2}{1+x_2} W_2 + \lambda (x_1 + x_2 - B)$ while the first order conditions are $\frac{W_1}{(1+x_1)^2} - \lambda = 0$, $\frac{W_2}{(1+x_2)^2} - \lambda = 0$, and $x_1 + x_2 = B$. The first two conditions imply that

$$\frac{W_2}{W_1} = \frac{(1+x_2)^2}{(1+x_1)^2} \text{ or } \sqrt{\frac{W_2}{W_1}} = \frac{(1+x_2)}{(1+x_1)}$$

Together with the budget constraint, we have two equations and two unknowns. Using the positive roots, we have $+\sqrt{\frac{W_2}{W_1}} = \frac{(1+B-x_1)}{(1+x_1)}$ which imply that $x_1^* = \frac{1+B-\sqrt{\frac{W_2}{W_1}}}{\sqrt{\frac{W_2}{W_1}}+1}$ and $x_2^* = \frac{1+\sqrt{\frac{W_2}{W_1}}(B-1)}{\sqrt{\frac{W_2}{W_1}}+1}$.

Inequality Constraints. The problem considered above requires that the constraints are of the form $g_j(\mathbf{x}) = 0$. A larger class of optimization problems require only that a system of inequality or equality constraints be satisfied. The general problem is then

$$\max f(\mathbf{x}) \text{ subject to } g_1(\mathbf{x}) = 0$$
$$g_j(\mathbf{x}) = 0 \text{ for } j = 1, .., k$$
$$g_t(\mathbf{x}) \leq 0 \text{ for } t = 1, .., w$$

Again we assume that all of the relevant functions are differentiable. The Kuhn-Tucker conditions are similar to the Lagrangian conditions save how the inequality constraints are treated. The relevant translated first order conditions are

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \sum_{j=1}^{k} \lambda_j \frac{\partial g_j(\mathbf{x})}{\partial x_i} + \sum_{t=1}^{w} \lambda_t \frac{\partial g_t(\mathbf{x})}{\partial x_i} \text{for each } i = 1, ..., n$$

$$g_j(\mathbf{x}) = 0 \text{ for each } j = 1, ..., k$$

$$\lambda_t g_t(\mathbf{x}) = 0 \text{ for each } t = 1, ..., w$$

The difference is that for inequality constraints, either the constraint binds (in the sense that $g_t(\mathbf{x}) = 0$ or the multiplier $\lambda_t$ is zero.

The Envelope Theorem*. In applications the objective function or the constraints may also depend on exogenous variables $\mathbf{y} = (\mathbf{y}_1, ..., y_l, ..y_z) \in \mathbb{R}^z$. Consider the problem

$$\max f(\mathbf{x}; \mathbf{y}) \text{ subject to } g_1(\mathbf{x}; \mathbf{y}) = 0$$

$$g_j(\mathbf{x}; \mathbf{y}) = 0 \text{ for } j = 1, .., k$$

$$g_t(\mathbf{x}; \mathbf{y}) \leq 0 \text{ for } t = 1, .., w$$

By $v(\mathbf{y})$ we denote the value function which is a mapping $v : \mathbb{R}^z \to \mathbb{R}^1$ with $v(\mathbf{y}) = f(\mathbf{x}^*(\mathbf{y}); \mathbf{y})$ where $\mathbf{x}^*(\mathbf{y})$ is a solution to the optimization problem above. The theorem of the maximum indicated that under suitable conditions the value function is continuous. We can use calculus to gain more insight into the dependence of the value function on the exogenous parameters.

PROPOSITION 12.13. *Assume that $v(\mathbf{y}')$ is differentiable at $\mathbf{y}'$ and that $(\mathbf{x}^*(\mathbf{y}'), \lambda(y'))$ solves the above problem and on some open set $A$ containing $\mathbf{x}^*(\mathbf{y}')$ and some open set $B$ containing $\mathbf{y}'$ the set of constraints which bind on the solution $x^* : B \to A$ is constant , then for each $i = 1, .., n$*

$$\frac{\partial v(\mathbf{y}')}{\partial y_l} = \frac{\partial f(\mathbf{x}^*(\mathbf{y}'); \mathbf{y})}{\partial y_l} - \sum_{j=1}^{k} \lambda_j \frac{\partial g_j(\mathbf{x}^*(\mathbf{y}'); \mathbf{y})}{\partial y_l} - \sum_{t=1}^{w} \lambda_t \frac{\partial g_t(\mathbf{x}^*(\mathbf{y}'); \mathbf{y})}{\partial y_l}.$$

The novelty of this result is that in characterizing $\frac{\partial v(\mathbf{y}')}{\partial y_l}$ we do not need to worry about $D_{y_l}\mathbf{x}^*(\mathbf{y})$.

6.2.6. *Multivariate Integrals.* We can also calculate the area under multivariate functions with the use of the multivariate definite integral

$$\int_{a_1}^{b_1} ... \int_{a_n}^{b_n} f(x_1, ..., x_n) \, dx_1...dx_n$$

Multivariate integrals are calculated by taking sequentially integrating with respect to one variable while holding the remaining constant.

Suppose that we integrated with respect to $x_1$. Let $F_1(x_1, ..., x_n)$ be the partial anti-derivative with respect to $x_1$, then

$$\int\limits_{a_1}^{b_1} ... \int\limits_{a_n}^{b_n} f(x_1, ..., x_n)\, dx_1...dx_n$$

$$= \int\limits_{a_2}^{b_2} ... \int\limits_{a_n}^{b_n} F_1(b_1, ..., x_n)\, dx_2...dx_n - \int\limits_{a_2}^{b_2} ... \int\limits_{a_n}^{b_n} F_1(a_1, ..., x_n)\, dx_2...dx_n$$

We can continue this iterative process by taking the partial anti-derivative of $F_1$ with respect to $x_2$ and so on. It does not matter which definite partial integral that we compute first.

EXAMPLE 12.8. *Consider* $\int\limits_1^2 \int\limits_{\frac{1}{2}}^1 x^2 y\, dx\, dy$. *Now we begin by computing* $F_1(x, y) = \frac{1}{3}x^3 y$. *Then* $\int\limits_1^2 \int\limits_{\frac{1}{2}}^1 x^2 y\, dx\, dy = \int\limits_{\frac{1}{2}}^1 \frac{8}{3}y\, dy - \int\limits_{\frac{1}{2}}^1 \frac{1}{3}y\, dy = \frac{8}{3}\left[\frac{1}{2} - \frac{1}{8}\right] - \frac{1}{3}\left[\frac{1}{2} - \frac{1}{8}\right] = \frac{7}{8}$.

## 7. Probability Theory

As we saw in chapter 3, models of decision making under uncertainty are heavily dependent upon probability theory. In this section, we outlines the basics of probability and review some of the basic results.

**7.1. Outcomes and Events.** The building blocks of probability theory are outcomes and events. Let $S$ is the set of all possible outcomes that can be generated by a random process. Such a set is known as a *sample space*. A generic element $s \in S$ is called an *outcome*.

EXAMPLE 12.9. *Flipping two coins.* $S = \{HH, HT, TH, TT\}$.

EXAMPLE 12.10. *Unemployment rates:* $S = [0, 100]$

The first example is that of a discrete sample space because the number of outcomes is finite while the latter is a continuous sample space as the number of outcomes is infinite.

Given a sample space, we define an *event* as a subset $A \subseteq S$. Thus, an event is any combination of outcomes.

EXAMPLE 12.11. $A = \{TH, HT\}$ *"flip two is different that flip one"*

EXAMPLE 12.12. $A = [4, 13]$ *"unemployment is between 4 and 13%"*

**7.2. The Axioms of Probability Theory.** Probability theory concerns itself with the likelihood that various events occur. We let $Pr(A)$ denote be the probability of event $A$ occurs. Classical probability theory is based one the following axiomatic statements about $\Pr(A)$.

AXIOM 12.1. *For any event $A$, $\Pr(A) \geq 0$.*

AXIOM 12.2. $\Pr(S) = 1$

AXIOM 12.3. *Let $A_1$, $A_2$, ... be an infinite sequence of disjoint events then* $\Pr\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \Pr(A_i)$.

Axiom 1 says simply that the probability of any event is non-negative while axiom 2 says that the probability that some event occurs is 1. Axiom 3 concerns the probability of mutually exclusive or *disjoint* events. It states that the probability of one of an infinite set of mutually exclusive events is equal to the sum of the individual events. These axioms lead directly to a number of useful properties of probabilities..

The probability of the null event is zero.

THEOREM 12.9. $\Pr(\phi) = 0$.

Axiom 3 extends directly to the case of a finite number of disjoint events.

THEOREM 12.10. *Let $A_1$, $A_2$, ..., $A_n$ be a finite sequence of disjoint events then* $\Pr\left(\bigcup_{i=1}^{n} A_i\right) = \sum_{i=1}^{n} \Pr(A_i)$.

The previous theorem plus axioms 1 and 2 imply that the probabilities of an event and its complement sum to one.

THEOREM 12.11. *Let $S|A$ be the complement of $A$, then $\Pr(A) + \Pr(S|A) = 1$.*

A direct implication of the previous theorem is that the probability of any event is less than one.

THEOREM 12.12. *For any event $A$, $0 \leq \Pr(A) \leq 1$.*

If the outcomes associated with event $B$ are a proper subset of those associated with event $A$, the probability of event $A$ have to be at least as large as the probability of $B$.

THEOREM 12.13. *If $A \subset B$, then $\Pr(A) \geq \Pr(B)$.*

The next two theorems concern the probability of a union of events. With just two events $A$ and $B$, we can decompose $A \cup B$ into three disjoint sets $A|B$, $B|A$, and $A \cap B$. Thus, $Pr(A \cup B) = \Pr(A|B) + \Pr(B|A) + \Pr(A \cap B)$ from Theorem 12.10. Theorem 12.10 also suggests that $\Pr(A) = \Pr(A|B) + \Pr(A \cap B)$ and $\Pr(B) = \Pr(B|A) + \Pr(A \cap B)$. These results produce Theorem 12.14.

THEOREM 12.14. *For any two events $A$ and $B$, $Pr(A \cup B) = Pr(A) + Pr(B) - Pr(A \cap B)$.*

Theorem 12.15 is a straightforward generalization of Theorem 12.14.

THEOREM 12.15. *For any n events $A_1$, $A_2$, ..., $A_n$,*

$$\Pr\left(\bigcup_{i=1}^{n} A_i\right) = \sum_{i=1}^{n}\left[\Pr\left(A_i\right) - \sum_{j>i}^{n}\Pr\left(A_i \cap A_j\right) + \sum_{k>j>i}^{n}\Pr\left(A_i \cap A_j \cap A_k\right) - ...\right].$$

7.2.1. *Dependence and Conditional Probability.* We now turn to the question of how the likelihood of distinct events are related. We are concerned with whether the occurrence of one event effects the probability of another. Consider two events $A$ and $B$. Suppose we know that event $B$ has occurred what is the probability that event $A$ will occur?

One obvious possibility is that the likelihoods of the events are unrelated. We say that two events $A$ and $B$ are *independent* if $\Pr(A \cap B) = \Pr(A)\Pr(B)$. When events are independent, the realization of one event has no effect on the probability of the other. Suppose that event $B$ occurs, then $\Pr(A \cap B)$ simply $\Pr(A)$. Thus, the occurrence of $A$ is not effected by the occurrence of $B$. This logic extends to a general definition of independence.

DEFINITION 12.32. *Let $A_1$, ..., $A_n$ be a set of events. They are independent if* $\Pr\left(\bigcap_{i=1}^{n} A_i\right) = \prod_{i=1}^{n}\Pr\left(A_i\right).$

For this to be true, any subset of the events must also be independent. See DeGroot and Schervish. (2001) for an example where 3 events are pairwise independent but not independent.

Now we turn to cases where there is dependency among events. A key concept for analyzing such relationships is that of *conditional probability.* Given two events $A$ and $B$, the conditional probability of $A$ given $B$ is the probability that $A$ occurs given that $B$ has occurred. We denote the conditional probability of $A$ given event $B$ as:

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)} \text{ assuming } \Pr(B) > 0$$

Note that if $A$ and $B$ are independent, $\Pr(A|B) = \Pr(A)$. It is also easy to see that conditional probabilities must satisfy $\Pr(A \cap B) = \Pr(A|B)\Pr(B) = \Pr(B|A)\Pr(A)$. This multiplication rule generalizes to

$$\Pr\left(\bigcap_{i=1}^{n} A_i\right) = \Pr(A_1) \prod_{i=2}^{n} \Pr\left(A_i | \bigcap_{j=1}^{i=1} A_i\right)$$

**7.3. Bayes' Theorem.** One of the most important uses of probability theory in political game theory is its predictions about how agents use observed events to make assessments about the probability of unobserved events. Bayes' Theorem specifies exactly how such assessments are formed. Before stating and proving this theorem, we need an additional definition. A partition is simply a group of mutually exclusive events that cover the entire sample space.

DEFINITION 12.33. *A partition of a sample space $S$ is a set of disjoint events $A_1,..,A_k$ such that* $\Pr\left(\bigcup_{i=1}^{k} A_i\right) = S$.

We can now state Bayes' Theorem.

THEOREM 12.16. *If $A_1,..,A_k$ form a partition of $S$, $\Pr(B) > 0$ $\Pr(A_i) > 0$ for all $i$*

$$\Pr(A_j|B) = \frac{\Pr(A_j)\Pr(B|A_j)}{\sum\limits_{i=1}^{k} \Pr(A_i)\Pr(B|A_i)}$$

PROOF. The proof proceeds in a number of steps.

Claim 1: If $A_1,..,A_k$ is a partition of $S$ and $B$ is a subset of $S$, the sets $A_1 \cap B, ..., A_k \cap B$ form a partition of $B$. This follows from the fact that $\bigcup\limits_{i=1}^{k} (A_i \cap B) = \bigcup\limits_{i=1}^{k} A_i \cap B = S \cap B = B$.

Claim 2: If $A_1,..,A_k$ form a partition of $S$, $\Pr(B) = \sum\limits_{i=1}^{k} \Pr(A_i \cap B)$. This is a direct application of Claim 1 and Theorem 12.10.

Claim 3: If $A_1,..,A_k$ form a partition of $S$ and $\Pr(A_i) > 0$ for all $i$, then $\Pr(B) = \sum\limits_{i=1}^{k} \Pr(B|A_i)\Pr(A_i)$. This is an application of Claim 2 and the multiplication rule for conditional probabilities.

Now we can prove the main result. Note that if $\Pr(B) > 0$, the definition of conditional probability implies that

$$\Pr(A_j|B) = \frac{\Pr(A_j \cap B)}{\Pr(B)}$$

Bayes' Law follows from the substitution of $\Pr(A_j \cap B) = \Pr(B|A_j)\Pr(A_j)$ for the numerator and $\Pr(B) = \sum_{i=1}^{k} \Pr(B|A_i)\Pr(A_i)$ for the denominator. $\square$

**7.4. Random Variables and Distributions.** It is often convenient to use numerical representations of outcomes and events. Such representations are known as *random variables*. A random variable is simply a function that maps all possible outcomes into real numbers.

DEFINITION 12.34. *Let* $X : S \rightarrow \mathbb{R}$ *for some sample space* $S$. *Then* $X$ *is a **random variable** that assigns a real number* $X(s)$ *to each possible outcome* $s \in S$.

Given this definition of random variables, it is straightforward to define events as sets of real numbers and to define probability *distributions* over the random variables. A distribution is simply an assignment of probabilities to such events.

DEFINITION 12.35. *Let* $A$ *be any subset of* $\mathbb{R}$ *and let* $\Pr(X \in A)$ *denote the probability that* $X$ *is in* $A$. *Then* $\Pr(X \in A) = Pr\{s : X(s) \in A\}$. *A probability distribution of* $X$ *is an specification of* $\Pr(X \in A)$ *for all* $A \subset \mathbb{R}$.

7.4.1. *Discrete Distributions.* A random variable $X$ has a discrete distribution if it can take on only a finite number of outcomes: $x_1, x_2, \ldots, x_k$. We call this set of possible outcomes the support of the distribution.

DEFINITION 12.36. *If a random variable has a discrete distribution, the **probability function** (pf.)* $f$ *of* $X$ *is defined as* $f(x) = \Pr(X = x)$ *for any real number* $x$.

If $x$ is not equal to one of the points in the support of $X$, then $f(x) = 0$. By the axioms of probability theory, we know that $0 \leq \sum_{i=1}^{k} f(x_i) \leq 1$ and $0 \leq f(x_i) \leq 1$.

EXAMPLE 12.13. *The Uniform Distribution over Integers: Suppose that the value of* $X$ *is equally likely to be one of* $k$ *integers* $1, 2, 3, \ldots k$. *Then the pf is*

$$f(x) = \begin{cases} \frac{1}{k} \ for \ x = 1, ..., k \\ 0 \ otherwise \end{cases}$$

EXAMPLE 12.14. *The Binomial Distribution: Suppose an experiment succeeds with probability* $p$ *and fails with probability* $1 - p$. *The*

*pf for x successes out of n trials is given by:*

$$f(x) = \begin{cases} \binom{n}{x}p^x(1-p)^{n-x} \text{ for } x = 0, ..., n \\ 0 \text{ otherwise} \end{cases}$$

*where* $\binom{n}{x} = \frac{n!}{x!(n-x)!}$.

7.4.2. *Continuous Distributions.* Suppose that $X$ can take on an infinite number of values. Then we say that $X$ is a continuous random variable. If $X$ is a continuous random variable, then there exists a non-negative function $f$ such that for any interval $A = [a, b]$

$$\Pr(X \in A) = \int_A f(x)dx = \int_a^b f(x)dx$$

The function $f$ is known as the probability density function (or pdf). It does not tell us the $Pr(X = x)$ (which is 0) but the $\lim \Pr(X \in [x - \epsilon, x + \epsilon])$ as $\epsilon$ goes to 0. Every pdf must satisfy the following:

$$f(x) > 0 \text{ for all } x$$

$$\int_{-\infty}^{\infty} f(x)dx = 1$$

The set $X^s = \{x : f(x) > 0\}$ is known as the support of $X$.

EXAMPLE 12.15. *The Uniform Distribution on an Interval: Let a and b be two real numbers. Consider an experiment where in which a point X is chosen from $S = [a, b]$ where the probability that X belongs to any subinterval is proportional to the length of that subinterval. This implies that the pdf must be the same on any point in S and 0 otherwise. Thus,* $\int_{-\infty}^{\infty} f(x)dx = \int_a^b f(x)dx = \int_a^b cdx = 1$. *Solving the integral for c, we obtain that* $c = \frac{1}{b-a}$. *Thus, the pdf for the uniform distribution is*

$$f(x) = \begin{cases} \frac{1}{b-a} \text{ for } x \in [a, b] \\ 0 \text{ otherwise} \end{cases}$$

7.4.3. *The Cumulative Distribution Function.* The *cumulative distribution function* (cdf) is a real-valued function that relates for any real number the probability that $X$ takes on a lower value:

$$F(x) = \Pr(X \le x)$$

For discrete distributions, $F(x) = \sum\limits_{\{i:x_i < x\}} f(x_i)$.   For continuous dis-

tributions, $F(x) = \int\limits_{-\infty}^{x} f(\xi)d\xi$. Note that at any point in which $F$ is

differentiable, we have $F'(x) = f(x)$.   Thus, continuous random vari-
able can be represented either by the pdf of the cdf.

The following are some important properties of the cdf.

(1) $\Pr(X > x) = 1 - F(x)$.
(2) $\Pr(x_2 > X \geq x_1) = F(x_2) - F(x_1)$.
(3) $F$ is non-decreasing i.e. if $x_2 > x_1$ then $F(x_2) \geq F(x_1)$.
(4) $\lim\limits_{x \to -\infty} F(x) = 0$ and $\lim\limits_{x \to \infty} F(x) = 1$.
(5) $F$ is always continuous from the right.  It may be discontinu-
    ous from the left at $x$ if $x$ occurs with a positive probability.
    See Figure 12.8.

### Insert Figure 12.8 Here

7.4.4. *Bivariate Distributions.* Sometimes we will be concerned with
the probabilities of two or more random variables, say $X$ and $Y$, si-
multaneously.   One useful tool for analyzing two random variables is
the joint distribution which characterizes the probability of pairs of
realizations of $X$ and $Y$.

Discrete Bivariate Probability Functions. For the case of two dis-
crete random variables, the bivariate probability function is given by

$$f(x, y) = \Pr\{X = x \text{ and } Y = y\}.$$

Let $x_1, \ldots, x_k$ and $y_1, \ldots, y_m$ be the support of $X$ and $Y$ respectively,
then $f(x, y)$ must satisfy the following properties:

(1) $\sum\limits_{i=1}^{k} \sum\limits_{j=1}^{m} f(x_i, y_j) = 1$.
(2) Let $A$ be any set of combinations of $\{x_1, \ldots, x_k\}$ and $\{y_1, \ldots, y_m\}$
    then $\Pr\{(x, y) \in A\} = \sum\limits_{(x_i, y_j) \in A} f(x_i, y_j)$.

Continuous Bivariate Density Functions. If $X$ and $Y$ are continuous
random variables, the bivariate density is defined by

$$\Pr\{(x, y) \in A\} = \iint\limits_{A} f(x, y)dxdy$$

for any $A \subset \mathbb{R}^2$.  The bivariate pdf must satisfy the following properties:

(1) For any $(x, y) \in \mathbb{R}^2$, $f(x, y) > 0$.
(2) $\iint\limits_{\mathbb{R}^2} f(x, y) = 1$.

Bivariate Distribution Function. We can also generalize the notion of the cdf to bivariate distributions. The joint distribution function can be denoted by

$$F(x, y) = \Pr\{x \leq X \text{ and } y \leq Y\}$$

We can use the joint distribution function to determine the probability that $(x, y)$ lies in a rectangle $[a, b]x[c, d]$

$$Pr(a < X < b \ \& \ c < Y < d) = Pr(a < X < b \ \& \ Y < d) - Pr(a < X < b \& Y < c)$$

$$= Pr(X < b \ \& \ Y < d) - Pr(X < a \ \& \ Y < d) -$$
$$Pr(aX < b \ \& \ Y < c) - Pr(X < a \ \& \ Y < c)$$

$$= F(b, d) - F(a, d) - F(b, c) + F(a, c)$$

We can derive the distribution functions of $X$ and $Y$ ($F_x$ and $F_y$) from the joint distribution function:

$$F_x(x) = \lim_{y \to \infty} F(x, y)$$
$$F_y(y) = \lim_{x \to \infty} F(x, y)$$

7.4.5. *Marginal Distributions.* Suppose we know the joint pdf of $X$ and $Y$. We can get the probability density of each of them individually. These distributions of the individual random variables are know as the marginal distributions. For discrete distributions, the marginal probability functions $f_x$ and $f_y$ are defined by

$$\Pr(X = x) = f_x(x) = \sum_y \Pr(x = X \text{ and } y = Y) = \sum_y f(x, y)$$
$$\Pr(Y = y) = f_y(y) = \sum_x \Pr(x = X \text{ and } y = Y) = \sum_x f(x, y)$$

For continuous random variables, the marginal density functions are

$$f_x(x) = \int_{-\infty}^{\infty} f(x, y) dy$$

$$f_y(y) = \int_{-\infty}^{\infty} f(x, y) dx$$

**7.5. Independent Random Variables.** We know that if $X$ and $Y$ are independent random variables, the $\Pr(X = x \text{ and } Y = y) = \Pr(X = x)\Pr(Y = y)$. This implies that

$$F(x, y) = F_x(x)F_y(y)$$

where $F_x$ and $F_y$ are the marginal cdfs. It is also true that

$$f(x, y) = f_x(x)f_y(x)$$

where $f_x$ and $f_y$ are marginal probability (density) functions.

**7.6. Conditional Distributions.** Suppose that $X$ and $Y$ are not independent. Then we can define conditional distributions of $X$ given $Y$ and $Y$ given $X$. The derivation of the conditional distributions follow directly from the definition of conditional probability given above. For the the discrete case, the conditional probability functions $g_x(x|y)$ and $g_y(y|x)$ are defined as follows.

$$g_x(x|y) = \Pr(X = x|Y = y) = \frac{\Pr(X = x \text{ and } Y = y)}{\Pr(Y = y)} = \frac{f(x, y)}{f_y(y)}$$

$$g_y(y|x) = \Pr(Y = y|X = x) = \frac{\Pr(X = x \text{ and } Y = y)}{\Pr(X = x)} = \frac{f(x, y)}{f_x(x)}$$

For continuous random variables,

$$g_x(x|y) = \frac{f(x, y)}{f_y(y)}$$

$$g_y(y|x) = \frac{f(x, y)}{f_x(x)}$$

**7.7. The Expectation of a Random Variable.** One of the most important features of any probability distribution is its *expectation* or central tendency. The expectation of a random variable is the average over all of its realizations weighted by its probability. For discrete distributions, the expectation of $X$ or $E(X)$ is defined as

$$E(X) = \sum_x x f(x).$$

For continuous distributions,

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

Often in this book, we will be interested in expectations of functions of a random variable. Let $Y = r(X)$, then $E(Y) = \int_{-\infty}^{\infty} r(x)f(x)dx$.

The expectation functions must satisfy the following properties:

(1) If $Y = a + bX$, then $E(Y) = a + bE(X)$.
(2) If there exists $a$ such that $Pr(X \geq a) = 1$, then $E(X) \geq a$. If there exists $b$ such that $Pr(X \leq b) = 1$, $E(X) \leq b$.
(3) If $X_1, \ldots, X_n$ are random variables, $E(X_1 + \ldots + X_n) = E(X_1) + \ldots + E(X_n)$.
(4) If $X_1, \ldots, X_n$ are INDEPENDENT random variables, then
$$E\left(\prod_{i=1}^{n} X_i\right) = \prod_{i=1}^{n} E(X_i).$$

**7.8. The Variance of a Random Variable.** Another important property of a random variable is the extent to which it deviates from its expected value. One such measure is the variance defined as

$$var(X) = \sigma_x^2 = E\left[(X - E(X))^2\right].$$

The variance function must satisfy a number of properties.

(1) If there exists $c$ such that $\Pr(X = c) = 1$, $var(X) = 0$.
(2) For any constants $a$ and $b$, $var(a + bX) = b^2 var(X)$
(3) For any random variable $X$, $var(X) = E(X^2) - [E(X)]^2$
(4) If $X_1, \ldots, X_n$ are INDEPENDENT random variables, then
$$var\left(\sum_{i=1}^{n} X_i\right) = \sum_{i=1}^{n} var(X_i).$$

**7.9. The Median and the Mode.** Two other important function that help to summarize random variables are the median and the mode.

The Median: Let $F$ be the cdf of $X$. A point $m$ is the median of $X$ if and only if $Pr(X \leq m) \geq .5$ and $Pr(X \geq m) \leq .5$ or (for continuous distributions) $F(m) = .5$.

The Mode: Let $f$ be the pf or pdf of $X$. Then a number $m$ is a mode of $X$ if and only if $m \in \arg\max f(x)$.

**7.10. Covariance and Correlation.** Given a joint distribution over $(X, Y)$, we are often interested in describing the relationship between $X$ and $Y$. In particular, we would like to know the extent to which they move together or covary. To this end, the covariance is defined as

$$cov(X, Y) = \sigma_{xy} = E\left[(X - \mu_x)(Y - \mu_y)\right]$$

where $\mu_x = E(X)$ and $\mu_y = E(Y)$.

If $X$ and $Y$ move "together", the covariance is the expectation of a positive function and therefore positive. If $X$ and $Y$ move "against one another", the covariance is the expectation of a negative function and is therefore negative.

The covariance has a scaling problem. Consider the covariance of $X$ and $Z = aY + b$. It is straightforward to show that $\sigma_{xz} = a\sigma_{xy}$. Thus, the covariance depends on how variables are scaled. The correlation adjusts for this problem. Formally, the correlation between $X$ and $Y$ is

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

Covariances and correlation coefficients have to satisfy the following properties.

(1) For any random variables $X$ and $Y$ with finite variances, $1 \geq \rho_{xy} \geq -1$.
(2) For any random variables $X$ and $Y$, $\sigma_{xy} = E(XY) - \mu_x \mu_y$.
(3) For independent random variables $X$ and $Y$ with finite variances, $\sigma_{xy} = \rho_{xy} = 0$.
(4) For random variable $X$ with a finite variance and $Y = aX + b$, $\rho_{xy} = 1$ if $a > 0$ and $\rho_{xy} = -1$ if $a < 0$.
(5) For any random variables $X$ and $Y$ with finite variances, $var(X + Y) = \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}$.
(6) If $X_1, \ldots, X_n$ are random variables with finite variances, $var(\sum_{i=1}^{n} X_i) = \sum_{i=1}^{n} \sigma_i^2 + 2\sum_{i=1}^{n} \sum_{j=i+1}^{n} \sigma_{ij}$.

**7.11. Conditional Expectation.** Often we are interested in computing expectations of random variables conditioned on the outcomes of other random variables. The conditional expectation function is defined as:

$$E(Y|x) = \sum_y y g_y(y|x)$$

$$E(Y|x) = \int_{-\infty}^{\infty} y g_y(y|x) dy$$

The conditional expectation is function of $X$ and has a distribution derived from the distribution of $X$. An important property of conditional expectations is the Law of Iterated Expectations.

$$E(E(Y|X)) = E(Y)$$

# Bibliography

[1] Acemoglu, Daron and James Robinson. 2004. *The Economic Origins of Democracy and Dictatorship.* Book manuscript.

[2] Arrow, Kenneth J. 1951. *Social Choice and Individual Values.* New York: Wiley.

[3] Ashworth, Scott and Ethan Bueno de Mesquita 2004. "Monotone Comparative Statics: An Introduction for Political Scientists." Typescript, Princeton University.

[4] Austen-Smith, David and Jeffrey Banks. 1988. "Elections, Coalitions, and Legislative Outcomes." *American Political Science Review,* 82(2): 405-422.

[5] Austen-Smith, David and Jeffrey S. Banks. 1999. *Positive Political Theory I: Collective Preference.* Ann Arbor, MI: University of Michigan Press.

[6] Austen-Smith, David, and John R. Wright. 1992. "Competitive Lobbying for a Legislator's Vote." *Social Choice and Welfare* 9:229-257.

[7] Austen-Smith, David and John R. Wright. 1994. "Counteractive Lobbying." *American Journal of Political Science* 38:25-44.

[8] Axelrod, Robert. 1970. *Conflict of Interest.* Chicago: Marham.

[9] Axelrod, Robert. 1984 *The Evolution of Cooperation.* New York: Basic Books.

[10] Banks, Jeffrey A. 1991. *Signaling Games in Political Science.* Harwood Academic Publishers.

[11] Banks, Jeffrey S. 1990. "Equilibrium Behavior in Bargaining Games." *American Journal of Political Science* 34(3):599-614.

[12] Banks, Jeffrey and Joel Sobel. 1987. "Equilibrium Selection in Signaling Games" *Econometrica* 55:647-661.

[13] Baron, David P. 1991. "Majoritarian Incentives, Pork Barrel Programs, and Procedural Control." *American Journal of Political Science* 35(1):57–90.

[14] Baron, David P. and John A. Ferejohn. 1989. "Bargaining in Legislatures." *American Political Science Review,* 83(4):1181-1206.

[15] Baron, David P. and Adam Meirowitz. 2004. "Fully-revealing Equilibria of Multiple-sender Signaling and Screening Models." Typescript, Princeton University.

[16] Battaglini, Marco. 2002 "Multiple Referrals and Multidimensional Cheap Talk." *Econometrica.* 70:1379-1401.

[17] Bellman, Richard. 1957. *Dynamic Programming.* Princeton, NJ: Princeton University Press

[18] Bendor, Jonathan and Adam Meirowitz. 2004. "Spatial Models of Delegation." *American Political Science Review* 98(2):293-310.

[19] Berge, Claude. 1997. *Topological Spaces.* Dover.

[20] Black, Duncan. 1958. *The Theory of Committees and Elections.* London: Cambridge University Press.

[21] Blau, Julian H. 1971. "A Direct Proof of Arrow's Theorem." *Econometrica* 40(1):61-67.

[22] Border, Kim C. 1989. *Fixed Point Theorems with Applications to Economics and Game Theory* New York: Cambridge University Press.

[23] Bredon, Glen E. 1993. *Topology and Geometry,* New York: Springer.

[24] Brouwer, L. E. J. 1910. "Über Abbildung von Mannigfaltigkeiten." *Mathematische Annalen* 71: 97-115.

[25] Calvert, Randall L. 1985. "Robustness of the Multidimensional Voting Model: Candidate Motivations, Uncertainty, and Convergence." 29(1):69-95.

[26] Cameron, Charles M. 2000. *Veto Bargaining: The Politics of Negative Power.* New York: Cambridge University Press.

[27] Cameron, Charles M. and Nolan McCarty. 2004. "Models of Vetoes and Veto Bargaining." *Annual Review of Political Science.*

[28] Chiang, Alpha. 2004. *Fundamental Methods of Mathematical Economics* McGraw-Hill.

[29] Clarke, Edward H. 1971. "Multi-part Pricing of Public Goods." *Public Choice* 2:19-33.

[30] Cho, In-Koo, and David Kreps, 1987. "Signaling Games and Stable Equilibria" *Quarterly Journal of Economics* 102:179-221.

[31] Cox, Gary and Mathew D. McCubbins. 1994. *Legislative Leviathan* University of California Press.

[32] Debreu, Gerard. 1959. *The Theory of Value.* New Haven, Conn: Yale University Press.

[33] DeGroot, Morris and Mark J. Schervish. 2001. *Probability and Statistics.* Pearson Addison Wesley.

[34] Downs, Anthony. 1957. *An Economic Theory of Democracy* New York: Harper and Row.

[35] Duggan, John and Cesar Martinelli. 2001. "A Bayesian Model of Voting in Juries." *Games and Economic Behavior*, 37: 259-294.

[36] Echenique, Federico, 2002. "A Characterization of Strategic Complementarities." Working Papers 1142, California Institute of Technology, Division of the Humanities and Social Sciences.

[37] Ellsberg, Daniel. 1961. "Risk, Ambiguity, and the Savage Axioms," *Quarterly Journal of Economics* 75:643-669.

[38] Epstein, David and Sharyn O'Halloran. 1994. "Administrative Procedures, Information, and Agency Discretion" *American Journal of Political Science* 38(3): 697-722.

[39] Epstein, David and Peter Zemsky. 1995. "Money Talks: Deterring Quality Challengers in Congressional Elections." *American Political Science Review*, 89(2):295-308.

[40] Huber, John D. and Nolan McCarty. 2004. "Bureaucratic Capacity, Delegation, and Political Reform." *American Political Science Review*

[41] Fearon, James. D. 1994. "Domestic Political Audiences and the Escalation of International Disputes." *American Political Science Review* 88(3):577-592.

[42] Fearon, James. D. 1995. "Rationalist Explanations for War." *International Organization*, 49(3):379-414.

[43] Fearon, James D. and David D. Laitin 1996. "Explaining Interethnic Cooperation." *American Political Science Review* 90(4):715-735.

[44] Fedderson, Timothy and Wolfgang Pessendorfer. 1998. "Convicting the Innocent: The Inferiority of Unanimous Jury Verdicts under Strategic Voting." *American Political Science Review*, 92(1):23-35.

[45] Ferejohn, John. 1986. "Incumbent Performance and Electoral Control." *Public Choice* 50: 5-26.

[46] Gaughan, Edward. 1993. *Introduction to Analysis, 4th ed.* Pacific Grove, CA: Brooks Cole Publishing.

[47] Gibbard, Alan. 1973. "Manipulation of Voting Schemes: A General Result." *Econometrica* 41(4): 587–602.

[48] Gill, Jeff. 2004. *Essential Mathematics for Political and Social Research.* New York: Cambridge University Press.

[49] Gilligan, Thomas and Keith Krehbiel. 1987. "Collective Decision-Making and Standing Committees: An Informational Rationale for Restrictive Amendment Procedures." *Journal of Law, Economics, and Organization*, 3(2):287-335.

[50] Gilligan, Thomas W. and Keith Krehbiel. 1989. "Asymmetric Information and Legislative Rules with a Heterogenous Committee." *American Journal of Political Science*, 33(2):459-90.

[51] Green, Edward and Rob Porter. 1984. "Noncooperative Collusion Under Imperfect Price Information." *Econometrica* 52(January):87-100.

[52] Groseclose, Timothy. 1999. "An Examination of the Market for Favors and Votes in Congress." *Economic Inquiry* 34:320-40.

[53] Groseclose, Timothy and Keith Krehbiel. 2002."Gatekeeping." paper presented at the Conference on Political Parties and Legislative Organization in Parliamentary and Presidential Regimes, Yale University, March, 2002.

[54] Groseclose, Timothy and Nolan McCarty. 2000. "The Politics of Blame: Bargaining Before an Audience." *American Journal of Political Science* 45(1):100-119.

[55] Groves, Theodore. 1973. "Incentives in Teams." *Econometrica* 45:.617–631.

[56] Harsanyi, John C. 1967-68. "Games with Incomplete Information Played by Bayesian Players." *Management Science* 14: 159-182, 320-334, 486-502.

[57] Hotelling, Harold. 1929. "Stability in Competition." *The Economic Journal* 39(153): 41-57.

[58] Kahn, Kim F. and Patrick J. Kenney. 1999. *The Spectacle of U.S. Senate Campaigns.* Princeton, NJ: Princeton University Press.

[59] Kahneman, Daniel and Amos Tversky. 1979. "Prospect Theory: An Analysis of Decision Under Risk." *Econometrica* 47(2):263-291.

[60] Kakutani, Shizuo. 1941. "A Generalization of Brouwer's Fixed Point Theorem." *Duke Mathematical Journal* 8: 457-459.

[61] Knight, Frank H. 1921. *Risk, Uncertainty, and Profit.* Boston: Houghton Mifflin.

[62] Kolmogorov. A.N., and S.V. Fomin. 1970. *Introductory Real Analysis.* Dover.

[63] Krehbiel, Keith. 2001. "Plausibility of Signals by Heterogeneous Committees." *American Political Science Review*, 95: 453-8.

[64] Krishna. Vijay. 2002. *Auction Theory* Academic Press.

[65] Krishna, Vijay and John Morgan. 1997. "An Analysis of the War of Attrition and the All-Pay Auction." *Journal of Economic Theory* 72: 343-362.

[66] Krishna, Vijay and John Morgan. 2001. "Asymmetric Information and Legislative Rules: Some Amendments." *American Political Science Review*, 95(2):435-452.

[67] Lasswell, Harold D. 1936. *Politics: Who Gets What, When and How.* New York: McGraw Hill Book Company.

[68] Matthews, Steven A. 1989. "Veto Threats: Rhetoric in a Bargaining Game." *Quarterly Journal of Economics* 104:347-369.

[69] McCarty, Nolan 1997. "Presidential Reputation and the Veto." *Economics and Politics* 9:1-27.

[70] McCarty, Nolan. 2000a. "Proposal Rights, Veto Rights, and Political Bargaining." *American Journal of Political Science*, 44(3):506-522.

[71] McCarty, Nolan. 2000b. "Presidential Pork: Executive Veto Power and Distributive Politics." *American Political Science Review*, 94(1):117-129.

[72] McCarty, Nolan 2002. "Vetoes in the Early Republic." Working paper, Woodrow Wilson School, Princeton University.

[73] McKelvey, Richard D. 1976. "Intransitivities in Multidimensional Voting Models and Some Implications for Agenda Control," *Journal of Economic Theory* 12:472-482.

[74] Meirowitz, Adam. 2002. "Informative Voting and Condorcet Jury Theorems with a Continuum of Types." *Social Choice and Welfare* 19:219-236.

[75] Meirowitz, Adam. 2004. "Costly Action in Electoral Contests." Typescript, Princeton University.

[76] Mertens, Jean-Francois, and Shmuel Zamir. 1985. "Formulation of Bayesian analysis for games with incomplete information." *International Journal of Game Theory* 14(1):1–29.

[77] Milgrom and Weber. 1985

[78] Muthoo, Abhinay. 1999. *Bargaining Theory with Applications* New York: Cambridge University Press.

[79] Myerson, Roger 1981

[80] Moulin, Herve. 1980. "On Strategy-proofness and Single Peakedness." *Public Choice* 35:437-455.

[81] Nash, John C. 1950a. "The Bargaining Problem." *Econometrica* 18(2): 155-162.

[82] Nash, John F. 1950b. "Equilibrium points in N-Person Games." *Proceedings of the National Academy of Science* 48-49.

[83] Palfrey, Thomas R. and Howard Rosenthal. 1984. "Participation and the Provision of Discrete Public Goods: A Strategic Analysis." *Journal of Public Economics* 24(2):171-193.

[84] Palfrey, Thomas R. and Howard Rosenthal. 1988. "Private Incentives in Social Dilemmas: The Effects of Incomplete Information and Altruism." *Journal of Public Economics* 35:309-332.

[85] Plott, Charles R. 1967. "A Notion of Equilibrium and Its Possibility under Majority Rule. *American Economic Review* 57:787-806.

[86] Powell, Robert. 1999. *In the Shadow of Power.* Princeton, NJ: Princeton Univ. Press

[87] Primo, David. 2004. "Open vs. Closed Rules in Budget Legislation: A Result and an Application." Typescript, University of Rochester.

[88] Olson, Mancur. 1965. *The Logic of Collective Action: Public Goods and the Theory of Groups* Cambridge: Harvard University Press.

[89] Riker, William. 1962. *The Theory of Coalitions.*

[90] Riley and Samuelson 1981

[91] Romer, Thomas and Howard Rosenthal. 1978. "Political Resource Allocation, Controlled Agendas, and the Status Quo." *Public Choice* 33(1):27-44.

[92] Royden. H.L. 1988. *Real Analysis,* 3ed. Prentice Hall.

[93] Rubinstein, Ariel. 1982. "Perfect Equilibrium in a Bargaining Model." *Econometrica* 50:97-110.

[94] Sattherwaite, Mark. 1975. "Strategy-proofness and Arrow's Conditions: Existence and Correspondence Theorems for Voting Prodedures and Social Welfare Functions." *Journal of Economics Theory* 10: 187-217.

[95] Savage, Leonard. 1954. *The Foundations of Statistics.* New York: John Wiley.

[96] Selten, Reinhard. 1965. "Spieltheoretische Behandlung eines Oligopolmodells mit Nachfragentragheit." *Zeitschrift fur die gesamte Staatswissenschaft* 12:201-324.

[97] Shepsle, Kenneth. 1979. "Institutional Arrangements and Equilibrium in Multidimensional Voting Models." *American Journal of Political Science*, 23(1):27-59.

[98] Shepsle, Kenneth A. and Barry R. Weingast. 1987. "The Institutional Foundations of Committee Power." *American Political Science Review*, 81(1): 85-104.

[99] Simon, Carl P. and Lawrence Blume. 1994. *Mathematics for Economists.* New York: W. W. Norton & Company.

[100] Taylor, Michael. 1976. *Anarchy and Cooperation* London: Wiley.

[101] Topkis, Donald. M. 1998. *Supermodularity and Complementarity.* Princeton University Press.

[102] Von Neumann, John and Oscar Morgenstern. 1944. *Theory of Games and Economic Behavior.* Princeton N.J.:Princeton University Press.

[103] Weingast, Barry R. 1997. "The Political Foundations of Democracy and the Rule of Law." *American Political Science Review*, 91(2): 245-263.

[104] Weingast, Barry R. and William J. Marshall. 1988. "The Industrial Organization of Congress; or, Why Legislatures, Like Firms, Are Not Organized as Markets." *Journal of Political Economy*, 96(1): 132-163.

[105] Whittman, Donald. 1977. "Candidates with Policy Preferences: A Dynamic Model." *Journal of Economic Theory* 14:180-189.

[106] Zhou, Lin. 1994. "The Set of Nash Equilibria of a Supermodular Game is a Complete Lattice." *Games and Economic Behavior* 7:295-300.

**Figure 2.1**

**Satiable and Non-Satiable Utility Functions**



Non-satiable                                  Satiable

# Figure 2.2

## Linear and Quadratic Preferences



Linear

Quadratic

**Figure 2.3**

**Indifference Curves for Two-Dimensional Quadratic Preferences**



$$u\left(x^{*}\right) > u\left(w\right) > u\left(y\right) > u\left(z\right)$$

# Figure 3.1

## The Simplex



Two-Dimensional Simplex



Three Dimensional Simplex

**Figure 3.2**

**Tree Representations of Lotteries**

# Figure 3.3

## Compound Lotteries

$$\tfrac{1}{4}\,p + \tfrac{3}{4}\,q = r$$

**Figure 3.4**

**Risk Averse Preferences**



Utilities

$u(x_2)$

$u(w)$

$pu(x_1)+(1-p)u(x_2)$

$u(x_1)$

$x_1$     $w = px_1 + (1-p)x_2$     $x_2$

**Outcomes**

**Figure 3.5**

**Risk Acceptant Preferences**



Utilities

$u(x_2)$

$pu(x_1)+(1-p)u(x_2)$

$u(w)$

$u(x_1)$

$x_1$

$w=px_1+(1-p)x_2$

$x_2$

Outcomes

**Figure 3.6**

**Risk Neutral Preferences**



Utilities

$u(x_2)$

$u(w) = pu(x_1) + (1-p)u(x_2)$

$u(x_1)$

$x_1$

$w = px_1 + (1-p)x_2$

$x_2$

Outcomes

# Figure 3.7

## Risk and Spatial Preferences

**Utilities**

$u(w_1)$

$pu(x_1)+(1-p)u(x_2)$

$u(w_2)$

$qu(x_2)+(1-q)u(x_3)$

$x_1 \qquad w_1 \qquad x_2 \qquad w_2 \qquad x_3$

**Outcomes**

$$w_1 = px_1 + (1-p)x_2 \qquad\qquad w_2 = qx_2 + (1-q)x_3$$

# Figure 3.8

# Figure 3.9

## Prospect Theoretical Value Functions

# Figure 3.10

## Decision Weights

# Figure 4.1

## Preferences over Sub-Fields

# Figure 4.2

## Single-Peaked Preferences over
## Sub-Fields

**Rankings**



**Alternatives**

A ——— C ·········· I —··—·· T —————— F — — — —

# Figure 4.3

## Condorcet Winner in Two Dimensions

# Figure 4.4

## No Condorcet Winner

# Figure 4.5

## McKelvey's Theorem

# Figure 5.1

### Mixed Strategy Nash Equilibrium to Colonel Blotto Game

**Figure 5.2**

Mixed Strategy Nash Equilibrium to Terrorist Hunt

# Figure 5.3

Mixed Strategy Nash Equilibrium to Terrorist Hunt
(Modified Payoffs)

# Figure 5.4

## Country 2 Payoffs as a Function of Country 1's Investment



$$s_1''' > s_1'' > s_1'$$

**Figure 5.5**

Best Response Functions for
Externality Game

# Figure 5.6:  Strategic Complementarity

# Figure 5.7: Fixed Point

**Figure 5.8:  Comparative Statics on Fixed Points**



*f(x)*

*x*

**Figure 5.9: Mixed Strategy Equilibria of Palfrey-Rosenthal Game**

# Figure 7.1

Escalation Game

A

Do  Not
Initiate

Initiate

(0,0)

B

Acquiesce

Escalate

(4,-4)

(-8,-8)

**Figure 7.2**

The Centipede Game

| 1 | | 2 | | 1 | | 2 | | 1 | |
|---|---|---|---|---|---|---|---|---|---|
| Left | | L | | L | | L | | L | |

(10,-10)

Down        D        D        D        D

(1,1)     (-2,2)     (3,3)     (-4,4)     (5,5)

# Figure 7.3

Regulatory Enforcement Game



$B$

$H$

$L$

$P$

$P$

Oversight

No oversight

Oversight

No oversight

$-c, 1-k$

$-c, 1$

$-f, 1-k$

$0, 0$

# Figure 7.4

Prisoner's Dilemma in Extensive Form

# Figure 7.5

Complex Information Sets



In stage 1, player 1 has one information set that is a singleton.

In stage 2, player 2 has two information sets

In stage 3, player 3 has four information sets.

# Figure 7.6

### Sequential Voting Game



1

x        y

2        2

x    y    x    y

3     3     3     3

x   y   x   y   x   y   x   y

*x wins*   x wins   *x wins*   y wins   *x wins*   y wins   y wins   y wins

### *X* Wins  Subgame



1

x        z

2        2

x    z    x    z

3     3     3     3

x   z   x   z   x   z   x   z

*x wins*   x wins   *x wins*   z wins   *x wins*   z wins   z wins   z wins

*Y* Wins Subgame

# Figure 7.7

Rule of Law Game

# Figure 7.8

## Equilibrium Policies from Romer-Rosenthal Game

# Figure 7.9

Equilibrium Outcomes from Veto Bargaining

# Figure 7.10

The Effects of Veto Overrides

# Figure 7.11

The Democratization Game



Rich

$D$       $N$

Poor

Rich

$\tau$

$\tau$

$R$     $NR$     $R$     $NR$

$v^p\left(R,\mu^s\right), v^r\left(R,\mu^s\right)$     $v^p\left(\tau^p\right), v^r\left(\tau^p\right)$     $v^p\left(R,\mu^s\right), v^r\left(R,\mu^s\right)$     $v^p\left(N,\hat{\tau}\right), v^r\left(N,\hat{\tau}\right)$

# Figure 8.1

Deterrence Game A

A

Do  Not
Initiate

Initiate

0,0

B

Acquiesce

Escalate

4,-4

-8,-8

# Figure 8.2

Deterrence Game (B)

```
                    A
   Do  Not        / \        Initiate
   Initiate      /   \
                /     \
                       B
   0,0                / \
              Acquiesce/   \ Escalate
                      /     \
                     /       \
                  4,-4       -8,-3
```

# Figure 8.3

# Figure 8.4

# Figure 8.5

**Figure 8.6**

**Figure 8.7**

# Figure 8.8

5,3      4,0      0,0      2,7

*l*      *h*      *l*      *h*

2

*a*      *a*

*A*      *B*

1    ( *N* )    1

*p*      1-*p*

*b*      *b*

2

*l*      *h*      *l*      *h*

1,9      0,2      3,1      4,4

# Figure 8.9

*Nature*

$p_0$       $1-p_0$

*S*       *W*

Incumbent

*WC*    *~WC*    *WC*    *~WC*

Challenger

*E*   *~E*   *E*   *~E*   *E*   *~E*   *E*   *~E*

$1-\pi_s-c_s,$    $1-c_s,$    $1-\pi_s,$    $1,$    $1-\pi_w-c_w,$    $1-c_w,$    $1-\pi_w,$    $1,$

$\pi_s-k$    $0$    $\pi_s-k$    $0$    $\pi_w-k$    $0$    $\pi_w-k$    $0$

**Figure 8.10**



C

run

don't run

m

(0,0,2)

endorse

don't
endorse

p=0

I

p=1

high
effort

low
effort

high
effort

low
effort

(-1,0,0)

(2,0,1)

(-2,1,1)

(-1,-1,0)

**Figure 8.11**



0,1    2,0        1,0    3,1

~a

a      Bush    a    ~a

y    y

w    ~w

Hussein   N   Hussein

$p = 1/4$    $1-p = . 3/4$

~y    ~y

Bush

a    ~a    a    ~a

1,1    3,0        2,1

**Figure 8.12**

Packages $f$

Indifference
curve for
type 1

$f = 2r$

Indifference curve
for type 2

$f(r^p)$

Region of
deviations
$r'$

Messages $r$

$r^p$

**Figure 8.13**



Packages $f$

$f = 3r$

$f = 2r$

Indifference curve for type 1

Indifference curve for type 2

$f(r^p)$

Messages $r$

$r^p$

$r^{min}$

**Figure 8.14**

**Figure 10.1:  Nash Bargaining Solution**

# Figure 10.2:  Veto Bargaining with Incomplete Information

## Panel a

| | Preferred by $e$ and $m$ | Preferred by $m$ only | |
|---|---|---|---|

$q$          $e$     $m$       $c$

$b_e$        $b_m$

## Panel b

| | Preferred by $e$ and $m$ | Preferred by $m$ only | |
|---|---|---|---|

$q$         $e$     $m$

$b_e$      $\begin{matrix} c \\ = \\ b_m \\ = \\ b_e \end{matrix}$

## Panel c

| | Preferred by $e$ and $m$ | Preferred by $m$ only | |
|---|---|---|---|

$q$        $e$     $m$

$\begin{matrix} c \\ = \\ b_m \\ = \\ b_e \end{matrix}$

**Figure 10.3:  Proposals in Cheap-talk Game**

$$r \qquad q \qquad e \qquad m \qquad\qquad\qquad c \qquad\qquad a$$

$$b_r \qquad\qquad\qquad b_e \qquad\qquad b_m \quad b_a$$

# Figure 10.4
## Proposals in "Babbling Equilibrium"



## Proposals in "Two Message Equilibrium"
### following Compromising Message

# Figure 10.5

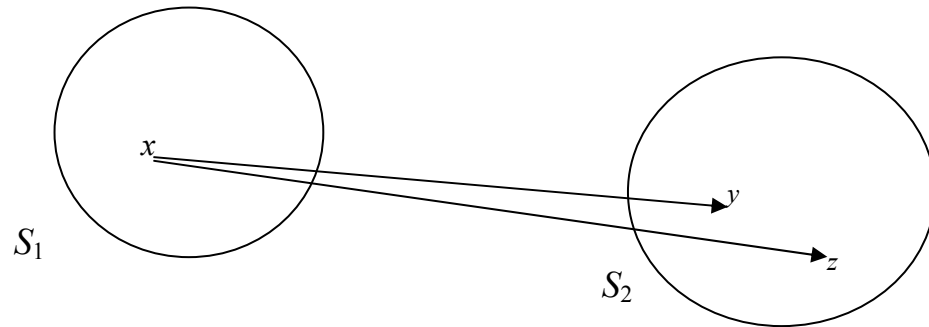## Conditions for Equilibrium Vetoes in the Blame Game Model

# Figure 11.1: Policy Outcomes from the Epstein-O'Halloran Model
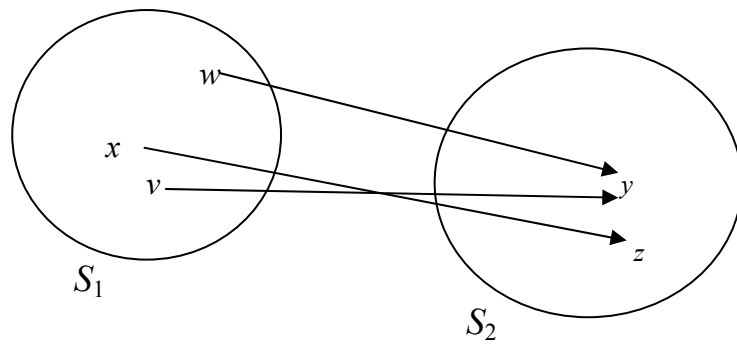
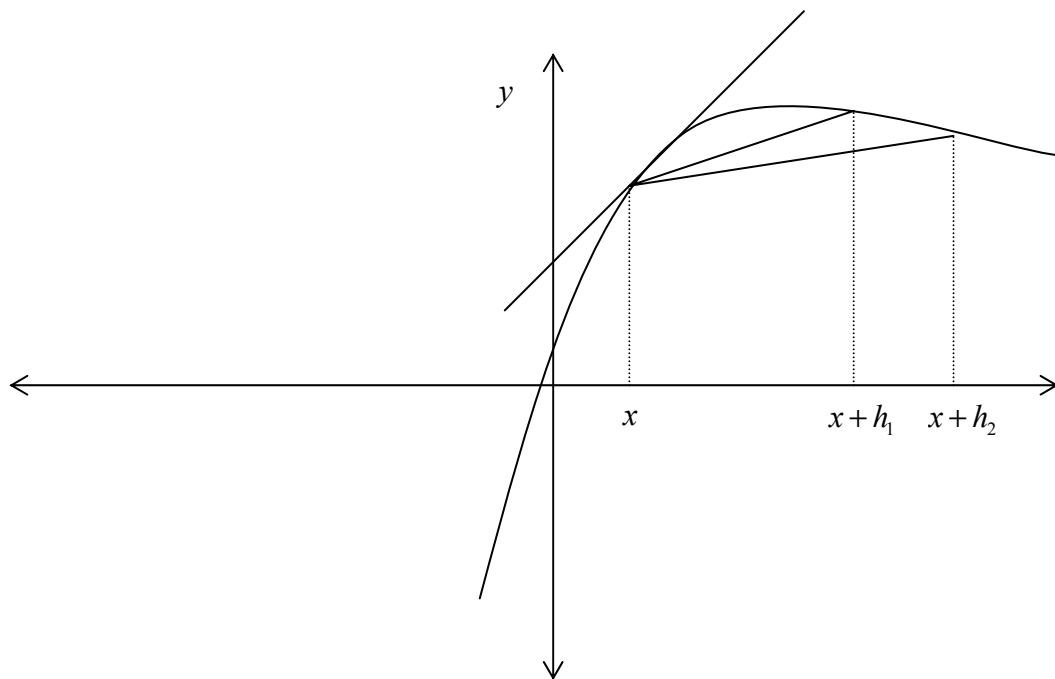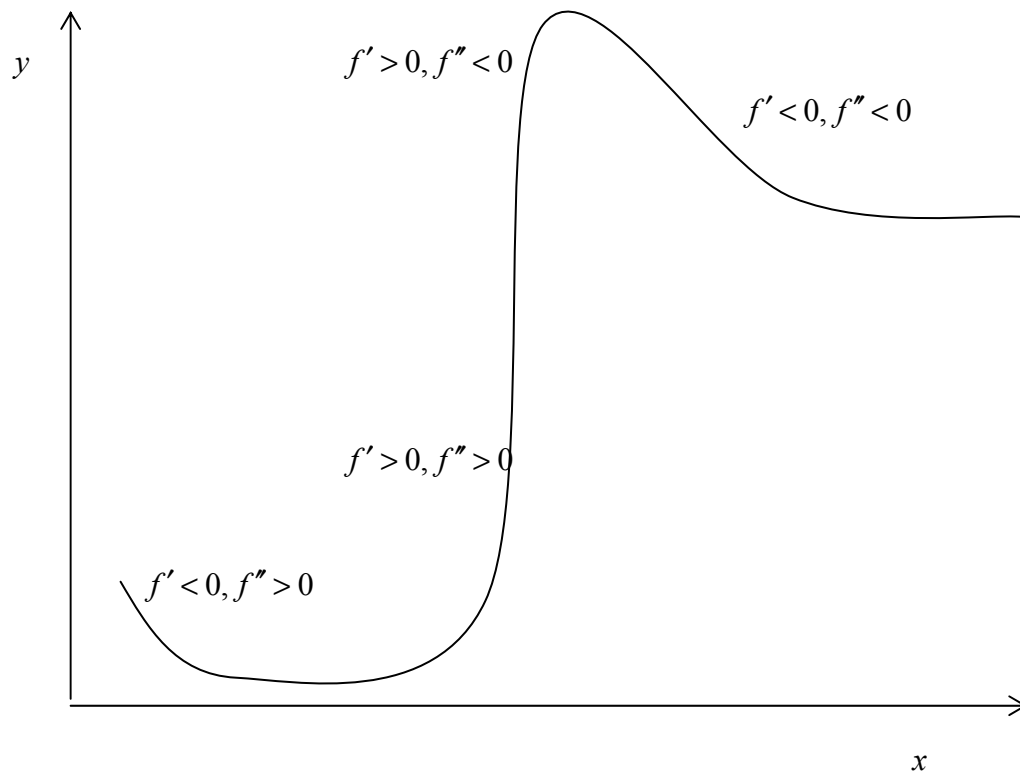**Figure 11.2: Policy Choice in Huber-McCarty Model**

$a + \varepsilon$

Marginal benefit of increasing $p$
$a - p + \varepsilon$

Marginal compliance
cost of increasing $p$
$\delta f\left(p - \bar{p}\right)$

$l_m = 0$

$\bar{p}_1$

$\bar{p}_2 = p_2^*$

$p_1^* = p_3^*$

$a + \varepsilon$

$\bar{p}_3$

**Figure 12.1**
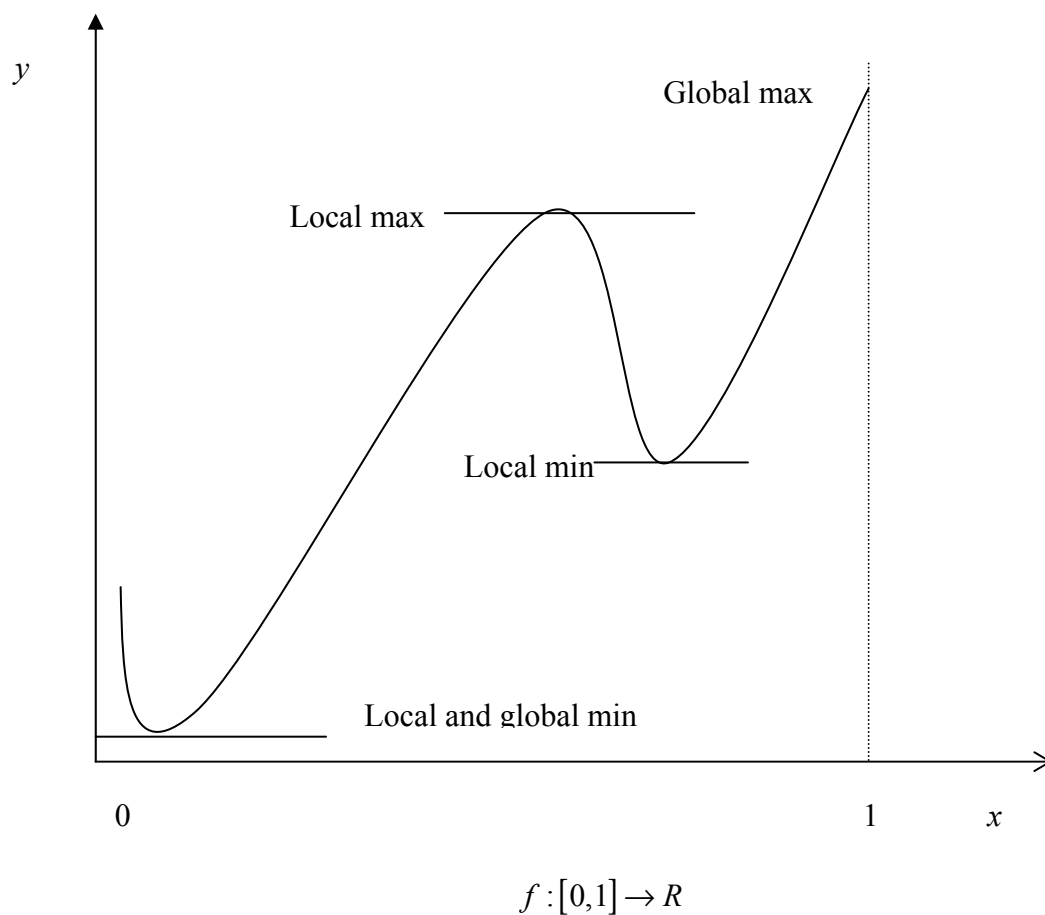**Correspondences**

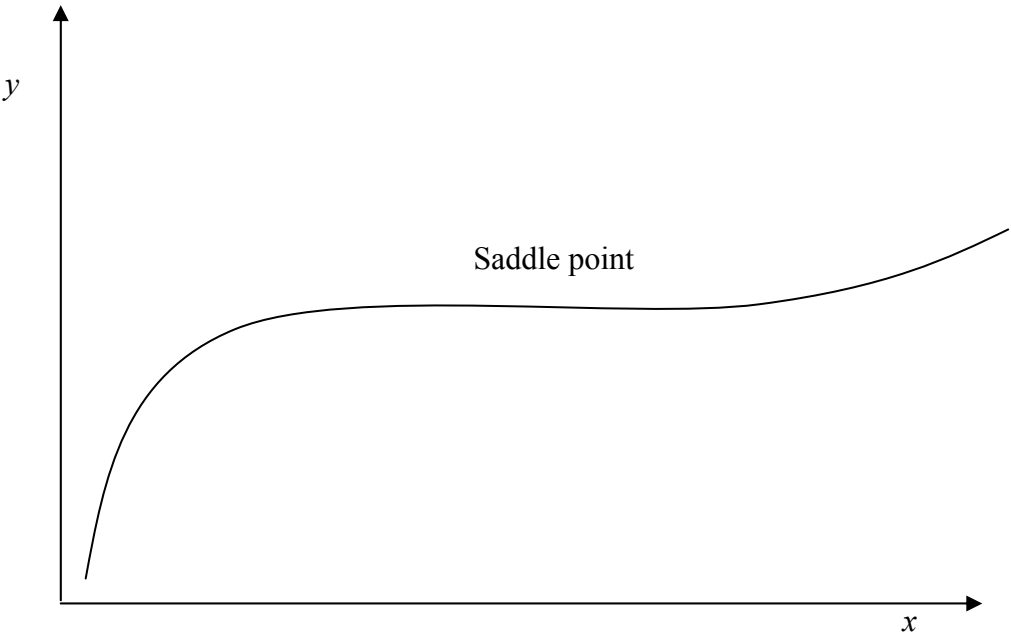**Figure 12.2**
**Functions**

**Figure 12.3**
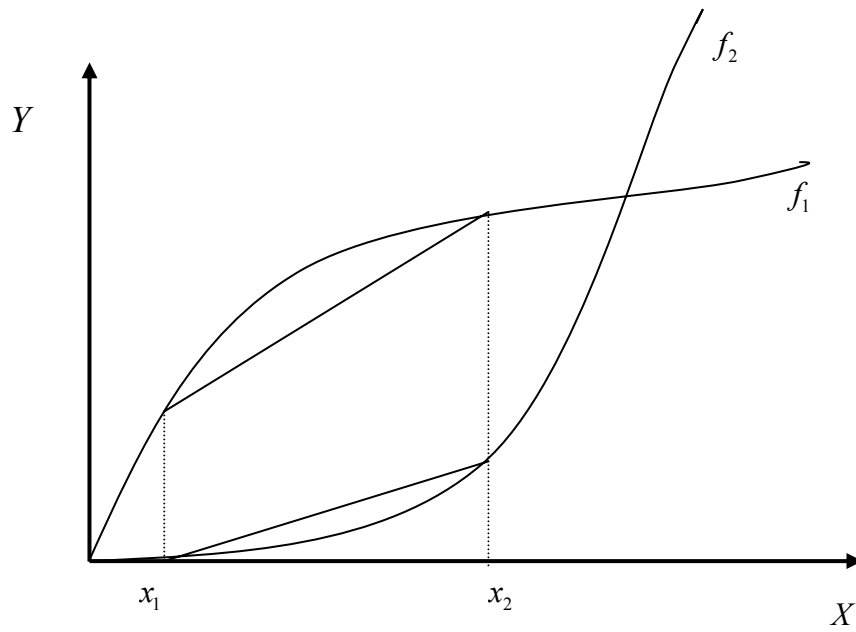**Derivatives**

# Figure 12.4
## Second Derivatives



$f' > 0, f'' < 0$

$f' < 0, f'' < 0$

$f' > 0, f'' > 0$

$f' < 0, f'' > 0$

$y$

$x$

**Figure 12.5**
**Extremum Points**

$f:[0,1] \rightarrow R$

**Figure 12.6**
**Saddle Point**

**Figure 12.7:**
**Convexity and Concavity**

# Figure 12.8

**Discontinuous Cumulative Density Functions**