

Notes on game theory

Sam Wycherley*

November 24, 2023

Latest version available [here](#).

These are a collection of notes covering topics in game theory. They are drawn from a number of sources and I claim no originality to any of the content.¹ The notes are not particularly well-organized, alas. They are already quite long but still work-in-progress. Any errors are my responsibility alone – please do let me know when you find them. I write notes like this because I find it is an effective way for me to learn. However, that means errors are more likely here than in notes written by someone with years of accumulated wisdom. Nonetheless, if you are reading these, I hope you find them useful – Sam.

I also plan on writing up some decision theory notes at some point (game theory is just decision theory when you have to worry about other people).

Some housekeeping: I usually write “wlog” for “without loss of generality” and “iff” for “if and only if”. If I write “s.t.” then I mean “such that” and if I write “wrt”, I mean “with respect to”.

Contents

| | | |
|-------|---------------------------------|----|
| 1 | The structure of games | 5 |
| 1.1 | Extensive form games | 5 |
| 1.2 | Normal form games | 10 |
| 1.3 | Representational equivalence | 10 |
| 1.4 | Mixed strategies | 13 |
| 1.5 | Knowledge | 18 |
| 1.5.1 | The standard model of knowledge | 19 |
| 1.5.2 | Common knowledge | 21 |
| 1.5.3 | Aumann’s agreement theorem | 22 |
| 1.5.4 | No trade? | 24 |

*Email: wycherley@stanford.edu.

¹A few books and figures worth mentioning here: Matt Jackson (Stanford lectures), Ben Brooks (Stanford lectures, particularly on mechanism and information design), Paul Milgrom (his book *Putting auction theory to work*, and some of his Stanford lectures); Hans Peters’ book, particularly for cooperative game theory; Drew Fudenberg & Jean Tirole’s book, of course; notes by Mihai Manea, Jon Levin, Ilya Segal, Arunava Sen, Debasis Mishra and Federico Echenique; many papers, far too numerous to list.

| | | |
|--------|--|-----|
| 2 | Games of complete information | 26 |
| 2.1 | Welfare and efficiency in games | 27 |
| 2.2 | Zero sum and matrix games | 29 |
| 2.3 | Maxmin and minmax | 30 |
| 2.4 | Strict dominance | 34 |
| 2.5 | Correlated rationalizability | 36 |
| 2.6 | Weak dominance | 39 |
| 2.7 | Nash equilibrium | 41 |
| 2.7.1 | Nash equilibrium inefficiency | 47 |
| 2.7.2 | Existence of Nash equilibrium | 47 |
| 2.7.3 | Upper hemicontinuity and Nash equilibria | 50 |
| 2.7.4 | Uniqueness of Nash equilibrium | 51 |
| 2.8 | (Trembling hand) perfect equilibrium | 54 |
| 2.9 | Proper equilibrium | 59 |
| 2.10 | Focal points | 61 |
| 2.11 | Payoff dominance and risk dominance | 62 |
| 2.12 | Correlated equilibrium | 63 |
| 2.13 | Coalition-proof Nash equilibrium | 72 |
| 2.14 | Epistemic foundations of equilibrium | 74 |
| 2.15 | Learning equilibrium | 76 |
| 2.15.1 | Best response dynamics | 76 |
| 2.15.2 | Fictitious play | 78 |
| 2.15.3 | Self-confirming equilibrium | 81 |
| 2.16 | Potential games | 83 |
| 2.16.1 | Congestion games | 89 |
| 2.17 | Supermodular and submodular games | 90 |
| 2.17.1 | Supermodular games and strategic complements | 90 |
| 2.17.2 | Submodular games | 93 |
| 3 | Games of incomplete information | 93 |
| 3.1 | Bayesian games | 93 |
| 3.2 | Dominance in Bayesian games | 98 |
| 3.3 | Ex post equilibrium | 99 |
| 3.4 | Purification | 100 |
| 3.5 | Beliefs, type spaces and the universal type space | 105 |
| 3.6 | Global games | 113 |
| 3.6.1 | Symmetric binary action global games | 113 |
| 4 | Mechanism design | 116 |
| 4.1 | Incentive compatibility and the revelation principle | 117 |
| 4.1.1 | Dominant strategy incentive compatibility | 117 |
| 4.1.2 | Ex post incentive compatibility | 118 |
| 4.1.3 | Bayesian incentive compatibility | 119 |
| 4.1.4 | Limits to the revelation principle | 120 |

| | | |
|-------|---|-----|
| 4.2 | Quasilinear preferences, transfers and private values | 121 |
| 4.2.1 | Groves mechanisms | 123 |
| 5 | Social choice theory | 129 |
| 5.1 | Arrow's impossibility theorem | 132 |
| 5.1.1 | How much of a problem is Arrow's theorem? | 136 |
| 6 | Sequential games | 139 |
| 6.1 | Sequential rationality | 139 |
| 6.1.1 | Backward induction | 140 |
| 6.2 | Subgame perfection | 145 |
| 6.3 | One-shot deviation principle | 150 |
| 6.4 | Perfect Bayesian equilibrium | 153 |
| 6.5 | Sequential equilibrium | 158 |
| 6.6 | Forward induction | 160 |
| 6.7 | Stable equilibria | 164 |
| 6.8 | Noncooperative theory of bargaining | 167 |
| 6.8.1 | Finite period alternating bargaining | 167 |
| 6.8.2 | Infinite period alternating bargaining | 170 |
| 7 | Communication games | 173 |
| 7.1 | Signalling games | 173 |
| 7.1.1 | Job market signalling | 176 |
| 7.1.2 | Forward induction in signalling games | 180 |
| 7.2 | Cheap talk | 184 |
| 7.3 | Bayesian persuasion | 186 |
| 8 | Repeated games | 191 |
| 8.1 | Finitely repeated games | 194 |
| 8.2 | Discounted infinitely repeated games | 197 |
| 8.3 | Undiscounted infinitely repeated games | 199 |
| 8.4 | The folk theorems | 200 |
| 8.4.1 | In discounted infinitely repeated games | 202 |
| 9 | Cooperative game theory | 206 |
| 9.1 | Cooperative games | 206 |
| 9.2 | Cooperative theory of bargaining | 207 |
| 9.2.1 | Nash bargaining solution | 209 |
| 9.2.2 | Raiffa-Kalai-Smorodinsky bargaining solution | 214 |
| 9.3 | Transferable utility games | 215 |
| 9.4 | The core and stable sets | 218 |
| 9.4.1 | The core | 219 |
| 9.4.2 | Stable sets | 223 |
| 9.4.3 | Balanced games | 226 |

| | | |
|--------|--|-----|
| 9.4.4 | Implementing the core | 227 |
| 9.4.5 | Competitive equilibrium and the core | 230 |
| 9.5 | Shapley value | 234 |
| 9.5.1 | Some axiomatic characterizations | 236 |
| 9.5.2 | Harsanyi dividends | 242 |
| 9.5.3 | Multilinear extensions | 243 |
| 9.5.4 | The potential approach | 245 |
| 9.5.5 | Reduced games | 248 |
| 9.5.6 | Myerson value | 249 |
| 10 | Mathematical appendix | 251 |
| 10.1 | Correspondences | 251 |
| 10.2 | Linear programming | 254 |
| 10.2.1 | Hyperplanes, convex sets and extreme points | 254 |
| 10.2.2 | Lemmas of the alternative | 260 |
| 10.2.3 | Duality theorems | 262 |
| 10.3 | Binary relations, ordered sets and lattices | 264 |
| 10.3.1 | Strong set order | 267 |
| 10.4 | Fixed point theorems | 268 |
| 10.5 | Envelope theorems | 273 |
| 10.6 | Monotone comparative statics | 278 |
| 10.6.1 | Supermodularity and increasing differences | 278 |
| 10.6.2 | Supermodular capacities and probability measures | 280 |
| 10.6.3 | Quasi-supermodularity and single crossing properties | 283 |
| 10.6.4 | Infinite supermodularity | 285 |

I use the following scheme for numbering axioms and assumptions:

- (**K***) Knowledge axioms.
- (**G***) Global game assumptions.
- (**P***) Social choice axioms.
- (**A***) One-shot deviation principle assumptions.
- (**B***) Bargaining axioms.
- (**S***) Shapley value axioms.

1 The structure of games

1.1 Extensive form games

Formally defining extensive form games is cumbersome but worth doing. The formal definition is a pretty ugly object, because it has so many components. We first introduce the definition of a finite extensive form game. Finite extensive form games are those most extensively studied in the literature. We then give the more general case.

We first require a very small amount graph theory so we can talk about game trees:

Definition 1 (Graphs).

- (a) *Graph*. A *graph* (V, E) is a pair consisting of a set of *vertices* (or *nodes*) V and a set of *edges* E . An *edge* is a pair of vertices (u, v) , with $u, v \in V$. We write this edge more compactly as uv . We say (V, E) is *directed* if edges are ordered pairs of vertices (the first entry is the *source* and the second is the *destination*), and *undirected* if edges are unordered pairs. For any directed graph (V, E) , the corresponding undirected graph is obtained by treating each edge pair as unordered.

We say that two vertices u, v are *adjacent* if $uv \in E$. We say that edges $e_1 = u_1u_2$ and $e_2 = v_1v_2$ are *adjacent* if $u_2 = v_1$ or $u_1 = v_2$. In an undirected graph, two edges are adjacent if they share a common vertex. In a directed graph, two edges are adjacent if the destination of one of the edges is the source of the other.

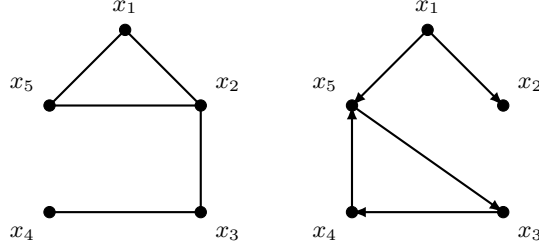
We may sometimes write $uv \in G$ to denote that the edge uv belongs to $G = (V, E)$, i.e. $uv \in E$. We write $G - uv$ for the graph obtained from G by removing the edge uv , and $G + uv$ for the graph obtained from G by adding the edge uv .

- (b) *Trails and paths*. In a graph (V, E) , a *trail* $(v_1v_2, \dots, v_{n-1}v_n)$ is a sequence of distinct adjacent edges such that $v_iv_{i+1} \in E$ for each $i = 1, \dots, n-1$. If, moreover, each vertex that is visited by the trail is visited precisely once, then the trail is called a *path*.
- (c) *Cycles*. A *cycle* is a trail $(v_1v_2, \dots, v_{n-1}v_n)$ such that $v_1 = v_n$. We say that a graph is *acyclic* if it contains no cycles.
- (d) *Connectedness*. An undirected graph is *connected* if for any pair of vertices $u, v \in V$, there exists a path from u to v . A directed graph is connected if the corresponding undirected graph is connected.
- (e) *Trees*. A graph (V, E) is a *tree* if it is connected and acyclic. A *rooted tree* (V, E, v) is a tree equipped with a *root*, which is some node $v \in V$. An *arborescence* is a directed rooted tree in which all edges are directed away from the root. That is, there is precisely one path between the root and any other vertex.

I will also sometimes call graphs *networks*.

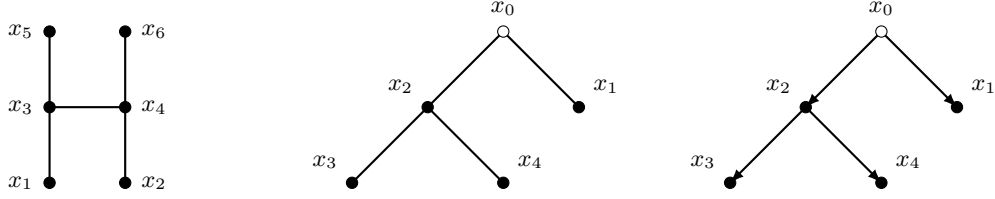
Example 1.

- (a) An undirected graph (left) and a directed graph (right) on five vertices.



Both graphs are connected. The left graph contains one cycle, (x_1x_2, x_2x_5, x_5x_1) , and there is a path from any node to any other. The right graph also contains a cycle: (x_5x_3, x_3x_4, x_4x_5) . There is a path from x_1 to any other node in the right graph, but for all other nodes, this is not true (e.g. there is no path to x_2 from any node but x_1).

- (b) An undirected tree (left), a rooted tree (centre) and an arborescence (right).



Both the rooted tree and the arborescence have root x_0 . We will always denote the root by a hollow node in our diagrams. Note that in the arborescence, any node can be reached via a unique path from x_0 , and there is no other node from which all the other nodes can be reached.

Definition 2 (Finite extensive form game). A *finite extensive form game* Γ is a tuple

$$\Gamma = (\mathcal{I}, T, P, \Phi, \mathcal{A}, (u_i)_{i \in \mathcal{I}}, \eta),$$

where:

- (i) *Players*. \mathcal{I} is a finite index set of n players, $i = 1, \dots, n$, and possibly including nature, conventionally indexed as player 0 or player N .²
- (ii) *Game tree*. The game tree $T = (X, E, x_0)$ is a finite arborescence with root node x_0 , which we refer to as the *initial node*.³ The set of nodes, X , is partitioned as

²Nature randomizes in a given way over a set of ‘states of nature’.

³Kreps (2023) allows the game tree to have multiple initial nodes, so it is not quite an arborescence. There is some probability distribution over the initial nodes but Nature can be a player at subsequent chance nodes. It’s not obvious to me what is gained by defining game trees in this way – we can just make Nature the player at a single initial node.

$X = X_\tau \cup X_d$, where X_τ is the set of *terminal nodes*, the nodes that are not the source of any edge in T . The set X_d is the set of *decision nodes*. Any decision node is the source of some edge in T .

It is convenient to define a partial order on X . We write $x > x'$ if there is a path from x to x' in T , and we say that x *precedes* x' , or x is a *predecessor* of x' . Conversely, we say that x' *succeeds* x or x' is a *successor* of x . If $x > x'$ and there is no node x'' such that $x > x'' > x'$, then we say that x is the *immediate predecessor* of x' , and x' is an *immediate successor* of x . Clearly, every node in the game tree except the initial node has a unique immediate predecessor, and any node that has no successors is a terminal node.

- (iii) *Player allocation function.* The *player allocation function* $P : X_d \rightarrow \mathcal{I}$ is a function assigning each decision node x to a player $P(x) \in \mathcal{I}$. We define the set of decision nodes of player i as $X_i := P^{-1}(i)$. At a decision node, a player makes a choice. Recall we potentially include nature as a player in \mathcal{I} . We refer to decision nodes of nature as *chance nodes*. At a chance node, nature plays randomly according to a given probability distribution.
- (iv) *Information sets.* $\Phi = (\Phi_i)_{i \in \mathcal{I}}$ is the set of information sets. For each player $i \in \mathcal{I}$, each X_i is partitioned into *information sets* $\phi_i \in \Phi_i$. The information sets ϕ_i of player i are disjoint nonempty subsets of X_i and $X_i = \bigcup_{\phi_i \in \Phi_i} \phi_i$. Say a path in T *intersects* an information set ϕ if it visits a node in ϕ . For each information set $\phi_i \in \Phi_i$, we assume every path in T intersects ϕ_i at most once, and every node in ϕ_i is the source of the same number of edges. At an information set ϕ_i , player i cannot distinguish between the nodes in ϕ_i . That is, if $x, x' \in \phi_i$ are distinct nodes then player i can determine that they are at information set ϕ_i , but cannot determine whether they are at node x as opposed to x' . For nature, we assume that $\Phi_0 = \{\{x\} \mid x \in X_0\}$.
- (v) *Actions.* $\mathcal{A} = (A_i)_{i \in \mathcal{I}}$ gives an action assignment correspondence $A_i : \Phi_i \rightarrow 2^A$ for each player i , where A denotes the set of all possible actions in the game. At each information set ϕ_i of i , $A_i(\phi_i)$ is a partition of the set of edges in T that originate at decision nodes in ϕ_i . An *action* $a \in A_i(\phi_i)$ is a set of edges such that for each decision node $x \in \phi_i$, a contains precisely one edge that has x as its source. We assume $|A_i(\phi_i)| \geq 2$, or else there is no choice to be made by i at ϕ_i .
- (vi) *History.* A *history* is any path from the initial node to a decision node or terminal node. A history that terminates at a terminal node is called a *terminal history* or an *outcome path*. An *outcome* is the unique terminal node associated with an outcome path.
- (vii) *Payoff functions.* Payoffs for each player i are defined over each terminal node by the von Neumann-Morgenstern expected utility function $u_i : X_\tau \rightarrow \mathbb{R}$.

- (viii) *Randomization over states of nature.* At each chance node, η assigns a probability distribution over the set of outgoing edges. We assume the probabilities assigned over all edges is positive.

Finite extensive form games have finite action spaces at each decision node. But in many circumstances, we might instead want to model infinitely many actions at a decision node. Often one sees such games represented in a “game tree like” way, with cones used to denote a continuum of possible actions at a decision node. Second, we often want to model interactions that are potentially repeated indefinitely, as in repeated games or in certain bargaining models. A more general formulation of extensive form games, applicable to such circumstances, is as follows.

Definition 3 (Extensive form game). An *extensive form game* Γ is a tuple

$$\Gamma = (\mathcal{I}, \mathcal{H}, P, \Phi, \mathcal{A}, (u_i)_{i \in \mathcal{I}}, \eta),$$

where

- (i) *Players.* \mathcal{I} is a finite index set of n players, $i = 1, \dots, n$, and possibly including nature, conventionally indexed as player 0 or player N .⁴
- (ii) *Histories.* \mathcal{H} is a set of histories. A *history* $h \in \mathcal{H}$ is a sequence consisting of actions. A *subhistory* of a history h is any right-truncation of h . That is, if $h = (a_0, a_1, \dots, a_k)$ then h' is a proper subhistory of h if it is of the form $h' = (a_0, a_1, \dots, a_{k'})$ for $k' \leq k$. If this holds with strict inequality then we say h' is a *proper subhistory* of h . I write $h' \leq h$ if h' is a subhistory of h and $h' < h$ if h' is a proper subhistory of h . If $h' < h$ and there is no h'' such that $h' < h'' < h$, then say that h' is an *immediate predecessor* of h , and h is an *immediate successor* of h' .

A valid set of histories is such that

- (a) $\emptyset \in \mathcal{H}$, where \emptyset denotes the *empty history*;
- (b) If $h \in \mathcal{H}$ then for any proper subhistory $h' \subset h$, we have that $h' \in \mathcal{H}$.

A history $h \in \mathcal{H}$ is a *terminal history* either if it is not a proper subhistory of any other history in \mathcal{H} or if it is an infinite sequence. We denote the set of terminal histories by \mathcal{H}_τ . If h is not a terminal history then it is called a *decision history*. We denote the set of decision histories by \mathcal{H}_d .

- (iii) *Player allocation function.* The *player allocation function* $P : \mathcal{H}_d \rightarrow \mathcal{I}$ is a function assigning each decision history h to a player $P(h) \in \mathcal{I}$. We say that player i is *called to play* at a history h if $i = P(h)$. We denote the set of histories at which i is called to play by $\mathcal{H}_i := P^{-1}(i)$. At each history at which i is called to play, i makes a choice. Recall we potentially include nature as a player in \mathcal{I} . We refer to histories at which nature is called to play as *chance histories*. At a chance history, nature plays randomly according to a given probability distribution.

⁴Again, nature randomizes over a set of ‘states of nature’.

- (iv) *Information sets.* $\Phi = (\Phi_i)_{i \in \mathcal{I}}$ is the set of information sets. For each player $i \in \mathcal{I}$, each \mathcal{H}_i is partitioned into *information sets* $\phi_i \in \Phi_i$. The information sets ϕ_i of player i are disjoint nonempty subsets of \mathcal{H}_i and $\mathcal{H}_i = \bigcup_{\phi_i \in \Phi_i} \phi_i$. Say that a history $h \in \mathcal{H}$ intersects an information set ϕ if there is a subhistory of h that lies in ϕ . For each information set $\phi_i \in \Phi_i$, we assume every history $h \in \mathcal{H}$ intersects ϕ_i at most once. At an information set ϕ_i , player i cannot distinguish between the histories in ϕ_i . That is, if $h, h' \in \phi_i$ are distinct nodes then player i can determine that they are at information set ϕ_i , but cannot determine whether they are being called to play at history h as opposed to h' . For nature, we assume that $\phi_0 = \{\{h\} \mid h \in \mathcal{H}_0\}$.
- (v) *Actions.* $\mathcal{A} = (A_i)_{i \in \mathcal{I}}$ gives an action assignment correspondence $A_i : \Phi_i \rightarrow 2^A$ for each player i , where A is the set of all possible actions in the game. At each information set ϕ_i belonging to i , $A_i(\phi_i)$ is the set of *actions* available to player i at ϕ_i . For any history $h \in \phi_i$, if h' is an immediate successor of h then the last entry in h' is an action $a \in A_i(\phi_i)$.
- (vi) *Payoff functions.* Payoffs for each player i are defined over each terminal history by the von Neumann-Morgenstern expected utility function $u_i : \mathcal{H}_\tau \rightarrow \mathbb{R}$.
- (vii) *Randomization over states of nature.* At each chance history h , η assigns a probability distribution over the set of actions $A_0(\{h\})$. We assume the probabilities assigned over all actions in $A_0(\{h\})$ is positive.

It should be clear that for finite games, this more general definition defines an object equivalent to the finite extensive form definition given in Definition 2, since each node in the finite game tree of Definition 2 uniquely corresponds to a history.

The definition of a history in Definition 3 can be very rich. For example, while we state that a history is a sequence of actions, since nature is potentially a player, we can incorporate e.g. signal realizations. This will come up when we discuss infinitely repeated games, where it is convenient to suppose there exists a public randomization device.

Definition 4 (Pure strategy). In an extensive form game Γ , a pure *strategy* $s_i : \Phi_i \rightarrow A_i$ for player i is a mapping from the set of information sets Φ_i into the set of i 's possible choices $A_i = \bigcup_{\phi_i \in \Phi_i} A_i(\phi_i)$, assigning to each information set ϕ_i an action $s_i(\phi_i) \in A_i(\phi_i)$. The set of pure strategies is denoted $S_i = \times_{\phi_i \in \Phi_i} A_i(\phi_i)$.⁵⁶

Informally, a strategy is a *complete contingent plan of action*.

⁵Given a collection of sets $\{E_k\}$, $\times_k E_k$ denotes the Cartesian product of the sets E_k .

⁶Sometimes, it is said that a strategy s_i must be *measurable* with respect to the information sets. Because we defined the strategy on the information sets, this is trivial in our setting. We could have instead defined strategies on nodes, i.e. so that each s_i is a mapping $s_i : X_i \rightarrow A_i$. Then the measurability requirement really boils down to the requirement that $s_i(x) = s_i(x')$ whenever nodes $x, x' \in \phi_i$ for some information set $\phi_i \in \Phi_i$. For a technical definition, let \mathcal{F}_i to be the σ -algebra generated by Φ_i . Then s_i is measurable with respect to i 's information sets if $s_i^{-1}(a_i) \in \mathcal{F}_i$ for every action $a_i \in A_i$, i.e. if s_i is measurable with respect to \mathcal{F}_i .

1.2 Normal form games

Definition 5 (Normal form game). A *normal form game*, also called a *strategic form game*, is a tuple $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, where:

- (i) \mathcal{I} is a nonempty finite set of players, of cardinality n .
- (ii) S_i is a nonempty set of pure strategies of player $i \in \mathcal{I}$. We denote the typical strategy of player i by $s_i \in S_i$.
- (iii) $u_i : S_i \times S_{-i} \rightarrow \mathbb{R}$ is a von Neumann-Morgenstern expected utility (payoff) function, where $S_{-i} = \times_{j \neq i} S_j$, the Cartesian product of the strategy sets of the j other players with typical element $s_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n) \in S_{-i}$.

Let $S = \times_{i \in \mathcal{I}} S_i$ denote the set of all strategy profiles with typical element $s = (s_1, \dots, s_n)$ and let $u(s) = (u_1(s), \dots, u_n(s))$ be the payoff profile of all players given strategy profile s .

In many cases, we can represent a normal form two-player game by means of a *payoff matrix*, which encodes all the relevant information defining the normal form game. This can be extended beyond two players. Note that the payoff matrix for a normal form game exists only if there is a finite strategy set for each player and a finite number of players.

1.3 Representational equivalence

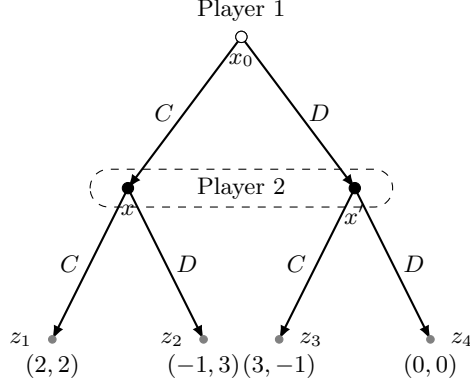
Any extensive form game can be represented in payoff matrix form.

Example 2 (Prisoner's dilemma).

- (a) *Classic prisoner's dilemma*. The classic *prisoner's dilemma* has the following motivation: suppose there are two prisoners (Player 1 and Player 2) who can either confess ("defect" D) or stay silent ("cooperate" C). If Player i confesses and the other stays quiet, then Player i is released (receiving immunity for providing the prosecutor with evidence) and the other player faces a full prison sentence. If both stay silent, then both face a short jail sentence (given the prosecutor's lack of evidence to convict on more serious charges). If both confess then they both face lengthy but reduced sentences. The set of players is $\mathcal{I} = \{1, 2\}$ and the strategy set for each player i is $S_i = \{C_i, D_i\}$. A payoff matrix for this game is:

| | | |
|-------|-------|-------|
| | C_2 | D_2 |
| C_1 | 3, 3 | 0, 4 |
| D_1 | 4, 0 | 1, 1 |

Equivalently, we can represent the prisoner's dilemma in extensive form:



The dotted rectangle denotes Player 2's information set.

- (b) *Repeated prisoner's dilemma.* Suppose the prisoner's dilemma game is played twice, and at the second round, each player observes and remembers the actions of both players in the first round. Suppose the payoffs in each round are as above except that the second round is discounted at some discount rate $\delta \in (0, 1)$. Again, the set of players is $\mathcal{I} = \{1, 2\}$. A strategy for player i now consists of a round 1 action and a round 2 action for each round 1 action profile that could be observed. For example, Player i could have as a strategy $CCDCD$ (“ C in round 1, C in round 2 if (C, C) observed in round 1, D in round 2 if (C, D) observed in round 1, C in round 2 if (D, C) observed in round 1, D in round 2 if (D, D) observed in round 1”). For Player 1, this is a “tit-for-tat” strategy. If both players played tit-for-tat strategies (i.e. Player 1 plays $CCDCD$ and Player 2 plays $CCCDD$) then the payoff profile would be $(2 + 2\delta, 2 + 2\delta)$.

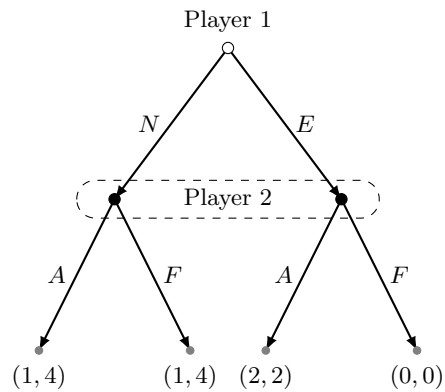
The number of strategies available to each player is $2^5 = 32$. We could, rather tediously, write out a 32×32 payoff matrix for this game. As the number of times the game is repeated increases, the number of strategies will blow up very quickly.

Note that the extensive form representation of a normal form game need not be unique. Consider the following example.

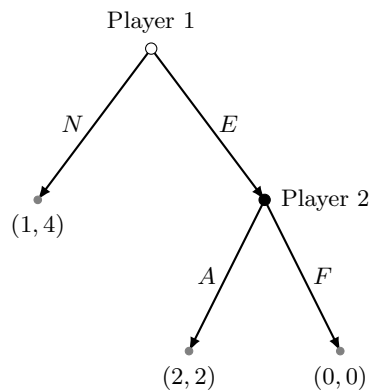
Example 3 (Market entry game I: non-unique extensive form representation). Suppose Player 1 is a potential market entrant and Player 2 is an incumbent monopolist ($\mathcal{I} = \{1, 2\}$). Player 1 chooses between entering the market (E) and not entering (N), so has pure strategy set $S_1 = \{N, E\}$. Simultaneously, Player 2 chooses between acquiescing to the entrant (A) or fighting to try to deter the entrant (F), so Player 2 has strategy set $S_2 = \{A, F\}$. There are thus four potential strategy profiles: $S = \{(E, A), (E, F), (N, A), (N, F)\}$. Whether Player 2 acquiesces or fights does not matter for payoffs if Player 1 chooses not to enter. Suppose we have payoff matrix,

| | | |
|-----|------|------|
| | A | F |
| N | 1, 4 | 1, 4 |
| E | 2, 2 | 0, 0 |

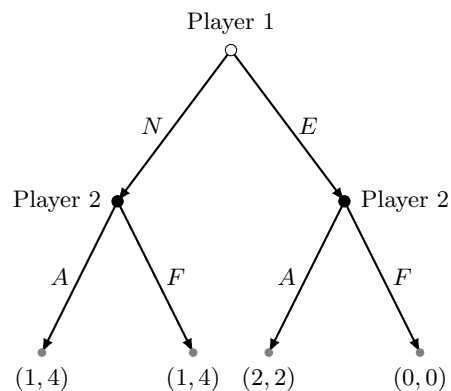
There are two equivalent ways to represent this in extensive form. We can represent this as (a) a *simultaneous move game*,



or (b) as a *chain store game*,



These are equivalent because the payoff if Player 1 plays *N* do not depend on Player 2's choice. Note the strategy set is the same in all the above representations (normal, simultaneous-move extensive and chain store extensive). Indeed, were we to instead have the *sequential move game*,



this would *not* be equivalent to the normal form game, for the strategy set of Player 2, $S_2 = \{(A, A), (F, A), (A, F), (F, F)\}$ is no longer that of the normal form game.

The extensive form representation of a game typically conveys more information about the game than the normal form representation. For dynamic games, such as repeated games or games with multiple stages, the extensive form representation is usually easier to interpret.⁷

1.4 Mixed strategies

Definition 6 (Randomized strategies). For any topological space A , let \mathcal{B} denote the Borel σ -algebra on A . Let $\Delta(A)$ be the set of all probability measures on \mathcal{A} :

$$\Delta(A) := \{\sigma : \mathcal{B} \rightarrow \mathbb{R} \mid \sigma \text{ is a measure and } \sigma(A) = 1\}.$$

In general, we will assume sets are endowed with their natural topology. If A is countable, this is the discrete topology, so $\mathcal{B} = 2^A$, the power set. If A is a connected subset of \mathbb{R}^k (such as an interval in the real line, or a box in \mathbb{R}^2), then the natural topology is induced by the Euclidean metric and the Borel σ -algebra is the Lebesgue σ -algebra restricted to A .

Recall that for σ to be a measure, we need that $\sigma \geq 0$, $\sigma(\emptyset) = 0$, and for any disjoint countable collection $\{E_n\}_{n=1}^\infty \subseteq \mathcal{B}$, we have $\sigma(\bigcup_{n=1}^\infty E_n) = \sum_n \sigma(E_n)$ (countable additivity).

If A is finite with $|A| = k$, we can more straightforwardly define $\Delta(A)$ as a simplex in \mathbb{R}^k . Then each $\sigma \in \Delta(A)$ is a probability vector, and we can interpret the value of the i th entry in σ as the probability that the i th action in A is played.

Now consider a game Γ with finite set of players $\mathcal{I} = \{1, \dots, n\}$ and where each player $i \in \mathcal{I}$ has strategy set S_i consisting of pure strategies s_i .

- (a) *Mixed strategy.* A *mixed strategy* is a probability measure $\sigma_i \in \Delta(S_i)$.

If S_i is countable, it is more convenient to focus on the probability measure of the singletons in the σ -algebra we associate with S_i . That is, abusing notation, we define $\sigma_i(s_i) = \sigma_i(\{s_i\})$, so we have $\sigma_i : S_i \rightarrow [0, 1]$ satisfying $\sum_{s_i \in S_i} \sigma_i(s_i) = 1$. This is of course a probability mass function $\sigma_i \in \Delta(S_i)$.

A mixed strategy profile is a profile $(\sigma_1, \dots, \sigma_n) \in \Delta(S_1) \times \dots \times \Delta(S_n)$.

We define the space of mixed strategies of opponents by $\Delta_{-i}(S_{-i}) = \times_{j \neq i} \Delta(S_j)$.

- (b) *Mixed strategy game.* Given a pure strategy game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a mixed strategy game is (in general) the game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$.

⁷These differences in information and the fact that the same normal form game can have non-unique extensive form representations can pose a problem. For some solution concepts, making seemingly irrelevant changes to the extensive form representation of a normal form game can lead to differing solutions. See e.g. Example 49.

- (c) *Correlated strategy*. In general, a *correlated strategy* is a strategy profile in $\Delta(S_1 \times \cdots \times S_n) = \Delta(S)$.

This is a generalization of a mixed strategy profile, for it allows correlation between players' strategies. Indeed, $\Delta(S_1) \times \cdots \times \Delta(S_n) \subseteq \Delta(S)$.

- (d) *Behavioural strategy*. In general, a *behavioural strategy* $\pi_i : \Phi_i \rightarrow \bigcup_{\phi_i \in \Phi_i} \Delta(A_i(\phi_i))$ of player i is a function assigning to each information set $\phi_i \in \Phi_i$ a probability measure over the action set at that information set, $\pi_i(\phi_i) \in \Delta(A_i(\phi_i))$.

A mixed strategy σ_i and a behavioural strategy π_i are called (outcome) *equivalent* if for every (mixed or behavioural) strategy profile σ_{-i} and any node x in the extensive form game Γ , it follows that

$$P_{\sigma_i, \sigma_{-i}}(x) = P_{\pi_i, \sigma_{-i}}(x),$$

where $P(x)$ denotes the probability that node x is reached.

We can equivalently define a behavioural strategy of player i in terms of the set of histories \mathcal{H}_i that terminate at an information set of player i . For each $h_i \in \mathcal{H}_i$, let $\phi_i(h_i)$ denote the information set at which h_i terminates. Then a behavioural strategy is a function $\pi_i : \mathcal{H}_i \rightarrow \bigcup_{\phi_i \in \Phi_i} \Delta(A_i(\phi_i))$ that assigns to each history $h_i \in \mathcal{H}_i$ a probability measure $\pi_i(h_i) \in \Delta(A_i(\phi_i(h_i)))$, with the condition that if $\phi_i(h_i) = \phi_i(h'_i)$ then $\pi_i(h_i) = \pi_i(h'_i)$. We will stick to the definition given in (d) unless otherwise stated (when it comes to repeated games we switch to the history version).

Note that any pure strategy is a degenerate mixed strategy (that is, a mixed strategy with all probability mass placed on a single pure strategy.) I abuse notation sometimes and write the pure strategy profile where I mean a degenerate mixed strategy profile.

It is common in treatments of mixed strategies to restrict exposition to mixed strategies over finite strategy sets. I have included the more general definition of a mixed, correlated and behavioural strategy above, defined in terms of a probability measure over a Borel σ -algebra. This immediately allows for mixed strategies over continuous strategy sets, where our simplified definition for countable sets is of little use.

Conversely, any mixed strategy over S_i is a pure strategy in $\Delta(S_i)$. Indeed, mixed strategy extensions of pure strategy games are equivalent to a certain class of pure strategy games with continuous strategy spaces.

For any correlated strategy profile $\sigma \in \Delta(S)$, the expected payoff to player i is $\int_S u_i d\sigma$. Since a mixed strategy profile is a product measure, for any mixed strategy profile $\sigma = (\sigma_1, \dots, \sigma_n) \in \times_{i \in \mathcal{I}} \Delta(S_i)$ the expected payoff is

$$\begin{aligned} u_i(\sigma_i, \sigma_{-i}) &= \int_S u_i(s) d\sigma \\ &= \int_{S_i \times S_{-i}} u_i(s_i, s_{-i}) d(\sigma_i \times \sigma_{-i}) \\ &= \int_{S_1 \times \cdots \times S_n} u_i(s_i, s_{-i}) d(\sigma_1 \times \cdots \times \sigma_n). \end{aligned}$$

Specializing to the countable case, the expected payoff to player i of correlated strategy profile $\sigma = (\sigma_1, \dots, \sigma_n) \in \Delta(S)$ is $u_i(s) = \sum_{s \in S} u_i(s) \sigma(s)$, and for any mixed strategy profile $\sigma = (\sigma_1, \dots, \sigma_n) \in \prod_{i \in \mathcal{I}} \Delta(S_i)$, the expected payoff is

$$\begin{aligned} u_i(\sigma_i, \sigma_{-i}) &= \sum_{s \in S} u_i(s) \sigma(s) \\ &= \sum_{(s_i, s_{-i}) \in S_i \times S_{-i}} u_i(s_i, s_{-i}) \sigma_i(s_i) \sigma_{-i}(s_{-i}) \\ &= \sum_{(s_1, \dots, s_n) \in S_1 \times \dots \times S_n} u_i(s_i, s_{-i}) \prod_{i \in \mathcal{I}} \sigma_i(s_i). \end{aligned}$$

Lemma 1. *For any mixed or behavioural strategy profile σ_{-i} , if mixed strategy σ_i is equivalent to behavioural strategy π_i then*

$$u_i(\sigma_i, \sigma_{-i}) = u_i(\pi_i, \sigma_{-i}).$$

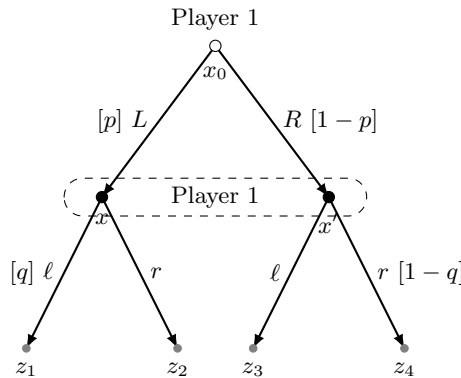
Proof. Follows immediately from the definition of equivalence and the expected payoff. \square

Definition 7 (Perfect recall). A game Γ is a game of *perfect recall* if every player in every stage of the game recalls all their previous information and every previous action they have taken.

Given a decision node x , let $\mathcal{E}_i(x)$ denote the the chronologically-ordered list of information sets encountered by player i on the path from initial node x_0 to x and which action i has taken at each such information sets. We call $\mathcal{E}_i(x)$ the *experience* of player i at x . Formally, a game is one of perfect recall if for every player i and every information set ϕ_i of i , we have that $\mathcal{E}_i(x) = \mathcal{E}_i(x')$ for all $x, x' \in \phi_i$.

Example 4 (Games without perfect recall).

(a) Player 1 cannot remember their first move:



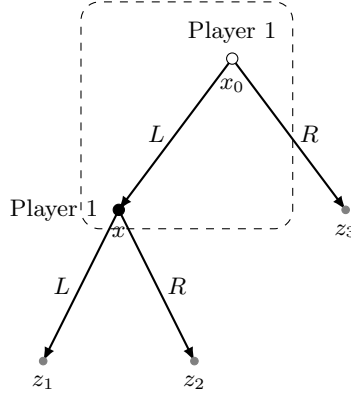
We have $\mathcal{E}_i(x) = \{x_0, L\} \neq \mathcal{E}_i(x') = \{x_0, R\}$, yet $x, x' \in \phi_1$.

The strategy set for Player 1 is $S_1 = \{L\ell, Lr, R\ell, Rr\}$. Now consider the mixed strategy $\sigma_1 = (1/2, 0, 0, 1/2)$ and consider an outcome-equivalent behavioural strategy σ_1^b s.t. Player 1 randomizes $(p, 1-p)$ over $\{L, R\}$ at node x_0 and $(q, 1-q)$ over $\{\ell, r\}$ at information set ϕ_1 . We have

$$\begin{aligned} P_{\sigma_1^b}(z_1) &= pq = \frac{1}{2} = P_{\sigma_1}(z_1), \\ P_{\sigma_1^b}(z_2) &= p(1-q) = 0 = P_{\sigma_1}(z_2), \\ P_{\sigma_1^b}(z_3) &= (1-p)q = 0 = P_{\sigma_1}(z_3), \\ P_{\sigma_1^b}(z_4) &= (1-p)(1-q) = \frac{1}{2} = P_{\sigma_1}(z_4). \end{aligned}$$

We have $q = \frac{1}{2p}$ and hence $p(1-q) = p(1 - 1/2p) = p - \frac{1}{2} = 0$ so $p = \frac{1}{2}$, $q = 1$. Yet this implies $(1-p)(1-q) = 0 \neq \frac{1}{2}$, yielding a contradiction. Hence there is no behavioural strategy σ_1^b that is equivalent to σ_1 .

(b) Suppose Player 1 forgets they make a first move:



We have $\mathcal{E}_i(x_0) = \emptyset \neq \mathcal{E}_i(x) = \{x_0, L\}$, but $x_0, x \in \phi_1$. The strategy set is $S_1 = \{L, R\}$. Now, any behavioural strategy assigning positive probability to both L and R can attain outcomes z_1, z_2, z_3 , each with positive probability. However, there is no mixed strategy that can attain z_2 with positive probability.

In games of perfect recall, to every mixed strategy there exists an outcome-equivalent behavioural strategy. Kuhn (1950, 1953) first proved this result for finite strategy sets, but it extends to (countably or uncountably) infinite strategy sets also.

Theorem 1 (Kuhn, 1953). *In every game of perfect recall, every mixed strategy has an equivalent behavioural strategy.*

Proof. We prove this for finite strategy sets, per Kuhn (1953). Aumann (1964) proved that the theorem also holds for infinite strategy sets.

We require that for any mixed strategy σ_i , there exists a behavioural strategy π_i s.t. for any (mixed or behavioural) strategy profile σ_{-i} ,

$$P_{\sigma_i, \sigma_{-i}}(x) = P_{\pi_i, \sigma_{-i}}(x)$$

for all nodes x . Now, for any x that lies at an information set irrelevant for σ_i (that is, not reached with positive probability given strategy σ_i for some σ_{-i}) then both sides are zero. Hence consider only those information sets ϕ_i that are relevant for σ_i . For each node x , let $R_i(\phi_i) = \{s_i \in S_i \mid \phi_i \text{ is on the path of } (s_i, s_{-i}) \text{ for some } s_{-i} \in S_{-i}\}$ be the set of relevant information sets for ϕ_i . Now for each relevant ϕ_i for σ_i , define the probability of playing move $a \in \phi_i$ under π_i by

$$\pi_i(a \mid \phi_i) = \frac{\sum_{s_i \in R_i(\phi_i): s_i(\phi_i)=a} \sigma_i(s_i)}{\sum_{s_i \in R_i(\phi_i)} \sigma_i(s_i)}.$$

Let $\phi_i^1, \dots, \phi_i^{\bar{k}}$ denote player i 's information sets preceding ϕ_i . Under perfect recall, reaching ϕ_i requires that i takes the appropriate action a^k at each ϕ_i^k :

$$R_i(\phi_i) = \{s_i \mid s_i(\phi_i^k) = a^k \text{ for all } k = 1, \dots, \bar{k}\}.$$

Now, conditional on reaching ϕ_i , the distribution of continuation play is

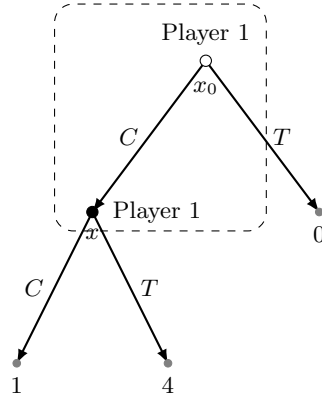
$$\pi_i(a \mid \phi_i) = \frac{\sum_{s_i \mid s_i(\phi_i^k)=a^k \text{ for all } k=1, \dots, \bar{k} \text{ and } s_i(\phi_i)=x} \sigma_i(s_i)}{\sum_{s_i \mid s_i(\phi_i^k)=a^k \text{ for all } k=1, \dots, \bar{k}} \sigma_i(s_i)},$$

which is precisely the probability of playing action a under σ_i conditional on reaching information set ϕ_i .

Since this holds for all ϕ_i , it follows that $P_{\sigma_i, \sigma_{-i}}(x) = P_{\pi_i, \sigma_{-i}}(x)$ for all nodes x . \square

Unless otherwise specified, we assume all games are games of perfect recall. Without perfect recall, one has to be very careful in interpreting information sets, strategies etc. The optimal strategy ceases to be obvious, because different principles of optimality can conflict, as the following example shows.

Example 5 (Paradox of the absent-minded driver; Piccione and Rubinstein, 1995). A weary traveller is sat in a pub planning their journey home. If they take the first exit on the way home, they will end up in a dangerous area (payoff 0). Ideally, they would to take the second exit, which leads them to their home (payoff 4). If they miss the second exit, they cannot turn around and must continue to the end of the road, where they can stay the night at a hotel (payoff 1). The driver is very absentminded, so when arriving at an exit, cannot tell if it is the first one or how many they have already passed. The driver knows this fact when planning the trip. The game can be represented as follows, where action T is taking the exit and C is continuing:



Suppose first that the driver cannot randomize. Then when planning the trip, they should plan to continue to the end of the road, receiving payoff 1 for if they instead choose to turn off when encountering an exit, they would turn off at the first exit and receive payoff 0. Yet on the road, when encountering an exit, the driver reasons there is probability $\frac{1}{2}$ that the exit is the second exit given their strategy. If they choose to turn off, their expected payoff is 2, and so it is optimal to take the exit. Hence the optimal strategy is time inconsistent.

There is a paradox here: on the one hand, the optimality of the ex ante optimal strategy should not require verification at execution if tastes or information have not changed; on the other, maximizing the driver's expected payoff given their beliefs at each stage leads them to deviate from the ex ante optimal strategy when at any exit. If pursuing this, the driver will take the first exit, and thus receive payoff 0.

The paradox persists if we allow the driver to randomize. The optimal behavioural strategy involves staying on the road with probability $\frac{2}{3}$ whenever encountering an exit. Now suppose p is the probability the driver does not exit, and α is the probability the driver assigns to being at the first exit. Then the expected payoff is $\alpha(p^2 + 4(1-p)p) + (1-\alpha)(p + 4(1-p))$, and so the optimal choice of p is $\max\left\{0, \frac{(7\alpha-3)}{6\alpha}\right\}$, which only gives $p = \frac{2}{3}$ if $\alpha = 1$, i.e. if the driver is certain they have not passed the first exit, which is implausible given the driver has no capacity to remember whether they have done so.

1.5 Knowledge

“Reports that say that something hasn’t happened are always interesting to me, because as we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns – the ones we don’t know we don’t know.” –

Donald Rumsfeld, 2004, on failing to find Iraq’s nonexistent WMD.

Reasoning about games requires developing some epistemic foundations. For us to model how players reason, we need to carefully set out what they know about the structure of the game, about what they know about each other, and so on. How much information agents have will alter their behaviour. For example, a monopolist is able to practice first-degree price discrimination if she knows her customers’ individual demand

schedules, whereas if she only knows the aggregate demand schedule, she must set a price uniformly.

“Higher order” knowledge also often matters. For example, suppose Firm 1 and Firm 2 are both profit-maximizing duopolists. Firm 1 can produce at a marginal cost known to both firms but Firm 2 has two possible types – either its costs are high, and so it produces a relatively smaller level of output, or its costs are low, in which case it produces a higher level of output. If Firm 1 does not know Firm 2’s type, then it has to set its output based only on its estimate of Firm 2’s output, based on its own beliefs about Firm 2’s type. Suppose instead that Firm 1 learns Firm 2’s output before setting its own. If Firm 2 does not know that Firm 1 knows its own output, then this puts Firm 1 in a much better position – whatever Firm 2’s output, Firm 1 can tailor its own output to maximize its own profits. However, if Firm 2 knows that Firm 1 knows its output, then Firm 2 is in a much better position, because it becomes a Stackelberg leader – it can set higher output itself knowing that Firm 1 must optimally choose a lower output in response.

1.5.1 The standard model of knowledge

The standard model of knowledge does not claim to accurately describe knowledge as we would understand it in everyday use, but since we are not philosophers, this is unimportant – the model works well enough to capture some interesting insights.⁸

Definition 8 (States, information and knowledge).

- (a) *States and events.* Let (Ω, \mathcal{F}, p) be a probability space. The elements $\omega \in \Omega$ are called *states* and the sets $E \in \mathcal{F}$ are called *events*.
- (b) *Information.* An *information function* for Ω is a function $h : \Omega \rightarrow \mathcal{F}$, associating to each state $\omega \in \Omega$ a nonempty set $h(\omega) \in \mathcal{F}$.

We call any partition \mathcal{P} of Ω such that each member of \mathcal{P} lies in \mathcal{F} an *information partition*.

We say that h is *partitional* if $\{h(\omega) \mid \omega \in \Omega\}$ is an information partition and $\omega \in h(\omega)$ for each $\omega \in \Omega$. We say h corresponds to information partition \mathcal{P} if $P = \{h(\omega) \mid \omega \in \Omega\}$.

Equivalently, h is partitional if

- (I1) $\omega \in h(\omega)$ for all $\omega \in \Omega$ and
- (I2) if $\omega' \in h(\omega)$ implies $h(\omega') = h(\omega)$.

- (c) *Knowledge.* A *knowledge operator* for an agent is a function $K : \mathcal{F} \rightarrow \mathcal{F}$ defined by

$$K(E) = \{\omega \in \Omega \mid h(\omega) \subseteq E\}.$$

⁸The analytic philosophy literature attempting to define knowledge is, like most things in analytic philosophy, a Sisyphean nightmare. I recommend avoiding.

We say that K is the knowledge operator *induced* by information function h or information partition \mathcal{P} .

(I1) has the interpretation that an agent in state ω cannot rule out being in state ω . Under (I2), if ω' were also deemed possible by the agent in state ω , then it must be that the set of states the agent would consider possible in states ω and ω' are the same.

The knowledge operator captures under which states an agent knows a given event. We say that the agent *knows* event E in state ω if $\omega \in K(E)$.

Example 6. Suppose $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$, and that we have partition $\{\{\omega_1, \omega_2\}, \{\omega_3\}, \{\omega_4\}\}$ associated with the partitional information function h . We have e.g. $h(\omega_1) = h(\omega_2) = \{\omega_1, \omega_2\}$ and $h(\omega_3) = \{\omega_3\}$. The corresponding knowledge operator K has $K(\{\omega_3, \omega_4\}) = \{\omega_3, \omega_4\}$, $K(\{\omega_1, \omega_3\}) = \{\omega_3\}$, and $K(\{\omega_1\}) = \emptyset$.

Suppose instead that $\Omega = \{\omega_1, \omega_2\}$, that $h(\omega_1) = \{\omega_1\}$ but $h(\omega_2) = \{\omega_1, \omega_2\}$. Then h is not partitional. We would have $K(\{\omega_1\}) = \{\omega_1\}$, $K(\{\omega_2\}) = \emptyset$ and $K(\{\omega_1, \omega_2\}) = \{\omega_1, \omega_2\}$.

Knowledge operators derived from partitional information functions satisfy the following axioms:

Axioms.

- (K1) *Awareness.* $K(\Omega) = \Omega$.
- (K2) *Omniscience.* $K(E \cap F) = K(E) \cap K(F)$ for all $E, F \in \mathcal{F}$.
- (K3) *Knowledge.* $K(E) \subseteq E$ for all $E \in \mathcal{F}$.
- (K4) *Transparency.* $K(K(E)) = K(E)$ for all $E \in \mathcal{F}$.
- (K5) *Wisdom.* $\neg K(E) = K(\neg K(E))$ for all $E \in \mathcal{F}$.

This axiomatic characterization was introduced by Milgrom (1981). All of them are somewhat contentious – see Samuelson (2004) who discusses this in some detail.

The axiom of awareness (K1), one of the more innocuous axioms, states that the agent always knows she is in some state, or equivalently knows the set of possible states.

The axiom of omniscience (K2) is so-named because despite its appeal – knowing E and F implies knowing E and knowing F – it carries quite severe implications. Note that under the axiom, $E \subseteq F$ implies that $K(E) \subseteq K(F)$. Now $E \subseteq F$ carries the interpretation that E implies F , and so if the agent knows E , she also knows F . It follows that if we were to, say, explain the rules of chess to the agent, she would know the optimal strategy to win chess.

The next two axioms are less contentious. The axiom of knowledge (K3) states that the agent knows E only if E happens. The axiom of transparency (K4) states that if the agent knows E then she also knows that she knows E . The axiom implies that she also knows that she knows that she knows E and so on. Indeed, $K(E)$ implies $K^n(E)$ for all $n \in \mathbb{N}$.

Finally, the axiom of wisdom (**K5**) implies that if the agent does not know an event then she knows she does not know it. This rules out that the agent can be unaware of any possibilities.

Proposition 1 (Bacharach, 1985). *A knowledge operator K satisfies the axioms (**K1**)-(**K5**) iff it is induced by a partitional information function.*

Proof. First suppose K is induced by a partitional information function h . Since $h(\omega) \subseteq \Omega$ for all $\omega \in \Omega$, and thus $K(\Omega) = \Omega$, (**K1**) holds. Suppose $E, F \in \mathcal{F}$. If $h(\omega) \subseteq E \cap F$, then $h(\omega) \subseteq E$ and $h(\omega) \subseteq F$. Conversely, if $h(\omega) \subseteq E$ and $h(\omega) \subseteq F$, then $h(\omega) \subseteq E \cap F$. Thus $\omega \in K(E \cap F)$ iff $\omega \in K(E) \cap K(F)$, so (**K2**) holds. Third, if $h(\omega) \subseteq E$ then $\omega \in E$ by (I1), and so $K(E) \subseteq E$, i.e. (**K3**) holds. Fourth, by (I2), for any $\omega' \in h(\omega)$, we have that $h(\omega') \subseteq E$ also, and so $h(\omega) \subseteq K(E)$. Thus if $\omega \in K(E)$ then $\omega \in K(K(E))$, so (**K4**) holds. Finally, $\neg K(E) = \{\omega \in \Omega \mid h(\omega) \not\subseteq E\}^c$. If $h(\omega) \cap E = \emptyset$, then $h(\omega) \subseteq E^c \subseteq \neg K(E)$, and so $\omega \in K(\neg K(E))$. If $h(\omega) \cap E \neq \emptyset$ but $h(\omega) \not\subseteq E$, then $K(E) \cap h(\omega) = \emptyset$ by (I1)-(I2), and hence $h(\omega) \subseteq \neg K(E)$, so $\omega \in K(\neg K(E))$. It follows that (**K5**) holds.

See Bacharach (1985) for proof of the converse, though a caution that his notation is messy. \square

1.5.2 Common knowledge

Definition 9 (Common knowledge). Consider a finite set \mathcal{I} of n agents with partitional information functions h_1, \dots, h_n and corresponding knowledge operators K_1, \dots, K_n .

- (a) *Mutual knowledge.* We say an event $E \in \mathcal{F}$ is *mutual knowledge* in state $\omega \in \Omega$ if $\omega \in \bigcap_{i=1}^n K_i(E)$, that is, if E is known to all agents in state ω . If E is mutual knowledge, we write $K^1(E)$.

Recursively, let $K^k(E) := K^1(K^{k-1}(E))$. Hence $K^2(E)$ denotes that mutual knowledge of E is mutual knowledge, and so on

- (b) *Common knowledge.* We say an event $E \in \mathcal{F}$ is *common knowledge* in state $\omega \in \Omega$ if $\omega \in \bigcap_{k=1}^{\infty} K^k(E)$. That is, if E is known to all agents, all agents know E is known to all agents, and so on.
- (c) *Self-evident events.* We say an event $E \in \mathcal{F}$ is *self-evident* if for all $\omega \in E$ and all $i \in \mathcal{I}$, we have $h_i(\omega) \subseteq E$.

If Ω is finite, an equivalent definition of common knowledge can be stated in terms of self-evident events.

Lemma 2. *The following statements are equivalent for any $E \in \mathcal{F}$:*

- (i) $K_i(E) = E$ for all $i \in \mathcal{I}$;
- (ii) E is a self-evident event;

(iii) for all $i \in \mathcal{I}$, E is a union of members of the partition of Ω induced by h_i .

Proof. E is self-evident iff $E \subseteq K_i(E)$ for all $i \in \mathcal{I}$, and by **(K3)**, $K_i(E) \subseteq E$. Hence (i) and (ii) are equivalent. Now, if E is self-evident then by definition, $\omega \in E$ implies $h(\omega) \subseteq E$, so $E = \bigcup_{\omega \in E} h_i(\omega)$ for each $i \in \mathcal{I}$. Thus (ii) implies (iii). If E is a union of the members of the partition induced by h_i , then $h_i(\omega) \subseteq E$ iff $\omega \in E$, and so (iii) implies (i). \square

Proposition 2. Suppose Ω is a finite set. An event $E \in \mathcal{F}$ is common knowledge in a state $\omega \in \Omega$ iff there exists a self-evident event $F \subseteq E$ such that $\omega \in F$.

Proof. Suppose E is common knowledge in state ω . By the axiom of transparency **(K4)**, $E \supseteq K^1(E) \supseteq K^2(E) \supseteq \dots$, and $\omega \in E$, $\omega \in K^k(E)$ for all $k \in \mathbb{N}$. Given Ω is finite, there exists a set $F = K^k(E)$ for some k such that $K_i(F) = F$ for all $i \in \mathcal{I}$. Now $F \subseteq E$ and $\omega \in F$, and F is self-evident since $F \subseteq K_i(F)$ for all i .

Conversely, suppose $E \in \mathcal{F}$ is such that there exists some self-evident event $F \subseteq E$ with $\omega \in F$. Then $K_i(F) = F$ for all $i \in \mathcal{I}$, so $K^1(F) = F$ and thus $K^1(F)$ is self-evident. Iterating, it follows that $K^k(F) = F$ is self-evident for all $k \in \mathbb{N}$. Since $F \subseteq E$, by the axiom of omniscience **(K2)**, we have that $F \subseteq K^k(E)$ for any k , and $\omega \in F$. Thus $\omega \in K^k(E)$ for all k , so $\omega \in \bigcup_{k=1}^{\infty} K^k(E)$, and hence E is common knowledge in state ω . \square

Common knowledge plays an important role in classifying games:

Definition 10 (Games and information).

- (a) *Perfect and imperfect information.* An extensive form game Γ is a game of *perfect information* if all information sets are singletons; that is, if all past moves are common knowledge. Otherwise, Γ is a game of *imperfect information*.
- (b) *Structure.* The *structure* of a game G (Γ) consists of all elements listed in the normal form tuple G or in the extensive form tuple Γ respectively.
- (c) *Complete and incomplete information.* A (normal or extensive form) game Γ is a game of *complete information* if the structure of the game is common knowledge. Otherwise, Γ is a game of *incomplete information*.

1.5.3 Aumann's agreement theorem

The formalization of common knowledge has spawned several famous, possibly surprising theorems. Famously, *Aumann's agreement theorem* says that no two agents with a common prior, whose posterior beliefs are common knowledge, will ever disagree about an event.

Definition 11. If \mathcal{P}_1 and \mathcal{P}_2 are information partitions of Ω , define the *join* $\mathcal{P}_1 \vee \mathcal{P}_2$ to be the coarsest common refinement of \mathcal{P}_1 and \mathcal{P}_2 , and define the *meet* $\mathcal{P}_1 \wedge \mathcal{P}_2$ to be the finest common coarsening of \mathcal{P}_1 and \mathcal{P}_2 .⁹

Let (Ω, \mathcal{F}, p) be a probability space, let \mathcal{P}_1 and \mathcal{P}_2 be information partitions, for agents 1 and 2 respectively, such that $\mathcal{P}_1 \vee \mathcal{P}_2$ consists of only non-null events, i.e. $p(E) > 0$ for all $E \in \mathcal{P}_1 \vee \mathcal{P}_2$. We interpret p as the *common prior* of agents 1 and 2.

Corollary 1. *An event E is common knowledge in state $\omega \in \Omega$ if the member $F \in \mathcal{P}_1 \wedge \mathcal{P}_2$ for which $\omega \in F$ is such that $F \subseteq E$.*

Proof. F is a self-evident event, and the ‘if’ direction in Proposition 2 holds generally (note that part of the proof does not rely on finiteness of Ω). \square

Fixing an event $E \in \mathcal{F}$, we denote the posterior probability for agent i by

$$q_i(\omega) = \frac{p(E \cap h_i(\omega))}{p(h_i(\omega))}.$$

Theorem 2 (Aumann, 1976). *Let $\omega \in \Omega$. If it is common knowledge in state ω that $q_1(\omega) = \bar{q}_1$ and $q_2(\omega) = \bar{q}_2$ for numbers \bar{q}_1 and \bar{q}_2 , then $\bar{q}_1 = \bar{q}_2$.*

Proof. Let $P \in \mathcal{P}_1 \wedge \mathcal{P}_2$ be such that $\omega \in P$. We can write $P = \bigcup_j P^j$ for a disjoint collection $\{P^j\} \subseteq \mathcal{P}_1$. Since $q_1 = \bar{q}_1$ on P , we have that $\bar{q}_1 = \frac{p(E \cap P^j)}{p(P^j)}$ for all j , giving $p(E \cap P^j) = \bar{q}_1 p(P^j)$. By countable additivity, $p(E \cap P) = \bar{q}_1 p(P)$.

Repeating the argument, we can write $P = \bigcup_k P^k$ for a disjoint collection $\{P^k\} \subseteq \mathcal{P}_2$ and follow the same steps to obtain $p(E \cap P) = \bar{q}_2 p(P)$. Hence $\bar{q}_1 = \bar{q}_2$. \square

Note the theorem generalizes beyond two agents (Bacharach, 1985; Rubinstein & Wolinsky, 1990, and Samet, 1990). The conclusion is that if agents share the same priors, they cannot “agree to disagree”. Of course, people disagree a lot in real life, and we could put this down to differences in individuals’ subjective priors. Yet Aumann’s agreement poses a challenge to Harsanyi’s (1968) argument that differences in subjective priors ought to only come about from differences in information. If this were true, then two reasonable people with the same information who have common knowledge

⁹A *refinement* \mathcal{P}' of a partition \mathcal{P} is a partition such that every element of \mathcal{P}' is the subset of some element of \mathcal{P} . We say that \mathcal{P}' is a *coarsening* of \mathcal{P} if \mathcal{P} is a refinement of \mathcal{P}' . We say that \mathcal{P} is a *common refinement* (common coarsening) of \mathcal{P}_1 and \mathcal{P}_2 if it is a refinement (coarsening) of both \mathcal{P}_1 and \mathcal{P}_2 . A partition \mathcal{P} is the *coarsest common refinement* of \mathcal{P}_1 and \mathcal{P}_2 if \mathcal{P} is the common refinement of \mathcal{P}_1 and \mathcal{P}_2 such that any other common refinement of \mathcal{P}_1 and \mathcal{P}_2 is a refinement of \mathcal{P} . Likewise, \mathcal{P} is the *finest common coarsening* of \mathcal{P}_1 and \mathcal{P}_2 if it is a common coarsening such that it is the refinement of any other common coarsening of \mathcal{P}_1 and \mathcal{P}_2 .

To give an example, consider $\Omega = \{1, 2, 3, 4\}$. Let $\mathcal{P}_1 = \{\{1, 2, 3\}, \{4\}\}$ and $\mathcal{P}_2 = \{\{1, 2\}, \{3, 4\}\}$. A coarsening of both partitions is $\{\{1, 2, 3, 4\}\}$, which is the only (and thus the finest) common coarsening of \mathcal{P}_1 and \mathcal{P}_2 . A refinement of both partitions is $\mathcal{P}_3 = \{\{1, 2\}, \{3\}, \{4\}\}$. There is no refinement of both partitions that is not a refinement of \mathcal{P}_3 , and so this is the coarsest common refinement of \mathcal{P}_1 and \mathcal{P}_2 .

of each others posterior beliefs should never disagree – but in practice, disagreements among expert colleagues with access to the same evidence is very common. Of course, it could be that reasonable people nevertheless ascribe errors to the calculation of each others' posterior distributions, arising due to systematic biases (see e.g. the behavioural literature).

1.5.4 No trade?

The agreement theorem undergirds a group of results that are known as the no-trade theorems – among which are Kreps (1977), Milgrom & Stokey (1982), Tirole (1982) and Rubinstein & Wolinsky (1990). These effectively rule out that rational risk-averse traders with common knowledge of rationality can engage in speculative trades unless they have different priors. We focus on the Milgrom-Stokey no-trade theorem, and follow the simple version presented in Levin's notes.

Suppose there are two agents. Let (Ω, \mathcal{F}, p) be a probability space with Ω finite, let X be a set of trading outcomes, and assume each agent's information function is partitionial.

We define a *contingent contract* to be a measurable function $a : \Omega \rightarrow X$. We let A be the space of contingent contracts. Each agent i has a utility function $u_i : X \times \Omega \rightarrow \mathbb{R}$, and the agent's utility from a contract a is the random variable $U_i(a)(\omega) = u_i(a(\omega), \omega)$. We denote i 's expectation of $U_i(a)$ given her information H_i by $\mathbb{E}[U_i(a) \mid H_i]$.

Proposition 3. *Let ϕ be a random variable on (Ω, \mathcal{F}, p) . If p is the common prior of i and j , and the posterior distributions of both agents are common knowledge, then it cannot be common knowledge that i 's expectation of ϕ is strictly greater than j 's expectation of ϕ .*

Proof. Fix $\omega \in \Omega$. Let $E(t) = \{\phi \leq t\}$ for any $t \in \mathbb{R}$. Define $q_i(\omega)(E) = \frac{p(E \cap h_i(\omega))}{p(h_i(\omega))}$ for each $E \in \mathcal{F}$ and $i = 1, 2$. Suppose these posterior distributions are common knowledge. Then by Aumann's agreement theorem (Theorem 2), $q_1(\omega) = q_2(\omega)$ everywhere. Now, $\mathbb{E}[\phi \mid h_i(\omega)] = \int_{\Omega} \phi \, dq_i(\omega)$ for each i ,¹⁰ and thus $\mathbb{E}[\phi \mid h_1(\omega)] = \mathbb{E}[\phi \mid h_2(\omega)]$. \square

Call a contingent contract b *ex ante efficient* if there does not exist a contract $a \in A$ such that $\mathbb{E}[U_i(a)] > \mathbb{E}[U_i(b)]$ for both agents i .

Theorem 3 (Milgrom-Stokey, 1982). *If a contingent contract b is ex ante efficient, then it cannot be common knowledge that every agent prefers some contract a to b .*

Proof. Let $E = \{\omega \in \Omega \mid \mathbb{E}[U_i(a) \mid h_i(\omega)] > \mathbb{E}[U_i(b) \mid h_i(\omega)] \text{ for all } i\}$. The theorem states there is no state ω in which E is common knowledge. Suppose otherwise. By Proposition 2 there is a self-evident set $F \subseteq E$ with $\omega \in F$. By definition, for all $\omega' \in F$, $h_i(\omega') \subseteq F$, and thus for all $\omega' \in F$ and all i , we have $\mathbb{E}[U_i(a) - U_i(b) \mid h_i(\omega')] > 0$.

Since h_i is partitionial and Ω is finite, we can write F as a disjoint union $F = \bigcup_{k=1}^n h_i(\omega_k)$ for some set of states $\{\omega_k\}_{k=1}^n \subseteq F$. Thus $\mathbb{E}[U_i(a) - U_i(b) \mid F] > 0$ for

¹⁰Note this integral is taken with respect to the measure $q_i(\omega)$, *not* with respect to ω .

all i . That is, a strictly Pareto dominates b on F (c.f. Definition 12). Now consider the contract c defined so that $c(\omega') = a(\omega')$ for each $\omega' \in F$ and $c(\omega') = b(\omega')$ for each $\omega' \notin F$. Then c yields the same payoff as b outside F and strictly better payoff than b on F , so b cannot be *ex ante* efficient, yielding a contradiction. \square

The Milgrom-Stokey no-trade theorem rules out the possibility of speculative trades under conditions that are (or were) standard in models of financial markets – namely that traders are rational Bayesian agents, traders have common priors, it is common knowledge when a trade takes place that it is feasible and mutually acceptable to both parties, and that markets are in an efficient equilibrium. An interesting implication of the theorem is that receipt of private information by a trader is not helpful, because they cannot find a counterparty to agree on a favourable trade – the only motive for trading in the Milgrom-Stokey setting is speculative, and a trader only places a speculative bet if they have private information, but then the counterparty can infer that the trader has better information and so will refuse to meet the bet.

In practice, the predictions of the no-trade theorem obviously do not hold – trade volumes in financial markets are very high. There are several ways we might accommodate this fact. The first is to suppose that some traders are noise traders, that is, non-rational traders who trade anyway. These traders create trade possibilities via two channels – first, rational traders can make gains from the losses of the noise traders, and secondly, rational traders may be willing to place trades between themselves because there is uncertainty over whether their counterparty is rational or a noise trader (thus a rational trader with private information can profit by trading with a rational trader without this information, for example). We could instead relax more fundamental assumptions, such as rationality more generally. For example, as with Aumann’s agreement theorem, the no-trade theorem fails if agents have systematic biases. Thirdly, we can weaken the assumption of common priors. Feinberg (2000) shows that under heterogeneous priors, there always exists some purely speculative trade that all agents would be willing to make.

Example 7 (Harrison & Kreps, 1978). Harrison & Kreps (1978) present a nice example (and a formal model) to illustrate how speculative trades become possible when agents have heterogeneous priors. Suppose there are two types of risk-neutral investor, $i = 1, 2$, and both have a common discount factor $\delta = \frac{3}{4}$. Investors can purchase a stock that pays a dividend, and in every period t , the dividend is either $d_t = 0$ or $d_t = 1$. They cannot short the stock. This is the state, and so the state space is $D = \{0, 1\}$. The dividend process is perceived by both types to follow a stationary Markov process, but the types disagree on transition probabilities, with type i believing the transition matrix is Q_i , with

$$Q_1 = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{2}{3} & \frac{1}{3} \end{pmatrix} \quad \text{and} \quad Q_2 = \begin{pmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{1}{4} & \frac{3}{4} \end{pmatrix}.$$

First, consider the value $v_i(d)$ to each type i of investor of buying a unit of the stock

when the current state is d and holding it forever. We have that

$$\begin{aligned} p_1(0) &= \frac{4}{3} = 1.33, & p_1(1) &= \frac{11}{9} = 1.22, \\ p_2(0) &= \frac{16}{11} = 1.45, & p_2(1) &= \frac{21}{11} = 1.91. \end{aligned}$$

The second type of investor always values holding the stock more than the first type, regardless of the current state. Nevertheless, there are opportunities for trade. In state 1, type 2 investors are optimistic that they will receive dividends in the future, since they assess that $d_{t+1} = 1$ with probability $\frac{3}{4}$. Type 1 investors are pessimistic about future dividends in state 1, but are unable to short-sell the stock. In state 0, however, type 1 investors are optimistic about a transition to state 1 relative to type 2 investors. Type 1 investors thus have an opportunity to trade by purchasing stock in state 0 and selling it to type 2 investors in state 1, realizing capital gains. We can thus expect that in equilibrium, the stock changes hands between investors: when there is a transition to state 1, type 1 investors sell the stock to type 2 investors, and when there is a transition to state 0, type 2 investors sell to type 1 investors.

Harrison & Kreps (1978) use the notion of *consistent equilibrium*, which imply the price $p_t(d_t)$ of the stock at time t should satisfy

$$p_t(d_t) = \max_k \delta \sum_{d_{t+1} \in D} [d_{t+1} + p_t(d_{t+1})] Q_k(d_t, d_{t+1}).$$

That is, the price of the stock in state d is the maximum discounted expected return across investors from buying and holding the asset for a single period. Note the price of the stock will be stationary, since the dividend process is stationary, and thus $p_t(d) = p_{t+k}(d)$ for any $k \in \mathbb{Z}$ and any $d \in D$. We have

$$\begin{aligned} p(0) &= \max \left\{ \frac{3}{4} \left[\frac{1}{2}p(0) + \frac{1}{2}(1 + p(1)) \right], \frac{3}{4} \left[\frac{2}{3}p(0) + \frac{1}{3}(1 + p(1)) \right] \right\}, \\ p(1) &= \max \left\{ \frac{3}{4} \left[\frac{2}{3}p(0) + \frac{1}{3}(1 + p(1)) \right], \frac{3}{4} \left[\frac{1}{4}p(0) + \frac{3}{4}(1 + p(1)) \right] \right\}. \end{aligned}$$

Solving this system gives $p(0) = \frac{24}{13} = 1.85$ and $p(1) = \frac{27}{13} = 2.04$. The price in both states is greater than the valuation “on fundamentals” of both types of investor! The type 2 investor is willing to pay in excess of their “fundamental” valuation in state 1 because they know they can sell the stock to type 1 investors in state 0 at a higher price than her valuation in that state. The type 1 investor, meanwhile, is willing to pay a high price in state 0 knowing she can sell it for a much higher price to type 2 investors in state 1.

2 Games of complete information

First, a bit about solution concepts. A solution concept that applies to single strategy profiles is often called an *equilibrium*. Other solution concepts only generally isolate sets.

Solution concepts serve several purposes:

- *Descriptive.* A solution concept to a game can aim to predict strategies that players will play in practice. A large empirical literature looks at behaviour of players in experimental settings (e.g. Güth et al (1982) study experimental evidence of behaviour in the ultimatum game, Forsythe et al. (1994) likewise analyse behaviour in practice for the dictator game) and non-experimental empirical settings (e.g. the empirical IO literature)
- *Normative.* A solution concept may aim to prescribe which strategies a rational player would play. They can thus provide a guide to action.
- *Theoretical.* Given certain assumptions about players' behaviour, a solution concept can aim to predict behaviour under those assumptions (without making a broader claim that those assumptions are accurate in practice)

The desirability of a solution concept lies in several properties:

- A solution concept should be reasonable in the sense that the assumptions about behaviour on which it relies are reasonable assumptions about agents (descriptively or normatively). While it might not always seem like it, we are trying to build a theory of how people actually interact strategically in the real world here.¹¹
- A solution concept should apply to a sufficiently broad class of games, or else it seems ad hoc.
- A solution concept ideally gives a clear-cut – and thus falsifiable – prediction about behaviour. This doesn't need to be a unique prediction necessarily, but some sufficiently small set to be useful. In practice, multiplicity of equilibria is the norm, and typically we might want to make equilibrium selection arguments or rely on refinements to narrow down the range of “reasonable” equilibria.

2.1 Welfare and efficiency in games

One often asks whether the solution to a game is efficient, in a Paretian sense or in the sense of maximizing social welfare. Non-cooperative games usually involve externalities and so we cannot usually expect reasonable solution concepts to yield efficient outcomes.

Definition 12.

- (a) *Pareto dominance.* In the game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a strategy profile s (weakly) *Pareto dominates* a strategy profile s' if

$$u_i(s) \geq u_i(s') \quad \text{for all } i \in \mathcal{I}, u_i(s) > u_i(s') \quad \text{for some } i \in \mathcal{I}.$$

¹¹And often successfully so! See https://twitter.com/ben_golub/status/1611917935399796743

In this case, we say s is a (weak) *Pareto improvement* on s' .

We say s *strongly Pareto dominates* s' if $u_i(s) > u_i(s')$ for all $i \in \mathcal{I}$. Then s is a *strong Pareto improvement* on s' .

- (b) *Pareto efficiency*. A strategy s is (strongly) *Pareto efficient* if there is no strategy profile s' such that s' (weakly) Pareto dominates s .

We say s is *weakly Pareto efficient* if there is no strategy s' that strongly Pareto dominates s .

- (c) *Strong efficiency*. A strategy s is *strongly efficient* if it is the solution to

$$\max_{s \in S} \sum_{i \in \mathcal{I}} u_i(s),$$

that is, if it maximizes total payoff.

The definition of strong efficiency is taken from Jackson & Wolinsky (1996). Clearly, it is equivalent to the statement that s maximizes a utilitarian social welfare function.

As mentioned, in general, we cannot expect our solution concepts to give us Pareto efficient – let alone strongly efficient – solutions. One exception is games with transferable utility (i.e. where players can costlessly transfer portions of their payoffs to other players) – here, appropriate solution concepts will typically give us strongly efficient solutions.

Often, we are interested in how large the gap is between welfare in an equilibrium and welfare at the social optimum. We could compare the ratio between welfare in a given equilibrium with welfare at the optimum. However, there are often multiple equilibria and computing equilibria explicitly is difficult. A more popular approach is to look at the worst- and best-case scenarios:

Definition 13. Consider a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ and a social welfare function $W : S \rightarrow \mathbb{R}$. Let S^* be the set of equilibria of G .

- (a) *Price of anarchy*. The *price of anarchy* is defined by

$$\text{PoA}(G, W) = \frac{\inf_{s \in S^*} W(s)}{\sup_{s \in S} W(s)}.$$

- (b) *Price of stability*. The *price of stability* is defined by

$$\text{PoS}(G, W) = \frac{\sup_{s \in S^*} W(s)}{\sup_{s \in S} W(s)}.$$

In words, the price of anarchy is the welfare ratio between the worst equilibrium outcome and the socially optimal outcome, and the price of stability is the welfare ratio between the best equilibrium outcome and the socially optimal outcome. I have kept “equilibrium outcome” vague here. In the definition, the set S^* is usually taken to be

the set of Nash equilibria of G , but in some contexts, other equilibrium concepts will be more appropriate – for example, we might want to look at the price of anarchy and price of stability where S^* is the set of subgame perfect equilibria, or Wardrop equilibria, and so on.

2.2 Zero sum and matrix games

Definition 14 (Zero sum game). A *zero sum game* is a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ such that

$$\sum_{i \in \mathcal{I}} u_i(s) = 0 \quad \text{for all } s \in S = \prod_{i \in \mathcal{I}} S_i.$$

If instead $\sum_{i \in \mathcal{I}} u_i(s) = c$ for all $s \in S$ and some constant c , then we call G a *constant sum game*.

The zero sum games constitute a large class of games. Zero sum games are games of pure competition – comparing any two strategies, any increase in the payoff of some player is matched by a decrease in the aggregate payoff of the other players. Any constant sum game is also a game of pure competition, and for any constant sum game G , we can define a best response equivalent zero sum game $\tilde{G} = (\mathcal{I}, (S_i, \tilde{u}_i)_{i \in \mathcal{I}})$ with $\tilde{u}_i(s) = u_i(s) - c$.

Corollary 2. If $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a zero sum game, the corresponding mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ is a zero sum game. That is,

$$\sum_{i \in \mathcal{I}} u_i(\sigma) = 0 \quad \text{for all } \sigma \in \prod_{i \in \mathcal{I}} \Delta(S_i).$$

Proof. By linearity of the Lebesgue integral,

$$\sum_{i \in \mathcal{I}} u_i(s) = \sum_{i \in \mathcal{I}} \int_S u_i(s) d\sigma = \int_S \left[\sum_{i \in \mathcal{I}} u_i(s) \right] d\sigma = 0.$$

□

Any outcome of a zero sum game is clearly Pareto optimal.

Finite two-person zero sum games are representable as *matrix games*.

Definition 15 (Matrix game). A *matrix game* is a real $m \times n$ matrix A , where m is the number of actions for Player 1 and n is the number of actions for Player 2. A mixed strategy of Player 1 is an m -dimensional probability vector p , and the set of mixed strategies of Player 1 is

$$\Delta^m := \left\{ p \in \mathbb{R}_+^m \mid \sum_{i=1}^m p_i = 1 \right\}.$$

A mixed strategy of Player 2 is an n -dimensional probability vector q , and the set of mixed strategies of Player 2 is

$$\Delta^n := \left\{ q \in \mathbb{R}_+^n \mid \sum_{i=1}^n q_i = 1 \right\}.$$

We call a strategy p or q of a matrix game a *pure strategy* if there is an entry of p or q with value 1. We denote the vector with i th entry 1 and all other entries 0 by e^i .

Because one player's gain in a zero-sum game is always another player's loss, any strategy profile in a zero-sum game is Pareto efficient. Moreover, any strategy profile is strongly efficient:

Proposition 4. *If $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a zero-sum game, then any strategy profile $s \in S$ is strongly efficient.*

Proof. We have $\sum_{i \in \mathcal{I}} u_i(s) = 0$ for all $s \in S$, and hence any s maximizes this sum. \square

2.3 Maxmin and minmax

Definition 16.

- (a) *Maximin.* Given a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, the pure strategy *maximin payoff* for player i is given by

$$w_i = \max_{s_i \in S_i} \min_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}).$$

A pure strategy $\alpha^i = (\alpha_i^i, \alpha_{-i}^i) \in S$ is a *maximin solution* for player i if

$$(\alpha_i^i, \alpha_{-i}^i) = \arg \max_{s_i \in S_i} \min_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}).$$

That is, the maximin payoff is the payoff attained at the maximin solution, $w_i = u_i(\alpha_i^i, \alpha_{-i}^i)$.

Given the mixed game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, the mixed strategy maximin payoff for player i is given by

$$w_i^m = \max_{\sigma_i \in \Delta(S_i)} \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} u_i(\sigma_i, \sigma_{-i}).$$

- (b) *Minimax.* Given a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, the pure strategy *minimax payoff* for player i is given by

$$v_i = \min_{s_{-i} \in S_{-i}} \max_{s_i \in S_i} u_i(s_i, s_{-i}).$$

A pure strategy $\gamma^i = (\gamma_i^i, \gamma_{-i}^i) \in S$ is a *minimax solution* for player i if

$$(\gamma_i^i, \gamma_{-i}^i) = \arg \min_{s_{-i} \in S_{-i}} \max_{s_i \in S_i} u_i(s_i, s_{-i}).$$

That is, the minimax payoff is the payoff attained at the minimax solution, $v_i = u_i(\gamma_i^i, \gamma_{-i}^i)$.

Given the mixed game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, the mixed strategy minimax payoff for player i is given by

$$v_i^m = \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} \max_{\sigma_i \in \Delta(S_i)} u_i(\sigma_i, \sigma_{-i}).$$

The maximin payoff for i is the largest payoff player i can guarantee themselves in the absence of any knowledge about their opponent's strategy. The minimax payoff for i , by contrast, is the least payoff that opponents $-i$ can enforce on player i .

Lemma 3. *Consider the mixed game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, and let w_i^m and v_i^m respectively denote the mixed strategy maximin and minimax payoffs for player i . Then*

$$v_i^m \geq w_i^m.$$

Proof. For any σ_i, σ'_{-i} ,

$$u_i(\sigma_i, \sigma'_{-i}) \geq \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} u_i(\sigma_i, \sigma_{-i}).$$

Hence

$$\max_{\sigma_i \in \Delta(S_i)} u_i(\sigma_i, \sigma'_{-i}) \geq \max_{\sigma_i \in \Delta(S_i)} \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} u_i(\sigma_i, \sigma_{-i}),$$

and so

$$\min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} \max_{\sigma_i \in \Delta(S_i)} u_i(\sigma_i, \sigma_{-i}) \geq \max_{\sigma_i \in \Delta(S_i)} \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} u_i(\sigma_i, \sigma_{-i})$$

□

In finite-strategy zero sum games, the mixed strategy maximin and minimax solutions for each player are equivalent. Indeed, mixed strategy maximin is a natural solution concept in this setting:

Theorem 4 (Minimax theorem; von Neumann, 1928). *In any finite two player zero sum game $G = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$,*

$$v_1^m = w_1^m = v_2^m = w_2^m.$$

Proof. Recall we can write any finite two player zero sum game as an $m \times n$ matrix game A . First we claim that $w_1^m = v_2^m$. Let p_1, q_1 be the choices of $p \in \Delta^m$ and $q \in \Delta^n$ that solve $w_1^m = \max_{p \in \Delta^m} \min_{q \in \Delta^n} p' A q$ and let p_2, q_2 be the choices of $p \in \Delta^m$ and $q \in \Delta^n$ that solve $v_2^m = \min_{q \in \Delta^n} \max_{p \in \Delta^m} p' A q$. Now, we have

$$w_1^m = p_1' A q_1 \leq p_2' A q_1 \leq p_2' A q_2 = v_2^m,$$

so $w_1^m \leq v_2^m$.

Suppose $w_1^m > v_2^m$. We prove this yields a contradiction using the lemma of the alternative for matrices, Lemma 32. Let B be an arbitrary $m \times n$ matrix game. Let $w_1(p) = \min_{q \in \Delta^n} p' B q$, let $w_1(B) = \max_{p \in \Delta^m} w_1(p)$, let $v_2(q) = \max_{p \in \Delta^m} p' B q$, and let $v_2(B) = \min_{q \in \Delta^n} v_2(q)$.

Recall from the lemma that precisely one of the following must hold:

- (a) There exist $y \in \mathbb{R}^n$ and $z \in \mathbb{R}^m$ s.t. $(y, z) \geq 0$, $(y, z) \neq 0$ and $By + z = 0$;
- (b) There is an $x \in \mathbb{R}^m$ s.t. $x > 0$ and $x'B > 0$.

Suppose (a) holds, so there exists $y \in \mathbb{R}^n$ and $z \in \mathbb{R}^m$, at least one of which nonzero, such that $(y, z) \geq 0$ and $By + z = 0$. If $y = 0$ then $z = 0$, yielding a contradiction, so $y \neq 0$ and $\sum_{k=1}^n y_k > 0$. Define $q \in \Delta^n$ so that $q_j = \frac{y_j}{\sum_{k=1}^n y_k}$ for each $j = 1, \dots, n$. It follows that $Bq = -\frac{z}{\sum_{k=1}^n y_k} \leq 0$. Thus $v_2(q) \leq 0$, so $v_2(B) \leq 0$.

Suppose (b) holds. Then there exists $x \in \mathbb{R}^m$ such that $x > 0$ and $x'B > 0$. Define $p \in \Delta^m$ so that $p = \frac{x}{\sum_{k=1}^m x_k}$. Then $w_1(p) > 0$ and so $w_1(B) > 0$. It follows that we cannot have $w_1(B) \leq 0 < v_2(B)$, since at least one of (a) and (b) must hold.

Now define B so that each ij th entry of B is $B_{ij} = A_{ij} - w_1(A)$. Then $w_1(B) = v_1(A) - v_1(A) = 0$ and $v_2(B) = v_2(A) - v_1(A) > 0$. Hence $w_1(B) \leq 0 < v_2(B)$, yielding a contradiction.

Hence we conclude that $w_1^m = v_2^m$. An identical argument shows that $v_1^m = w_2^m$. The equality $v_1^m = w_1^m = v_2^m = w_2^m$ now follows from Lemma 3. \square

We call $v := v_1^m = v_2^m = w_1^m = w_2^m$ the *value of the game* G . Von Neumann's minimax theorem as stated above is a special case of the following, more general version:

Theorem 5 (Minimax theorem). *Let $X \subseteq \mathbb{R}^n$ and $Y \subseteq \mathbb{R}^m$ be nonempty, convex, compact sets. If $f : X \times Y \rightarrow \mathbb{R}$ is a continuous function such that*

- (i) *f is concave in its first argument, that is, for each $y \in Y$, $g(x) = f(x, y)$ is concave, and*
- (ii) *f is convex in its second argument, that is, for each $x \in X$, $h(y) = f(x, y)$ is convex,*

then

$$\max_{x \in X} \min_{y \in Y} f(x, y) = \min_{y \in Y} \max_{x \in X} f(x, y).$$

There are a number of generalizations – Sion's minimax theorem and Parthasarthy's theorem, to name the two most famous examples.

The value of the game need not exist once we allow for infinite strategy spaces, as the following example shows.

Example 8 (Game without a value, Sion & Wolfe, 1957). Suppose Player 1 chooses a number $x \in [0, 1]$ and Player 2 chooses a number $y \in [0, 1]$. Player 1 receives payoff

$$\pi(x, y) = \begin{cases} -1 & \text{if } x < y < x + \frac{1}{2}, \\ 0 & \text{if } x = y \text{ or } y = x + \frac{1}{2}, \\ 1 & \text{otherwise,} \end{cases}$$

and Player 2 receives payoff $-\pi(x, y)$.

This is a kind of continuous *Colonel Blotto game*: imagine that Player 1 assigns a fraction x of their forces to attack one mountain pass and $1 - x$ to attack the other. Player 2 assigns a fraction y to defend the first mountain pass and $1 - y$ to defend the second, at which a second permanent garrison of $\frac{1}{2}$ is also located. A player receives payment 1 from their opponent at each pass if their forces are larger than their opponent.

If the value of this game exists, then it is the value

$$\sup_f \inf_g \iint \pi \, df \, dg = \inf_g \sup_f \iint \pi \, df \, dg.$$

However, it can be shown that $\sup_f \inf_g \iint \pi \, df \, dg = \frac{1}{3} \neq \frac{3}{7} = \inf_g \sup_f \iint \pi \, df \, dg$.

Example 9 (Matching pennies). Matching pennies is a zero sum game in which each of two players simultaneously announce heads (H) or tails (T). Player 1 receives a payment from Player 2 if the two announcements match, and Player 2 receives a payment from Player 1 if the two announcements differ.

| | H_2 | T_2 |
|-------|-------|-------|
| H_1 | 1, -1 | -1, 1 |
| T_1 | -1, 1 | 1, -1 |

Consider the mixed strategies $\sigma_1 = (p, 1 - p)$ and $\sigma_2 = (q, 1 - q)$. Consider the maximin strategy for Player 1. Player 1's expected payoff is

$$u_1(p, q) = p[q - (1 - q)] + (1 - p)[(1 - q) - q] = (1 - 2p)(1 - 2q).$$

Given $p \in [0, 1]$, Player 2's problem (in order to minimize Player 1's payoff) is

$$\begin{aligned} \min_{q \in [0, 1]} u_1(p, q) &= \begin{cases} 1 - 2p & \text{if } p \leq \frac{1}{2}, \\ 2p - 1 & \text{if } p > \frac{1}{2}, \end{cases} \\ &= \min\{u_1(p, 0), u_1(p, 1)\} \\ &= \min\{1 - 2p, 2p - 1\} \\ &= -|1 - 2p|. \end{aligned}$$

Player 1's mixed maximin strategy thus solves

$$w_1^m = \max_{p \in [0, 1]} \min_{q \in [0, 1]} u_1((p, 1 - p), (q, 1 - q)) = \max_{p \in [0, 1]} -|1 - 2p| = 0,$$

for $(p, q) = (\frac{1}{2}, \frac{1}{2})$. The game is symmetric, so Player 2's maxmin payoff is also $w_2^m = 0$.
Likewise, given $q \in [0, 1]$, Player 1's problem is

$$\begin{aligned} \max_{p \in [0, 1]} u_1(p, q) &= \begin{cases} 1 - 2q & \text{if } q \leq \frac{1}{2}, \\ 2q - 1 & \text{if } q > \frac{1}{2}, \end{cases} \\ &= \max\{u_1(0, q), u_1(1, q)\} \\ &= \max\{1 - 2q, 2q - 1\} \\ &= |1 - 2q|. \end{aligned}$$

Player 2's minmax strategy solves

$$w_1^m = \min_{q \in [0, 1]} \max_{p \in [0, 1]} u_1((p, 1 - p), (q, 1 - q)) = \min_{q \in [0, 1]} |1 - 2q| = 0,$$

for $(p, q) = (\frac{1}{2}, \frac{1}{2})$. The game is symmetric, so we also have $w_2^m = 0$.

We see that here, $v_i^m = w_i^m$. This is of course as predicted by von Neumann's minmax theorem.

2.4 Strict dominance

Strict dominance captures the idea that rational (i.e. expected payoff-maximizing) players will never play strategies that perform uniformly worse than some alternative strategy.

Definition 17 (Strict dominance).

- (a) *Strict dominance in pure strategies.* In a pure strategy game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a pure strategy $s_i \in S_i$ for player i is said to *strictly dominate* a strategy $s'_i \in S_i$ if

$$u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i}.$$

We say that a strategy s'_i is *strictly dominated* if there exists some strategy s_i that strictly dominates s'_i .

We say that a strategy $s_i \in S_i$ is *strictly dominant* if s_i strictly dominates every $s'_i \in S_i - \{s_i\}$.

- (b) *Strict dominance in mixed strategies.* In a mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, a mixed strategy $\sigma_i \in \Delta(S_i)$ is said to *strictly dominate* a strategy $\sigma'_i \in \Delta(S_i)$ if

$$u_i(\sigma_i, \sigma_{-i}) > u_i(\sigma'_i, \sigma_{-i}) \quad \text{for all } \sigma_{-i} \in \Delta_{-i}(S_{-i}).$$

We say that σ'_i is *strictly dominated* if there exists some strategy σ_i that strictly dominates σ'_i .

Equivalently, by linearity and monotonicity of expectations, σ_i strictly dominates σ'_i if

$$u_i(\sigma_i, s_{-i}) > u_i(\sigma'_i, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i}.$$

- (c) *Strict dominance of a pure strategy by a mixed strategy.* In a mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, a mixed strategy $\sigma_i \in \Delta(S_i)$ is said to *strictly dominate* pure strategy s_i if

$$u_i(\sigma_i, s_{-i}) > u_i(s_i, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i}.$$

- (d) *Strictly dominant strategy equilibrium.* A strategy profile $s^* \in S$ ($\sigma^* \in \times_{i \in \mathcal{I}} \Delta(S_i)$) is a *strictly dominant strategy equilibrium* if, for each player i , $s_i^* \in S_i$ ($\sigma_i^* \in \Delta(S_i)$) is a strictly dominant strategy.

Proposition 5. *In a game G , if s^* (σ^*) is a strictly dominant strategy equilibrium, then it is the unique strictly dominant strategy equilibrium.*

Proof. Wlog, consider a pure strategy game. Suppose s^* and s' are strictly dominant strategy equilibria and that $s^* \neq s'$. Then for some i , $s_i^* \neq s_i'$. Since s_i^* and s_i' are both strictly dominant strategies for i , we have

$$u_i(s_i^*, s_{-i}) > u_i(s_i', s_{-i}) \quad \text{and} \quad u_i(s_i', s_{-i}) > u_i(s_i^*, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i},$$

a clear contradiction. \square

A strictly dominant strategy equilibrium requires only that each player is rational. Beliefs about other players are irrelevant. However, it applies only to a very small class of games: those in which every player has a strictly dominant strategy. In general, a strictly dominant strategy equilibrium need not exist. The prisoner's dilemma (Example 2(a)) has a strictly dominant strategy equilibrium (D, D) , for example, but the market entry game in Example 3 does not.

Definition 18 (Iterated strict dominance).

- (a) *Level- k rationality.* An assumption of *level-1 rationality* is that all players are rational (i.e. expected payoff-maximizing). We say that players \mathcal{I} are *level-2 rational* if they are all rational (level-1) and know that all other players are rational. Iteratively, we say players are *level- k rational* if they know that all players are level- $(k-1)$ rational.
- (b) *Iterated strict dominance.* Let $\mathcal{D}_i^0 = S_i$ for all i , and define

$$\mathcal{D}_i^k = \left\{ s_i \in \mathcal{D}_i^{k-1} \mid \nexists \sigma_i \in \Delta(\mathcal{D}_i^{k-1}) \text{ s.t. } u_i(\sigma_i, s_{-i}) > u_i(s_i, s_{-i}) \text{ for all } s_{-i} \in \mathcal{D}_{-i}^{k-1} \right\},$$

where $\mathcal{D}_{-i}^{k-1} = \times_{j \neq i} \mathcal{D}_j^{k-1}$. Clearly, $\mathcal{D}_i^k \subseteq \mathcal{D}_i^{k-1}$ for all $k \geq 1$.

We call \mathcal{D}_i^k the set of pure strategies that survive k rounds of iterated strict dominance for player i . We might also refer to this as the level- k iterated strict dominance set.

We call $\mathcal{D}_i = \bigcap_{k=0}^{\infty} \mathcal{D}_i^k$ the *set of pure strategies that survive iterated strict dominance* for player i , and $\mathcal{D} = \mathcal{D}_1 \times \cdots \times \mathcal{D}_n$ the *set of pure strategy profiles that survive iterated strict dominance*.

The set of mixed strategies that survive iterated strict dominance is defined analogously.

If \mathcal{D} contains only one strategy profile s , we call s the *iterated strict dominance solution* and say the game is *iterated strict dominance solvable*.

Iterated strict dominance is often known by the more verbose but more descriptive name “iterated deletion of strictly dominated strategies”. Shorter names are preferable so we don’t call it that.

If the structure of the game and rationality are common knowledge, then no player i will choose a strategy that is not contained in \mathcal{D}_i , i.e. that does not survive iterated strict dominance. The predictive power of iterated strict dominance is limited – indeed, in some games $\mathcal{D}_i = S_i$ – though not as poor as that of strictly dominant strategy equilibrium.

If an iterated strict dominance solution isolates a unique strategy profile, then we have a sharp prediction of equilibrium play. Even if \mathcal{D}_i is not a singleton for all players i , iterated strict dominance may allow us to considerably reduce the size of the strategy space that we need to consider.

Example 10. Consider the game

| | A_2 | B_2 | C_2 | D_2 |
|-------|-------|-------|-------|--------|
| A_1 | 0, 7 | 2, 5 | 7, 0 | 0, 1 |
| B_1 | 5, 2 | 3, 3 | 5, 2 | 0, 1 |
| C_1 | 7, 0 | 2, 5 | 0, 7 | 0, 1 |
| D_1 | 0, 0 | 0, -2 | 0, 0 | 10, -1 |

The strategy D_2 is strictly dominated by mixed strategy $(\frac{1}{2}, 0, \frac{1}{2}, 0)$, and no other strategies of either player are strictly dominated. Hence $\mathcal{D}_1^1 = \{A_1, B_1, C_1, D_1\}$ and $\mathcal{D}_1^2 = \{A_2, B_2, C_2\}$. In the second round, D_1 is strictly dominated by B_1 . No other strategies of either player are strictly dominated. Hence $\mathcal{D}_1 = \mathcal{D}_1^2 = \{A_1, B_1, C_1\}$ and $\mathcal{D}_2 = \mathcal{D}_2^2 = \{A_2, B_2, C_2\}$.

2.5 Correlated rationalizability

The notion of *correlated rationalizability* is closely related to iterated strict dominance. We define a player i ’s *belief* μ_{-i} about opponents’ play to be a subjective probability distribution over S_{-i} . That is, $\mu_{-i} \in \Delta(S_{-i})$. Intuitively, μ_{-i} gives the probabilities i attaches to the (possibly correlated) strategy profile of i ’s opponents.

Definition 19 (Best response). A strategy $s_i \in S_i$ ($\sigma_i \in \Delta(S_i)$) is a *best response*, or *best reply*, to $s_{-i} \in S_{-i}$ ($\sigma_{-i} \in \Delta(S_{-i})$) if

$$\begin{aligned} u_i(s_i, s_{-i}) &\geq u_i(s'_i, s_{-i}) && \text{for all } s_i \in S_i \\ (u_i(\sigma_i, \sigma_{-i}) &\geq u_i(\sigma'_i, \sigma_{-i}) && \text{for all } \sigma'_i \in \Delta(S_i). \end{aligned}$$

We write $B_i(s_{-i})$ ($B_i(\sigma_{-i})$) for the set of i 's best responses to s_{-i} (σ_{-i}).

A subset $B_1 \times \cdots \times B_n \subseteq S = S_1 \times \cdots \times S_n$ is called a *best response set* if, for all i and all $s_i \in B_i$, there exists a $\sigma_{-i} \in \Delta(B_{-i})$ such that s_i is a best response to σ_{-i} .

Call a subset $B_1^I \times \cdots \times B_n^I \subseteq S$ a *best response set to independent strategies* if, for all i and all $s_i \in B_i^I$, there exists a $\sigma_{-i} \in \Delta_{-i}(B_{-i}^I)$ such that s_i is a best response to σ_{-i} .

To check that a mixed strategy σ_i is a best response to some opponents' strategy profile σ_{-i} , we need only compare it to pure strategies s_i :

Lemma 4. *Consider a finite mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$. A mixed strategy $\sigma_i \in \Delta(S_i)$ is a best response to $\sigma_{-i} \in \Delta(S_{-i})$ iff $u_i(\sigma_i, \sigma_{-i}) \geq u_i(s_i, \sigma_{-i})$ for all $s_i \in S_i$.*

Proof. Suppose $u_i(\sigma_i, \sigma_{-i}) \geq u_i(s_i, \sigma_{-i})$ for all $s_i \in S_i$. Consider any $\sigma'_i \in \Delta_i(S_i)$. Then

$$u(\sigma_i, \sigma_{-i}) = \sum_{s_i \in S_i} u(s_i, \sigma_{-i}) \sigma_i(s_i) \geq \sum_{s_i \in S_i} u(s_i, \sigma_{-i}) \sigma'_i(s_i) = u_i(\sigma'_i, \sigma_{-i}),$$

so $\sigma_i \in B_i(\sigma_{-i})$. The converse is immediate by definition. \square

Since player i 's belief μ_{-i} is a correlated strategy profile, we use $B_i(\mu_{-i})$ to denote i 's best response to belief μ_{-i} .

Definition 20 (Rationalizability).

- (a) *Correlated rationalizability.* Let $\mathcal{R}_i^0 = S_i$ for all i . Given \mathcal{R}_i^{k-1} , the set of k -correlated rationalizable strategies is defined by

$$\mathcal{R}_i^k = \{s_i \in \mathcal{R}_i^{k-1} \mid s_i \in B_i(\mu_{-i}) \text{ for some } \mu_{-i} \in \Delta(\mathcal{R}_{-i}^{k-1})\},$$

where $\mathcal{R}_{-i}^{k-1} = \times_{j \neq i} \mathcal{R}_j^{k-1}$ and $\mathcal{R}^k = \times_{i \in \mathcal{I}} \mathcal{R}_i^k$.

The set of correlated rationalizable strategies R is defined as

$$\mathcal{R} = \bigcap_{k=0}^{\infty} \mathcal{R}^k.$$

- (b) *Independent rationalizability.* Let $\mathcal{R}_i^{I,0} = S_i$ for all i . Given $\mathcal{R}_i^{I,k-1}$, the set of k -independent rationalizable strategies is defined by

$$\mathcal{R}_i^{I,k} = \{s_i \in \mathcal{R}_i^{I,k-1} \mid s_i \in B_i^I(\mu_{-i}) \text{ for some } \mu_{-i} \in \Delta_{-i}(\mathcal{R}_{-i}^{I,k-1})\},$$

where $\mathcal{R}_{-i}^{I,k-1} = \times_{j \neq i} \mathcal{R}_j^{I,k-1}$ and $\mathcal{R}^{I,k} = \times_{i \in \mathcal{I}} \mathcal{R}_i^{I,k}$.

The set of independent rationalizable strategies R^I is defined as

$$\mathcal{R}^I = \bigcap_{k=0}^{\infty} \mathcal{R}^{I,k}.$$

Lemma 5 (Myerson, 1991). *In any finite game, a strategy s_i of player i is a best response to some belief μ_{-i} iff it is not strictly dominated (by a mixed strategy).*

Proof. Fix any pure strategy $s_i \in S_i$. Consider the linear program,

$$\delta^* := \min_{\delta \in \mathbb{R}, \sigma_{-i} \in \Delta(S_{-i})} \delta,$$

subject to

$$\delta + \sum_{s_{-i} \in S_{-i}} \sigma_{-i}(s_{-i}) [u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i})] \geq 0 \text{ for all } s'_i \in S_i.$$

Suppose $\delta^* > 0$. Then for any correlated profile σ_{-i} , we have that there is some strategy profile $s'_i \in S_i$ such that $u_i(s'_i, \sigma_{-i}) > u_i(s_i, \sigma_{-i})$, and so s_i is never a best response to any $\sigma_{-i} \in \Delta(S_{-i})$. Conversely, if $\delta^* \leq 0$, then for some $\sigma_{-i} \in \Delta(S_{-i})$ we have $u_i(s_i, \sigma_{-i}) \geq u_i(s'_i, \sigma_{-i})$ for all $s'_i \in S_i$, so s_i is a best response to $\sigma_{-i} =: \mu_{-i}$. Hence s_i is a best response to some belief over $\Delta(S_{-i})$ iff $\delta^* \leq 0$.

Next consider the linear program

$$\epsilon^* := \max_{\epsilon \in \mathbb{R}, \sigma_i \in \Delta(S_i), \{\eta_{s_{-i}} \in \mathbb{R}_+ | s_{-i} \in S_{-i}\}} \epsilon$$

subject to

$$\eta_{s_{-i}} + \epsilon + \sum_{s'_i \in S_i} \sigma_i(s'_i) [u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i})] = 0 \text{ for all } s_{-i} \in S_{-i}.$$

Suppose $\epsilon^* > 0$. Then i has a mixed strategy σ_i with $\sum_{s'_i \in S_i} \sigma_i(s'_i) [u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i})] < 0$ for all s_{-i} , which gives $u_i(s_i, s_{-i}) < u_i(\sigma_i, s_{-i})$ for all s_{-i} , and so σ_i strictly dominates s_i . Conversely, suppose $\epsilon^* \leq 0$. Then for any σ_i there is some s_{-i} for which we have $\sum_{s'_i \in S_i} \sigma_i(s'_i) [u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i})] \geq 0$ implying $u_i(s_i, s_{-i}) \geq u_i(\sigma_i, s_{-i})$, and so s_i is not strictly dominated. Hence s_i is not strictly dominated iff $\epsilon^* \leq 0$.

Now, the second linear program is the dual of the first. Applying the strong duality theorem (Theorem 54 together with Proposition 84), we have that $\delta^* = \epsilon^*$. Hence s_i is a best response to some belief μ_{-i} iff it is not strictly dominated. \square

Proposition 6. *In any finite game, the set of strategies surviving iterated strict dominance and the set of correlated rationalizable strategies coincide, i.e. $\mathcal{D}_i = \mathcal{R}_i$ for each player i .*

Proof. If $s_i \in \mathcal{R}_i$, then s_i is the best response to some belief about opponents' play $\mu_{-i} \in \Delta(\mathcal{R}_{-i})$. Since $\mathcal{R}_{-i} \subseteq S_{-i}$, the lemma implies s_i is not strictly dominated. Hence $\mathcal{R}_i \subseteq \mathcal{D}_i^1$. Applying the argument iteratively gives $\mathcal{R}_i \subseteq \mathcal{D}_i^k$ for all $k \in \mathbb{N}$ and hence $\mathcal{R}_i \subseteq \mathcal{D}_i$. Since this holds for all i , $\mathcal{R} \subseteq \mathcal{D}$.

Conversely, by definition, no strategy profile in \mathcal{D} is strictly dominated in the reduced game in which the strategy set is \mathcal{D} . Hence any $s_i \in \mathcal{D}_i$ is the best response to some beliefs μ_{-i} over \mathcal{D}_{-i} and hence over S_{-i} . It follows that \mathcal{D} is a best response set. Thus $\mathcal{D} \subseteq \mathcal{R}$.

Since $\mathcal{D} \subseteq \mathcal{R}$ and $\mathcal{R} \subseteq \mathcal{D}$, we have $\mathcal{R} = \mathcal{D}$. \square

Corollary 3. *A strategy profile σ is a profile of correlated rationalizable strategies, i.e. $\sigma \in \mathcal{R}$, iff for each i , σ_i is a best response to some belief μ_{-i} consistent with common knowledge of rationality and the structure of the game.*

Proof. Follows from the coincidence of \mathcal{R} and \mathcal{D} and the fact that no player i will play $s_i \notin \mathcal{D}_i$ if rationality and the structure of the game are common knowledge, while every $s_i \in \mathcal{D}_i$ is the best response to some $\sigma_{-i} \in \Delta(\mathcal{D}_{-i})$ or else it is strictly dominated. \square

Bernheim (1984) and Pierce (1984) use the concept of independent rationalizability rather than correlated rationalizability. The set of strategies surviving iterated strict dominance coincides with the set of independent rationalizable strategies for two player finite games, but for $n \geq 3$, independent rationalizability is a refinement of iterated strict dominance, i.e. $\mathcal{R}^I \subseteq \mathcal{D}$ but in general $\mathcal{R}^I \neq \mathcal{D}$.

2.6 Weak dominance

Definition 21 (Weak dominance).

- (a) *Weak dominance.* In a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a pure strategy $s_i \in S_i$ weakly dominates $s'_i \in S_i$ if

$$u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i}, \text{ with strict inequality for some } s_{-i}.$$

We call s'_i a *weakly dominated strategy* if there is some s_i that weakly dominates s'_i .

We call s_i a *weakly dominant strategy* if it weakly dominates all $s'_i \neq s_i$.

In a mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, a mixed strategy $\sigma_i \in \Delta(S_i)$ is said to *weakly dominate* a strategy $\sigma'_i \in \Delta(S_i)$ if

$$u_i(\sigma_i, s_{-i}) > u_i(\sigma'_i, s_{-i}) \quad \text{for all } s_{-i} \in S_{-i}.$$

- (b) *Weakly dominant strategy equilibrium.* A strategy profile $s^* \in S$ ($\sigma^* \in \times_{i \in \mathcal{I}} \Delta(S_i)$) is a *weakly dominant strategy equilibrium* if, for each player i , s_i^* (σ_i^*) is a weakly dominant strategy.
- (c) *Iterated weak dominance.* Let $\mathcal{W}_i^0 = S_i$ for all i , and define

$$\mathcal{W}_i^k = \left\{ s_i \in \mathcal{W}_i^{k-1} \mid \nexists \sigma_i \in \Delta(\mathcal{W}_i^{k-1}) \text{ s.t. } u_i(\sigma_i, s_{-i}) \geq u_i(s_i, s_{-i}) \quad \forall s_{-i} \in S_{-i}, \right. \\ \left. \text{with strict inequality for some } s_{-i} \right\}$$

We call \mathcal{W}_i^k the set of pure strategies that survive k rounds of iterated weak dominance for player i . The set of strategies that *survive iterated weak dominance* for player i is

$$\mathcal{W}_i = \bigcap_{k=0}^{\infty} \mathcal{W}_i^k.$$

Note that weakly dominant strategy equilibrium need not be unique.

Iterated weak dominance (and indeed weakly dominant strategy equilibrium) is a bit more difficult to justify than iterated strict dominance. Strategies that are weakly but not strictly dominated are rationalizable – any such strategy is a best response to some belief about opponents’ play. This said, such strategies are not robust, in the sense that if a player has any doubt about their beliefs, or believes their opponents can ‘tremble’ (make a mistake) with positive probability, then they are better off choosing a strategy that weakly dominates a weakly dominated strategy. Indeed, there is never any loss to choosing a strategy that weakly dominates a weakly dominated strategy, only potential gain. Luce & Raiffa (1957) thus have as an axiom of decision theory that no player will ever choose a weakly dominated strategy. We could push this further to argue that players will never play a strategy profile that does not survive iterated weak dominance (see section 6.6).

Example 11. Sometimes the process of iterated strict dominance or iterated weak dominance is conceived of differently. Rather than eliminating all (strictly/weakly) dominated strategies at each step, players take turns to remove a single strategy from their strategy set. This conception of iterated dominance has a few pitfalls. First, it only works if players’ strategy sets are finite. Second, under the assumption that the strategy sets are finite, for iterated weak dominance the set of strategies we arrive at can differ depending on the order in which we delete strategies (for iterated strict dominance, we will end up at the same set). Thus this version of the process is conceptually flawed.

Consider the game

| | A_2 | B_2 | C_2 |
|-------|-------|-------|-------|
| A_1 | 1, 2 | 2, 3 | 0, 3 |
| B_1 | 2, 2 | 2, 1 | 3, 2 |
| C_1 | 2, 1 | 0, 0 | 1, 0 |

B_1 is a weakly dominant strategy for Player 1. However, if Player 1 is sure that Player 2 will play A_2 , then C_1 is reasonable.

We have $\mathcal{W}_1^1 = \{B_1\}$ and $\mathcal{W}_2^1 = \{A_2, C_2\}$. No further deletions are possible.

The order in which we delete strategies matters here. For example, suppose we proceed in rounds where, rather than finding level- k iterated weak dominance sets, we alternate deleting strategies between players.

First delete A_1 . Then we have $\{B_1, C_1\}$, and Player 2 should delete C_2 , leaving $\{A_2, B_2\}$. Now Player 1 deletes C_1 leaving $\{B_1\}$, and Player 2 deletes B_2 , leaving $\{A_2\}$. We are left with strategy profile (B_1, A_2) .

Now suppose instead we first delete C_1 , leaving $\{A_1, B_1\}$. Then Player 2 can delete A_2 or B_2 . If Player 2 deletes A_2 , this leaves $\{B_2, C_2\}$; Player 1 deletes A_1 leaving $\{B_1\}$, and Player 2 deletes B_2 , leaving us with strategy profile (B_1, C_2) . If Player 2 instead deletes B_2 , this leaves $\{A_2, C_2\}$; Player 1 deletes C_1 , leaving $\{B_1\}$, leaving us with strategy profiles (B_1, A_2) and (B_1, C_2) .

Example 12 (Beauty contest). Suppose there are n players and each player’s strategy is a guess $x_i \in [0, 1]$. A player i wins if x_i is closer to $2/3$ of the average guess than any other

player. Ties are broken at random and a player receives payoff 1 if winning, 0 otherwise. No guess is strictly dominated. However, any guess in $(2/3, 1]$ is weakly dominated by $200/3$, and $\mathcal{W}_i^1 = [0, 2/3]$ survives the first round of iterated weak dominance. Iterating, we see that $\mathcal{W}_i^k = [0, (2/3)^k]$. It follows that $\mathcal{W}_i = \bigcap_{k=0}^{\infty} \mathcal{W}_i^k = \{0\}$. The unique weakly dominant strategy equilibrium (which is also the unique Nash equilibrium) thus has every player guessing 0.

In experimental studies, players guess significantly higher numbers than 0. Nagel (1995) takes this as evidence for bounded rationality.

2.7 Nash equilibrium

Nash equilibrium captures the notion that, given their opponents' strategies, a rational player should play a strategy that does at least as well against those opponents' strategies as any other available strategy.

Definition 22 (Nash equilibrium).

- (a) *Best response correspondence.* The *best response correspondence* $B_i : S_{-i} \rightrightarrows S_i$ for player i is defined by

$$B_i(s_{-i}) = \{s_i \in S_i \mid u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i}) \text{ for all } s'_i \in S_i\},$$

for each $s_{-i} \in S_{-i}$. The *best response correspondence* $B : S \rightrightarrows S$ is defined by

$$B(s) = \{s' \in S \mid s'_i \in B_i(s_{-i}) \text{ for all } i \in \mathcal{I}\},$$

for each $s \in S$. The extension to mixed strategies is trivial.

- (b) *Nash equilibrium.* A strategy profile $s^* \in S$ is a (pure strategy) *Nash equilibrium* if

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

A strategy profile $\sigma^* \in \times_{i \in \mathcal{I}} \Delta(S_i)$ is a (mixed strategy) *Nash equilibrium* if

$$u_i(\sigma_i^*, \sigma_{-i}^*) \geq u_i(s_i, \sigma_{-i}^*) \quad \text{for all } s_i \in S_i.$$

- (c) *Strict Nash equilibrium.* A strategy profile s^* is a *strict Nash equilibrium* if

$$u_i(s_i^*, s_{-i}^*) > u_i(s_i, s_{-i}^*) \quad \text{for all } s_i \in S_i - \{s_i^*\}.$$

All strict Nash equilibria are Nash equilibria. Whenever we refer to Nash equilibrium, we refer to the definition in (b).

A Nash equilibrium requires that every player plays a best response to their opponents' strategies. Indeed, a Nash equilibrium is a fixed point of the best response correspondence:

Proposition 7. *A strategy profile $s^* \in S$ is a (pure strategy) Nash equilibrium iff*

$$s^* \in B(s^*).$$

Likewise, a mixed strategy profile $\sigma^ \in \times_{i \in \mathcal{I}} \Delta(S_i)$ is a (mixed strategy) Nash equilibrium iff*

$$\sigma^* \in B(\sigma^*).$$

Proof. If $s^* \in B(s^*)$ then $s_i^* \in B_i(s_{-i}^*)$. Hence $u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*)$ for all $s_i \in S_i$. Since this holds across all $i \in \mathcal{I}$, it follows that s^* is a Nash equilibrium. Conversely, suppose s^* is a Nash equilibrium. Then $u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*)$ for all $s_i \in S_i$. Thus $s_i^* \in B_i(s_{-i}^*)$. Since this holds for all i , $s^* \in B(s^*)$.

The extension to mixed strategies follows immediately, since a mixed strategy game is a special case of a pure strategy game. \square

Nash equilibrium is famously the ‘workhorse’ solution concept in game theory. The (mixed strategy) Nash existence result makes Nash equilibrium particularly appealing – allowing mixed strategies, any finite player game has at least one Nash equilibrium. On close reflection, however, the concept is not quite as intuitive as might be assumed. A Nash equilibrium is a point of mutual best response. But in general, to play the Nash equilibrium solution, all players must not only be rational but also have correct beliefs about the play of their opponents. In general, multiple Nash equilibria may exist, and it is not obvious which Nash equilibrium will be played, or even if a Nash equilibrium will be played at all.

How then do we get to Nash equilibrium? Here are several suggestions, though none are applicable to all settings:

- *Nash equilibrium as the outcome of introspection.* If players are rational and reason about the behaviour of the other players, then we might expect Nash equilibrium to be played. This usually requires that we have a point prediction of rational play – for example, if there is a unique iterated strict dominance solution and rationality is common knowledge, or if we have some reasonable equilibrium selection criterion that isolates a unique “reasonable” Nash equilibrium from a set of multiple ones.
- *Nash equilibrium as the outcome of a learning process.* Nash equilibrium can be thought of as the result of a learning process. There is some empirical evidence that in repeated games, agents learn to play a Nash equilibrium over time (Smith, 1990; McCabe et al., 1991; Prasnikar & Roth, 1991). Kalai & Lehrer (1993) developing a rational learning model in an infinitely repeated game, provide theoretical motivation for why Nash equilibrium play should emerge from learning. The obvious drawback is that a learning process requires repeated interaction. It does not help explain why we might expect Nash equilibrium in a one-shot game.
- *Nash equilibrium as the outcome of an evolutionary process.* An evolutionary interpretation of Nash equilibrium has large populations of types of agents playing pure strategies against each other. These types reproduce at different rates depending

on their evolutionary fitness. This process will converge a Nash equilibrium. It is particularly intuitive in biological settings – for example, in explaining the mix of predators and prey in an ecosystem, as in the Hawk-Dove game. This only applies in matrix games (i.e. in the usual interpretation, games of two players with symmetric strategy sets and symmetric payoffs.) However, evolutionary processes are myopic, and thus do not provide a compelling description of strategic behaviour in repeated games. If players value outcomes in future periods and realize that their current actions affect their opponents' future play, then we should not expect an evolutionary process to predict their behaviour well.

- *Nash equilibrium as a self-enforcing agreement.* One interpretation of Nash equilibrium is contractual. Suppose players get together and reach an agreement that specifies the strategy profile that players should play. If this strategy profile is a Nash equilibrium, then players have no incentive to break the agreement, and so the agreement is self-enforcing. We discuss this in more detail when we get to correlated equilibrium (section 2.12)

Example 13. Consider the game,

| | a_2 | b_2 | c_2 |
|-------|--------|--------|--------|
| a_1 | 4, 3 | -1, -1 | 0, 0 |
| b_1 | -1, -1 | -2, -2 | -1, -1 |
| c_1 | 0, 0 | -1, -1 | 5, 2 |

For Player 1, b_1 is strictly dominated by a_1 (or c_1). For Player 2, b_2 is strictly dominated by a_2 (or c_2). Player 1's best response to each of Player 2's pure strategies is coloured blue and Player 2's best response to each of Player 1's pure strategies is coloured red. Note that while, for compactness of representation, we have coloured the payoff of the best response by each player, the Nash equilibrium is a *strategy profile*, not a payoff profile. We see that there are two pure strategy Nash equilibria in this game: (a_1, a_2) and (c_1, c_2) .

There is also a mixed strategy Nash equilibrium. Suppose Player 1 plays a_1 with probability p , b_1 with probability 0 (since b_1 is strictly dominated we can rule it out of any Nash equilibrium) and c_1 with probability $1 - p$. That is, $\sigma_1 = (p, 0, 1 - p)$. Likewise, suppose Player 2 plays a_2 with probability q , b_2 with probability 0 and c_2 with probability $1 - q$. That is, $\sigma_2 = (q, 0, 1 - q)$.

Player 1's expected payoff is

$$u_1(\sigma_1, \sigma_2) = p(4q) + (1 - p)(5(1 - q)).$$

In Nash equilibrium, Player 2 randomizes such that Player 1 does not have a profitable deviation from randomizing. This requires that

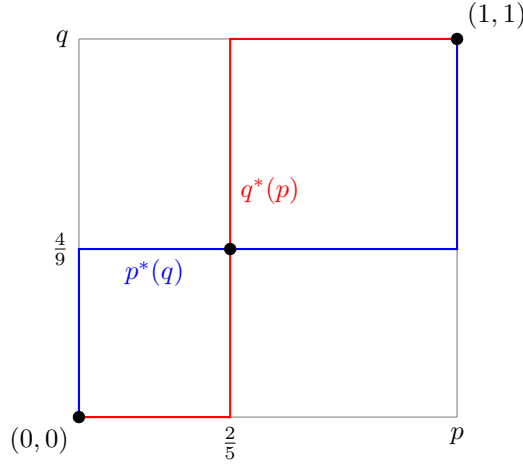
$$4q = 5(1 - q),$$

i.e. $q = \frac{5}{9}$. If $q > \frac{5}{9}$, then $u_1(a_1, \sigma_2) > u_1(\sigma_1, \sigma_2)$ for $p < 1$. If $q < \frac{5}{9}$ then $u_1(c_1, \sigma_2) > u_1(\sigma_1, \sigma_2)$ for $p > 0$.

Likewise, Player 1 randomizes so that

$$3p = 2(1 - p),$$

i.e. $p = \frac{2}{5}$. Hence we have mixed strategy Nash equilibrium $((\frac{2}{5}, 0, \frac{3}{5}), (\frac{5}{9}, 0, \frac{4}{9}))$. Graphically, we can plot the ‘best response probability’ $p^*(q)$ for Player 1 and $q^*(p)$ for Player 2 associated to the best response to σ_2 and σ_1 respectively. This is of course not the same as the best response itself, which is a (set of) strategy profile(s) over all the available pure strategies. However, given how we defined σ_1, σ_2 , it does fully characterize these best responses.



To see that correct beliefs are, in general, necessary for Nash equilibrium play, suppose Player 1 believes that Player 2 will play a_2 (i.e. Player 1’s belief over her opponent’s play is $\mu_{-1} = (1, 0, 0)$). Player 1’s best response to her belief is to play a_1 . Now suppose Player 2 believes that Player 1 will play c_1 (i.e. $\mu_{-2} = (0, 0, 1)$). Player 2’s best response to his belief is to play c_2 . The result is the profile (a_1, c_2) , and both players receive a payoff of 0. This is not a Nash equilibrium since both players have profitable deviations (c_1 and a_2 respectively).

Proposition 8. *In any finite game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, if a unique pure strategy profile survives iterated strict dominance then it is the unique pure strategy Nash equilibrium.*

Proof. Let (s_i^*, s_{-i}^*) be the unique strategy profile surviving iterated strict dominance. Suppose $u_i(s_i, s_{-i}^*) \geq u_i(s_i^*, s_{-i}^*)$ for some $s_i \neq s_i^*$.

Define a strict order $>$ on S_i by $s_i > s'_i$ iff $u_i(s_i, s_{-i}^*) > u_i(s'_i, s_{-i}^*)$.

Lemma 6. $u_i(s_i^*, s_{-i}^*) > u_i(s_i, s_{-i}^*)$ for all $s_i \in S_i$ such that $s_i \neq s_i^*$.

Proof. Since S_i is finite, there is some k_1 -level iterated strict dominance set $\mathcal{D}_i^{k_1}$ s.t. $s_i \in \mathcal{D}_i^{k_1}$ and $s_i \notin \mathcal{D}_i^{k_1+1}$. Thus we have that there is some $s'_i \in \mathcal{D}_i^{k_1}$ that strictly dominates s_i . Since $s_{-i}^* \in \mathcal{D}_{-i}^j$ for all $j \in \mathbb{N}$, it follows that

$$u_i(s'_i, s_{-i}^*) > u_i(s_i, s_{-i}^*).$$

If $u_i(s'_i, s_{-i}^*) \geq u_i(s_i^*, s_{-i}^*)$, then applying the same argument as above gives us some k_2 s.t. $s'_i \in \mathcal{D}_i^{k_2}$ but $s'_i \notin \mathcal{D}_i^{k_2+1}$. We have some $s_i^{(2)}$ that strictly dominates s'_i , implying $u_i(s_i^{(2)}, s_{-i}^*) > u_i(s'_i, s_{-i}^*)$. Since there are only finitely many members of S_i , reapplying this argument finitely many times, we can always find a chain $s_i^* > s_i^{(n)} > \dots > s'_i > s_i$. The lemma follows. \square

Since i is arbitrary, it follows from the lemma that (s_i^*, s_{-i}^*) is a pure strategy Nash equilibrium.

Now suppose there is a second pure strategy Nash equilibrium $(s'_i, s'_{-i}) \neq (s_i^*, s_{-i}^*)$. Then $u_i(s'_i, s'_{-i}) \geq u_i(s_i, s'_{-i})$ for all $s_i \in S_i$. Since $s'_i \notin \mathcal{D}_i$, there is some k -level iterated strict dominance set s.t. $s'_i \in \mathcal{D}_i^k$ but $s'_i \notin \mathcal{D}_i^{k+1}$. If $s'_{-i} \in \mathcal{D}_{-i}^{k+1}$, then we must have some $s''_i \in \mathcal{D}_i^k$ s.t. s''_i strictly dominates s'_i , implying $u_i(s''_i, s'_{-i}) > u_i(s'_i, s'_{-i})$. Given (s'_i, s'_{-i}) is a Nash equilibrium, this would yield a contradiction.

if $s'_{-i} \notin \mathcal{D}_{-i}^{k+1}$, then there is some $j \neq i$ s.t. s'_j is strictly dominated in \mathcal{D}_j^k by some $s''_j \in \mathcal{D}_j^k$, implying $u_j(s''_j, s'_{-j}) > u_j(s'_j, s'_{-j})$. Since (s'_j, s'_{-j}) is Nash, this again yields a contradiction. \square

Proposition 9. *In a finite mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, a strategy profile σ^* is a mixed strategy Nash equilibrium only if each player i chooses $\sigma_i^*(s_i) > 0$ for some $s_i \in S_i$ only if s_i is a best response to σ_{-i}^* .*

Proof. Given a strategy σ_i ,

$$u_i(\sigma_i, \sigma_{-i}^*) = \int_{S_i} u_i(s_i, \sigma_{-i}^*) d\sigma_i = \sum_{s_i \in S_i} u_i(s_i, \sigma_{-i}^*) \sigma_i(s_i).$$

Suppose σ^* is s.t. for each i , $\sigma_i^*(s_i) > 0$ for $s_i \in S_i$ only if s_i is a best response to σ_{-i}^* . We mean to show this is a sufficient condition for σ^* to be a Nash equilibrium. Let

$$B_i^s(\sigma_{-i}) = \{s_i \in S_i : u_i(s_i, \sigma_{-i}) \geq u_i(s'_i, \sigma_{-i}) \text{ for all } s'_i \in S_i\},$$

the set of pure strategies that are best responses to σ_{-i} . Clearly, $u_i(s_i, \sigma_{-i}^*) = u_i(s'_i, \sigma_{-i}^*) =: \bar{u}_i(\sigma_{-i}^*)$ for all $s_i, s'_i \in B_i^s(\sigma_{-i}^*)$, for if $u_i(s_i, \sigma_{-i}^*) > u_i(s'_i, \sigma_{-i}^*)$ then $s'_i \notin B_i^s(\sigma_{-i}^*)$, and the converse if $u_i(s'_i, \sigma_{-i}^*) > u_i(s_i, \sigma_{-i}^*)$. Now, by the hypothesis, σ_i^* assigns $\sigma_i^*(s_i) > 0$ only to $s_i \in B_i^s(\sigma_{-i}^*)$ and $\sigma_i^*(s_i) = 0$ otherwise. Hence

$$u_i(\sigma_i^*, \sigma_{-i}^*) = \sum_{s_i \in B_i^s(\sigma_{-i}^*)} u_i(s_i, \sigma_{-i}^*) \sigma_i^*(s_i) = \sum_{s_i \in B_i^s(\sigma_{-i}^*)} \bar{u}_i(\sigma_{-i}^*) \sigma_i^*(s_i) = \bar{u}_i(\sigma_{-i}^*).$$

Thus $u_i(\sigma_i^*, \sigma_{-i}^*) \geq u_i(\sigma_i, \sigma_{-i}^*)$ for all $\sigma_i \in \Delta(S_i)$. Since this holds for all i , it follows that σ^* is a Nash equilibrium.

Conversely, suppose σ^* is a Nash equilibrium. If s'_i is not a best response to σ_{-i}^* then there is some strategy s''_i s.t. $u_i(s''_i, \sigma_{-i}^*) > u_i(s'_i, \sigma_{-i}^*)$. If $\sigma_i^*(s_i) > 0$, then the strategy

$$\sigma'_i(s_i) = \begin{cases} \sigma_i^*(s_i) & \text{if } s_i \neq s'_i, s''_i, \\ 0 & \text{if } s_i = s'_i, \\ \sigma_i^*(s'_i) + \sigma_i^*(s''_i) & \text{if } s_i = s''_i, \end{cases}$$

yields expected payoff

$$\begin{aligned}
u_i(\sigma'_i, \sigma_{-i}^*) &= \sum_{s_i \neq s'_i, s''_i} u_i(s_i) \sigma_i^*(s_i) + u_i(s''_i, \sigma_{-i}^*) \sigma_i^*(s''_i) + u_i(s'_i, \sigma_{-i}^*) \sigma_i^*(s'_i) \\
&> \sum_{s_i \neq s'_i, s''_i} u_i(s_i) \sigma_i^*(s_i) + u_i(s''_i, \sigma_{-i}^*) \sigma_i^*(s''_i) + u_i(s'_i, \sigma_{-i}^*) \sigma_i^*(s'_i) \\
&= u_i(\sigma_i^*, \sigma_{-i}^*).
\end{aligned}$$

Since σ^* is a Nash equilibrium, this yields a contradiction. \square

Definition 23 (Never-best response). A pure strategy $s_i \in S_i$ of i is a *never-best response* if for all profiles $\sigma_{-i} \in \Delta_{-i}(S_{-i})$, there exists a $\sigma_i \in \Delta(S_i)$ such that

$$u_i(\sigma_i, \sigma_{-i}) > u_i(s_i, \sigma_{-i}).$$

A strictly dominated strategy is of course a never-best response. An obvious corollary to Proposition 9 is:

Corollary 4. *Suppose s_i is a never-best response. Then if σ^* is a mixed strategy Nash equilibrium $\sigma_i^*(s_i) = 0$.*

In general, some intuitively undesirable solutions can be Nash equilibria. For example:

Proposition 10. *A Nash equilibrium can be weakly dominated.*

Proof. Proof follows from the following example. \square

Example 14 (Weakly dominated Nash equilibrium). Consider the game with payoff matrix:

| | | |
|-------|-------|-------|
| | a_2 | b_2 |
| a_1 | 1, 1 | 0, 1 |
| b_1 | 1, 0 | 2, 2 |

Now, (a_1, a_2) is a Nash equilibrium since $u_1(a_1, a_2) \geq u_1(b_1, a_2)$ and $u_2(a_1, a_2) \geq u_2(a_1, b_2)$. However, b_1 is a weakly dominant strategy for Player 1 and b_2 is a weakly dominant strategy for Player 2, so (a_1, a_2) is weakly dominated.

Such a weakly dominated equilibrium is fragile in the sense that if either player's beliefs involve any degree of uncertainty over the actions of the other player, then that player would optimally choose the weakly dominant action. This is not the only kind of case where a Nash equilibrium might seem unreasonable. To isolate the more reasonable equilibria, we typically need some kind of refinement – in this case, the appropriate refinement is trembling hand perfect equilibrium.

2.7.1 Nash equilibrium inefficiency

Nash equilibria are not generally efficient. We show this by means of a counterexample.

Proposition 11. *Nash equilibria need not be Pareto efficient.*

Proof. See the following example. □

Example 2 (continued). In the prisoner's dilemma, we have payoff matrix

| | | |
|-------|-------|-------|
| | C_2 | D_2 |
| C_1 | 3, 3 | 0, 4 |
| D_1 | 4, 0 | 1, 1 |

The best responses are highlighted in blue and red for Player 1 and Player 2 respectively. We see that the (unique) Nash equilibrium is (D_1, D_2) [also, a strict Nash equilibrium and a strictly dominant strategy equilibrium]. This is Pareto inefficient since

$$u_1(C_1, C_2) = u_2(C_1, C_2) = 3 > 1 = u_1(D_1, D_2) = u_2(D_1, D_2).$$

2.7.2 Existence of Nash equilibrium

In finite games of pure strategies, Nash equilibrium need not exist. However, Nash (1950,1951) proves the existence of mixed strategy Nash equilibrium in games of finite players. Debreu (1952), Glicksberg (1952) and Fan (1952) prove a more general existence result.

Example 9 (continued). Recall the payoff matrix for matching pennies is

| | | |
|-------|-------|-------|
| | H_2 | T_2 |
| H_1 | 1, -1 | -1, 1 |
| T_1 | -1, 1 | 1, -1 |

The best responses for Player 1 is highlighted in blue and for Player 2 in red. We see that there is no pure strategy Nash equilibrium.

The proof of the existence results for the Debreu-Glicksberg-Fan and Nash existence results rely on Kakutani's fixed point theorem (Theorem 6). For a discussion of correspondences and upper hemicontinuity, see the Mathematical Appendix, section 10.1.

Theorem 6 (Kakutani's fixed point theorem). *Let $K \subset \mathbb{R}^n$ be a nonempty, compact, convex set and suppose $F : K \rightrightarrows K$ is a correspondence satisfying*

- (i) $F(x)$ is nonempty valued;
- (ii) $F(x)$ is convex valued;
- (iii) $F(x)$ is upper hemicontinuous.

Then F has a fixed point.

We discuss fixed point theorems, including this one, in more detail in section 10.4.

Theorem 7 (Debreu-Glicksberg-Fan, 1952). *Let $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ be a game of $n \in \mathbb{N}$ players such that for each player $i \in \mathcal{I}$, $S_i \subset \mathbb{R}^{k_i}$ is nonempty, convex and compact. If, for each $i \in \mathcal{I}$, u_i is continuous in s and quasiconcave in s_i , then the game has a pure strategy Nash equilibrium.*

Proof. Since each S_i is convex and compact, it follows that $S = \times_{i \in \mathcal{I}} S_i$ is convex and compact. Furthermore, since u_i is continuous in S_i and S_i is compact, to each s_{-i} there is some $s_i \in S_i$ s.t. $s_i \in \arg \max_{s'_i \in S_i} u_i(s'_i, s_{-i})$, by the extreme value theorem. Hence the best response correspondence $B_i(s_{-i})$ is nonempty for all $s_{-i} \in S_{-i}$.

Recall that a function f defined on a convex set X is quasiconcave if $f(\lambda x + (1-\lambda)y) \geq \min\{f(x), f(y)\}$ for all $\lambda \in [0, 1]$ and all $x, y \in X$.

Since u_i is quasiconcave in s_i , if $s_i, s'_i \in B_i(s_{-i})$ then

$$u_i(\lambda s_i + (1-\lambda)s'_i) \geq \min\{u_i(s_i), u_i(s'_i)\}$$

for all $\lambda \in [0, 1]$. Hence any convex combination of s_i, s'_i achieves payoff against s_{-i} that is at least as great as that of s_i or s'_i . Now, $u_i(s_i, s_{-i}) = u_i(s'_i, s_{-i})$, for if $u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$ then $s'_i \notin B_i(s_{-i})$ and likewise if $u_i(s'_i, s_{-i}) > u_i(s_i, s_{-i})$ then $s_i \notin B_i(s_{-i})$. Applying the same argument to the convex combination $\lambda s_i + (1-\lambda)s'_i$, we have $u_i(\lambda s_i + (1-\lambda)s'_i, s_{-i}) = u_i(s_i, s_{-i}) = u_i(s'_i, s_{-i})$ and hence $\lambda s_i + (1-\lambda)s'_i \in B_i(s_{-i})$ for all $\lambda \in [0, 1]$. It follows that $B_i(s_{-i})$ is convex for all $s_{-i} \in S_{-i}$. Since this holds for all i , we have that $B(s)$ is convex for all $s \in S$, i.e. B is convex valued.

Since u_i is continuous, u^{-1} maps closed sets into closed sets. Define $\bar{u}_i(s_{-i}) := u_i(s_i, s_{-i})$ for any $s_i \in B_i(s_{-i})$. By definition of the best response correspondence, $B_i(s_{-i}) = u^{-1}(\bar{u}_i(s_{-i}))$, and the latter is a single point and thus closed. Hence $B_i(s_{-i})$ is closed for all $s_{-i} \in S_{-i}$ and all i . Since $B_i(s_{-i}) \subseteq S_i$, a compact set, it follows that $B_i(s_{-i})$ is therefore also compact. Hence B is compact valued.

This allows us to use the sequence definition of upper hemicontinuity. For any $\{s^n\}$ in S s.t. $s^n \rightarrow s$ and $\{\hat{s}^n\}$ s.t. $\hat{s}^n \in B(s^n)$ for all n and $\hat{s}^n \rightarrow \hat{s}$, we mean to prove that $\hat{s} \in B(s)$. Suppose otherwise. Then for some player i , $\hat{s}_i \notin B_i(s_{-i})$ and so there exists s'_i s.t.

$$u_i(\hat{s}_i, s_{-i}) < u_i(s'_i, s_{-i}).$$

Yet by continuity of u_i , this implies that for some $N \in \mathbb{N}$,

$$u_i(\hat{s}_i^n, s_{-i}^n) < u_i(s'_i, s_{-i}^n)$$

for all $n \geq N$, and hence $\hat{s}_i^n \notin B_i(s_{-i}^n)$ for some n , yielding a contradiction. Hence B_i is upper hemicontinuous for all i and so B is upper hemicontinuous. We could have instead concluded that B is upper hemicontinuous by using Berge's theorem of the maximum (Theorem 45) to conclude that B has a closed graph, and then applying the closed graph theorem (Theorem 44).

Applying Kakutani's fixed point theorem, we have that B has a fixed point, i.e. there exists some $s^* \in S$ s.t. $s^* \in B(s^*)$. \square

Nash's existence theorem came earlier, and was a remarkable achievement in its own right, but armed with the Debreu-Glicksberg-Fan theorem, it becomes a corollary:

Theorem 8 (Nash's existence theorem, 1950). *Any game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ with a finite number of players and such that S_i is finite and nonempty for all $i \in \mathcal{I}$ has a mixed strategy Nash equilibrium.*

Proof. Clearly, $\Delta(S_i)$ is a nonempty, compact and convex subset of $\mathbb{R}^{|S_i|}$.

By linearity of the Lebesgue integral, u_i is continuous.

For any $\sigma_i, \sigma'_i \in \Delta(S_i)$ and any $\lambda \in [0, 1]$, we have

$$\begin{aligned} u_i(\lambda\sigma_i + (1-\lambda)\sigma'_i, \sigma_{-i}) &= \sum_{s \in S} u_i(s) [\lambda\sigma_i(s) + (1-\lambda)\sigma'_i(s)] \sigma_{-i}(s_{-i}) \\ &= \sum_{s \in S} u_i(s) [\lambda\sigma_i(s) + (1-\lambda)\sigma'_i(s)] \sigma_{-i}(s_{-i}) \\ &= \lambda \sum_{s \in S} u_i(s) \sigma_i(s) \sigma_{-i}(s_{-i}) + (1-\lambda) \sum_{s \in S} u_i(s) \sigma'_i(s) \sigma_{-i}(s_{-i}) \\ &\geq \min \left\{ \sum_{s \in S} u_i(s) \sigma_i(s) \sigma_{-i}(s_{-i}), \sum_{s \in S} u_i(s) \sigma'_i(s) \sigma_{-i}(s_{-i}) \right\} \\ &= \min \{ u_i(\sigma_i, \sigma_{-i}), u_i(\sigma'_i, \sigma_{-i}) \}. \end{aligned}$$

Hence u_i is quasiconcave in σ_i .

It follows that the hypotheses of the Debreu-Glicksberg-Fan theorem are satisfied. \square

Nash (1951) gives a direct proof using Brouwer's fixed point theorem:

Proof. With slight abuse of notation, take s_i to represent player i 's degenerate mixed strategy in which they play s_i with probability 1. For each pure strategy $s_i \in S_i$ of player i , define the function $f_i^{s_i} : \times_{i=1}^n \Delta(S_i) \rightarrow \mathbb{R}$ by

$$f_i^{s_i}(\sigma) = \max \{ 0, u_i(s_i, \sigma_{-i}) - u_i(\sigma_i, \sigma_{-i}) \}.$$

Note each $f_i^{s_i}$ is continuous, since u_i is continuous in σ and s_i is held fixed. We now define a mapping T of $\times_{i=1}^n \Delta(S_i)$ into itself as follows. For each σ , define

$$\sigma'_i = \frac{\sigma_i + \sum_{s_i \in S_i} s_i f_i^{s_i}(\sigma)}{1 + \sum_{s_i \in S_i} f_i^{s_i}(\sigma)}.$$

Let $T(\sigma) := (\sigma'_1, \dots, \sigma'_n)$.

Now, if σ is a Nash equilibrium of G^m , then if $s_i \in \text{supp}(\sigma_i)$, we must have $u_i(\sigma_i, \sigma_{-i}) = u_i(s_i, \sigma_{-i}) \geq u_i(s'_i, \sigma_{-i})$ for all $s'_i \in S_i$ and so $f_i^{s_i}(\sigma) = 0$ for all $s_i \in S_i$. Conversely, if $f_i^{s_i}(\sigma) = 0$ for all $s_i \in S_i$ and every player i , then for those $s_i \in \text{supp}(\sigma)$, we have that there is no $s'_i \in S_i$ with $u_i(s'_i, \sigma_{-i}) > u_i(s_i, \sigma_{-i})$, and so each pure strategy s_i in the support of σ_i is a best response to σ_{-i} , so by linearity, σ_i is a best response to σ_{-i} . Hence σ is a Nash equilibrium.

Now, T is continuous in σ , by continuity of the functions $f_i^{s_i}$. Clearly, $\times_{i=1}^n \Delta(S_i)$ is nonempty, convex and compact, and so applying Brouwer's fixed point theorem (Theorem 57) gives us that T has a fixed point. Hence G^m has a Nash equilibrium. \square

2.7.3 Upper hemicontinuity and Nash equilibria

We can make use of the fact that best response correspondences are upper hemicontinuous (under the appropriate assumptions on payoffs and strategy spaces) to find Nash equilibria in limits of games.

Definition 24. Let Λ be a parameter space, which we assume to be a compact metric space. Consider a finite set of players \mathcal{I} and let $S = \times_{i \in \mathcal{I}} S_i$ be the set of strategy profiles. Suppose each player $i \in \mathcal{I}$ has a payoff function $u_i : S \times \Lambda \rightarrow \mathbb{R}$ that is continuous in both strategy profiles S and parameters Λ .

1. *Family of games.* For each $\lambda \in \Lambda$, define the game $G(\lambda) := (\mathcal{I}, (S_i, u_i(\cdot, \lambda))_{i \in \mathcal{I}})$. We call $\mathcal{G}_\Lambda = \{G(\lambda) \mid \lambda \in \Lambda\}$ a *family of games* parameterized by Λ .
2. *Nash equilibrium correspondence.* Given a family of games \mathcal{G}_Λ , the *Nash equilibrium correspondence* $\text{NE} : \Lambda \rightrightarrows S$ is defined as

$$\text{NE}(\lambda) = \{s \in S \mid s \text{ is a Nash equilibrium of } G(\lambda)\}$$

for all $\lambda \in \Lambda$.

We assume the value of λ is common knowledge in game $G(\lambda)$.

Proposition 12. Under the assumptions of Definition 24, the Nash equilibrium correspondence of a family of games \mathcal{G} has a closed graph.

Proof. Consider any sequence $\{(s^n, \lambda^n)\}$ such that $(s^n, \lambda^n) \rightarrow (s, \lambda)$ with $s^n \in \text{NE}(\lambda^n)$ for each $n \in \mathbb{N}$. Suppose that $s \notin \text{NE}(\lambda)$. Then, for some player $i \in \mathcal{I}$, $u_i(s'_i, s_{-i}, \lambda) > u_i(s_i, s_{-i}, \lambda)$ for some $s'_i \in S_i$. Since u_i is continuous and $(s^n, \lambda^n) \rightarrow (s, \lambda)$, it follows that $u_i(s'_i, s_{-i}^n, \lambda^n) > u_i(s_i^n, s_{-i}^n, \lambda^n)$ for sufficiently large n , but then $s^n \notin \text{NE}(\lambda^n)$, yielding a contradiction. \square

By the closed graph theorem (Theorem 44), the Nash equilibrium correspondence is thus upper hemicontinuous. This implies that if we take a sequence of parameters $\{\lambda_n\}$ such that $\lambda_n \rightarrow \lambda$, then $\lim_{n \rightarrow \infty} \text{NE}(\lambda_n) \subseteq \text{NE}(\lambda)$. Thus we can find Nash equilibria of games by looking at approximations of those games and taking limits. Note however that the set inclusion relation only goes in one direction: because the Nash equilibrium correspondence is not necessarily lower hemicontinuous, there may be Nash equilibria of the limit game that are not limits of Nash equilibria of the approximations.

Since any mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ can be thought of as effectively a pure strategy game with strategy sets $\Delta(S_i)$ for each i , Proposition 12 clearly holds for mixed strategy games.

2.7.4 Uniqueness of Nash equilibrium

In most games, multiplicity of Nash equilibria is the norm – indeed, many games have infinitely many Nash equilibria. This fact is very annoying. We can rely on refinements or other selection criteria to eliminate “unreasonable” equilibria, but even so, often this does not necessarily leave us with a point prediction.

It is therefore interesting to ask under what conditions we get unique Nash equilibria. In general, this is a futile task, but we have some results for specific classes of games. Probably the most well-known result is due to Rosen (1965), who gives a sufficient condition for uniqueness in concave games.

Definition 25 (Concave game). Consider an n -player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, and let $S = S_1 \times \cdots \times S_n$. We say that G is *concave* if

- (i) S_i is convex and compact for each player $i \in \mathcal{I}$;
- (ii) $u_i(s)$ is continuous in $s \in S$ for each $i \in \mathcal{I}$, and
- (iii) for each player i and each strategy profile $s_{-i} \in S_{-i}$, $u_i(s_i, s_{-i})$ is concave in $s_i \in S_i$.

A pure strategy Nash equilibrium s^* of this game by definition must solve, for each player i , the maximization problem $\max_{s_i} u_i(s_i, s_{-i}^*)$ subject to $(s_i, s_{-i}^*) \in S$.

Immediately from the Debreu-Glicksberg-Fan theorem, any concave game has a pure strategy Nash equilibrium.

To discuss uniqueness, we need some machinery to describe the strategy sets. Suppose each player i 's strategy set $S_i \subset X^{m_i}$ where X is an Euclidean space. Then $S \subset X^m$ with $m = \sum_{i=1}^n m_i$. We assume S is described by some function $h : X^m \rightarrow X^k$ where each j th component $h_j : X^m \rightarrow X$ of h is a concave function. We define

$$S = \{s \in X \mid h(s) \geq 0\},$$

and we assume S is nonempty and bounded. Since it is closed and each h_j is concave, it follows that each S_i is compact and convex.

Some more technical conditions are needed: assume there is some point $\hat{s} \in S$ such that $h_j(\hat{s}) > 0$ for each $j = 1, \dots, k$, and assume each h_j is continuously differentiable. Also assume that for each $s \in S$, $u_i(s)$ has continuous first derivatives wrt the components of s_i . These conditions are all necessary to apply the Kuhn-Tucker conditions for a constrained maximum.

For any function $f : X \rightarrow \mathbb{R}$, let $\nabla_i f$ denote the gradient of f with respect to the components of s_i . The Kuhn-Tucker conditions for a solution s^* to the maximization problem are $h(s^*) \geq 0$, and for each player i , there exists some nonnegative $\lambda_i^* \in X^{k_i}$ such that

- (i) $\lambda_i^* \cdot h(s^*) = 0$, and

- (ii) $u_i(s^*) \geq u_i(s_i, s_{-i}^*) + \lambda_i^* \cdot h(s_i, s_{-i}^*)$ for all strategies $s_i \in S_i$. Since ϕ_i and h_j are concave and differentiable, this can equivalently be stated as

$$\nabla_i u_i(s^*) + \sum_{j=1}^k \lambda_{ij}^* \nabla_i h_j(s^*) = 0.$$

We define a weighted sum of the players' payoff,

$$w(s, \lambda) = \sum_{i=1}^n \lambda_i u_i(s)$$

for each payoff profile $s \in S$ and each nonnegative $\lambda \in X^n$.

Definition 26 (Diagonal strict concavity).

- (a) *Pseudogradient*. The *pseudogradient* g of the weighted sum w is defined as

$$g(s, \lambda) = \begin{pmatrix} \lambda_1 \nabla_1 u_1(s) \\ \lambda_2 \nabla_2 u_2(s) \\ \vdots \\ \lambda_n \nabla_n u_n(s) \end{pmatrix}.$$

- (b) *Diagonal strict concavity*. We call $w(s, \lambda)$ *diagonally strictly concave* for fixed $\bar{\lambda} \geq 0$ if for every $s, s' \in S$, we have that

$$(s' - s) \cdot g(s, \bar{\lambda}) + (s - s') \cdot g(s', \bar{\lambda}) > 0.$$

To check that w is diagonally strictly concave, we can check the Jacobian of the pseudogradient, because of the following sufficient condition:

Proposition 13. *Consider the weighted sum $w(s, \lambda)$ and let $J(s, \lambda)$ be the Jacobian of the corresponding pseudogradient $g(s, \lambda)$. Then w is diagonally strictly concave for some fixed $\bar{\lambda} > 0$ if the symmetric matrix $J(s, \bar{\lambda}) + J(s, \bar{\lambda})^\top$ is negative definite for all $s \in S$.*

Proof. Let $s^1, s^2 \in S$ be distinct and define $s(\theta) = \theta s^1 + (1 - \theta) s^2$ for $\theta \in [0, 1]$. Since S is convex, $s(\theta) \in S$. We have

$$\frac{dg(s(\theta), \bar{\lambda})}{d\theta} = J(s(\theta), \bar{\lambda}) \frac{ds(\theta)}{d\theta} = J(s(\theta), \bar{\lambda})(s^2 - s^1),$$

given $J(s, \bar{\lambda})$ is the Jacobian of $g(s, \bar{\lambda})$. Equivalently, we have

$$g(s^2, \bar{\lambda}) - g(s^1, \bar{\lambda}) = \int_0^1 J(s(\theta), \bar{\lambda})(s^2 - s^1) d\theta.$$

Left-multiplying by $(s^1 - s^2)^\top$ gives us

$$\begin{aligned} (s^1 - s^2)^\top g(s^2, \bar{\lambda}) + (s^1 - s^2)^\top g(s^1, \bar{\lambda}) &= - \int_0^1 (s^2 - s^1)^\top J(s(\theta), \bar{\lambda}) (s^2 - s^1) d\theta \\ &= - \frac{1}{2} \int_0^1 (s^2 - s^1)^\top (J(s(\theta), \bar{\lambda}) + J(s(\theta), \bar{\lambda})^\top) (s^2 - s^1) d\theta \\ &> 0, \end{aligned}$$

and so w is strictly diagonally concave. \square

Now we can state Rosen's theorem:

Theorem 9 (Rosen's uniqueness theorem). *Consider an n -player concave game G . If $w(s, \lambda)$ is diagonally strictly concave for some $\lambda > 0$ and s^* is a Nash equilibrium of G , then s^* is the unique Nash equilibrium of G .*

Proof. Suppose s^1 and s^2 are two Nash equilibria of G . The Kuhn-Tucker conditions tell us that for each player i and each of the two strategy profiles s^ℓ , $h_i(s_i^\ell) \geq 0$ and there exists $\lambda_i^\ell \in X^{k_i}$ so that $\lambda_i^\ell \cdot h_i(s_i^\ell) = 0$ and $\nabla_i u_i(s^\ell) + \sum_{j=1}^{k_i} \nabla_i h_{ij}(s_i^\ell) = 0$. Left-multiplying this final condition through by $\bar{\lambda}_i(s_i^2 - s_i^1)^\top$ for $\ell = 1$ and $\bar{\lambda}_i(s_i^1 - s_i^2)^\top$ for $\ell = 2$, and taking the sum gives $\beta + \gamma = 0$ for

$$\begin{aligned} \beta &= (s^2 - s^1)^\top g(s^1, \bar{\lambda}) + (s^1 - s^2)^\top g(s^2, \bar{\lambda}), \\ \gamma &= \sum_{i=1}^n \sum_{j=1}^{k_i} \bar{\lambda}_i \left[\lambda_{ij}^1 (s_i^2 - s_i^1)^\top \nabla_i h_{ij}(s_i^1) + \lambda_{ij}^2 (s_i^1 - s_i^2)^\top \nabla_i h_{ij}(s_i^2) \right] \\ &\geq \sum_{i=1}^n \sum_{j=1}^{k_i} \bar{\lambda}_i \left[\lambda_{ij}^1 (h_{ij}(s_i^2) - h_{ij}(s_i^1)) + \lambda_{ij}^2 (h_{ij}(s_i^1) - h_{ij}(s_i^2)) \right] \\ &= \sum_{i=1}^n \bar{\lambda}_i \left[\lambda_i^1 \cdot h_i(s_i^2) + \lambda_i^2 \cdot h_i(s_i^1) \right], \end{aligned}$$

where the inequality follows because h_i is concave and

$$h_j(s^2) - h_j(s^1) \leq (s^2 - s^1)^\top \nabla h_j(s^1) = \sum_{i=1}^n (s_i^2 - s_i^1) \cdot \nabla_i h_j(s^1)$$

for all j . Now, $h_i(s_i^\ell) \geq 0$ implies that $\gamma \geq 0$, and by strict diagonal concavity, we have that $\beta > 0$. Thus $\gamma + \beta > 0$, yielding a contradiction. Therefore $s^1 = s^2$. We conclude then that the Nash equilibrium is unique. \square

Because strict diagonal concavity in concave games gives us uniqueness and the functions associated with these games have nice differentiability properties, it is also very easy to compute the unique Nash equilibrium using gradient descent methods. This is very desirable, because computing Nash equilibria in general is computationally complex – formally speaking, it is PPAD-complete (Daskalakis, Goldberg & Papadimitrou, 2009; Chen & Deng, 2006).

2.8 (Trembling hand) perfect equilibrium

In retrospect the earlier use of the word “perfect” was premature. Therefore a perfect equilibrium point in the old sense will be called “subgame perfect”. The new definition of perfectness has the property that a perfect equilibrium point is always subgame perfect but a subgame perfect equilibrium point may not be perfect.

“Let’s all aspire to the bravery of Selten (1975), renaming the original “perfect equilibrium” to make way for the new perfect equilibrium, confident that the concept has now reached perfection.” – Shengwu Li (@ShengwuLi), Twitter, 31 January 2022

A serious drawback of Nash equilibrium in general is multiplicity, and so we may consider selection criteria and refinements to rule out certain Nash equilibria that are less plausible under certain criteria.

Trembling hand perfect equilibrium (or just, *perfect equilibrium*), introduced by Selten (1975), is a refinement of Nash equilibrium that excludes Nash equilibria that are fragile to noise in players’ beliefs.

Example 14 (continued). Recall that the game in Example 14 has payoff matrix

$$\begin{array}{cc} & \begin{array}{cc} a_2 & b_2 \end{array} \\ \begin{array}{c} a_1 \\ b_1 \end{array} & \begin{array}{cc} 1, 1 & 0, 1 \\ 1, 0 & 2, 2 \end{array} \end{array}$$

This game has two pure strategy Nash equilibria, (a_1, a_2) and (b_1, b_2) . Recall that b_1 and b_2 are weakly dominant strategies for Players 1 and 2 respectively. It follows that (a_1, a_2) can only be sustained if each player i has sure belief that the other player j will play a_j . The equilibrium is thus not robust – any very small doubt on the part of either player will ensure that the (rational) player plays b_i .

Definition 27 (Trembling hand perfect equilibrium). Consider a mixed strategy normal form game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ with n players.

- (a) *Totally mixed strategy*. A mixed strategy $\sigma_i \in \Delta(S_i)$ is called *totally mixed* or *completely mixed* if $\sigma_i(s_i) > 0$ for all $s_i \in S_i$.
- (b) *Perturbed game*. A *perturbation* for a player i is a function $\epsilon_i : S_i \rightarrow (0, 1)$ such that

$$\sum_{s_i \in S_i} \epsilon_i(s_i) < 1.$$

A *perturbation* $\epsilon : S \rightarrow (0, 1)^n$ is a function defined by $\epsilon(s) = (\epsilon_i(s_i)_{i \in \mathcal{I}})$ where ϵ_i is a perturbation for player i .

An ϵ -*perturbed game* G_ϵ of G^m is a game $G_\epsilon^m = (\mathcal{I}, (\Delta_{\epsilon_i}(S_i), u_i)_{i \in \mathcal{I}})$ with

$$\Delta_{\epsilon_i}(S_i) = \{\sigma_i \in \Delta(S_i) \mid \sigma_i(s_i) \geq \epsilon_i(s_i) \text{ for all } s_i \in S_i\}$$

for each $i \in \mathcal{I}$, where ϵ_i is a perturbation. That is, $\Delta_{\epsilon_i}(S_i)$ is the set of totally mixed strategies of player i bounded below by the perturbation ϵ_i .

- (c) *Trembling hand perfect equilibrium* (normal form). A mixed strategy profile σ^* of $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ is a (*trembling hand*) *perfect equilibrium* if
- (i) there exists a sequence of perturbations $\{\epsilon^k\}$ such that $\epsilon_i^k(s_i) \rightarrow 0$ for all $s_i \in S_i$ and all $i \in \mathcal{I}$, and
 - (ii) for each $k \in \mathbb{N}$, there exists a (totally mixed) Nash equilibrium σ^k of the ϵ^k -perturbed game $G_{\epsilon^k}^m$ such that $\sigma^k \rightarrow \sigma^*$.¹²
- (d) *Trembling hand perfection in extensive form games*. The *agent-normal form* of an extensive form game Γ with information sets $(\Phi_i)_{i \in \mathcal{I}}$ is the corresponding normal form game $G = \left((\Phi_i)_{i \in \mathcal{I}}, (A_i(\phi_i), u_i)_{(\Phi_i)_{i \in \mathcal{I}}} \right)$, that is, assigning to each information set ϕ_i of player i a new ‘player’ with payoff function u_i and strategy set $A_i(\phi_i)$. A (*trembling hand*) *perfect equilibrium* for Γ is a perfect equilibrium of the corresponding agent-normal form game.

Note that a perfect equilibrium σ^* need not itself be totally mixed. We also do not need to limit the definition to mixed strategy games – it makes sense to talk about perfect equilibria of pure strategy games but we must consider its mixed strategy counterpart when it comes to modelling trembles.

Note also that we can equivalently define perfect equilibrium in terms of best responses:

Proposition 14. *A mixed strategy profile σ^* of $G^m(\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ is a perfect equilibrium iff there exists a sequence of totally mixed strategy profiles $\{\sigma^k\}$ such that*

- (i) $\sigma^k \rightarrow \sigma^*$, and
- (ii) for all $k \in \mathbb{N}$ and for each $i \in \mathcal{I}$, $\sigma_i^* \in B_i(\sigma_{-i}^k)$, that is,

$$u_i(\sigma_i^*, \sigma_{-i}^k) \geq u_i(s_i, \sigma_{-i}^k) \quad \text{for all } s_i \in S_i.$$

Proof. The proof is quite longwinded. See Selten (1975), pp. 49-51. □

Perfect equilibrium excludes weakly dominated Nash equilibria:

Proposition 15. *Consider any mixed strategy finite player game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, where S_i is finite for all $i \in \mathcal{I}$.*

- (i) *Any perfect equilibrium σ^* of G^m is a Nash equilibrium in weakly undominated strategies.*
- (ii) *If G^m is a two player game, if σ^* is a Nash equilibrium and if σ^* is not weakly dominated, then σ^* is trembling hand perfect.*

¹²We say that a Nash equilibrium of an ϵ -perturbed game is ϵ -perfect.

Proof. (i) Assume G^m has a perfect equilibrium σ^* . Then there exists a sequence of perturbations $\{\epsilon^k\}$ satisfying Definition 27(c)(i) and for each k , there exists a totally mixed Nash equilibrium σ^k of the ϵ^k -perturbed game $G_{\epsilon^k}^m$ s.t. $\sigma^k \rightarrow \sigma^*$. First we show σ^* is a Nash equilibrium. Suppose otherwise. Then for some player $i \in \mathcal{I}$, there exists a strategy $s'_i \in S_i$ s.t. $u_i(s'_i, \sigma_{-i}^*) > u_i(\sigma_i^*, \sigma_{-i}^*)$. It follows that there exists some $K \in \mathbb{N}$,

$$u_i(s'_i, \sigma_{-i}^k) > u_i(\sigma_i^k, \sigma_{-i}^k)$$

for all $k \geq K$, since $\sigma^k \rightarrow \sigma^*$. Hence for some perturbed game σ^k is not a Nash equilibrium of the perturbed game $G_{\epsilon^k}^m$, yielding a contradiction.

Next we show that if s'_i is a weakly dominated strategy, then $\sigma_i^*(s'_i) = 0$. Suppose s'_i is weakly dominated, necessarily by some strategy $\hat{\sigma}_i$ with $\hat{\sigma}_i(s'_i) = 0$. Then there is some profile s'_{-i} s.t. $u_i(s'_i, s'_{-i}) < u_i(\hat{\sigma}_i, s'_{-i})$. Now, for any totally mixed strategy profile σ_{-i} ,

$$\begin{aligned} u_i(s'_i, \sigma_{-i}) &= \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) \sigma_{-i}(s_{-i}) \\ &= u_i(s'_i, s'_{-i}) \sigma_{-i}(s'_{-i}) + \sum_{s_{-i} \neq s'_{-i}} u_i(s'_i, s_{-i}) \sigma_{-i}(s_{-i}) \\ &< u_i(\hat{\sigma}_i, s'_{-i}) \sigma_{-i}(s'_{-i}) + \sum_{s_{-i} \neq s'_{-i}} u_i(\hat{\sigma}_i, s_{-i}) \sigma_{-i}(s_{-i}) \\ &= u_i(\hat{\sigma}_i, \sigma_{-i}). \end{aligned}$$

Suppose $\sigma_i^*(s'_i) > 0$ and define $\sigma'_i = \sigma_i^*(s_i) + \sigma_i^*(s'_i) \hat{\sigma}_i(s_i)$ for all $s_i \in S_i$. Then for any k , and any totally mixed profile σ_{-i}^k

$$\begin{aligned} u_i(\sigma_i^*, \sigma_{-i}^k) &= \sum_{s_i \in S_i} u_i(s_i, \sigma_{-i}^k) \sigma_i^*(s_i) \\ &= u_i(s'_i, \sigma_{-i}^k) \sigma_i^*(s'_i) + \sum_{s_i \neq s'_i} u_i(s_i, \sigma_{-i}^k) \sigma_i^*(s_i) \\ &< u_i(\hat{\sigma}_i, \sigma_{-i}^k) \sigma_i^*(s'_i) + \sum_{s_i \neq s'_i} u_i(s_i, \sigma_{-i}^k) \sigma_i^*(s_i) \\ &= u_i(\sigma'_i, \sigma_{-i}^k). \end{aligned}$$

Hence if σ_i^* plays a weakly dominated strategy with positive probability, it is not a best response to any totally mixed strategy σ_{-i}^k , and thus σ^* cannot be perfect.

(ii) Proof omitted. □

Proposition 16. *In a mixed strategy finite player game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$, if a player i has a strictly dominant strategy s_i then in any perfect equilibrium σ^* of G^m , $\sigma_i^*(s_i) = 1$.*

Proof. Suppose otherwise for some i . Let s'_i be a strictly dominant strategy. Since σ^* is a perfect equilibrium, there is some sequence of totally mixed profiles $\{\sigma^k\}$ s.t. $\sigma^k \rightarrow \sigma^*$ and $u_i(\sigma_i^*, \sigma_{-i}^k) \geq u_i(s_i, \sigma_{-i}^k)$ for all k and all $s_i \in S_i$, yet $\sigma_i^*(s'_i) < 1$. But, since s'_i is strictly dominant, it is the unique best response, i.e.

$$u_i(s'_i, \sigma_{-i}^k) > u_i(s_i, \sigma_{-i}^k)$$

for any σ_i s.t. $\sigma_i(s_i) \neq 1$. This yields a contradiction. \square

Lemma 7. *Any finite player ϵ -perturbed game $G_\epsilon^m = (\mathcal{I}, (\Delta_{\epsilon_i}(S_i), u_i)_{i \in \mathcal{I}})$, where S_i is finite, has a (totally) mixed strategy Nash equilibrium.*

Proof. $\Delta_{\epsilon_i}(S_i)$ is convex and compact, since $\Delta_{\epsilon_i}(S_i) = \Delta(S_i) \cap [\epsilon_i(s_1), 1] \times \cdots \times [\epsilon_i(s_{k_i}), 1]$, where $k_i = |S_i|$. Since this is the intersection of two convex sets, it is convex, and since it is the finite intersection of two compact (and thus closed) sets, it is a closed subset of a compact set in \mathbb{R}^{k_i} , and thus compact.

By definition of ϵ_i , $\sum_{s_i \in S_i} \epsilon_i(s_i) < 1$. Let $\delta_i = 1 - \sum_{s_i \in S_i} \epsilon_i(s_i)$. Then $\delta_i > 0$. Define $\sigma_i(s_i) = \epsilon_i(s_i) + \frac{\delta_i}{k_i}$ for all $s_i \in S_i$. Then $\sum_{s_i \in S_i} \sigma_i(s_i) = 1$, so $\sigma_i \in \Delta(S_i)$. Furthermore, $\sigma_i(s_i) \geq \epsilon_i(s_i)$ for all $s_i \in S_i$, and so $\sigma_i \in \Delta_{\epsilon_i}(S_i)$. Hence $\Delta_{\epsilon_i}(S_i)$ is nonempty. Since these properties hold for all i , it follows that $\Delta_\epsilon(S) := \times_{i \in \mathcal{I}} \Delta_{\epsilon_i}(S_i)$ is compact, convex and nonempty.

Now u_i is continuous in σ , by linearity, and quasiconcave in σ_i , since for any $\sigma_i, \sigma'_i \in \Delta_{\epsilon_i}(S_i)$ and any $\lambda \in [0, 1]$,

$$\begin{aligned} u_i(\lambda \sigma_i + [1 - \lambda] \sigma'_i, \sigma_{-i}) &= \sum_{s \in S} u_i(s) [\lambda \sigma_i + (1 - \lambda) \sigma'_i](s) \sigma_{-i}(s_{-i}) \\ &= \sum_{s \in S} u_i(s) [\lambda \sigma_i(s_i) + (1 - \lambda) \sigma'_i(s_i)] \sigma_{-i}(s_{-i}) \\ &= \lambda \sum_{s \in S} u_i(s) \sigma_i(s_i) \sigma_{-i}(s_{-i}) + (1 - \lambda) \sum_{s \in S} u_i(s) \sigma'_i(s_i) \sigma_{-i}(s_{-i}) \\ &\geq \min \left\{ \sum_{s \in S} u_i(s) \sigma_i(s_i) \sigma_{-i}(s_{-i}), \sum_{s \in S} u_i(s) \sigma'_i(s_i) \sigma_{-i}(s_{-i}) \right\} \\ &= \min\{u_i(\sigma_i, \sigma_{-i}), u_i(\sigma'_i, \sigma_{-i})\}. \end{aligned}$$

The hypotheses of the Debreu-Glicksberg-Fan theorem (Theorem 7) thus hold, and hence there exists a (necessarily totally mixed) Nash equilibrium of the ϵ -perturbed game. \square

Theorem 10 (Existence of perfect equilibria). *Any finite-player mixed strategy game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ for which each S_i is finite has at least one perfect equilibrium.*

Proof. Given perturbations ϵ_i for player i , call the function $\epsilon : S \rightarrow (0, 1)^n$ defined by $\epsilon(s) = (\epsilon_i(s_i)_{i \in \mathcal{I}})$ a perturbation.

By the lemma, for any sequence of perturbations $\{\epsilon^k\}$ there exists a sequence $\{\sigma^k\}$ s.t. each σ^k is a Nash equilibrium of the ϵ^k -perturbed game. Suppose $\epsilon^k \rightarrow 0$. Since

$\times_{i \in \mathcal{I}} \Delta(S_i)$ is a compact set, the sequence $\{\sigma^k\}$ must have some accumulation point σ^* , and hence there exists some subsequence $\{\sigma^{k_j}\}$ s.t. $\sigma^{k_j} \rightarrow \sigma^*$. Taking this subsequence, and the corresponding subsequence $\{\epsilon^{k_j}\}$ of $\{\epsilon^k\}$, there exists a sequence of perturbations $\{\epsilon^{k_j}\}$ s.t. $\epsilon^{k_j} \rightarrow 0$ and a sequence $\{\sigma^{k_j}\}$ s.t. σ^{k_j} is a Nash equilibrium of the ϵ^{k_j} -perturbed game $G_{\epsilon^k}^m$ and $\sigma^{k_j} \rightarrow \sigma^*$. Hence σ^* is a perfect equilibrium. This completes the proof. \square

Given a family of games \mathcal{G}_Λ parameterized by Λ , Proposition 12 stated that the Nash equilibrium correspondence on \mathcal{G}_Λ has a closed graph. This is quite useful for finding Nash equilibria via approximations of a game. Define the *perfect equilibrium correspondence* $\text{PE} : \Lambda \rightrightarrows \times_{i \in \mathcal{I}} \Delta(S_i)$ by

$$\text{PE}(\lambda) = \{\sigma \in \times_{i \in \mathcal{I}} \Delta(S_i) \mid \sigma \text{ is a perfect equilibrium of } G^m(\lambda)\}.$$

Unlike the Nash equilibrium correspondence, the perfect equilibrium correspondence does not have a closed graph, as the following counterexample shows.

Example 15. Consider the following family of games $G(n)$ parameterized by \mathbb{N} :

| | | |
|-------|--|--|
| | L_2 | R_2 |
| U_1 | 1 , 1 | 0, 0 |
| D_1 | 0, 0 | 1/n , 1/n |

Best responses are highlighted in blue (Player 1) and red (Player 2). For each $n \in \mathbb{N}$, we claim that (D_1, R_2) is a perfect equilibrium. To see this, define

$$\begin{aligned} \sigma_1^\epsilon(U_1) &= \epsilon, & \sigma_1^\epsilon(D_1) &= 1 - \epsilon, \\ \sigma_2^\epsilon(L_2) &= \epsilon, & \sigma_2^\epsilon(R_2) &= 1 - \epsilon. \end{aligned}$$

For any $\epsilon < \frac{1}{n+1}$, we have

$$u_1(D_1, \sigma_2^\epsilon) = \frac{1 - \epsilon}{n} > \epsilon = u_1(U_1, \sigma_2^\epsilon)$$

and so D_1 is a best response for Player 1. A symmetric calculation shows R_1 is a best response for Player 2 to σ_1^ϵ . Taking $\epsilon \rightarrow 0$ gives $(\sigma_1^\epsilon, \sigma_2^\epsilon) \rightarrow (D_1, R_2)$, and thus (D_1, R_2) is a perfect equilibrium of $G(n)$.

Now consider the limit game $G(\infty) = \lim_{n \rightarrow \infty} G(n)$. This game is:

| | | |
|-------|--|--|
| | L_2 | R_2 |
| U_1 | 1 , 1 | 0, 0 |
| D_1 | 0, 0 | 0 , 0 |

Now (D_1, R_2) is a weakly dominated Nash equilibrium, so by Proposition 15(i), (D_1, R_2) cannot be perfect. Hence we have found a sequence of games where the limit of the corresponding perfect equilibria is not a perfect equilibrium of the limit game.

2.9 Proper equilibrium

Perfect equilibrium is not quite as perfect as Selten had so confidently hoped. Consider the following example.

Example 16 (Myerson, 1978). Consider the two player game with payoff matrix

| | a_2 | b_2 | c_2 |
|-------|--------|--------|--------|
| a_1 | 1, 1 | 0, 0 | -9, -9 |
| b_1 | 0, 0 | 0, 0 | -7, -7 |
| c_1 | -9, -9 | -7, -7 | -7, -7 |

Again, best responses are highlighted in blue (Player 1) and red (Player 2). Clearly, there are three pure strategy Nash equilibria: (a_1, a_2) , (b_1, b_2) and (c_1, c_2) . It can be shown there are no non-degenerate mixed strategy Nash equilibria.

(c_1, c_2) is weakly dominated and thus not a perfect equilibrium. However, both (a_1, a_2) and (b_1, b_2) are trembling hand perfect. To see that (b_1, b_2) is a perfect equilibrium, define

$$\begin{aligned}\sigma_1^\epsilon(a_1) &= \epsilon, & \sigma_1^\epsilon(b_1) &= 1 - 2\epsilon, & \sigma_1^\epsilon(c_1) &= \epsilon, \\ \sigma_2^\epsilon(a_2) &= \epsilon, & \sigma_2^\epsilon(b_2) &= 1 - 2\epsilon, & \sigma_2^\epsilon(c_2) &= \epsilon,\end{aligned}$$

Then for any $\epsilon < \frac{1}{2}$,

$$u_1(b_1, \sigma_2^\epsilon) = -7\epsilon = u_1(c_1, \sigma_2^\epsilon) > -9\epsilon = u_1(a_1, \sigma_2^\epsilon),$$

and hence b_1 is a best response to σ_2^ϵ . A symmetric calculation for Player 2 shows that b_2 is a best response to any σ_1^ϵ . Taking $\epsilon \rightarrow 0$, we have $(\sigma_1^\epsilon, \sigma_2^\epsilon) \rightarrow (b_1, b_2)$. Hence (b_1, b_2) is a perfect equilibrium.

Likewise, defining

$$\begin{aligned}\sigma_1^\epsilon(a_1) &= 1 - 2\epsilon, & \sigma_1^\epsilon(b_1) &= \epsilon, & \sigma_1^\epsilon(c_1) &= \epsilon, \\ \sigma_2^\epsilon(a_2) &= 1 - 2\epsilon, & \sigma_2^\epsilon(b_2) &= \epsilon, & \sigma_2^\epsilon(c_2) &= \epsilon,\end{aligned}$$

we have, for any $\epsilon < \frac{1}{4}$, that

$$\begin{aligned}u_1(a_1, \sigma_2^\epsilon) &= 1 - 11\epsilon \\ &> -7\epsilon = u_1(b_1, \sigma_2^\epsilon) \\ &> -9 + 4\epsilon = u_1(c_1, \sigma_2^\epsilon),\end{aligned}$$

so a_1 is a best response to σ_2^ϵ . Again, a similar calculation shows a_2 is a best response to σ_1^ϵ for all $\epsilon < \frac{1}{4}$. As $\epsilon \rightarrow 0$, $(\sigma_1^\epsilon, \sigma_2^\epsilon) \rightarrow (a_1, a_2)$, and thus (a_1, a_2) is a perfect equilibrium.

Yet there is a clear sense in which the equilibrium (b_1, b_2) seems less compelling than (a_1, a_2) . In particular, we had to assume that costly trembles (c_i) happened with the same probability as less costly trembles (a_i).

Myerson thus (1978) proposes a further refinement of perfect equilibrium: *proper equilibrium*. Proper equilibrium imposes a restriction on the relative likelihood of trembles. Namely, given a strategy profile σ , consider ordering a player i 's pure strategies $s_i \in S_i$ by their payoff to i . In an ϵ -proper equilibrium, we assume i plays the second-highest payoff strategy with at most ϵ times the probability of the highest payoff strategy, i plays the third-highest payoff strategy with at most ϵ times the probability of the second-highest payoff strategy, and so forth. Thus proper equilibrium captures the notion that a player will tremble to play a higher-payoff strategy relatively more often than they tremble to play a lower-payoff strategy.

Definition 28 (Proper equilibrium). Consider a finite mixed strategy normal form game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ with n players.

- (a) *ϵ -proper equilibrium*. A totally mixed strategy profile σ is an *ϵ -proper equilibrium* if, for all $i \in \mathcal{I}$, if

$$u_i(s_i, \sigma_{-i}) < u_i(s'_i, \sigma_{-i})$$

then

$$\sigma_i(s_i) \leq \epsilon \cdot \sigma_i(s'_i).$$

That is, a totally mixed strategy profile σ is an ϵ -proper equilibrium if every player places greater probability weight on their better responses than their worse responses, by a factor of $1/\epsilon$.

- (b) *Proper equilibrium*. A strategy profile σ^* is a *proper equilibrium* if

- (i) there exists a sequence $\{\epsilon^k\}$ with $\epsilon^k \geq 0$ and $\epsilon^k \rightarrow 0$, and
- (ii) there exists a sequence of profiles $\{\sigma^k\}$ such that each σ^k is an ϵ^k -proper equilibrium and $\sigma^k \rightarrow \sigma^*$.

Proposition 17. Any proper equilibrium σ of a finite mixed strategy game G^m is a perfect equilibrium.

Proof. Note that σ is a proper equilibrium if it is the limit of ϵ -proper equilibrium. Now, any ϵ -proper equilibrium is a Nash equilibrium of the $(\epsilon, \dots, \epsilon)$ -perturbed game G_ϵ^m . Thus σ is the limit of ϵ -perfect equilibria as $\epsilon \rightarrow 0$, and so σ is perfect. \square

Theorem 11. Every finite player mixed strategy normal form game $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$ such that S_i is finite for each $i \in \mathcal{I}$ has a proper equilibrium.

Proof. The proof is almost identical to Theorem 10. \square

Example 16 (continued). Returning to Myerson's (1978) example, we show that (a_1, a_2) is the only proper equilibrium. By Theorem 11, a proper equilibrium must exist. Suppose $\epsilon \in (0, 1)$, and let (σ_1, σ_2) be an ϵ -proper equilibrium. Since b_1 dominates c_1 and σ_2 is totally mixed, $u_1(b_1, \sigma_2) > u_1(c_1, \sigma_2)$, and thus $\sigma_1(c_1) \leq \epsilon \cdot \sigma_1(b_1)$. But then $u_2(c_2, \sigma_1) < u_2(a_2, \sigma_1)$, and so $\sigma_2(c_2) \leq \epsilon \cdot \sigma_2(a_2)$. It follows that $u_1(b_1, \sigma_2) < u_1(a_1, \sigma_2)$ and so

$\sigma_1(b_1) \leq \epsilon \cdot \sigma_1(a_1) \leq \epsilon$, and $\sigma_1(c_1) \leq \epsilon \cdot \sigma_1(b_1) \leq \epsilon^2$. By a similar argument, $\sigma_2(b_2) \leq \epsilon \cdot \sigma_2(a_2) \leq \epsilon$ and $\sigma_2(c_2) \leq \epsilon \sigma_2(b_2) \leq \epsilon^2$. These imply that $\sigma_1(a_1) \geq 1 - \epsilon - \epsilon^2$ and $\sigma_2(a_2) \geq 1 - \epsilon - \epsilon^2$. Taking $\epsilon \rightarrow 0$, we have that $(\sigma_1, \sigma_2) \rightarrow (a_1, a_2)$, so the only proper equilibrium is (a_1, a_2) .

2.10 Focal points

In coordination games, Schelling (1960) introduced a particular refinement of Nash equilibrium relying on information external to the description of the game – in certain games, there may be a particular strategy that can act as a default in the absence of communication between players. Such a strategy is a *focal point* or *Schelling point*.

Example 17 (Schelling, 1960). Consider a coordination game motivated as follows: two strangers are to meet in New York, and they cannot communicate with each other. If they choose the same location, they receive a positive payoff, and zero otherwise. Experimentally, Schelling (1960) finds that the most common strategy is to choose Grand Central Terminal, a central New York landmark, rather than some less obvious location. Yet there is nothing in the formal structure of the game that favours this particular location over any other.

Example 18. Consider the following scenario.¹³ Player 1 and Player 2 are in two separate rooms. In front of Player 1 is a display that shows a paragraph of text. Player 1 can only message Player 2 by sending a string of real numbers. If Player 2 correctly inputs letters into their computer to replicate the paragraph shown to Player 1, then both players win \$1 million, and if Player 2 makes a mistake or otherwise fails to replicate the paragraph within a reasonable time limit, neither player receives anything. Assume getting the order of the letters correct is all that matters, not punctuation or spaces. Also assume players were not able to confer to agree a strategy before learning the task.

When we formalize this as a game, any strategies involving Player 1 choosing a unique number for each letter, generating an accurate sequence corresponding to the order of the letters in the paragraph using these numbers, and Player 2 decoding the sequence perfectly constitute a Nash equilibrium (as do all sorts of other strategies). However, we assumed players could not confer, and it is very implausible that if Player 1 chooses arbitrary real numbers, then Player 2 can perfectly guess which number corresponds to which letter of the alphabet and decode the message.¹⁴ However, there is a focal point in this setting – suppose Player 1 chooses 1 for *A*, 2 for *B* and so on. Then it is plausible that Player 2 can infer how Player 1 has encoded the paragraph in her message.

Like many of Schelling’s insights, focal points are not really something that can be easily formalized:

¹³This example came up in a characteristically heated discussion with Filip Tokarski.

¹⁴Of course, if we allowed players to communicate before the task, this would be possible. The scenario would then be closer to the situations that often arise in cryptography, where encryption/decryption keys are known to the appropriate parties.

What exactly is the formal definition of a focal point or a focal-point equilibrium? There is no formal answer. U.S. Supreme Court Justice Potter Stewart famously said about pornography, “I know it when I see it.” That pretty much captures what can be said about focal points.
– David M. Kreps (2023).

2.11 Payoff dominance and risk dominance

In games with coordination, equilibrium selection concepts refine Nash equilibrium – when there are multiple equilibria, determining which equilibrium will be played requires specifying additional criteria. The classic Harsanyi-Selten approach involves specifying desirable properties of equilibria to determine which agents may arrive at.

Definition 29 (Harsanyi & Selten, 1988). Consider a finite player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$.

- (a) *Payoff dominance*. A Nash equilibrium s^* is said to *payoff dominate* a Nash equilibrium s' if

$$u_i(s^*) > u_i(s') \quad \text{for all } i \in \mathcal{I}.$$

A Nash equilibrium s^* is said to be a *payoff dominant equilibrium* of G if it payoff dominates all other Nash equilibrium s' of G .

- (b) *Bilateral risk dominance*. In a two-player game, the *Nash product* at a Nash equilibrium s^* relative to an alternative equilibrium s' is given by the product of deviation losses of both players, that is

$$(u_1(s_1^*, s_2^*) - u_1(s_1', s_2^*))(u_2(s_1^*, s_2^*) - u_2(s_1^*, s_2')).$$

A Nash equilibrium s^* is said to (bilaterally) *risk dominate* a Nash equilibrium s' if s^* has strictly greater Nash product than s' , that is, if

$$(u_1(s_1^*, s_2^*) - u_1(s_1', s_2^*))(u_2(s_1^*, s_2^*) - u_2(s_1^*, s_2')) > (u_1(s_1', s_2') - u_1(s_1^*, s_2'))(u_2(s_1', s_2') - u_2(s_1', s_2')).$$

If a Nash equilibrium s^* of G risk dominates all other Nash equilibria of G then it is called a *risk dominant equilibrium*.

Example 19. Consider the coordination game:

| | | |
|-------|-------|-------|
| | a_2 | b_2 |
| a_1 | 6, 6 | 3, 4 |
| b_1 | 4, 3 | 4, 4 |

There are two pure strategy Nash equilibria, (a_1, a_2) and (b_1, b_2) [there is also a mixed strategy Nash equilibrium]. Since $u_i(a_1, a_2) = 6 > 3 = u_i(b_1, b_2)$ for $i = 1, 2$, we have that (a_1, a_2) payoff dominates (b_1, b_2) .

However, we have Nash product relation (a_1, a_2) vs (b_1, b_2) of

$$\begin{aligned} (u_1(a_1, a_2) - u_1(b_1, a_2))(u_2(a_1, a_2) - u_2(a_1, b_2)) &= (6 - 4)(6 - 4) = 4 \\ &> 1 = (4 - 3)(4 - 3) \\ &= (u_1(b_1, b_2) - u_1(a_1, b_2))(u_2(b_1, b_2) - u_2(b_1, a_2)). \end{aligned}$$

Hence (a_1, a_2) risk dominates (b_1, b_2) .

2.12 Correlated equilibrium

Nash equilibrium can be envisaged as a kind of self-enforcing agreement, in the sense that in a Nash equilibrium, no player can do any better by unilaterally deviating and doing something else. This is not to say that a Nash equilibrium is necessarily a credible agreement – some Nash equilibria only exist due to threats that a player would be better off reneging on.¹⁵ Yet Nash equilibrium is somewhat less general than the notion of a self-enforcing agreement, as the following examples show. In some situations, Nash equilibrium is the inappropriate solution concept because it is inconsistent with certain kinds of communication.

Example 20 (Bach-or-Stravinsky). Suppose two friends are, for reasons unknown to us, classical music aficionados. There are two concerts taking place at the same time. One is a Bach recital (b) whereas the other is an orchestral concert centred around Stravinsky's *Feu d'artifice* (f). One of the friends much prefers Bach to Stravinsky, and the other much prefers Stravinsky to Bach, but neither wants to go to a concert alone.¹⁶ They cannot communicate, and only receive a positive payoff if they choose the same event. Assigning the woman as Player 1, the payoff matrix is

| | | |
|-------|-------|-------|
| | b_2 | f_2 |
| b_1 | 3, 2 | 0, 0 |
| f_1 | 0, 0 | 2, 3 |

The pure strategy Nash equilibria are (b_1, b_2) and (f_1, f_2) . There is also a mixed Nash equilibrium. Suppose Player 1 plays b_1 with probability p . Then Player 2 is indifferent over b_2 and f_2 iff $2p = 3(1 - p) \Rightarrow p = \frac{3}{5}$. Likewise, if Player 2 plays b_2 with probability q , then Player 1 is indifferent over b_1 and f_1 iff $3q = 2(1 - q) \Rightarrow q = \frac{2}{5}$. Hence we have mixed strategy Nash equilibrium $((\frac{3}{5}, \frac{2}{5}), (\frac{2}{5}, \frac{3}{5}))$.

Now to see that Nash equilibrium can be interpreted as a self-enforcing agreement, we might imagine that players jointly agree to play (b_1, b_2) . The profile (b_1, b_2) is self-enforcing, in the sense that neither player, if rational, gains from deviating. The same argument of course applies to (f_1, f_2) and to the mixed strategy Nash equilibrium.

Yet there are self enforcing agreements that are not Nash equilibria. Suppose the couple agree to flip a fair coin. If it lands on Heads, they will both attend the ballet, i.e. play the Nash equilibrium (b_1, b_2) , and if it lands on Tails, they both attend the prize fight, i.e. play the Nash equilibrium (f_1, f_2) . This achieves an expected payoff of $\frac{5}{2}$ for

¹⁵In the literature, you sometimes see it argued that Nash equilibria are *not* self-enforcing, in the sense that they are not *strategically stable*, e.g. Kohlberg & Mertens (1986) say this. However, we are not using “self-enforcing” in their sense here (which is one of credibility – for example, Nash equilibria can involve threats off-path that a player would prefer not to follow through with).

¹⁶This game is also commonly known as *Battle of the Sexes*. The now rather dated motivation for this game, introduced by Luce and Raiffa (1957), imagines a dating couple (man and woman) hoping to meet for the evening. They have a choice between a prize fight F (preferred by the man) and a ballet B (preferred by the woman). This plays heavily on gendered stereotypes but the motivation for the payoffs makes slightly more sense. You might still enjoy going to the concert alone if the music is good, whereas you can't enjoy a one-person date – or maybe you can, dear reader, who am I to judge your preferences.

each player, which is not a Nash equilibrium payoff – it is greater than the expected payoff in the mixed strategy Nash equilibrium, in which (b_1, f_2) and (f_1, b_2) are played with positive probability. Furthermore, this is a self-enforcing agreement. If the coin lands on Heads, neither player i can profitably deviate by playing f_i . We see then that Nash equilibrium might not allow sufficient room to communicate.

Example 21. The coin flip story is applicable in the battle-of-the-sexes, but it is a primitive way to communicate prior to play. More generally, we might imagine that there is a mediator who can perform randomizations and who tells the players what to play depending on the outcome.

Consider the following game:

| | | |
|-----|------|------|
| | L | R |
| U | 5, 1 | 0, 0 |
| D | 4, 4 | 1, 5 |

There are three Nash equilibria: pure strategy Nash equilibria (U, L) and (D, R) and a mixed strategy Nash equilibrium $((\frac{1}{2}, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2}))$.

Suppose players find a mediator who chooses $\omega \in \{1, 2, 3\}$ uniformly at random, so each with probability $\frac{1}{3}$. Players do not observe ω directly. The mediator proposes the following:

- (i) if $\omega = 1$, tell Row to play U and Column to play L ;
- (ii) if $\omega = 2$, tell Row to play D and Column to play L ;
- (iii) if $\omega = 3$, tell Row to play D and Column to play R .

It is a perfect Bayesian equilibrium for players to follow the mediator's advice here.¹⁷

The notion of *correlated equilibrium*, first introduced by Aumann (1974), fully captures the notion of a self-enforcing agreement.

Definition 30 (Correlated equilibrium).

- (a) *Correlating mechanism.* A *correlating mechanism* is a tuple $(\Omega, \{\mathcal{H}_i\}_{i \in \mathcal{I}}, p)$ where
 - (i) Ω is a finite set of states of the world;
 - (ii) for each player i in finite player set \mathcal{I} , \mathcal{P}_i is a partition of Ω , with the function $h_i : \Omega \rightarrow \mathcal{P}_i$ assigning to each ω the element $h_i(\omega) = (P_i \in \mathcal{P}_i \mid \omega \in P_i)$, and
 - (iii) p is a probability distribution on Ω .¹⁸

¹⁷See section 6.4 for definition and discussion of perfect Bayesian equilibria. If Row hears U then Row believes Column will play L . Row's best response to this belief is to play U . If Row hears D , then Row believes Column will play L with probability $\frac{1}{2}$ and R with probability $\frac{1}{2}$, and so D is a best response for Row. Column's strategy can be checked similarly.

¹⁸More generally (i.e. with a continuum of states), we can define a correlating mechanism as a tuple $((\Omega, \mathcal{F}, p), \{\mathcal{P}_i\}_{i \in \mathcal{I}})$, where (Ω, \mathcal{F}, p) is a probability space, with \mathcal{F} being a σ -algebra, and $\{\mathcal{P}_i\}$ is defined as before.

- (b) *Correlated strategy*. A *correlated strategy* for i is a function $\sigma_i : \Omega \rightarrow S_i$ that is measurable with respect to information partition \mathcal{P}_i , i.e. if $h_i(\omega) = h_i(\omega')$ then $\sigma_i(\omega) = \sigma_i(\omega')$.¹⁹²⁰
- (c) *Correlated equilibrium*. A strategy profile $\sigma^* = (\sigma_1^*, \dots, \sigma_n^*)$ is an (objective) *correlated equilibrium* relative to the correlating mechanism $(\Omega, \{\mathcal{P}_i\}, p)$ if, for every i and every correlated strategy σ_i ,

$$\sum_{\omega \in \Omega} u_i(\sigma_i^*(\omega), \sigma_{-i}^*(\omega))p(\omega) \geq \sum_{\omega \in \Omega} u_i(\sigma_i(\omega), \sigma_{-i}^*(\omega))p(\omega).$$

Unpacking the definition, correlated equilibrium requires that σ_i^* maximizes i 's *ex ante* payoff. That is, the strategy can be considered a contingent plan to be implemented after learning the partition element ω . Note that this is equivalent to σ_i maximizing i 's *interim* payoff for each $P_i \in \mathcal{P}_i$ that is reached with positive probability: that is, for all $i \in \mathcal{I}$, for all $\omega \in \Omega$, and for all $s'_i \in S_i$,

$$\sum_{\omega' \in h_i(\omega)} u_i(\sigma_i^*(\omega), \sigma_{-i}^*(\omega'))p(\omega' | h_i(\omega)) \geq \sum_{\omega' \in h_i(\omega)} u_i(s'_i, \sigma_{-i}^*(\omega'))p(\omega' | h_i(\omega)),$$

where $p(\omega' | h_i(\omega))$ is the conditional probability of ω' given that the true state lies in $h_i(\omega)$. By Bayes' rule, this is

$$p(\omega' | h_i(\omega)) = \frac{\mathbb{P}\{h_i(\omega) | \omega'\}p(\omega')}{\sum_{\omega'' \in h_i(\omega)} \mathbb{P}\{h_i(\omega) | \omega''\}p(\omega'')} = \frac{p(\omega')}{p(h_i(\omega))}.$$

The definition is unwieldy, in that naively searching for all the correlated equilibria would require checking many arbitrary correlating mechanisms. Fortunately, the problem can be reduced to checking only a certain kind of correlating mechanism – the direct mechanisms:

Definition 31 (Direct mechanism). A *direct mechanism* is a correlating mechanism $(\Omega, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, p)$ such that $\Omega = S$, $h_i(s) = \{s' \in S : s'_i = s_i\}$, and where p is some probability distribution over pure strategy profiles.

In a direct mechanism, the state space is the set of pure strategy profiles.

For any correlated equilibrium relative to some correlating mechanism, there is an outcome-equivalent correlated equilibrium relative to some direct mechanism. This result is known as the *revelation principle*.²¹

¹⁹Note this differs from our definition of a correlated strategy *profile*, which is simply a point in $\Delta(S)$.

²⁰Again, the technical definition of measurable here is that $\sigma_i^{-1}(B)$ lies in the σ -algebra generated by \mathcal{P}_i for all Borel sets $B \subseteq S_i$. But this is unnecessary formalism in our setting – the condition boils down to $\sigma_i(\omega) = \sigma_i(\omega')$ whenever $h_i(\omega) = h_i(\omega')$.

²¹A version of this result is also widely applied in mechanism design.

Theorem 12 (Revelation principle). *Suppose σ^* is a correlated equilibrium relative to correlating mechanism $(\Omega, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, p)$. Define $q(s) = \mathbb{P}\{\sigma^*(\omega) = s\}$. Then the strategy profile $\tilde{\sigma}$ with $\tilde{\sigma}_i(s) = s_i$ for all $s_i \in S_i$ and all $i \in \mathcal{I}$ is a correlated equilibrium relative to the direct mechanism $(S, \{\tilde{\mathcal{P}}_i\}_{i \in \mathcal{I}}, q)$.*

Proof. Suppose s_i is recommended to i with positive probability, i.e. $p(s_i, s_{-i}) > 0$ for some s_{-i} . We require that under the direct mechanism, i does not benefit from choosing some $s'_i \neq s_i$ when s_i is suggested. If s_i is recommended, then i 's expected payoff from playing s'_i is

$$\sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) q(s_{-i} \mid s_i).$$

If there is only one information set $P_i \in \mathcal{P}_i$ s.t. $\sigma_i^*(P_i) = s_i$, then conditioning on s_i is equivalent to conditioning on H_i so the proposition holds trivially. More generally, substituting for q gives expected payoff to playing s'_i of

$$\frac{1}{\mathbb{P}\{\sigma_i^*(\omega) = s_i\}} \sum_{\omega \mid \sigma_i^*(\omega) = s_i} u_i(s'_i, \sigma_{-i}^*(\omega)) p(\omega).$$

Rearranging,

$$\frac{1}{\mathbb{P}\{\sigma_i^*(\omega) = s_i\}} \sum_{H_i \mid \sigma_i^*(P_i) = s_i} \mathbb{P}\{P_i\} \left[\sum_{\omega \in P_i} u_i(s'_i, \sigma_{-i}^*(\omega)) p(\omega \mid P_i) \right].$$

Given $(\Omega, \{\mathcal{P}_i\}, p, \sigma^*)$ is a correlated equilibrium, we have that each bracketed term for which $\mathbb{P}\{H_i\} > 0$ is maximized at $\sigma_i(P_i) = s_i$. Hence s_i is optimal given recommendation s_i , under the direct mechanism. \square

Hence to find all correlated equilibria, we need only consider the class of direct mechanisms. A direct mechanism is effectively the mechanism we considered in the mediator story we discussed in Example 21, where an impartial mediator recommends a strategy to each player.

The probability distribution induced over strategy profiles is all that matters for a correlated equilibrium. We call the probability distribution q over strategy profiles s a *correlated equilibrium distribution* if it is the distribution generated by some correlated equilibrium.

Proposition 18. *The distribution $q \in \Delta(S)$ is a correlated equilibrium distribution iff for all $i \in \mathcal{I}$, for all $s_i \in S_i$ with $q(s_i) > 0$, and for all $s'_i \in S_i$,*

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) q(s_{-i} \mid s_i) \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) q(s_{-i} \mid s_i).$$

Proof. Suppose q satisfies the inequality. Then the profile σ^* with $\sigma_i^*(s) = s_i$ is a correlated equilibrium given direct mechanism $(S, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, q)$, for the inequality states precisely that s_i is the best response for i on being recommended s_i .

Conversely, suppose q is a correlated equilibrium distribution. Then q corresponds to some correlated equilibrium σ^* relative to some direct mechanism $(S, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, q)$. Hence for all i and all recommendations s_i , the inequality must hold by optimality of the recommendation s_i . \square

We have the corollary:

Corollary 5 (Aumann, 1987). *The distribution $q \in \Delta(S)$ is a correlated equilibrium distribution iff for all $i \in \mathcal{I}$, for all $s_i \in S_i$ with $q(s_i) > 0$, and for all $s'_i \in S_i$,*

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i})q(s_i, s_{-i}) \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i})q(s_i, s_{-i}).$$

Proof. Let $q(s_i)$ denote the marginal probability of s_i under q . Then we have $q(s_{-i} | s_i) = \frac{q(s_i, s_{-i})}{q(s_i)}$. Substituting into the inequality in Proposition 18 and multiplying through by $q(s_i)$ yields the result. \square

As we might expect from the self-enforcing agreement motivation, Nash equilibrium is a refinement of correlated equilibrium:

Proposition 19. *Every Nash equilibrium is a correlated equilibrium.*

Proof. We provide a proof for finite player games with finite strategy sets.

We require that for all i and all s_i with $q(s_i) > 0$, that

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i})q(s_{-i} | s_i) \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i})q(s_{-i} | s_i)$$

for all $s'_i \in S_i$, under the distribution q induced by the Nash equilibrium. In case of a pure strategy Nash equilibrium s^* , we have

$$q(s_{-i} | s_i^*) = \begin{cases} 1 & \text{if } s_{-i} = s_{-i}^*, \\ 0 & \text{otherwise.} \end{cases}$$

Hence the inequality reduces to

$$u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*),$$

which is precisely the definition of pure strategy Nash equilibrium.

In case of a mixed strategy Nash equilibrium, since players mix independently, we have that $q(s_{-i} | s_i^*) = \sigma_{-i}^*(s_{-i})$ for any s_i^* in the support of σ_i^* . Hence, for all i , s_i^* in the support of σ_i^* , and $s_i \in S_i$, we have that

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i^*, s_{-i})\sigma_{-i}^*(s_{-i}) \geq \sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i})\sigma_{-i}^*(s_{-i}),$$

which is precisely the definition of mixed strategy Nash equilibrium. \square

Corollary 6. *In any finite game, there exists a correlated equilibrium.*

Proof. Follows immediately from the Nash existence theorem and the fact that Nash equilibria are correlated equilibria. See Hart and Schmeidler (1989) for a direct existence proof. \square

Proposition 20. *The sets of correlated equilibrium distributions and payoffs are convex.*

Proof. For any distribution $p \in \Delta(S)$, define $S_i^p := \{s_i \in S_i : p(s_i) > 0\}$. Note that if $p \in \Delta(S)$, then for any i and any $s_i \in S_i^p$, we have that $p(s_{-i} | s_i) = \frac{p(s_i, s_{-i})}{p(s_i)}$, where $p(s_i)$ is the marginal probability of s_i induced by the distribution p .

Suppose q and q' are correlated equilibrium distributions. For any pair $\lambda \in (0, 1)$, let $q''(s_i, s_{-i}) = \lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})$, so

$$q''(s_{-i} | s_i) = \frac{q''(s_i, s_{-i})}{q''(s_i)} = \frac{\lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})}{q''(s_i)}.$$

By Proposition 18, it is sufficient to show that for all $i \in \mathcal{I}$, for all $s'_i \in S_i$ and for all $s_i \in S_i^{q''} = S_i^q \cup S_i^{q'}$,

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) q''(s_{-i} | s_i) \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) q''(s_{-i} | s_i)$$

iff

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) \frac{\lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})}{q''(s_i)} \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) \frac{\lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})}{q''(s_i)}.$$

Multiplying both sides through by $q''(s_i)$, we have

$$\sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) [\lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})] \geq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) [\lambda q(s_i, s_{-i}) + (1 - \lambda)q'(s_i, s_{-i})],$$

which is equivalent to

$$\begin{aligned} \lambda \sum_{s_{-i}} u_i(s_i, s_{-i}) q(s_i, s_{-i}) + (1 - \lambda) \sum_{s_{-i}} u_i(s_i, s_{-i}) q'(s_i, s_{-i}) &\geq \lambda \sum_{s_{-i}} u_i(s_i, s_{-i}) q(s_i, s_{-i}) \\ &\quad + (1 - \lambda) \sum_{s_{-i}} u_i(s'_i, s_{-i}) q'(s_i, s_{-i}). \end{aligned}$$

By Corollary 5, this holds iff

$$\begin{aligned} \lambda \sum_{s_{-i}} u_i(s_i, s_{-i}) q(s_{-i} | s_i) + (1 - \lambda) \sum_{s_{-i}} u_i(s_i, s_{-i}) q'(s_{-i} | s_i) &\geq \lambda \sum_{s_{-i}} u_i(s_i, s_{-i}) q(s_{-i} | s_i) \\ &\quad + (1 - \lambda) \sum_{s_{-i}} u_i(s'_i, s_{-i}) q'(s_{-i} | s_i). \end{aligned}$$

This inequality holds, since q and q' are correlated equilibrium distributions. Hence q'' is a correlated equilibrium distribution. Since $\lambda \in (0, 1)$ was arbitrary, it follows that any convex combination of correlated equilibrium distributions is a correlated equilibrium distribution, and thus the set of correlated equilibrium distributions is convex. The statement wrt payoffs is a straightforward corollary. \square

From the inequality in Corollary 5 and the fact that the set of probability vectors is a simplex, we see that the set of correlated equilibrium distributions in any finite game is a *convex polytope*. Moreover, the set of correlated equilibrium payoffs in any finite game is a convex polytope.

Definition 32 (Public correlating mechanism). A correlating mechanism $(\Omega, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, p)$ is a *public correlating mechanism* if for all $i \in \mathcal{I}$, $\mathcal{P}_i = \mathcal{P}$ for some partition \mathcal{P} . That is, every player has the same partition with respect to states of the world. Since if \mathcal{P} contains some non-singleton set P , we can replace the subset $P \subseteq \Omega$ with some state ω' such that $p(\omega') = \sum_{\omega \in P} p(\omega)$. Hence, wlog, we can represent a public correlating mechanism by (Ω, p) .

We call σ^* a *public correlated equilibrium* if it is a correlated equilibrium relative to a public correlating mechanism.

Proposition 21. *The set of public correlated equilibrium payoffs is the convex hull of the set of Nash equilibrium payoffs.*

Since Proposition 20 shows the set of correlated equilibrium payoffs is convex, it is often convenient to look at payoffs graphically.

Definition 33 (Feasibility and individual rationality).

- (a) *Feasible payoffs.*²² In a finite game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a payoff profile $v = (v_1, \dots, v_n)$ is *feasible* if there is some probability distribution $p \in \Delta(S)$ such that

$$v_i = \sum_{s \in S} u_i(s) p(s) \quad \text{for all } i \in \mathcal{I}.$$

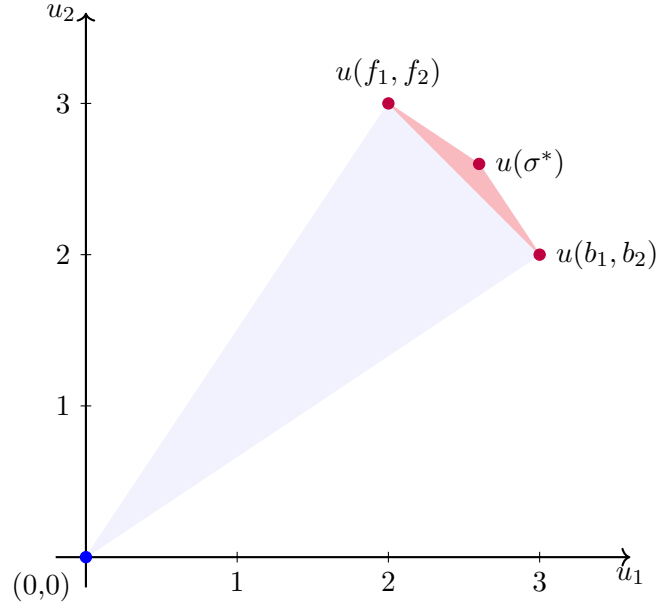
- (b) *Individually rational payoffs.* In a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a payoff profile $v = (v_1, \dots, v_n)$ is *individually rational* if for each $i \in \mathcal{I}$,

$$v_i \geq \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} \max_{\sigma_i \in \Delta(S_i)} u_i(\sigma_i, \sigma_{-i}) = v_i.$$

A payoff profile v is *strictly individually rational* if this holds with strict inequality for all i .

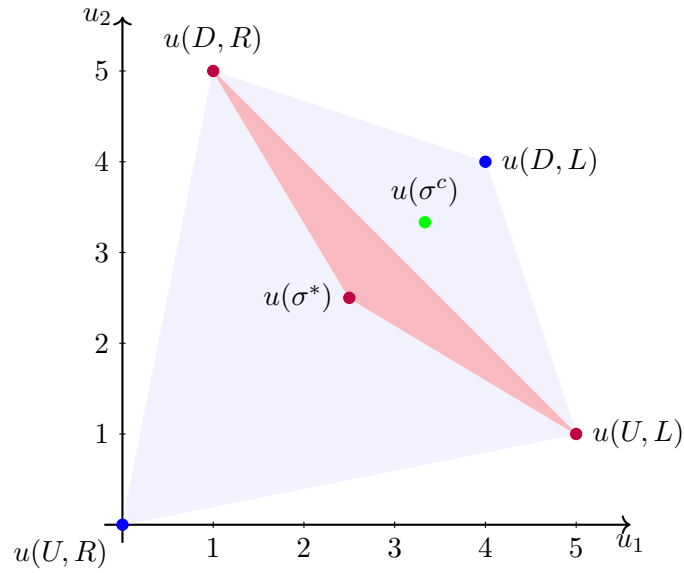
Example 20 (continued). We show, for the battle-of-the-sexes game, the set of feasible payoffs (blue region), the set of Nash equilibrium payoffs (red points), and correlated equilibrium payoffs (red region).

²²More generally, a payoff profile v is *feasible* if there is some probability distribution $p \in \Delta(S)$ such that $v_i = \int_S u_i dp$.



σ^* denotes the mixed strategy Nash equilibrium profile $\sigma^* = ((\frac{1}{2}, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2}))$. The convex hull of Nash equilibria is denoted in red. Any point in the convex hull of Nash equilibria is the payoff profile of some public correlated equilibrium.

Example 21 (continued). When the correlating mechanism is not public, it is possible to achieve correlated equilibrium payoffs that lie outside the convex hull of the set of Nash equilibrium payoffs.



The set of feasible payoffs is shaded blue and the convex hull of the set of Nash equilibrium payoffs is shaded red (σ^* is the mixed strategy Nash equilibrium $((\frac{1}{2}, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2}))$).

The correlated equilibrium we previously considered, where (U, L) , (D, L) and (D, R) were each played with probability $\frac{1}{3}$, is denoted σ^c . We see that $u(\sigma^c)$ lies outside the convex hull.

We motivated the correlated equilibrium σ^c by means of a mediator telling to each player privately which strategy to play. This is not a public correlating mechanism, since players have private information – if told to play D , for example, Row cannot infer whether the mediator told Column to play L or R . Were the mediator to announce which strategy she recommends to each player publicly, then σ^c could not be sustained. For example, if the mediator publicly announced that players should play (D, L) , then Row would have a profitable deviation U and Column would have a profitable deviation R .

In the (objective) correlating mechanisms $(\Omega, \{\mathcal{P}_i\}_{i \in \mathcal{I}}, p)$ previously introduced, we assume a *common prior* p over set of states Ω . Equivalently, the players share the same probability distribution over equilibrium play (by Theorem 12). Subjective correlated equilibrium weakens this assumption.

Definition 34 (Subjective correlated equilibrium).

- (a) *Subjective correlating mechanism.* A *subjective correlating mechanism* is a tuple $(\Omega, \{\mathcal{P}_i, p_i\}_{i \in \mathcal{I}})$ where Ω and \mathcal{P}_i are defined as in Definition 30 and for each player $i \in \mathcal{I}$, p_i is a probability distribution over Ω .
- (b) *Subjective correlated equilibrium.* A profile $\sigma^* = (\sigma_1^*, \dots, \sigma_n^*)$ of functions $\sigma_i^* : \Omega \rightarrow S_i$ is a *subjective correlated equilibrium* relative to subjective correlating mechanism $(\Omega, \{\mathcal{H}_i, p_i\}_{i \in \mathcal{I}})$ if for every i and every correlated strategy σ_i ,

$$\sum_{\omega \in \Omega} u_i(\sigma_i^*(\omega), \sigma_{-i}^*(\omega)) p_i(\omega) \geq \sum_{\omega \in \Omega} u_i(\sigma_i(\omega), \sigma_{-i}^*(\omega)) p_i(\omega).$$

Example 21 (continued). Subjective correlated equilibrium does not require that players are correct about the objective probability distribution over states (or about play). We can obtain (D, L) (with payoff $(4, 4)$) in subjective correlated equilibrium play. Consider the direct mechanism $p_1 = p_2 = \frac{1}{3}(U, L) + \frac{1}{3}(D, L) + \frac{1}{3}(D, R)$. This is a subjective correlated equilibrium, and since players need not be correct about beliefs, we can have (D, L) played with probability 1.

Example 20 (continued). Moreover, a subjective correlated equilibrium can achieve *ex ante* payoffs that are not even feasible (of course, such payoffs can never be achieved *ex post*). In the battle-of-the-sexes game, suppose $\Omega = \{\omega', \omega''\}$, $p_1(\omega') = 1$ and $p_2(\omega'') = 1$. Then each player i finds it optimal to play

$$\sigma_i(\omega) = \begin{cases} b_i & \text{if } \omega = \omega', \\ f_i & \text{if } \omega = \omega''. \end{cases}$$

Given each player's priors, we have *ex ante* expected payoffs,

$$\begin{aligned} u_1(\omega) &= \sum_{\omega \in \Omega} u_1(\sigma(\omega))p_1(\omega) = u_1(\sigma(\omega'))p_1(\omega') = 4, \\ u_2(\omega) &= \sum_{\omega \in \Omega} u_2(\sigma(\omega))p_2(\omega) = u_2(\sigma(\omega''))p_2(\omega'') = 4. \end{aligned}$$

So the *ex ante* payoff is $(4, 4)$, yet this is not feasible!

2.13 Coalition-proof Nash equilibrium

Correlated equilibrium can often be motivated via a communication story. In Example 20 (Bach or Stravinsky) we argued that if the players got together beforehand and agreed to flip a coin to decide where to go, they could generate payoffs that didn't lie within the set of Nash equilibrium payoffs. This expanded the set of solutions.

However, communication possibilities can also motivate us to refine the set of Nash equilibria. In the Bach or Stravinsky game, the non-degenerate mixed strategy Nash equilibrium is Pareto dominated by either of the two pure strategy equilibria. The non-degenerate mixed strategy equilibrium is pretty implausible in context since you can presumably communicate beforehand – if you make arrangements to go on a date or go to a concert with a friend, then it would be bizarre to instead go someplace else in the expectation that the other person might show up there. One obvious refinement, then, is simply to filter out Pareto-dominated equilibria, restricting to the *Pareto set of Nash equilibria*, i.e. the set of all Nash equilibria σ for which there is no other Nash equilibrium σ' with $u_i(\sigma') \geq u_i(\sigma)$ for all players $i \in \mathcal{I}$ with strict inequality for some player i .

However, in settings with more than two players, private communication may be possible among subsets of players, and so we need to consider deviations other than unilateral deviations – this restricts the set of credible equilibria even further than if we simply rule out Nash dominated equilibria. Example 22 illustrates this.

Call any (nonempty) subset $C \subseteq \mathcal{I}$ a *coalition*, let $S_C = \times_{i \in C} S_i$ be the set of strategy profiles $(s_i)_{i \in C}$ associated with coalition C . A stronger requirement proposed by Aumann (1959) is that equilibrium is robust to every joint deviation by every possible coalition:

Definition 35 (Strong Nash equilibrium). Given a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a strategy profile $s^* \in \times_{i \in \mathcal{I}} S_i$ is a *strong Nash equilibrium* if for every nonempty coalition $C \subseteq \mathcal{I}$ and for all $s_C \in S_C := \times_{j \in C} S_j$, there is some player $i \in C$ for which

$$u_i(s^*) \geq u_i(s_C, s_{-C}^*),$$

where $s_{-C}^* = (s_j^*)_{j \notin C}$.

This is far more stringent than Nash equilibrium – whereas Nash equilibrium is robust to unilateral deviations, strong Nash equilibrium is robust to all deviations. What strong Nash equilibrium rules out is that any nonempty coalition C can jointly deviate in a way such that all its members are strictly better off. This ensures *weak* Pareto optimality:

Proposition 22. *All strong Nash equilibria are weakly Pareto optimal.*

Proof. Fix a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$. Suppose $s \in \times_{i \in \mathcal{I}} S_i$ is strongly Pareto dominated by s' . Then if the grand coalition \mathcal{I} jointly deviates to s' , we have $u_i(s') > u_i(s)$ for all $i \in \mathcal{I}$, so there is no i for which $u_i(s) \geq u_i(s')$. Hence s is not a strong Nash equilibrium. \square

Note weak Pareto optimality implies there are no strong Nash improvements on a strong Nash equilibrium, but a strong Nash equilibrium can still be weakly Pareto dominated. Consider the following game.

| | | |
|-------|-------|-------|
| | A_2 | B_2 |
| A_1 | 3, 2 | 0, 0 |
| B_1 | 0, 0 | 2, 2 |

Here, (A_1, A_2) weakly Pareto dominates (B_1, B_2) , but does not strictly dominate it since Player 2's payoff is the same under both profiles. For any possible joint deviation s' from (B_1, B_2) involving Player 2, $u_2(s') = u_2(B_1, B_2)$, and if Player 1 deviates alone, her payoff is strictly lower. Hence (B_1, B_2) is a strong Nash equilibrium.²³

Strong Nash equilibrium has two main drawbacks. One drawback is inconvenience: because it is so stringent, it rarely exists. More importantly, it is not very credible, because it gives too much freedom to coalitions in how they choose joint deviations. In particular, fix some strategy profile s . Suppose there is a coalition C that has a profitable joint deviation s'_C for all its members, but that there is some subcoalition of C that has a further incentive to deviate from s'_C . Then in practice, we would not expect C to threaten the stability of s , because the members of C would never agree to deviate to s'_C since a subcoalition would form a side agreement to deviate from s'_C .

This limits the set of plausible joint deviations to those that are “self-enforcing” in the sense that players in the deviating coalition would all agree to follow through on the deviation. This is Bernheim, Peleg & Whinston's (1987) insight. The definition is recursive.

Definition 36 (Coalition-proof Nash equilibrium). For any n -player normal form game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, let $\mathcal{C} = 2^{\mathcal{I}} - \{\emptyset, \mathcal{I}\}$ be the set of *possible coalitions*. For each coalition $C \in \mathcal{C}$, write $S_C = \times_{i \in C} S_i$ and $S_{-C} = \times_{i \notin C} S_i$.

Given a coalition $C \in \mathcal{C}$, define $u_i(\cdot \mid s_{-C}) : S_C \rightarrow \mathbb{R}$ by $u_i(s_C \mid s_{-C}) = u_i(s_C, s_{-C})$ for all $s_C \in S_C$. We say that $G|_{s_{-C}} = (C, (S_i, u_i(\cdot \mid s_{-C}))_{i \in C})$ is the *subgame induced on subgroup C by strategy profile s_{-C}* .

- (i) In a one player game $G = (\{1\}, (S_1, u_1))$, a strategy $s^* \in S$ is a *coalition-proof Nash equilibrium* if $u_1(s^*) \geq u_1(s)$ for all $s \in S$.
- (ii) Fix an n -player game G for $n > 1$ and suppose we have defined coalition-proof Nash equilibrium for all games with $|\mathcal{I}| < n$.

²³We could make strong Nash equilibrium even more stringent by requiring that $u_i(s^*) \geq u_i(s_C, s_{-C}^*)$ for all $i \in C$ for each nonempty coalition C and all $s_C \in S_C$. Under this more stringent criterion, a strong Nash equilibrium is strongly Pareto optimal.

- (a) Call a strategy profile $s^* \in S$ *self-enforcing* if s_C^* is a coalition-proof Nash equilibrium of $G_{|s_C^*}$ for every coalition $C \in \mathcal{C}$. Let $\mathcal{S} \subseteq S$ denote the set of self-enforcing strategy profiles of G .
- (b) Call strategy profile $s^* \in S$ a *coalition-proof Nash equilibrium* of G if s^* is self-enforcing and there is no self-enforcing $s \in S$ such that $u_i(s) > u_i(s^*)$ for all $i \in \mathcal{I}$.

Proposition 23. *Every strong Nash equilibrium is a coalition-proof Nash equilibrium.*

Proof. Fix a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$. If $s^* \in S$ is a strong Nash equilibrium, then it is a strong Nash equilibrium in every subgame $G_{|s_C^*}$ for all $C \in \mathcal{C}$. Now, for $|\mathcal{I}| = 1$, coalition-proof Nash and strong Nash equilibrium coincide. For $n > 1$, suppose if s^* is a strong Nash equilibrium for all $|\mathcal{I}| < 1$ then it is coalition-proof for n . \square

Example 22 (Bernheim, Peleg & Whinston, 1987). Consider the following three player game.

| | | | | | |
|-------|-----------|-----------|-------|-----------|-----------|
| | A_3 | | | B_3 | |
| | A_2 | B_2 | | A_2 | B_2 |
| A_1 | 1, 1, -5 | -5, -5, 0 | A_1 | -1, -1, 5 | -5, -5, 0 |
| B_1 | -5, -5, 0 | 0, 0, 10 | B_1 | -5, -5, 0 | -2, -2, 0 |

This game has two pure Nash equilibria, (B_1, B_2, A_3) and (A_1, A_2, B_3) . Moreover, (B_1, B_2, A_3) strongly Pareto dominates (A_1, A_2, B_3) . Thus the Pareto dominance criterion selects the equilibrium (B_1, B_2, A_3) .

However, (B_1, B_2, A_3) is not coalition-proof since (A_1, A_2) is a profitable joint deviation for the coalition $\{1, 2\}$ in the subgame induced by $s_3 = A_3$ so (B_1, B_2, A_3) is not self-enforcing. On the other hand, (A_1, A_2, B_3) is the only self-enforcing Nash equilibrium so is a coalition-proof Nash equilibrium.

Since any strong Nash equilibrium is also a coalition-proof Nash equilibrium and is weakly Pareto efficient, there is no strong Nash equilibrium in this game, since the unique weakly Pareto efficient Nash equilibrium in this game is not coalition-proof.

2.14 Epistemic foundations of equilibrium

Nash equilibrium and its cousins are often treated as something of a black box. Indeed, much thought about what conditions are sufficient for equilibrium to be played has only been developed long after the initial solution concepts were developed. In this section, we try to clarify what rational agents need to know in order to play equilibrium.

Consider a normal form game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ of n players, and let (Ω, \mathcal{F}) be a set of states Ω equipped with a σ -algebra \mathcal{F} . We assume each player i has a partitioned information function $h_i : \Omega \rightarrow \mathcal{F}$. Each state $\omega \in \Omega$ specifies, for each player $i \in \mathcal{I}$, i 's knowledge $h_i(\omega) \in \mathcal{F}$, i 's pure strategy $s_i(\omega) \in S_i$, and i 's belief $\mu_i(\omega) \in \Delta(S_{-i})$ about opponents' play. We assume that if $\omega, \omega' \in h_i(\omega)$ then $\mu_i(\omega) = \mu_i(\omega')$. We use $s = (s_1, \dots, s_n)$ and $\mu = (\mu_1, \dots, \mu_n)$ to denote the strategy and belief profiles respectively.

Proposition 24. *Let G be a finite game. Suppose that in state $\omega \in \Omega$, for each player $i \in \mathcal{I}$*

- (i) $h_i(\omega) \subseteq \{\omega' \in \Omega \mid s_{-i}(\omega') = s_{-i}(\omega)\}$ (player i knows the actions of other players);
- (ii) the support of $\mu_i(\omega)$ lies in $\{s_{-i}(\omega') \in S_{-i} \mid \omega' \in h_i(\omega)\}$ (player i has a belief consistent with such knowledge);
- (iii) $s_i(\omega)$ is a best response to $\mu_i(\omega)$ (player i is rational).

Then $s(\omega)$ is a Nash equilibrium of G .

Proof. By (i) and (ii), each player i 's beliefs assign probability one to the profile $s_{-i}(\omega)$, and by (iii), $s_i(\omega)$ is a best response to $\mu_i(\omega) = s_{-i}(\omega)$. Hence $s(\omega)$ must be a Nash equilibrium. \square

That each player knows the strategy of their opponents is a very strong assumption. Aumann & Brandenburger (1995) ask how much we can relax this assumption. In the case of two-player games, we can replace the assumption that players know the actions of other players with the weaker assumption that players mutually know each others' beliefs and know they are both rational:

Proposition 25. *Let G be a two-player finite game. Suppose that in state $\omega \in \Omega$, for each player $i \in \mathcal{I}$,*

- (i) $h_i(\omega) \subseteq \{\omega' \in \Omega \mid \mu_{-i}(\omega') = \mu_{-i}(\omega)\}$ (player i knows the belief of her opponent);
- (ii) for any $\omega' \in h_i(\omega)$, an action $s_{-i}(\omega')$ is in the support of $\mu_i(\omega')$ only if $s_{-i}(\omega')$ is a best response for player $-i$ to $\mu_{-i}(\omega')$ (player i knows her opponent is rational and i 's beliefs are consistent with her knowledge).

Then the strategy profile $\sigma = (\mu_1(\omega), \mu_2(\omega))$ is a Nash equilibrium of G .

Proof. Let $s_{-i} \in S_{-i}$ lie in the support of $\mu_i(\omega)$. Then there is some state $\omega' \in h_i(\omega)$ such that $s_{-i}(\omega') = s_{-i}$, and so s_{-i} is a best response to $\mu_{-i}(\omega')$ by (ii). Now by (i), $\mu_{-i}(\omega) = \mu_{-i}(\omega')$. \square

The assumptions here are quite weak – we don't need beliefs to be derived from a common prior, we don't need that beliefs or rationality are common knowledge, and we don't actually need that the game is common knowledge (only mutual knowledge).

Unfortunately, this result does not extend to games with more than two players. In general, we need that beliefs are derived from a common prior, that rationality is mutual knowledge, and that beliefs are common knowledge.

Proposition 26. *Let G be a finite game of n players, and let $\mu(\omega) = (\mu_1(\omega), \dots, \mu_n(\omega))$. Suppose players have a common prior p on Ω that assigns positive probability to the structure of the game being mutually known. Suppose that in state $\omega \in \Omega$, the structure of G is mutual knowledge, that players' rationality is mutual knowledge, and that the belief profile $\mu(\omega)$ is common knowledge. Then $\mu(\omega)$ is a Nash equilibrium of G .*

Proof. See Aumann & Brandenburger (1995) for the proof. \square

Correlated equilibrium also has an epistemic foundation:

Proposition 27 (Aumann, 1987). *Let G be a finite game and let $(\Omega, \{\mathcal{H}_i\}_{i \in \mathcal{I}}, p)$ be a correlating mechanism. For each player i , let \mathcal{F}_i be the σ -algebra generated by \mathcal{H}_i . Let σ be a correlated strategy profile. For every player $i \in \mathcal{I}$, suppose that*

- (i) *player i is rational;*
- (ii) *player i 's belief μ_i is derived from the common prior p and $p(h_i(\omega)) > 0$ for all $\omega \in \Omega$, and*
- (iii) *player i 's strategy $\sigma_i : \Omega \rightarrow S_i$ is measurable with respect to \mathcal{F}_i .*

Then σ is a correlated equilibrium relative to the correlating mechanism $(\Omega, \{\mathcal{H}_i\}, p)$.

Proof. This is immediate from Definition 30. \square

2.15 Learning equilibrium

We mentioned that one hypothesis for how Nash equilibrium (or other related notions) might arise is as the outcome of a learning process. These are processes in which players myopically adapt their strategies in response to what they observe of their opponents' play. There is a large literature that attempts to model such learning processes, and it is now big in computer science because of the machine learning craze.

2.15.1 Best response dynamics

A particular simple myopic process is the *best response dynamic*. Suppose the same pure strategy n -player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is played repeatedly. Let $B_i(s_{-i})$ denote the set of best responses of player i to their opponents' strategy profile s_{-i} . The best response dynamic is an algorithm as follows:

- (I) Fix an initial strategy profile $s \in S = S_1 \times \cdots \times S_n$.
- (II) While s is not a Nash equilibrium, choose any player $i \in \mathcal{I}$ for whom $s_i \notin B_i(s_{-i})$. Set $s_i \in B_i(s_{-i})$.
- (III) Terminate when s is a Nash equilibrium, i.e. when there is no player left i for whom $s_i \notin B_i(s_{-i})$.

If the dynamic terminates, then clearly it arrives at a Nash equilibrium. The process has much in common with local search algorithms in computer science. In fact, it is easiest to think about the best response dynamic in terms of a *best response graph*. This is a graph $G = (V, E)$ with vertices $V = S$ being the set of strategy profiles and edges being pairs of strategies, with the edge $(s, (s'_i, s_{-i}))$ being included in E iff

- (i) s_i is not a best response for player i to s_{-i} , and

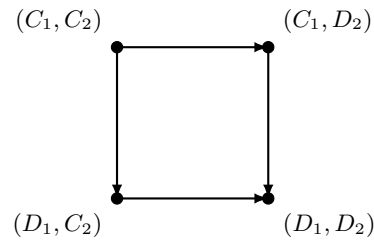
(ii) s'_i is a best response for i to s_{-i} .

If a strategy profile s has no outgoing edges in the best response graph (i.e. it is a *sink*), then it is a pure strategy Nash equilibrium.

Example 2 (continued). Consider again the prisoner's dilemma:

| | C_2 | D_2 |
|-------|-------|-------|
| C_1 | 3, 3 | 0, 4 |
| D_1 | 4, 0 | 1, 1 |

The best response graph for this game is:



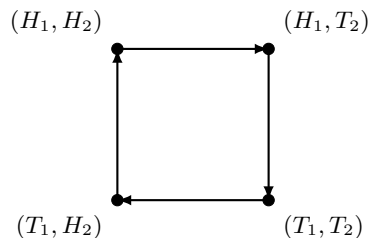
As we would expect, the only sink node is (D_1, D_2) , since this is the unique Nash equilibrium – every other node has an outgoing edge.

For games that do not have a pure strategy Nash equilibrium, the best response dynamic will never terminate. Even if a game does have a pure strategy Nash equilibrium, the best response dynamic still might not terminate depending on our initial choice of strategy profile.

Example 9 (continued). Recall the game of matching pennies:

| | H_2 | T_2 |
|-------|-------|-------|
| H_1 | 1, -1 | -1, 1 |
| T_1 | -1, 1 | 1, -1 |

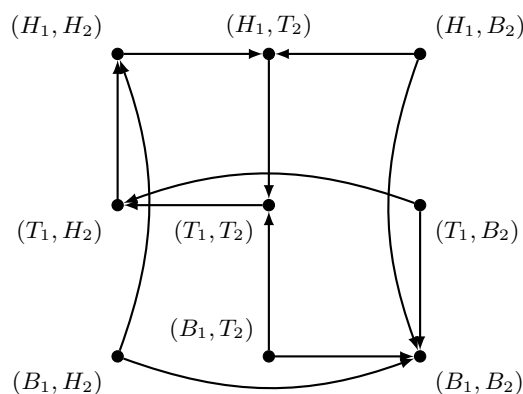
Drawing out the best response graph, we see that the best response dynamic will never terminate, as we would expect given no pure strategy Nash equilibrium exists in this game: The best response graph for this game is:



Now consider a modified version of matching pennies. Each player i can choose heads or tails as normal, but they also have the option to balk, B_i . The payoffs if both choose heads or tails are as usual, but if a player balks and their opponent does not, they have to pay their opponent 50 cents, but if they both balk, they both get paid a dollar. The normal form of this game is:

| | H_2 | T_2 | B_2 |
|-------|---------------------------|---------------------------|---|
| H_1 | $\textcolor{blue}{1}, -1$ | $-1, \textcolor{red}{1}$ | $1/2, -1/2$ |
| T_1 | $-1, \textcolor{red}{1}$ | $\textcolor{blue}{1}, -1$ | $1/2, -1/2$ |
| B_1 | $-1/2, 1/2$ | $-1/2, 1/2$ | $\textcolor{blue}{1}, \textcolor{red}{1}$ |

The best response graph is now:



As we would expect, (B_1, B_2) is the unique sink. However, if we start at (H_1, H_2) , there is no path to (B_1, B_2) , and the best response dynamic will never terminate.

The best response dynamic is only guaranteed to terminate for any initial strategy profile iff the best response graph is acyclic. If the graph contains some cycle, then we can always choose the best responses in (II) so that we stay in the cycle forever.

2.15.2 Fictitious play

Brown (1951) introduced *fictitious play*, a learning rule in which each player assumes their opponents play stationary mixed strategies and choose a best response based on the empirical frequency of opponents' play.²⁴ Suppose there are two players, $i = 1, 2$. Suppose the players play a finite game $G = (\{1, 2\}, \{S_i, u_i\}_{i=1,2})$ at times $t = 0, 1, 2, \dots$. For each t , let $\eta_i^t : S_{-i} \rightarrow \mathbb{N}$ be such that $\eta_i^t(s_{-i})$ gives the number of times at time t that player i has observed s_{-i} in the past. We take players' actions in period 0 as exogenously given. Since each player assumes their opponent follows a stationary mixed strategy, each

²⁴Fictitious play was big in the 50s and early 60s. Its popularity then waned, both because of the negative results of Shapley (1964) and because game theory moved in a less behavioural direction over the next couple of decades. There was a resurgence of interest in fictitious play and other myopic learning processes in the 1990s and 2000s, as the field swung back towards questioning some of the strong rationality assumptions it had developed.

player's beliefs are some distribution μ_i^t on $\Delta(S_{-i})$. Players update their beliefs using Bayesian updating. Since the distribution of $\eta_i^t(s_{-i})$ is multinomial, a typical assumption is that μ_i^0 is a Dirichlet distribution with $\mu_i^0(\sigma_{-i}) = \frac{1}{B} \prod_{s_{-i} \in S_{-i}} \sigma_{-i}(s_{-i})^{\eta_i^0(s_{-i})}$, where B is a normalizing constant.²⁵ In period t , the player expects play $\hat{\mu}_i^t(s_{-i}) = \mathbb{E}_{\mu_i^t} \sigma_{-i}(s_{-i})$. Because the prior is Dirichlet, Bayesian updating implies that $\hat{\mu}_i^t(s_{-i}) = \frac{\eta_i^t(s_{-i})}{\sum_{s'_{-i} \in S_{-i}} \eta_i^t(s'_{-i})}$. Note that each player's assumption that the other is playing a stationary mixed strategy is wrong, because both are instead following fictitious play. Thus even though players update their forecasts correctly, their priors are unreasonable. At each time t , player i chooses a best response to her forecast μ_i^t , that is, she chooses a strategy $s_i^t \in \arg \max_{s_i \in S_i} u_i(s_i, \mu_i^t)$. Under the assumption that her opponent plays a stationary mixed strategy – and thus is *not* updating their own play based on her actions – this myopic choice would be optimal.

To investigate convergence, we should also track the number of times a player has played a given strategy in the preceding periods. For each t , let $\alpha_i^t : S_i \rightarrow \mathbb{N}$ be such that $\alpha_i^t(s_i)$ gives the number of times at time t that player i has previously played s_i .

Definition 37 (Convergence of fictitious play).

- (a) *Convergence.* We say that the sequence of pure strategy profiles $\{s^t\}$ *converges* to a pure strategy profile s if there exists a time T such that $s^t = s$ for all $t \geq T$.
- (b) *Convergence in a time-average sense.* We say that a sequence of pure strategy profiles $\{s^t\}$ *converges in a time-average sense* to a mixed strategy profile σ if, for each $i \in \mathcal{I}$, $\sigma_i(s_i) = \lim_{t \rightarrow \infty} \frac{1}{t} \alpha_i^t(s_i)$ for all $s_i \in S_i$.

Proposition 28. *Given a finite game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, if a sequence of fictitious play $\{s^t\}$ converges in a time-average sense to a mixed strategy profile σ , then σ is a Nash equilibrium of $G^m = (\mathcal{I}, (\Delta(S_i), u_i)_{i \in \mathcal{I}})$.*

Proof. Suppose $s^t \rightarrow \sigma$ in a time-average sense but σ is not a Nash equilibrium of G^m . Then for some player i , there exists a pair of strategies $s_i, s'_i \in S_i$ such that $\sigma_i(s_i) > 0$ and $u_i(s'_i, \sigma_{-i}) > u_i(s_i, \sigma_{-i})$. Let $N = |S_{-i}|$. Choose $\epsilon > 0$ s.t. $\epsilon < \frac{1}{2}(u_i(s'_i, \sigma_{-i}) - u_i(s_i, \sigma_{-i}))$ and T sufficiently large that $|\hat{\mu}_i^t(s_{-i}) - \sigma_{-i}(s_{-i})| < \frac{\epsilon}{2N}$ for all $t \geq T$. Now, for any $t \geq T$,

$$\begin{aligned}
u_i(s_i, \hat{\mu}_i^t) &= \sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) \hat{\mu}_i^t(s_{-i}) \\
&\leq \sum_{s_{-i} \in S_{-i}} u_i(s_i, s_{-i}) \sigma_{-i}(s_{-i}) + \epsilon \\
&< \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) \sigma_{-i}(s_{-i}) - \epsilon \\
&\leq \sum_{s_{-i} \in S_{-i}} u_i(s'_i, s_{-i}) \hat{\mu}_{-i}^t(s_{-i}) = u_i(s'_i, \hat{\mu}_i^t).
\end{aligned}$$

²⁵The Dirichlet distribution is the conjugate prior distribution of the multinomial distribution.

Hence s_i is never a best response for $t \geq T$, and so $\lim_{t \rightarrow \infty} \frac{1}{t} \alpha_i^t(s_i) = 0$, yet $\sigma_i(s_i) > 0$, yielding a contradiction. \square

This is good news for fictitious play – if fictitious play converges, then agents learn a Nash equilibrium. Even better, we know that in some two-person games, convergence under fictitious play is guaranteed. For example, Robinson (1951) shows that in two-person finite zero-sum games, fictitious play converges in the time-average sense, and Miyasawa (1951) proves the same result for 2×2 games (two-person, two-strategy).

The bad news is that in general, we often do not get convergence even in the time-average sense. Even worse, even when there is convergence in the time-average, it is often very counterintuitive, as the following examples illustrate.

Example 23.

- (a) *Rock-Paper-Scissors*. Shapley (1964) shows that fictitious play can fail to converge in the following game of rock-paper-scissors:

| | L | C | R |
|-----|------|------|------|
| T | 0, 0 | 1, 0 | 0, 1 |
| M | 0, 1 | 0, 0 | 1, 0 |
| B | 1, 0 | 0, 1 | 0, 0 |

The unique Nash equilibrium in this game has both players play mixed strategy $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. Suppose in period 0, (T, M) is played. Then in period 1, Player 1 expects M and Player 2 expects T , so (T, R) is played. (T, R) continues to be played until Player 1 switches to M , after which (M, R) is played, until Player 2 switches to L . Then (M, L) is played until Player 1 switches to B , and (B, L) is played until Player 2 switches to M , and then (B, M) is played until Player 1 switches to T , and we are back where we started. The number of periods at which each strategy pair is played before switching grows exponentially over time.

- (b) *Matching pennies*. Consider matching pennies:

| | H_2 | T_2 |
|-------|-------|-------|
| H_1 | 1, -1 | -1, 1 |
| T_1 | -1, 1 | 1, -1 |

Suppose players begin at (H_1, H_2) . The sequence of play is then as follows:

| t | $\eta_1 t$ | η_2^t | s^t |
|-----|------------|------------|--------------|
| 0 | (0,0) | (0,0) | (H_1, T_2) |
| 1 | (0,1) | (1,0) | (T_1, T_2) |
| 2 | (0,2) | (1,1) | (T_1, H_2) |
| 3 | (1,2) | (1,2) | (T_1, H_2) |
| 4 | (2,2) | (1,3) | (T_1, H_2) |
| 5 | (3,2) | (1,4) | (H_1, H_2) |
| 6 | (4,2) | (2,4) | (H_1, H_2) |
| 7 | (5,2) | (3,4) | (H_1, H_2) |
| 8 | (6,2) | (4,4) | (H_1, T_2) |
| 9 | (6,3) | (5,4) | (H_1, T_2) |
| 10 | (6,4) | (6,4) | (H_1, T_2) |

Play converges in the time average sense to $((1/2, 1/2), (1/2, 1/2))$, yet neither player actually plays a mixed strategy – play here is deterministic except when $\eta_i^t = (k, k)$ for some k (in this case, either strategy is a best response). Because players wrongly believe each other to be playing stationary mixed strategies, a rational player who realizes that their opponent is following fictitious play can easily take advantage of this fact by predicting what their opponent will play next given their opponent's learning rule and past play. For example, at $t = 5$, if Player 2 knows that Player 1 is playing fictitious play, she knows Player 1 will play H_1 next, and so she is better off playing T_2 .

There are of course many other models of learning. Borgers & Sarin (1997) and Erev & Roth (1998) consider models of reinforcement learning, for example. Reinforcement learning is now very big in machine learning.

2.15.3 Self-confirming equilibrium

If learning is a foundation for (some) Nash equilibria, then it is worth asking which equilibria we can expect to be the product of some learning process. Fudenberg & Levine (1993) suggest *self-confirming equilibrium* as an answer to this question.

Definition 38 (Self-confirming equilibrium). Consider an extensive form game Γ . For each mixed strategy σ_i of player i , let $\pi_i(\phi_i \mid \sigma_i)$ denote the behavioural strategy induced by σ_i at information set ϕ_i . Let Π_i denote the set of player i 's behavioural strategies. For each player i , let μ_i denote player i 's belief over $\Pi_{-i} = \times_{j \neq i} \Pi_j$. Let Φ_{-i} denote the set of information sets not belonging to i , and let $\Phi(s_i, \sigma_{-i})$ denote the set of information sets that can be reached with non-zero probability if player i plays pure strategy s_i and i 's opponents play strategy profile σ_{-i} .

- (a) *Nash equilibrium*. A mixed strategy profile σ^* is a *Nash equilibrium* of Γ if for every player $i \in \mathcal{I}$ and each pure strategy $s_i \in S_i$ in the support of σ_i ,
 - (i) $u_i(s_i, \mu_i) \geq u_i(s'_i, \mu_i)$ for all $s'_i \in S_i$ (s_i is a best response to i 's beliefs), and

- (ii) for all information sets $\phi_j \in \Phi_{-i}$, we have $\mu_i(\{\pi_{-i} \mid \pi_j(\phi_j) = \pi_j(\phi_j \mid \sigma_j^*)\}) = 1$ (i 's beliefs are correct.)
- (b) *Self-confirming equilibrium.* A mixed strategy profile σ^* is a *self-confirming equilibrium* of Γ if for every player $i \in \mathcal{I}$ and each pure strategy $s_i \in S_i$ in the support of σ_i ,
- (i) $u_i(s_i, \mu_i) \geq u_i(s'_i, \mu_i)$ for all $s'_i \in S_i$ (s_i is a best response to i 's beliefs), and
 - (ii) for all information sets $\phi_j \in \Phi(s_i, \sigma_{-i}^*)$, we have $\mu_i(\{\pi_{-i} \mid \pi_j(\phi_j) = \pi_j(\phi_j \mid \sigma_j^*)\}) = 1$ (i 's beliefs are empirically correct.)

Unlike Nash equilibrium, in a self-confirming equilibrium, beliefs only need to be correct for histories that can be reached. In the definition, we assume players perfectly observe the actions of their opponents. In general, self-confirming equilibrium can be defined relaxing this assumption.

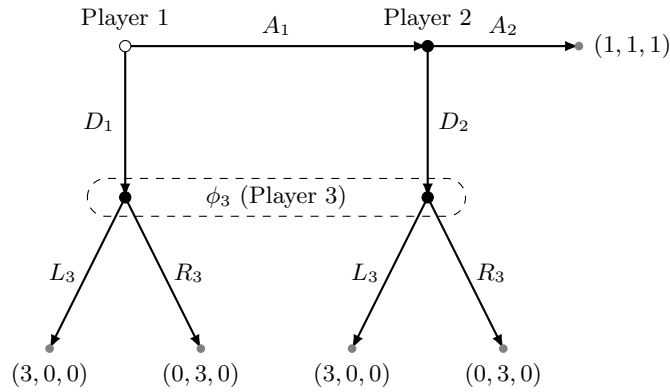
Example 9 (continued). Once again, consider the game of matching pennies:

| | | |
|-------|-------|-------|
| | H_2 | T_2 |
| H_1 | 1, -1 | -1, 1 |
| T_1 | -1, 1 | 1, -1 |

The mixed strategy profile $\sigma^* = ((1/2, 1/2), (1/2, 1/2))$ supported by beliefs $\mu^* = ((1/2, 1/2), (1/2, 1/2))$ is a self-confirming equilibrium. Any other self-confirming equilibrium must involve one of the players i playing a pure strategy s_i , but then j must believe that i will play this strategy since this is the only belief j can hold that is empirically correct. Now j 's best response to this belief and s_i are not mutual best responses, and hence this cannot be a self-confirming equilibrium.

A self-confirming equilibrium does not have to be a Nash equilibrium:

Example 24 (Fudenberg & Kreps, 1993). Consider the following three-player variant of Selten's horse:



Suppose μ_1 is such that Player 1 is certain that Player 3 will play R_3 and Player 2 will play A_2 and μ_2 is such that Player 2 is certain that Player 3 will play L_3 . Then it is optimal for Player 1 to play A_1 and for Player 2 to play A_2 . Since ϕ_3 is never reached, these beliefs can support a self-confirming equilibrium in which the equilibrium path involves (A_1, A_2) .

Yet there is no Nash equilibrium in which (A_1, A_2) is played. Consider any strategy $\sigma_3 = (p, 1 - p)$ of Player 3. If $p \geq \frac{1}{3}$, then Player 2's best response is A_2 , and Player 1's best response is thus D_1 . If $p \leq \frac{1}{3}$, then Player 2's best response is D_2 . In neither case is (A_1, A_2) a pair of best responses for Players 1 and 2.

2.16 Potential games

As Example 23(a) showed, fictitious play need not converge. An interesting question is the following: for what kinds of game can we be sure that learning processes will converge to the set of Nash equilibria? Monderer & Shapley (1996) show that in potential games, fictitious play will always converge. Likewise, the best response dynamic in a potential game will always terminate.

In physics, a *potential* for a vector function $(f_1, \dots, f_n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a function $P : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\frac{\partial f_i}{\partial x_i} = \frac{\partial P}{\partial x_i}$ for all $i = 1, \dots, n$. A potential game is a game that can be summarized by a similar type of function:

Definition 39 (Potential game). Consider an n -player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$.

- (a) *Ordinal potential game*. A function $P : S \rightarrow \mathbb{R}$ is called an *ordinal potential* for G if for every player i and every strategy profile $s_{-i} \in S_{-i}$,

$$u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}) > 0 \quad \text{iff} \quad P(s_i, s_{-i}) - P(s'_i, s_{-i}) > 0 \quad \text{for all } s_i, s'_i \in S_i.$$

We call G an *ordinal potential game* if there exists an ordinal potential for G .

- (b) *Potential game*. A function $P : S \rightarrow \mathbb{R}$ is called an (*exact*) *potential* for G if for every player i and every strategy profile $s_{-i} \in S_{-i}$,

$$u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}) = P(s_i, s_{-i}) - P(s'_i, s_{-i}) \quad \text{for all } s_i, s'_i \in S_i.$$

We call G a *potential game* if there exists a potential for G .

There are other varieties of potential function (weighted potentials, generalized ordinal potentials, and so on). See Monderer & Shapley (1996) for more details on these.

The set of all potential games is a subset of the set of all ordinal potential games:

Proposition 29. *If $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ has a potential then it has an ordinal potential.*

Proof. If P is a potential for G then $u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}) = P(s_i, s_{-i}) - P(s'_i, s_{-i})$ for every player i , every strategy profile $s_{-i} \in S_{-i}$, and every pair of strategies $s_i, s'_i \in S_i$. Clearly the left-hand side is positive iff the right-hand side is positive. \square

Proposition 30. *P is a potential for a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ iff there exists a collection of functions $h_i : S_{-i} \rightarrow \mathbb{R}$ such that $u_i(s) = P(s) + h_i(s_{-i})$ for every player i .*

Proof. Suppose we can write u_i as $u_i(s) = P(s) + h_i(s_{-i})$ for each $i \in \mathcal{I}$ and $s \in S$. Then for each player i ,

$$\begin{aligned} u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}) &= P(s_i, s_{-i}) + h_i(s_{-i}) - P(s'_i, s_{-i}) - h_i(s_{-i}) \\ &= P(s_i, s_{-i}) - P(s'_i, s_{-i}) \end{aligned}$$

for every $s_i, s'_i \in S_i$ and $s_{-i} \in S_{-i}$. Thus P is a potential.

Conversely, suppose P is a potential. Then

$$u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}) = P(s_i, s_{-i}) - P(s'_i, s_{-i})$$

for all $i \in \mathcal{I}$, $s_i, s'_i \in S_i$ and $s_{-i} \in S_{-i}$. Defining

$$h_i(s_{-i}) := u_i(s_i, s_{-i}) - P(s_i, s_{-i}) = u_i(s'_i, s_{-i}) - P(s'_i, s_{-i})$$

gives us $u_i(s) = P(s) + h_i(s_{-i})$. □

Potentials are unique up to a constant:

Proposition 31. *If P_1 and P_2 are potentials for a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ then there is some constant $c \in \mathbb{R}$ such that $P_1 = P_2 + c$. Moreover, if P is a potential for G then $P + c$ is a potential for G for any $c \in \mathbb{R}$.*

Proof. Suppose G has n players. Fix $\hat{s} \in S$ and define

$$H(s) = \sum_{i=1}^n [u_i(\tilde{s}^{i-1}) - u_i(\tilde{s}^i)]$$

where $\tilde{s}^0 = s$ and $\tilde{s}^i = (s_{-i}^{i-1}, \hat{s}_i)$, for all $s \in S$. Now, if P is a potential for Γ , then

$$\begin{aligned} P(s) - P(\hat{s}) &= \sum_{i=1}^n [P(\tilde{s}^{i-1}) - P(\tilde{s}^i)] \\ &= \sum_{i=1}^n [u_i(\tilde{s}^{i-1}) - u_i(\tilde{s}^i)] = H(s), \end{aligned}$$

for all $s \in S$. Since P_1 and P_2 are both potentials for G , we have that $P_1(s) - P_1(\hat{s}) = H(s) = P_2(s) - P_2(\hat{s})$, or $P_1(s) - P_2(s) = P_1(\hat{s}) - P_2(\hat{s})$, and the rhs is fixed. Thus $P_1(s) - P_2(s) = c$ for constant $c := P_1(\hat{s}) - P_2(\hat{s})$. □

Say that a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is *smooth* if each S_i is a closed interval of \mathbb{R} and each u_i is twice continuously differentiable. In the case of smooth potential games, we have:

Proposition 32. Suppose $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is smooth.

(a) P is a potential for G iff

$$\frac{\partial P(s)}{\partial s_i} = \frac{\partial u_i(s)}{\partial s_i} \quad \text{for all } s \in S \text{ and } i \in \mathcal{I}.$$

(b) G is a potential game iff

$$\frac{\partial^2 u_i}{\partial s_i \partial s_j} = \frac{\partial^2 u_j}{\partial s_j \partial s_i} \quad \text{for every pair of players } i \neq j.$$

Proof. (a) Suppose P is a potential for G , and fix a strategy profile $s \in S^\circ$, the interior of S . By definition,

$$P(s_i + t, s_{-i}) - P(s_i, s_{-i}) = u_i(s_i + t, s_{-i})$$

for all t s.t. $s_i + t \in S_i$ and each player i . Thus

$$\frac{P(s_i + t, s_{-i}) - P(s_i, s_{-i})}{t} = \frac{u_i(s_i + t, s_{-i})}{t}.$$

Taking limits gives us

$$\begin{aligned} \frac{\partial P(s)}{\partial s_i} &= \lim_{t \rightarrow 0} \frac{P(s_i + t, s_{-i}) - P(s_i, s_{-i})}{t} \\ &= \lim_{t \rightarrow 0} \frac{u_i(s_i + t, s_{-i}) - u_i(s_i, s_{-i})}{t} = \frac{\partial u_i(s)}{\partial s_i}. \end{aligned}$$

Conversely, suppose $\frac{\partial P(s)}{\partial s_i} = \frac{\partial u_i(s)}{\partial s_i}$ for all $s \in S^\circ$ and every player i . Fix a strategy profile $s \in S$. We have, by the fundamental theorem of calculus,

$$\begin{aligned} P(s_i, s_{-i}) - P(s'_i, s_{-i}) &= \int_{s'_i}^{s_i} \frac{\partial P(s_i, s_{-i})}{\partial s_i} ds_i \\ &= \int_{s'_i}^{s_i} \frac{\partial u(s_i, s_{-i})}{\partial s_i} ds_i \\ &= u_i(s_i, s_{-i}) - u_i(s'_i, s_{-i}), \end{aligned}$$

and so P is a potential for G .

(b) Suppose G is a potential game. Then taking (a) and differentiating,

$$\frac{\partial^2 u_i(s)}{\partial s_i \partial s_j} = \frac{\partial^2 P(s)}{\partial s_i \partial s_j} = \frac{\partial^2 P(s)}{\partial s_j \partial s_i} = \frac{\partial^2 u_j(s)}{\partial s_j \partial s_i}$$

for all $i \neq j$, where the second equality follows because u_i is twice continuously differentiable, and thus from (a), so must be P , and thus its second order partial derivatives are symmetric.

Conversely, suppose $\frac{\partial^2 u_i(s)}{\partial s_i \partial s_j} = \frac{\partial^2 u_j(s)}{\partial s_j \partial s_i}$ for all pairs $i \neq j$. Fix a strategy profile $\bar{s} \in S$. Let $\tau : [0, 1] \rightarrow S$ be an almost-everywhere continuously differentiable path from \bar{s} to s in S .²⁶ We claim

$$P(s) = \sum_{i \in \mathcal{I}} \int_0^1 \frac{\partial u_i(\tau(t))}{\partial s_i} \frac{\partial \tau_i(t)}{\partial t} dt$$

is a potential for G .

Fix $s_{-i} \in S_{-i}$ and take any $s_i, s'_i \in S_i$. Now,

$$P(s_i, s_{-i}) - P(s'_i, s_{-i}) = \int_0^a \frac{\partial u_i(\tau(t))}{\partial s_i} \frac{\partial \tau_i(t)}{\partial t} dt$$

□

Example 2 (continued). The prisoner's dilemma is a potential game:

| | | | | | |
|-------|-------|-------|--|-------|-------|
| | C_2 | D_2 | | C_2 | D_2 |
| C_1 | 3, 3 | 0, 4 | | 1 | 2 |
| D_1 | 4, 0 | 1, 1 | | 2 | 3 |

The payoff matrix is on the left and the potential function is represented on the right.

Example 25 (Potential in Cournot competition).

- (a) First consider a symmetric Cournot game G with n firms, where each firm has constant marginal cost $c > 0$. Let $Q := \sum_{i=1}^n q_i$ and assume the inverse demand function $p(Q)$ is positive. We make no further assumptions about $p(Q)$. Firm i 's profit function is $\pi_i(q) = (p(Q) - c)q_i$ for all $q = (q_1, \dots, q_n) \in (0, \infty)^n$. Define $P : (0, \infty)^n \rightarrow \mathbb{R}$ by

$$P(q) = (p(Q) - c) \prod_{i=1}^n q_i.$$

This is an ordinal potential for G . To see this, note $P(q_i, q_{-i}) = \pi_i(q_i, q_{-i}) \prod_{j \neq i} q_j$ for each firm i . Fixing i and q_{-i} , we have that if $\pi_i(q'_i, q_{-i}) > \pi_i(q_i, q_{-i})$ then multiplying both sides by $\prod_{j \neq i} q_j$ shows $P(q'_i, q_{-i}) > P(q_i, q_{-i})$, and conversely if $P(q'_i, q_{-i}) > P(q_i, q_{-i})$, then dividing both sides by $\prod_{j \neq i} q_j$ shows $\pi_i(q'_i, q_{-i}) > \pi_i(q_i, q_{-i})$.

²⁶That is, τ is almost everywhere continuously differentiable and has $\tau(0) = \bar{s}$ and $\tau(1) = s$.

- (b) Next, consider quasi-Cournot competition with n firms and linear inverse demand $p(Q) = a - bQ$ where $a, b > 0$ and suppose each firm i has an arbitrary differentiable cost function $c_i(q_i)$. Define

$$P(q) = a \sum_{j=1}^n q_j - b \sum_{j=1}^n q_j^2 - b \sum_{1 \leq i < j \leq n} q_i q_j - \sum_{j=1}^n c_j(q_j).$$

This is a potential: fixing i and q_{-i} , we have

$$\begin{aligned} P(q_i, q_{-i}) - P(q'_i, q_{-i}) &= a(q_i - q'_i) - b(q_i^2 - (q'_i)^2) - b \sum_{1 \leq i < j \leq n} (q_i - q'_i)q_j - c_i(q_i) + c_i(q'_i) \\ &= \left(a - bq_i - b \sum_{j \neq i} q_j \right) q_i - c_i(q_i) - \left(a - bq'_i - b \sum_{j \neq i} q_j \right) q'_i + c_i(q'_i) \\ &= \pi_i(q_i, q_{-i}) - \pi_i(q'_i, q_{-i}). \end{aligned}$$

Lemma 8. *If P is an ordinal potential for a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, then G is best response equivalent to the game $\tilde{G} = (\mathcal{I}, (S_i, P)_{i \in \mathcal{I}})$. Thus $s^* \in S$ is a Nash equilibrium for G iff*

$$P(s_i^*, s_{-i}^*) \geq P(s'_i, s_{-i}^*) \quad \text{for all } s'_i \in S_i.$$

Hence if P admits a maximal value then G has a pure strategy Nash equilibrium.

Proof. By definition of the ordinal potential, $u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$ for all $s'_i \in S_i$ iff $P(s_i, s_{-i}) > P(s'_i, s_{-i})$ for all $s'_i \in S_i$, and hence s_i is a best response to s_{-i} in G iff it is a best response to s_{-i} in \tilde{G} . The Nash equilibrium condition follows immediately. Finally, suppose P admits a maximal value, i.e. there is some $s \in S$ s.t. $P(s) \geq P(s')$ for all $s' \in S$. Then $P(s_i, s_{-i}) \geq P(s'_i, s_{-i})$ for all $s'_i \in S_i$ and all i , so s is a Nash equilibrium. \square

Note that while any strategy profile maximizing the potential is a Nash equilibrium, the lemma *does not* imply that all Nash equilibria need to maximize the potential function. Note that if G is a finite ordinal potential game then the lemma implies it has a pure strategy Nash equilibrium. In general, potential games can also have (non-degenerate) mixed strategy Nash equilibria.

Proposition 33. *If $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a finite potential game then it has an acyclic best response graph. Thus the best response dynamics always terminate.*

Proof. Suppose the best response graph for G contains a cycle $(s^0 s^1, s^1 s^2, \dots, s^k s^0)$. For each $\ell = 0, \dots, k$, s^ℓ and $s^{\ell+1}$ differ in the strategy of precisely one agent i_ℓ (let $s^{k+1} := s^0$). For this agent, we have $u_{i_\ell}(s^\ell) > u_{i_\ell}(s^{\ell+1})$. If a potential P were to exist for this game, we would have $P(s^0) < P(s^1) < \dots < P(s^k) < P(s^0)$, yielding a contradiction. Hence if P exists, the best response graph must be acyclic. \square

For the very simple learning dynamic captured by best response dynamics, we see then that the learning process always converges to a Nash equilibrium in any potential game.

We can get an even stronger result. Say that a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ has the *fictitious play property* if every possible sequence of fictitious play converges in a time-average sense to some mixed strategy Nash equilibrium.

Proposition 34 (Monderer-Shapley). *Every finite potential game has the fictitious play property.*

Proof. For any potential game G with potential function P consider $G' = (\mathcal{I}, (S_i, P)_{i \in \mathcal{I}})$.²⁷ The games G and G' are best-response-equivalent, in the sense that for any strategy profile s_{-i} in either game, player i 's set of best responses $B_i(s_{-i})$ is the same in both games for every player i . By Theorem A in Monderer & Shapley (1996), any finite game of identical interests has the fictitious play property (I will not prove Theorem A because it would require quite a lot of set-up). Hence G' has the fictitious play property and, by best-response-equivalence, so must G . \square

Finally, we might ask under what conditions we will have a unique Nash equilibrium in a potential game. This depends on the strict concavity of the potential and on some smoothness conditions.

Proposition 35. *Suppose $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a finite player potential game with potential P , such that S_i is a convex compact set for each $i \in \mathcal{I}$ and P is twice continuously differentiable and strictly concave. Then G has a unique Nash equilibrium (in pure strategies).*

Proof. By Lemma 8, G is best-response-equivalent to the game $\tilde{G} = (\mathcal{I}, (S_i, P)_{i \in \mathcal{I}})$, and if P admits a maximal value in S then G has a pure strategy Nash equilibrium. Since S is compact, P attains a maximum value on S provided P is concave and continuously differentiable on S , and so G has a pure strategy Nash equilibrium. Define weighted sum $w(s, \lambda) = \sum_{i \in \mathcal{I}} \lambda_i P(s)$ for all $s \in S$ and $\lambda \geq 0 \in \mathbb{R}^n$, and let $g(s, \lambda)$ be the corresponding pseudogradient (c.f. Definition 26). By Rosen's uniqueness theorem (Theorem 9), \tilde{G} (and thus G) has a unique pure strategy Nash equilibrium if (i) it has a pure strategy Nash equilibrium and (ii) w is diagonally strictly concave for some λ , i.e. if

$$(s' - s) \cdot g(s, \lambda) + (s - s') \cdot g(s', \lambda) > 0$$

for all $s, s' \in S$ and some λ . By Proposition 13, it is sufficient to show that given the Jacobian $J(s, \lambda)$ of g at (s, λ) , $J(s, \lambda) + J(s, \lambda)^\top$ is negative definite for all $s \in S$ and some λ .

If P is strictly concave, then we can take $\lambda = 1 \in \mathbb{R}^{|\mathcal{I}|}$, and so $J(s, 1) + J(s, 1)^\top$ is the Hessian $H(s)$ of $P(s)$. If P is twice continuously differentiable on S , then the Hessian is symmetric so $J(s, 1) + J(s, 1)^\top = 2H(s)$, and since P is strictly concave,

²⁷We call a game in which everyone has the same payoff function a *game of identical interests*.

we therefore have that $J(s, 1) + J(s, 1)^\top$ is negative definite and so $w(s, \lambda)$ is strictly diagonally concave for $\lambda = 1$. Applying Rosen's uniqueness theorem thus gives us that G has a unique pure strategy equilibrium if P is strictly concave and twice continuously differentiable. \square

2.16.1 Congestion games

Congestion games are models in which agents choose a path across a network and incur a cost based on the amount of traffic flowing through each edge. An obvious application is to model road traffic, where the agents are drivers who wish to travel from their origin to their destination in the shortest amount of time, and strategically choose which route to take to achieve this goal. The edges are road segments in this case. Communications networks can be modelled in a similar way, where the edges are communication channels and information is routed between nodes in the network. A final example: suppose n firms can produce via different production processes requiring different inputs. To use an input, they must pay a setup cost that depends on the number of other firms demanding to use that same input, but face no variable costs.

Definition 40 (Congestion game). An n -player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a *congestion game* if

- (i) for each player $i \in \mathcal{I}$, we have $S_i \in 2^E - \{\emptyset\}$ where E is a nonempty finite set of *facilities* (such as road segments);
- (ii) for each player i , we have $u_i(s) = -c_i(s)$ where $c_i(s)$ is a *cost function*. Let $t_e(s) \in \{0, 1, \dots, n\}$ denote the number of users of facility $e \in E$ under strategy profile s . Let $\kappa : E \times \{0, \dots, n\} \rightarrow \mathbb{R}$ be such that $\kappa(e, t)$ gives the cost of using facility e when it has t users. Then we define the cost function $c_i : S \rightarrow \mathbb{R}$ by

$$c_i(s) = \sum_{e \in S_i} \kappa(e, t_e(s)).$$

Finite potential games and congestion games are closely connected. Rosenthal (1973) proves the following:

Proposition 36 (Rosenthal, 1973). *Every congestion game is a potential game.*

Proof. Consider any n -player congestion game G . For each strategy profile $s \in S$, define

$$P(s) = \sum_{e \in \bigcup_{i=1}^n S_i} \sum_{k=1}^{t_e(s)} \kappa(e, k).$$

We claim P is a potential function for G . Define 1_A to be the indicator function for a set A , that is, $1_A(x) = 1$ if $x \in A$ and $1_A(x) = 0$ otherwise. For each player i , we have that

$$c_i(s) = \sum_{e \in E} \kappa(e, t_e(s)) 1_{S_i}(e).$$

Together with the definition of P , this implies

$$P(s_i, s_{-i}) - c_i(s_i, s_{-i}) = \sum_{e \in E} \sum_{k=0}^{t_e(s_i, s_{-i}) - 1_{s_i}(e)} \kappa(e, k),$$

$$P(s'_i, s_{-i}) - c_i(s'_i, s_{-i}) = \sum_{e \in E} \sum_{k=0}^{t_e(s'_i, s_{-i}) - 1_{s'_i}(e)} \kappa(e, k).$$

Now, $t_e(s_i, s_{-i}) - 1_{s_i}(e) = t_e(s'_i, s_{-i}) - 1_{s'_i}(e)$ for all i and $e \in E$. Thus P is a potential function. \square

Monderer & Shapley (1996) establish a converse:

Proposition 37. *For every potential game, there is a congestion game with the same potential function.*

Proof. The proof is quite longwinded. See the proof of Theorem 3.2 in Appendix B of Monderer & Shapley (1996). \square

2.17 Supermodular and submodular games

A particularly interesting class of games are supermodular games, which encompass many applied models in industrial organization, the theory of signalling, in the networks literature, and so on. These games exhibit strategic complements – a player’s best response is increasing in the actions of her opponents. Such games are very well-behaved – they have pure strategy Nash equilibria, bounded sets of Nash equilibria, and if they have a unique Nash equilibrium, it can be found by iterated strict dominance. An inverse notion is that of submodular games, such as public goods games. These exhibit strategic substitutes.

This section relies a lot on concepts that are detailed in the appendix on monotone comparative statics. It is worth reading section 2.17 first.

2.17.1 Supermodular games and strategic complements

Supermodular games capture the notion of *strategic complements*. In a game of strategic complements, each player’s best response is increasing in the actions of their opponents. This is common in many settings. For example, there is a big literature on peer effects in network games, and these games exhibit strategic complements (see e.g. Ballester, Calvó-Armengol & Zenou, 2006). As we will see, there are many examples in industrial organization and macroeconomics.

Definition 41 (Supermodular game). A game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a *supermodular game* if, for each player $i \in \mathcal{I}$,

- (i) the strategy set S_i is a complete lattice;

- (ii) the payoff function $u_i : S_i \times S_{-i} \rightarrow \mathbb{R}$ is order upper semicontinuous in s_i , order continuous in s_{-i} and upper bounded;
- (iii) u_i is supermodular in s_i ;
- (iv) u_i has increasing differences in s_i and s_{-i} .

If (iii) is replaced with the condition that u_i is strictly supermodular in s_i and (iv) is replaced with the condition that u_i has strictly increasing differences in s_i and s_{-i} , then we say that G is a *strictly supermodular game*.

If each S_i is a rectangle in \mathbb{R}^{k_i} and u_i is twice-differentiable with only nonnegative cross-derivatives, then Topkis' characterization theorem (Theorem 68) tells us the game is supermodular. We call such games *smooth supermodular games*.

It actually turns out that the set of Nash equilibria of any supermodular game is a complete lattice:²⁸

Theorem 13 (Zhou, 1994). *If $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is an n -player supermodular game with each S_i compact, then G has a Nash equilibrium and the set of Nash equilibria of G is a complete sublattice.*

Proof. Let $B : S \rightrightarrows S$ be the best response correspondence (see Definition 22). Since each S_i is a complete lattice for every player i , $S = \times_{i \in \mathcal{I}} S_i$ is also a complete lattice.²⁹ The set of Nash equilibria is the set of fixed points of B . To apply Zhou's fixed point theorem (Theorem 63), all that is left to check is that B is monotonically increasing with respect to the partial order of S . It is sufficient to show that $B_i : S \rightrightarrows S_i$ is monotonically increasing for each player i . Since S is compact and u_i is order upper semicontinuous, $B_i(s)$ is compact for each $s \in S$.

Now consider strategies $s_i, s'_i \in B_i(s)$. Since u_i is supermodular in s_i , $u_i(s_i, s_{-i}) + u_i(s'_i, s_{-i}) \leq u_i(s_i \vee s'_i, s_{-i}) + u_i(s_i \wedge s'_i, s_{-i})$. Thus $s_i \vee s'_i \in B_i(s)$ and $s_i \wedge s'_i \in B_i(s)$, so $B_i(s)$ is a sublattice of S_i .

Finally, consider any $s_{-i} \geq s'_{-i}$, $s_i \in B_i(s)$ and $s'_i \in B_i(s')$. Now,

$$\begin{aligned} 0 &\geq u_i(s_i \wedge s'_i, s_{-i}) - u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i}) - u_i(s_i \wedge s'_i, s_{-i}) \\ &\geq u_i(s'_i, s'_{-i}) - u_i(s_i \wedge s'_i, s'_{-i}) \geq 0, \end{aligned}$$

where the second inequality is by supermodularity, the third is by increasing differences, and the first and last are by definition of B_i . It follows that all the inequalities hold with equality, and thus $s_i \vee s'_i \in B_i(s)$ and $s_i \wedge s'_i \in B_i(s')$. Thus B_i is monotonically increasing. Applying Zhou's fixed point theorem gives the result. \square

Supermodular games have other nice properties. For one, they are often solvable by iterated strict dominance, and even if not, we can find Nash equilibria by applying iterated strict dominance:

²⁸This strengthens Vives' (1990) result that for strictly supermodular games, the set of Nash equilibria is a complete lattice.

²⁹If X, Y are partially ordered sets, we define the partial order on $X \times Y$ as follows: $(x, y) \leq (x', y')$ iff $x \leq x'$ and $y \leq y'$. This extends easily to finite products of sets.

Theorem 14 (Milgrom-Roberts, 1990). *Suppose $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a supermodular game with each strategy set S_i a bounded subset of an Euclidean space and each payoff profile u_i order continuous. Then the supremum \bar{s} and the infimum \underline{s} of the set of strategies surviving iterated strict dominance are both Nash equilibria. Moreover, they are the greatest and least elements of the set of Nash equilibria of G , and if $\bar{s} = \underline{s}$, then G is iterated strict dominance solvable.*

Proof. Let $s^0 = (s_1^0, \dots, s_n^0)$ be the supremum of S , and let $s_i, s'_i \in B_i(s_{-i}^0)$ be such that there is no $s''_i \in B_i(s_{-i}^0)$ with either $s''_i > s_i$ or $s''_i > s'_i$. Suppose $s_i \neq s'_i$. Suppose $u_i(s_i \wedge s'_i, s_{-i}^0) \leq u_i(s_i, s_{-i}^0)$. Supermodularity gives

$$u_i(s_i, s_{-i}^0) - u_i(s_i \wedge s'_i, s_{-i}^0) \leq u_i(s_i \vee s'_i, s_{-i}^0) - u_i(s'_i, s_{-i}^0) < 0,$$

where the final strict inequality is because $s_i \vee s'_i > s'_i$ so cannot be a best response to s_{-i}^0 . Hence $B_i(s_{-i}^0)$ has a supremum s_i^1 . We have $s^1 := (s_1^1, \dots, s_n^1) \leq s^0$. We proceed by induction. In the k th round of iterated strict dominance, if s_i does not satisfy $s_i \leq s_i^k$ then s_i is strictly dominated by $s_i \wedge s_i^k < s_i$. Since $s_{-i} < s_{-i}^{k-1}$ for all strategy profiles s_{-i} remaining after $k-1$ rounds of iterated strict dominance, we have

$$\begin{aligned} u_i(s_i, s_{-i}) - u_i(s_i \wedge s_i^k, s_{-i}) &\leq u_i(s_i, s_{-i}^{k-1}) - u_i(s_i \wedge s_i^k, s_{-i}^{k-1}) \\ &\leq u_i(s_i \vee s_i^k, s_{-i}^{k-1}) - u_i(s_i^k, s_{-i}^{k-1}) < 0. \end{aligned}$$

Since the sequence $\{s^k\}$ is monotonically decreasing and S is a complete lattice, $\bar{s} := \lim_{k \rightarrow \infty} s^k = \bigwedge \{s^k\} \in S$, and so $\{s^k\}$ converges to \bar{s} . To see \bar{s} is a Nash equilibrium, note that for arbitrary s_i we have $u_i(s_i^{k+1}, s_{-i}^k) \geq u_i(s_i, s_{-i}^k)$ at each round k , and by continuity, $u_i(s_i^{k+1}, s_{-i}^k) \rightarrow u_i(\bar{s}_i, \bar{s}_{-i})$ as $k \rightarrow \infty$. Hence it follows that $u_i(\bar{s}_i, \bar{s}_{-i}) \geq u_i(s_i, \bar{s}_{-i})$ for all $s_i \in S_i$.

Proof for \underline{s} is symmetric.

The final statement of the theorem is obvious because any Nash equilibrium is rationalizable and the set of Nash equilibria is a complete sublattice. \square

It is worth discussing some examples here. As we mentioned, supermodular games have a lot of applications in IO and in macroeconomics.

Example 26 (Linear Bertrand game with differentiated products). While the original Bertrand model where firms produce a homogeneous product is *not* a supermodular game, the linear Bertrand game in which firms produce differentiated products is supermodular. Suppose there are n firms, and each firm i sets a price $p_i \geq 0$ for its product. Each firm i faces a demand function

$$D_i(p_i, p_{-i}) = \alpha_i - \beta_i p_i + \sum_{j \neq i} \gamma_{ij} p_j,$$

where $\alpha_i, \beta_i > 0$ and $\gamma_{ij} > 0$ for each $j \neq i$. Firm i can produce units at a constant marginal cost c_i . Thus firm i 's profit function is

$$\pi_i(p_i, p_{-i}) = (p_i - c_i) \left(\alpha_i - \beta_i p_i + \sum_{j \neq i} \gamma_{ij} p_j \right).$$

We can see that for $j \neq i$,

$$\frac{\partial \pi_i}{\partial p_i \partial p_j} = \gamma_{ij} > 0,$$

and so by Topkis' characterization theorem (Theorem 68), each π_i is supermodular in (p_i, p_{-i}) . It follows that π_i is supermodular in p_i and has increasing differences in both p_i and p_{-i} .

2.17.2 Submodular games

Submodular games capture the notion of *strategic complements*. In a game of strategic complements, each player's best response is decreasing in the actions of their opponents. A classic example is public goods games.

Definition 42 (Submodular game). A game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ is a *submodular game* if, for each player $i \in \mathcal{I}$,

- (i) the strategy set S_i is a complete lattice;
- (ii) the payoff function $u_i : S_i \times S_{-i} \rightarrow \mathbb{R}$ is order lower semicontinuous in s_i , order continuous in s_{-i} and lower bounded;
- (iii) u_i is submodular in s_i ;
- (iv) u_i has decreasing differences in s_i and s_{-i} .

If (iii) is replaced with the condition that u_i is strictly submodular in s_i and (iv) is replaced with the condition that u_i has strictly decreasing differences in s_i and s_{-i} , then we say that G is a *strictly submodular game*.

Clearly, this definition is symmetric to the definition of a supermodular game. However, the results for supermodular games do not in general extend to submodular games. This is because of a rather annoying asymmetry in lattice theory – the fixed point theorems for lattices only apply for functions and correspondences that are order-preserving, i.e. monotonically increasing.

In the case of *two-player* submodular games, however, we can apply all of the results for supermodular games, thanks to the following result:

Proposition 38. *For every submodular two-player game $G = (\{1, 2\}, (S_i, u_i)_{i=1,2})$, the game $\tilde{G} = (\{1, 2\}, (\tilde{S}_i, \tilde{u}_i)_{i=1,2})$ where $\tilde{S}_1 = S_1$, $\tilde{S}_2 = -S_2$ and $\tilde{u}_i(\tilde{s}) = u_i(\tilde{s}_1, -\tilde{s}_2)$ for both players i is a supermodular game.*

3 Games of incomplete information

3.1 Bayesian games

Definition 43 (Bayesian game). A *Bayesian game* (or a *game of incomplete information*) is a tuple $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$, where \mathcal{I} is a set of players, S_i is a set of

strategies for player i , Θ_i is a set of types θ_i for player i , $u_i : S \times \Theta \rightarrow \mathbb{R}$ is a von Neumann-Morgenstern expected payoff function for i (where $S = \times_i S_i$ and $\Theta = \times_i \Theta_i$), and p is a joint probability distribution over Θ .

Note that the payoff to a player i depends not only on own type θ_i but also on the type profile θ_{-i} of i 's opponents. If u_i is independent of θ_{-i} (i.e. $u_i(s, \theta_i, \theta_{-i}) = u_i(s, \theta_i, \theta'_{-i})$ for all $s \in S$, $\theta_i \in \Theta_i$, and $\theta_{-i}, \theta'_{-i} \in \Theta_{-i}$, then we say that the game has *private values*.)

Note also the assumption of a common prior p . This can be relaxed, although one does not gain much by relaxing this assumption. A more elaborate version of a Bayesian game is a tuple $(\mathcal{I}, \Omega, p, (S_i, \Theta_i, u_i, \tau_i)_{i \in \mathcal{I}})$ where Ω is a set of states, p is a common prior over Ω and the functions $\tau_i : \Omega \rightarrow \Theta_i$ map states of the world into types for each i . In this case, the assumption of a common prior is arguably wlog, since we can define the state space Ω appropriately to 'shift' differences in prior into the state space.

As Harsanyi (1968) notes, any game of incomplete information is equivalent to a game of complete but imperfect information with an additional player *Nature* that randomly chooses a type for each player according to the probability distribution p . In Harsanyi's formulation, the complete but imperfect equivalent game proceeds in three stages:

1. *Ex ante stage*. Players know only the probabilities with which Nature assigns those elements of the game that are not common knowledge (i.e. in Bayesian games, types).
2. *Interim stage*. Players learn their private information (in Bayesian games, their own type θ_i) and, on the basis of this information, choose strategies simultaneously.
3. *Ex post stage*. Given realized structure of the game (types) and strategies played, payoffs are realized.

Note that in Bayesian games, player i learns only θ_i and not θ_{-i} . We assume the set of type profiles Θ is common knowledge.

In the interim stage, players' posterior beliefs p_i are generated by *Bayesian updating*, i.e.

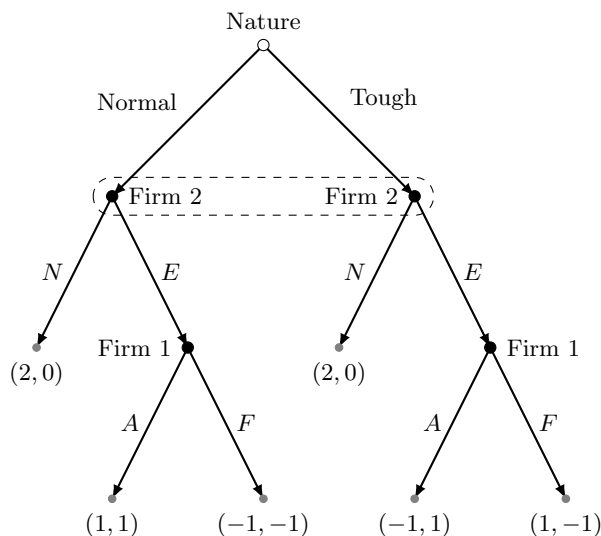
$$p_i(\theta_{-i} \mid \theta_i) = \frac{p(\theta_i, \theta_{-i})}{p(\theta_i)} = \frac{p(\theta_i, \theta_{-i})}{\int_{\Theta_{-i}} p(\theta_i, d\theta_{-i})}.$$

This requirement is called *consistency* of beliefs.

Example 27 (Market entry II). Suppose Firm 1 is an incumbent monopolist and Firm 2 is a potential entrant. Firm 2 can choose whether to enter (E) or not (N), and if choosing to enter, Firm 1 can choose whether to fight (F) or accommodate (A). With probability $1 - p$, the incumbent is "tough" and with probability p , the incumbent is "normal". If tough, Firm 1's payoff from fighting is greater than from accommodating, and the converse is true if normal:

| | "Normal" | | "Tough" | |
|-----|----------|--------|----------|-------|
| | N | E | N | E |
| A | 2, 0 | 1, 1 | A 2, 0 | -1, 1 |
| F | 2, 0 | -1, -1 | F 2, 0 | 1, -1 |

As a game of complete but imperfect information, we can envisage the game as follows:



Definition 44 (Bayes Nash equilibrium).

- (a) *Bayesian strategy.* A *Bayesian pure strategy* is a function $\sigma_i : \Theta_i \rightarrow S_i$. A *Bayesian mixed strategy* is a function $\sigma_i : \Theta_i \rightarrow \Delta(S_i)$. Let S_{Θ_i} denote the set of Bayesian pure strategies of player i and let $S_{\Theta} = \times_{i \in \mathcal{I}} S_{\Theta_i}$ denote the set of Bayesian strategy profiles.
- (b) *Expected payoffs.* Given a Bayesian strategy profile $\sigma \in S_{\Theta}$, the *ex ante expected payoff* $u_i(\sigma)$ of player i is defined by

$$u_i(\sigma) = \mathbb{E}_{\theta} u_i(\sigma, \theta) = \int_{\Theta} u_i(\sigma, \theta) p(d\theta),$$

and, further given type θ_i in the support of p , the *interim expected payoff* $u_i(\sigma | \theta_i)$ is defined by

$$u_i(\sigma | \theta_i) = \mathbb{E}_{\theta} [u_i(\sigma, \theta) | \theta_i] = \int_{\Theta_{-i}} u_i(\sigma, \theta_i, \theta_{-i}) p(d\theta_{-i} | \theta_i).$$

- (c) *Bayes Nash equilibrium.* A (pure strategy) *Bayes Nash equilibrium* of an n -player Bayesian game G_{Θ} is a strategy profile $\sigma^* = (\sigma_1^*, \dots, \sigma_n^*) \in S_{\Theta}$ such that

$$u_i(\sigma_i^*, \sigma_{-i}^*) \geq u_i(\sigma_i, \sigma_{-i}^*) \quad \text{for all } \sigma_i \in S_{\Theta_i} \text{ and } i \in \mathcal{I}.$$

Bayes Nash equilibrium is just the extension of Nash equilibrium to Bayesian games. Bayes Nash is sometimes just called *Bayesian equilibrium*, but we will later want to draw a distinction between Bayes Nash equilibria – the Bayesian counterpart of Nash

equilibria – and Bayes correlated equilibria – the Bayesian counterpart of correlated equilibria. Since Bayes Nash equilibria are extensions of Nash, their existence is just a corollary of Nash’s existence theorem (Theorem 8):

Corollary 7. *Any finite Bayesian game G_Θ has a Bayes Nash equilibrium.*

Proof. Immediate from Theorem 8. □

The definition we give above requires that the strategies in Bayes Nash equilibrium are *ex ante* optimal. This is equivalent to requiring that the strategies in Bayes Nash equilibrium are interim optimal:

Proposition 39. *A Bayesian strategy profile $\sigma^* \in S_\Theta$ is a (pure strategy) Bayes Nash equilibrium iff*

$$u_i(\sigma_i^*, \sigma_{-i}^* \mid \theta_i) \geq u_i(s_i, \sigma_{-i}^* \mid \theta_i) \quad \text{for all } s_i \in S_i, \theta_i \text{ s.t. } p(\theta_i) > 0, \text{ and } i \in \mathcal{I}.$$

Proof. Assume wlog that Θ is the support of p (otherwise we can simply redefine the type space.) Then

$$\begin{aligned} u_i(\sigma_i^*, \sigma_{-i}^*) &= \int_{\Theta} u_i(\sigma_i^*, \sigma_{-i}^*, \theta) p(d\theta) \\ &= \int_{\Theta_i \times \Theta_{-i}} u_i(\sigma_i^*, \sigma_{-i}^*, \theta_i, \theta_{-i}) p(d\theta_i, d\theta_{-i}) \\ &= \int_{\Theta_i} \left(\int_{\Theta_{-i}} u_i(\sigma_i^*, \sigma_{-i}^*, \theta_i, \theta_{-i}) p(d\theta_{-i} \mid \theta_i) \right) p(d\theta_i) \\ &= \int_{\Theta_i} u_i(\sigma_i^*, \sigma_{-i}^* \mid \theta_i) p(d\theta_i). \end{aligned}$$

□

More explicitly, if Θ is finite, σ^* is a Bayes Nash equilibrium if

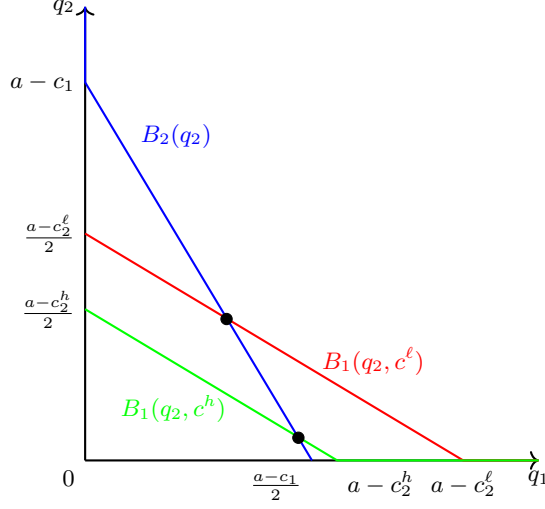
$$\sum_{\theta \in \Theta} u_i(\sigma_i^*(\theta_i), \sigma_{-i}^*(\theta_{-i}), \theta) p(\theta) \geq \sum_{\theta \in \Theta} u_i(s_i, \sigma_{-i}^*(\theta_{-i}), \theta) p(\theta)$$

for all $i \in \mathcal{I}$ and $s_i \in S_i$, or equivalently,

$$\sum_{\theta_{-i} \in \Theta_{-i}} u_i(\sigma_i^*(\theta_i), \sigma_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i}) p(\theta_{-i} \mid \theta_i) \geq \sum_{\theta_{-i} \in \Theta_{-i}} u_i(s_i, \sigma_{-i}^*(\theta_{-i}), \theta_i, \theta_{-i}) p(\theta_{-i} \mid \theta_i)$$

for all $i \in \mathcal{I}$, $s_i \in S_i$ and θ_i in the support of p .

Example 28 (Cournot duopoly with uncertain cost). Consider a Cournot duopoly with linear inverse demand function $P(Q) = \max\{a - Q, 0\}$ with $Q = q_1 + q_2$. Firm 1 has marginal cost c_1 , known to both firms, and Firm 2 has marginal cost $c_2 \in \{c_2^h, c_2^\ell\}$, with $c^h > c^\ell$. Firm 1 does not know Firm 2’s marginal cost, and believes it is c^h with probability p . We have type spaces $\Theta_1 = \{1\}$ and $\Theta_2 = \{h, \ell\}$.



A Bayes Nash equilibrium in this game is a strategy profile $(q_1^*, (q_2^*(h), q_2^*(\ell)))$ such that

$$q_2^*(\theta_2) = \arg \max_{q_2 \in \mathbb{R}_+} u_2(q_1^*, q_2, \theta_2) = \arg \max_{q_2 \in \mathbb{R}_+} (a - q_2 - q_1^* - c_2^{\theta_2})q_2,$$

for each $\theta_2 = h, \ell$, and

$$\begin{aligned} q_1^* &= \arg \max_{q_1 \in \mathbb{R}_+} \mathbb{E}_{\theta_2} u_1(q_1, q_2^*(\theta_2)) \\ &= \arg \max_{q_1 \in \mathbb{R}_+} p(a - q_2^*(h) - q_1 - c_1)q_1 + (1 - p)(a - q_2^*(\ell) - q_1 - c_1)q_1. \end{aligned}$$

We have best response correspondences

$$\begin{aligned} q_2^*(q_1, \theta_2) &= \max \left\{ \frac{a - q_1 - c_2^{\theta_2}}{2}, 0 \right\}, \\ q_1^*(q_2^*) &= \max \left\{ \frac{a - c_1 - pq_2^*(h) - (1 - p)q_2^*(\ell)}{2}, 0 \right\}. \end{aligned}$$

Assuming an interior solution, we obtain

$$\begin{aligned} q_2^*(h) &= \frac{a - 2c^h + c_1}{3} + \frac{(1 - p)(c_2^h - c_2^\ell)}{6}, \\ q_2^*(\ell) &= \frac{a - 2c_2^\ell + c_1}{3} - \frac{p(c_2^h - c_2^\ell)}{6}, \\ q_1^* &= \frac{a - 2c_1 + pc_2^h + (1 - p)c_2^\ell}{3}. \end{aligned}$$

We see that $q_2^*(\ell) > q_2^*(h) > q_2^{h*} = \frac{a - 2c_2^h + c_1}{3}$. Even in the high cost case, Firm 2 benefits from Firm 1's lack of information and so produces more than if its type were known to Firm 1.

3.2 Dominance in Bayesian games

Bayes Nash equilibrium requires very strict common knowledge assumptions. In particular, if we expect a particular Bayes Nash equilibrium to be played then we require common knowledge of the joint distribution of types of agents. Wilson (1987) famously critiques the dependence on such common knowledge assumptions. For example, in analysing trading processes, the design of institutions, and so on, analyses should not be too dependent on details.³⁰ Overdependence on details restricts the applicability of such analyses to the real world. While in some settings, it might be plausible that agents have sufficient knowledge about the distribution of types of other agents, often this will not be the case and the distribution of types might be subject to change.³¹

The notion of strict and weak dominance naturally extend to Bayesian games, and has some attractive properties in light of the Wilson doctrine. Of course, many Bayesian games lack a strictly or weakly dominant equilibrium. Nevertheless, when we come to mechanism design this will not turn out to be a problem.

Definition 45 (Dominance).

- (a) *Strict dominance.* In the Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$, a Bayesian strategy σ_i *strictly dominates* σ'_i if, for all $s_{-i} \in S_{-i}$ and all $\theta \in \Theta$,

$$u_i(\sigma_i(\theta), s_{-i}, \theta) > u_i(\sigma'_i(\theta_i), s_{-i}, \theta)$$

- (b) *Weak dominance.* In the Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$, a Bayesian strategy σ_i *weakly dominates* σ'_i if, for all $s_{-i} \in S_{-i}$ and all $\theta \in \Theta$,

$$u_i(\sigma_i(\theta), s_{-i}, \theta) \geq u_i(\sigma'_i(\theta_i), s_{-i}, \theta), \quad \text{with strict inequality for some } s_{-i} \in S_{-i}.$$

- (c) *Weakly dominant equilibrium.* A Bayesian strategy profile σ^* is a *weakly dominant equilibrium* of the Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$ if, for all $i \in \mathcal{I}$, $\theta \in \Theta$, $s_{-i} \in S_{-i}$ and any $\sigma'_i \in S_{\Theta_i}$,

$$u_i(\sigma_i^*(\theta_i), s_{-i}, \theta) \geq u_i(\sigma'_i(\theta_i), s_{-i}, \theta), \quad \text{with strict inequality for some } s_{-i} \in S_{-i}.$$

That is, σ^* is a weakly dominant equilibrium if, for all $i \in \mathcal{I}$, σ_i^* is a weakly dominant strategy.

Proposition 40. *In any Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$, any weakly dominant equilibrium is a Bayes Nash equilibrium.*

Proof. Assume G_Θ has a weakly dominant equilibrium σ^* . Fix any $\theta_i \in \Theta_i$ in the support of p and any $s_i \in S_i$. By definition of weakly dominant equilibrium, we have

$$u_i(\sigma_i^*(\theta_i), s_{-i}, \theta_i, \theta_{-i}) \geq u_i(s_i, s_{-i}, \theta_i, \theta_{-i})$$

³⁰This is the Wilson doctrine (no, not the British one).

³¹Repeated, anonymous interactions are a good example where we might expect that agents could know the joint distribution of types.

for all $\theta_{-i} \in \Theta_{-i}$ and $s_{-i} \in S_{-i}$. Hence

$$u_i(\sigma_i^*(\theta_i), \sigma_{-i}(\theta_{-i}), \theta_i, \theta_{-i}) \geq u_i(s_i, \sigma_{-i}(\theta_{-i}), \theta_i, \theta_{-i})$$

for all $\theta_{-i} \in \Theta_{-i}$ and $\sigma_{-i} \in S_{\Theta_{-i}}$. It follows that

$$\int_{\Theta_{-i}} u_i(\sigma_i^*(\theta_i), \sigma_{-i}(\theta_{-i}), \theta_i, \theta_{-i}) p(d\theta_{-i} | \theta_i) \geq \int_{\Theta_{-i}} u_i(s_i, \sigma_{-i}(\theta_{-i}), \theta_i, \theta_{-i}) p(d\theta_{-i} | \theta_i),$$

for all $\sigma_{-i} \in S_{\Theta_{-i}}$, which is to say,

$$u_i(\sigma_i^*, \sigma_{-i} | \theta_i) \geq u_i(s_i, \sigma_{-i} | \theta_i)$$

for all σ_{-i} , and in particular,

$$u_i(\sigma_i^*, \sigma_{-i}^* | \theta_i) \geq u_i(s_i, \sigma_{-i}^* | \theta_i).$$

Since this holds for all $s_i \in S_i$ all θ_i in the support of p , and all $i \in \mathcal{I}$, it follows that σ^* is a Bayes Nash equilibrium. \square

3.3 Ex post equilibrium

If we buy the Wilson doctrine, then we want to be looking at equilibrium concepts that are “robust” in the sense that they do not depend on the distribution of types. Ex post equilibrium is robust in this sense. Note we can associate each type profile realization in a Bayesian game with a state of the world. Imagine a setting as follows. Each player i chooses a strategy $\sigma_i(\theta_i)$ that depends only on their own type. Nature then realizes a state and agents learn the state, so they know the realized profile of types. In an ex post equilibrium, no rational agent would want to change their strategy upon learning the state.

Note that given a Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}}, p)$, each type profile $\theta \in \Theta$ determines a normal form game $G_\theta = (\mathcal{I}, (S_i, u_i(\cdot, \theta))_{i \in \mathcal{I}})$.

Definition 46 (Ex post equilibrium). In the Bayesian game $G_\Theta = (\mathcal{I}, (S_i, \Theta_i, u_i)_{i \in \mathcal{I}})$, a strategy profile σ^* is an *ex post equilibrium* if

$$u_i(\sigma_i^*(\theta_i), \sigma_{-i}^*(\theta_{-i}), \theta) \geq u_i(s_i, \sigma_{-i}^*(\theta_{-i}), \theta)$$

for every player $i \in \mathcal{I}$, $s_i \in S_i$, and $\theta \in \Theta$.

That is, an ex post equilibrium σ^* induces a Nash equilibrium in all the normal form games G_θ . Trivially:

Proposition 41. *If σ^* is a weakly dominant equilibrium of a Bayesian game G_Θ then is an ex post equilibrium of G_Θ .*

Proof. Immediately we have that for any $s_{-i} = \sigma_{-i}^*(\theta_{-i})$, $\sigma_i^*(\theta_i)$ is a best response. Hence $\sigma_i^*(\theta_i)$ is a best response to σ_{-i}^* in every normal form game G_θ . \square

The converse is not generally true, though it is true if each player’s payoff depends only on their own type.

3.4 Purification

In games of complete information, mixed strategies are sometimes controversial, because of skepticism that humans in practice actually randomize because they lack the psychological power to truly randomize actions. The skepticism does not make much sense, because people can actually play mixed strategies quite easily. Hunters in subsistence societies often randomize their choice of hunt location, and in many sports, we have good evidence that players play mixed strategies – see Walker & Wooders (2001) for evidence from tennis and Chiappori, Levitt & Groseclose (2002) for evidence from penalty kicks in association football.³²

A more compelling problem with any mixed strategy Nash equilibrium σ is that since each player i is indifferent between any of the pure strategies s_i in the support $\text{supp}(\sigma_i)$ of σ_i , then it is not obvious why i should mix between these pure strategies in such a way that each other player j is indifferent between the pure strategies in the support of σ_j .

Harsanyi (1973) provides a justification for (most) mixed strategy Nash equilibria in finite games. The idea is that if we consider the pure strategy Nash equilibria of a sequence of “payoff perturbed” games in which players have a modicum of uncertainty about their opponents’ payoffs, and the limit game of this sequence is a game of complete information, then almost all of the time, we get the mixed strategy equilibria of this complete information game as the limit of the pure strategy equilibria of the payoff perturbed games. Intuitively, imagine that a player’s payoffs are known to players but that player has some private inclination towards particular actions, independent of the payoffs and unknown to the other players. Then that player’s behaviour will appear as though she is randomizing her actions, even though she is not randomizing – actually she almost always choosing a strict best response.

The proof of this result is unfortunately rather technical and quite long, so I’d advise skipping through to Example 29 if you are not in the mood for a couple of pages of new equilibrium concepts, Jacobians, and a sprinkling of algebraic topology. Believe it or not, the proof described here is actually much shorter and more elegant than Harsanyi (1973). First, we need to introduce the notion of regular equilibrium, alongside some other equilibrium notions. Consider a mixed strategy finite normal form game $G^m = (\mathcal{I}, (\Delta(S_i), u_i))$. Let $X_i = \mathbb{R}^{|S_i|}$ so that $\Delta(S_i) \subset X_i$ for each $i \in \mathcal{I}$ (taking the elements of $\Delta(S_i)$ to be probability vectors), and let $X = \times_{i \in \mathcal{I}} X_i$. Clearly, the elements $\sigma_i \in X_i$ do not need to be probability vectors, unlike $\sigma_i \in \Delta(S_i)$. We extend u_i to X by

$$u_i(\sigma) := \sum_{s \in S} \sigma(s) u_i(s) \quad \text{for all } \sigma \in X.$$

This agrees with the interpretation of u_i as expected utility on $\Delta(S) \supset \times_{i=1}^n \Delta(S_i)$.

Now, a mixed strategy profile $\sigma \in \times_{i=1}^n \Delta(S_i)$ is a Nash equilibrium iff for every pure strategy $s_i \in S_i$ such that $\sigma_i(s_i) > 0$, we have that $u_i(s_i, \sigma_{-i}) = \max_{s'_i \in S_i} u_i(s'_i, \sigma_{-i})$, for

³²I remember hearing/reading somewhere about using the direction of bone spitting out from a fire as a randomization device for hunting in a hunter-gatherer society, but I can’t find a reference for this anywhere.

all players $i \in \mathcal{I}$. Thus σ is a Nash equilibrium iff it is a solution to the set of $|S| + |\mathcal{I}|$ equations:

$$\begin{aligned} \sigma_i(s_i)[u_i(s_i, \sigma_{-i}) - \max_{s'_i \in S_i} u_i(s'_i, \sigma_{-i})] &= 0 & \text{for all } i \in \mathcal{I} \text{ and } s_i \in S_i, \\ \sum_{s_i \in S_i} \sigma_i(s_i) - 1 &= 0 & \text{for all } i \in \mathcal{I}. \end{aligned}$$

This system defines a mapping, but unfortunately the mapping isn't necessarily differentiable. Instead, fix a pure strategy profile $\bar{s} \in S$. Consider the following system:

$$\begin{aligned} \sigma_i(s_i)[u_i(s_i, \sigma_{-i}) - u_i(\bar{s}_i, \sigma_{-i})] &= 0 & \text{for all } i \in \mathcal{I} \text{ and } s_i \in S_i - \{\bar{s}_i\}, \\ \sum_{s_i \in S_i} \sigma_i(s_i) - 1 &= 0 & \text{for all } i \in \mathcal{I}. \end{aligned}$$

If $\hat{\sigma}$ is an equilibrium such that $\hat{\sigma}(\bar{s}) > 0$ (i.e. $\hat{\sigma}_i(\bar{s}_i) > 0$ for all players i), then $\hat{\sigma}$ is a solution to this system. Define the mapping

$$\begin{aligned} F_i^{s_i}(\sigma \mid \bar{s}) &= \sigma_i(s_i)[u_i(s_i, \sigma_{-i}) - u_i(\bar{s}_i, \sigma_{-i})] & \text{for } s_i \in S_i - \{\bar{s}_i\}, \\ F_i^{\bar{s}_i}(\sigma \mid \bar{s}) &= \sum_{s_i \in S_i} \sigma_i(s_i) - 1 \end{aligned}$$

for $i \in \mathcal{I}$. This mapping is infinitely differentiable. We denote the Jacobian of $F(\cdot \mid \bar{s})$ evaluated at $\tilde{\sigma}$ by

$$J(\tilde{\sigma} \mid \bar{s}) := \frac{\partial F(\sigma \mid \bar{s})}{\partial \sigma}(\tilde{\sigma}).$$

Definition 47 (Regular equilibrium). Let $G^m = (\mathcal{I}, (\Delta(S_i), u_i))$ be a mixed strategy finite normal form game.

- (a) *Regular equilibrium*. A Nash equilibrium σ of G^m is called a *regular equilibrium* if $J(\sigma \mid \bar{s})$ is nonsingular for some $\bar{s} \in S$ such that $\sigma(\bar{s}) > 0$. We say that a Nash equilibrium σ is *irregular* if it is not regular.
- (b) *Quasi-strict equilibrium*. A Nash equilibrium σ of G^m is called a *quasi-strict equilibrium* if $\text{supp } \sigma = \tilde{B}(\sigma)$ where $\tilde{B} : \times_{i=1}^n \Delta(S_i) \rightrightarrows S$ is the best response correspondence mapping each mixed strategy profiles to the set of pure best response strategy profiles. That is, a Nash equilibrium σ is quasi-strict if every pure strategy best response for player i to σ is in the support of σ_i for every player i (i.e. σ_i is totally mixed over $\tilde{B}_i(\sigma)$).

Not all games have regular equilibria, but it turns out that almost all normal form games do:

Theorem 15. *For almost all finite mixed strategy normal form games G^m , all of the Nash equilibria of G^m are regular.*

Proof. See the (long) proof of Theorem 2.6.1 in Van Damme (1991).³³ Be warned that his notation is confusing (he uses s for the mixed strategy profile and φ for pure strategies – one day, economists might learn the benefits of standardization...) \square

We will follow the setting of Govindan, Reny & Robson (2003), who give a shorter proof of a generalization of Harsanyi's purification theorem.

Definition 48 (Payoff-perturbed game).

- (a) *Payoff-perturbed game.* Given a (finite) n -player game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, the corresponding *payoff-perturbed* game of incomplete information is $\tilde{G} = (\mathcal{I}, (S_i, \mathbb{R}^{|\mathcal{S}|}, u_i)_{i \in \mathcal{I}}, \nu)$ where
 - (i) $\nu = \nu_1 \times \cdots \times \nu_n$ is a product probability measure on $(\mathbb{R}^{|\mathcal{S}|}, \mathcal{B})^n$;
 - (ii) player i 's type $\eta_i \in \mathbb{R}^{|\mathcal{S}|}$ is a *payoff perturbation*, a random vector generated according to ν_i (necessarily, $\{\eta_i\}_{i \in \mathcal{I}}$ are independent random variables since ν is a product measure). Player i 's type is private information;
 - (iii) the *perturbed payoff* to player i is $u_i(s) + \eta_i(s)$ for all strategy profiles $s \in S$.
- (b) *Essentially strict equilibrium.* A Bayes Nash equilibrium σ of the payoff-perturbed game \tilde{G} is called an *essentially strict equilibrium* if for each player, $\sigma_i(\eta_i) \in S_i$ is a strict best response almost everywhere with respect to ν_i .

Some technical assumptions are needed for the proof. Assume ν_i is absolutely continuous with respect to the Lebesgue measure on $\mathbb{R}^{|\mathcal{S}|}$.³⁴ This ensures that the set of payoff perturbations η_i for which i is indifferent between strategies s_i and s'_i for any pair $s_i \neq s'_i$ is measure zero, and so i has a unique best response for almost every payoff perturbation η_i and opponents' strategy profile σ_{-i} .

Theorem 16 (Harsanyi's purification theorem; Govindan-Reny-Robson, 2003). *Consider a finite normal form game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$ and consider a family of corresponding perturbed games $\{\tilde{G}^k\}_{k \in \mathbb{N}}$, $\tilde{G}^k = (\mathcal{I}, (S_i, \mathbb{R}^{|\mathcal{S}|}, u_i)_{i \in \mathcal{I}}, \nu^k)$ such that each ν_i^k converges to a mass point at $0 \in \mathbb{R}^{|\mathcal{S}|}$ for each player $i \in \mathcal{I}$.*

If σ^ is a regular equilibrium of G , then for all $\epsilon > 0$ and sufficiently large k , there is an essentially strict Bayes Nash equilibrium $\tilde{\sigma}^k$ of \tilde{G}^k such that the distribution on S induced by $\tilde{\sigma}^k$ is within ϵ of σ^* .*

Proof. The proof relies on an algebraic topology result that we will not prove here:

Lemma 9. *Suppose $U \subset \mathbb{R}^m$ is a bounded open set and $f, g : \bar{U} \rightarrow \mathbb{R}^m$ are continuous functions on the closure of U . Suppose also that f is continuously differentiable on U , that f has a unique zero x_0 in U , and that the determinant of the Jacobian of f at x_0 is nonzero. If for all $\lambda \in [0, 1]$, $\lambda g + (1 - \lambda)f$ has no zero on the boundary of U then g has a zero in U .*

³³Specifically, Van Damme (1991), *Stability and Perfection of Nash Equilibria* (Second Ed.)

³⁴Actually, we need only the weaker assumption that for any given hyperplane $H(\sigma_{-i})$ with normal σ_{-i} in $\mathbb{R}^{|\mathcal{S}-i|}$, we have that $\nu_i\{\eta_i(s_i, \cdot) - \eta_i(s'_i, \cdot) \in H(\sigma_{-i})\} = 0$.

Let $m_i := |S_i|$ and $m := \sum_{i=1}^n m_i$. For each player i , order the pure strategies $s_i \in S_i$ so that $S_i = \{s_{i1}, \dots, s_{iK}\}$. As with when we defined regular equilibrium, for each $s_i \in S_i$ we extend $u_i(s_i, \cdot)$ to $\times_{j \neq i} \mathbb{R}^{|S_j|}$ so that $u_i(s_i, x_{-i}) = \sum_{s_{-i}} x_{-i}(s_{-i}) u_i(s_i, s_{-i})$. Write $B_i(\sigma \mid \eta_i)$ for i 's best response correspondence when i is of type η_i , and extend this to $B_i : \mathbb{R}^m \rightarrow \Delta(S_i)$. Let $b_i(\sigma \mid \eta_i)$ be a selection from $B_i(\sigma \mid \eta_i)$. Then $g_i^k(\sigma) = \int b_i(\sigma \mid \eta_i) d\nu_i^k(\eta_i)$ is a continuous function $g_i^k : \mathbb{R}^m \rightarrow \Delta(S_i)$ and $g = (g_1, \dots, g_n) : \mathbb{R}^m \rightarrow \times_{i=1}^n \Delta(S_i)$ is continuous. Note the fixed points of g are the distributions on S induced by the Bayes Nash equilibria of \tilde{G}^k , and for all $\sigma \in \times_{i=1}^n \Delta(S_i)$, we have that $b_i(\sigma \mid \eta_i)$ is unique almost everywhere with respect to ν_i . Hence all Bayes Nash equilibria of \tilde{G}^k are essentially strict.

Fix a pure strategy profile \bar{s} in the support of σ^* and define the mapping F as in the definition of regular equilibrium above. Since σ^* is a regular equilibrium, $\det J(\sigma \mid \bar{s}) \neq 0$. Applying the implicit function theorem at σ^* , we can choose an open set $U \subset \mathbb{R}^m$ s.t.

- (i) $\sigma^* \in U$,
- (ii) $\|\sigma^* - \sigma\| < \epsilon$ for all $\sigma \in U$,
- (iii) $\bar{U} \cap F^{-1}(0) = \{\sigma^*\}$,
- (iv) if $s_i \in \text{supp } \sigma_i^*$ then $s_i \in \text{supp } \sigma_i$ for all $\sigma \in \bar{U}$, and
- (v) if $u_i(s_i, \sigma_{-i}^*) - u_i(\bar{s}_i, \sigma_{-i}^*) < 0$ then $u_i(s_i, \sigma_{-i}) - u_i(\bar{s}_i, \sigma_{-i}) < 0$ for all $s_i \in S_i$, all players i , and all $\sigma \in \bar{U}$.

Now define the continuous function $h^k : \bar{U} \rightarrow \mathbb{R}^m$ by

$$\begin{aligned} h_{i,s_i}^k(\sigma \mid \bar{s}) &= g_i^k(\sigma)(s_i) - \sigma_i(s_i) \quad \text{for } s_i \in S_i - \{\bar{s}_i\}, \\ h_{i,\bar{s}_i}^k(\sigma \mid \bar{s}) &= \sum_{s_i \in S_i} \sigma_i(s_i) - 1, \end{aligned}$$

for all players $i \in \mathcal{I}$. It is sufficient to show that h^k has a zero in U for all sufficiently large k since the equilibria of \tilde{G}^k are essentially strict and all $\sigma \in U$ lie within an Euclidean distance ϵ of σ^* (since then g^k has a fixed point in U). Now, to apply Lemma 9, we need only establish the following:

Lemma 10. *For sufficiently large n and all $\lambda \in [0, 1]$, $\lambda h^k + (1 - \lambda)F$ has no zero on the boundary ∂U of U .*

Proof. Towards contradiction, assume there is some sequence $\{\lambda^k\}$, $\lambda^k \in [0, 1]$ with $\lambda^k \rightarrow \hat{\lambda}$ and a sequence $\{\sigma^k\}$, $\sigma^k \in \partial U$ with $\sigma^k \rightarrow \hat{\sigma} \in \partial U$ s.t. each σ^k is a zero of $\lambda^k h^k + (1 - \lambda^k)F$.

Now, if $u_i(s_i, \hat{\sigma}_{-i}) - u_i(\bar{s}_i, \hat{\sigma}_{-i}) < 0$ then $g_i^k(\sigma^k) \rightarrow 0$. Thus for all i and $s_i \in S_i$, we have

$$\begin{aligned} 0 &= \lim_{k \rightarrow \infty} \lambda^k h_{i,s_i}^k(\sigma^k \mid \bar{s}) + (1 - \lambda^k) F_i^{s_i}(\sigma^k \mid \bar{s}) = \hat{\lambda}(-\hat{\sigma}_i(s_i)) + (1 - \hat{\lambda}) F_i^{s_i}(\hat{\sigma} \mid \bar{s}) \\ &= -\hat{\sigma}_i(s_i) [\hat{\lambda} + (1 - \hat{\lambda})(u_i(\bar{s}_i, \hat{\sigma}_{-i}) - u_i(s_i, \hat{\sigma}_{-i}))], \end{aligned}$$

so for every player i , $\hat{\sigma}(s_i) = 0$ for all s_i such that $u_i(s_i, \hat{\sigma}_{-i}) - u_i(\bar{s}_i, \hat{\sigma}_i) < 0$.

For $\sigma \in \mathbb{R}^m$, let $B_i(\sigma) = \arg \max_{s_i \in S_i} u_i(s_i, \sigma_{-i})$. If $\hat{\sigma}_i(s_i) < 0$ then by (iv) we have $\sigma_i^*(s_i) = 0$. Since σ^* is quasi-strict, if $s_i \notin B_i(\sigma^*)$ then (v) implies $u_i(s_i, \hat{\sigma}_{-i}) - u_i(\bar{s}_i, \hat{\sigma}_i) > 0$ so $\hat{\sigma}_i(\bar{s}_i) = 0$, yielding a contradiction. Hence $\hat{\sigma}_i(s_i) \geq 0$ for all $s_i \in S_i$ and all i . Since $h_i^k(\cdot | \bar{s}) = F_i(\cdot | \bar{s})$, we have $F_i(\sigma^k | \bar{s}) = \lambda^k h_i^k(\sigma^k | \bar{s}) + (1 - \lambda^k) F_i(\sigma^k | \bar{s}) = 0$ for all i and k . Thus $\hat{\sigma} \in \times_{i=1}^n \Delta(S_i)$. Moreover, since $\hat{\sigma}$ is not a zero of F by (iii), $F_i(\hat{\sigma} | \bar{s}) \neq 0$ for some player i and some pure strategy profile $\bar{s} \in S$. Fix i and \bar{s} .

Now, we have that $u_i(\bar{s}_i, \hat{\sigma}_{-i}) - u_i(\bar{s}_i, \hat{\sigma}_i) > 0$, so $\bar{s} \notin B_i(\hat{\sigma})$. Moreover, $\hat{\sigma}_i(s_i) > 0$ for all $s_i \in B_i(\hat{\sigma})$, since (v) implies $B_i(\hat{\sigma}) \subseteq B_i(\sigma^*)$, σ^* is a quasi-strict equilibrium, and by (iv). Thus $\sum_{s_i \in B_i(\hat{\sigma})} F_i^{s_i}(\sigma^k | \bar{s}) > 0$ for all sufficiently large k . For $s_i \notin B_i(\hat{\sigma})$, we have from the definition of g_i^k that $g_i^k(\sigma^k)(s_i) \rightarrow 0$ so $\sum_{s_i \in B_i(\hat{\sigma})} g_i^k(\sigma^k)(s_i) \rightarrow 1$ as $n \rightarrow \infty$. Yet $\sum_{s_i \in B_i(\hat{\sigma})} \hat{\sigma}_i(s_i) < 1$, given $\sigma_i(s_i^*) > 0$ by (iv) and $\hat{\sigma} \in \times_{i=1}^n \Delta(S_i)$. Thus

$$\sum_{s_i \in B_i(\hat{\sigma})} h_{i,s_i}^k(\sigma^k | \bar{s}) > 0$$

for all sufficiently large k . Yet $\sum_{s_i \in B_i(\hat{\sigma})} F_i^{s_i}(\sigma^n | \bar{s}) > 0$ and $\sum_{s_i \in B_i(\hat{\sigma})} h_{i,s_i}^k(\sigma^k | \bar{s}) > 0$ together contradict that σ^k is a zero of $\lambda^k h^k + (1 - \lambda^k) F$ for all k . \square

\square

Theorem 16 was a lot of work to prove but has a fairly good payoff: we have a fairly plausible justification of mixed strategy Nash equilibria!

Example 29 (Harsanyi purification in a coordination game). Consider the following coordination game G :

$$\begin{array}{cc} & \begin{array}{cc} a_2 & b_2 \end{array} \\ \begin{array}{c} a_1 \\ b_1 \end{array} & \begin{array}{cc} 2, 2 & 0, 0 \\ 0, 0 & 3, 3 \end{array} \end{array}$$

This has two pure strategy Nash equilibria, (a_1, a_2) and (b_1, b_2) , and a mixed strategy equilibrium $((\frac{3}{5}, \frac{2}{5}), (\frac{3}{5}, \frac{2}{5}))$.

Now consider the corresponding payoff-perturbed games:

$$\begin{array}{cc} & \begin{array}{cc} a_2 & b_2 \end{array} \\ \begin{array}{c} a_1 \\ b_1 \end{array} & \begin{array}{cc} 2 + \eta_{1a}^k, 2 + \eta_{2a}^k & \eta_{1a}^k, \eta_{2b}^k \\ \eta_{1b}^k, \eta_{2a}^k & 3 + \eta_{1b}^k, 3 + \eta_{2b}^k \end{array} \end{array}$$

where k is common knowledge and the (privately known) payoff shocks $(\eta_{ia}^k, \eta_{ib}^k)$ are identically and independently distributed across players, and $\eta_i^k := \eta_{ia}^k - \eta_{ib}^k$ is a random variable uniformly distributed on $[-1/k, 1/k]$. Let F^k be the cdf of $U[-1/k, 1/k]$.

In \tilde{G}^k , suppose both players play pure strategies $\tilde{s}_i^k : \mathbb{R} \rightarrow \{a_i, b_i\}$ that each take the form of a threshold strategy:

$$\tilde{s}_i^k(\eta_i) = \begin{cases} a_i & \text{if } \eta_i \geq \bar{\eta}_i^k, \\ b_i & \text{if } \eta_i < \bar{\eta}_i^k, \end{cases}$$

for some threshold $\bar{\eta}_i^k$.

Given the threshold $\bar{\eta}_i^k$, player i chooses b_i with probability $\pi_i^k = F(\bar{\eta}_i^k)$, and a_i with complementary probability, that is,

$$\tilde{s}_i^k(\eta_i) = \begin{cases} a_i & \text{if } \eta_i \geq (F^k)^{-1}(\pi_i^k), \\ b_i & \text{if } \eta_i < (F^k)^{-1}(\pi_i^k). \end{cases}$$

Now, a_1 is optimal for Player 1 given Player 2 plays a_2 with probability $1 - \pi_2^k$ under \tilde{s}_2^k iff

$$\tilde{u}_1^k(a_1, \tilde{s}_2^k \mid \eta_1) - \tilde{u}_1^k(b_1, \tilde{s}_2^k \mid \eta_1) = 2(1 - \pi_2^k) + \eta_1^k - 3\pi_2^k \geq 0.$$

Thus, for \tilde{s}_1^k to be a best response to \tilde{s}_2^k , the threshold $\bar{\eta}_1^k = F^{-1}(\pi_1^k)$ must satisfy

$$(F^k)^{-1}(\pi_1^k) = 5\pi_2^k - 2,$$

and since the game is symmetric, if Player 2 is playing a best response, her threshold $\bar{\eta}_2^k = (F^k)^{-1}(\pi_2^k)$ must satisfy

$$(F^k)^{-1}(\pi_2^k) = 5\pi_1^k - 2.$$

Since the strategies are symmetric and η_1 and η_2 are identically distributed, we have that $\pi^k := \pi_1^k = \pi_2^k$ satisfies

$$(F^k)^{-1}(\pi^k) = 5\pi^k - 2,$$

or $\pi^k = F^k(5\pi^k - 2) = \frac{5k\pi^k - 2k + 1}{2}$, which gives us $\pi^k = \frac{2k-1}{5k-2}$. As $k \rightarrow \infty$, we see that $\pi^k \rightarrow \frac{2}{5}$, which is precisely the probability that i plays b_i in the mixed strategy Nash equilibrium of the original unperturbed game G .

3.5 Beliefs, type spaces and the universal type space

So far, we have been working with types happily enough. But introducing types was not, prior to their introduction by Harsanyi (1967), an ex ante obvious modelling device for situations of strategic interaction under uncertainty. It is worth spelling out the original problem before Harsanyi came along.³⁵

Savage's (1954) subjective expected utility theory had already provided a coherent and closed account of individual decisionmaking under uncertainty. Say an agent chooses from among a set of alternatives A . The value of an alternative $a \in A$ to the agent depends on the state of the world ω , which lies in a (measurable) state space Ω . The agent has a belief μ over the states of the world, which is a probability measure on Ω , and the value of the possible alternatives in each state are summarized by a utility function $u : A \times \Omega \rightarrow \mathbb{R}$. Then the agent's preferences given their belief μ are represented by the subjective expected utility function $u : A \rightarrow \mathbb{R}$ defined by

$$u(a) = \int_{\Omega} u(a, \omega) \mu(d\omega).$$

³⁵This is worth doing not just for historical value, but because we often take for granted many of the foundational assumptions and modelling practices that underlie economic theory. Thinking about these deeply is often worthwhile.

Since the agent is rational, she simply chooses the action that maximizes $u(a)$.

This works fine when agents are taking decisions in isolation. But in the real-world, agents are burdened with the misfortune of encountering other people. In that case, agents have preferences over the actions of other people as well as their own. Suppose there are agents $i \in \{1, \dots, n\}$. Each agent i has a set of possible actions A_i , but now cares about the profile of actions $a = (a_1, \dots, a_n) \in A := \times_{i=1}^n A_i$. Again, the value $u_i(a, \omega)$ to i of an action profile a depends on the state of the world $\omega \in \Omega$, and agent i has belief μ_i over Ω . Agent's preferences again take the form $\int_{\Omega} u_i(a, \omega) \mu_i(d\omega)$. Rational decisionmaking in this context is more subtle. Given a_{-i} , agent i should choose the action $a_i \in A_i$ that maximizes her expected payoff given ϕ_i . But how does she know which profile of actions a_{-i} her opponents will choose? Well if everyone is rational (and this is common knowledge), then to predict how each of her opponents j will play she needs to worry about their beliefs ϕ_j . But this requires i to have beliefs about her opponents' beliefs. Likewise, each of her opponents will need to have beliefs about their opponents' beliefs (including i 's), and so i needs to consider what others believe about her belief ϕ_i , and so on for higher order beliefs – we get an infinite regress.

Early game theorists struggled with the question of how to close this model. Modelling the hierarchy of beliefs directly poses problems. Not only is a hierarchy of belief infinite-dimensional (we don't usually want to be worrying about 1000th order beliefs), but it does not directly complete the model. For example, say we have several agents and two possible states of the world $\omega \in \{0, 1\}$. Suppose each agent is certain that $\omega = 1$ but believes that every other agent j is certain that $\omega = 0$. Then each agent i is best responding to “phantom” opponents, because i 's beliefs about how j will behave need not accord with how j actually behaves. Because there is just one “version” of each player the model is not closed.

Harsanyi (1967) instead proposed completing the model by introducing types. This avoids the need to explicitly model hierarchies of belief. Rather, each type corresponds to a different “version” of the player, capturing uncertainty over the preferences, beliefs, action sets, information, or other characteristics of the other players.

Definition 49 (Type space). A *Harsanyi type space* (also called an *information structure*) is a tuple $(\Omega, \{\Theta_i, \pi_i\}_{i=1}^n)$, where

- (i) Ω is a space of fundamentals, i.e. a state space capturing the underlying realm of uncertainty;
- (ii) Θ_i is a measurable set of possible types θ_i of player i ;
- (iii) each player i has belief $\pi_i : \Theta_i \rightarrow \Delta(\Omega \times \Theta_{-i})$.

Type spaces complete the model, because by construction, players only assign positive probability over types for the other players, which are internal to the model. There is no infinite regress problem with beliefs because each type pins down a belief hierarchy. Agent i 's belief $\pi_i(\theta_i)$ if type θ_i lives in $\Delta(\Omega \times \Theta_{-i})$, and so i 's first-order belief $\mu_i^1(\theta_i)$ is simply the marginal on Ω , i.e. $\mu_i^1(E \mid \theta_i) = \pi_i(E \times \Theta_{-i} \mid \theta_i)$ for all Borel sets E of Ω .

Agent i 's second-order belief $\mu_i^2(\theta_i)$ lives in $\Delta(\Omega \times (\Delta(\Omega))^{n-1})$, and is given by

$$\mu_i^2(E \mid \theta_i) = \pi_i(\{(\omega, \theta_{-i}) \mid (\omega, \mu_{-i}^1(\theta_{-i})) \in E\} \mid \theta_i)$$

for all Borel sets E in $\Omega \times (\Delta(\Omega))^{n-1}$. Higher order beliefs can all be pinned down similarly.

Example 30 (A very simple type space). Take $n = 2$, $\Omega = \{\omega\}$, take $\Theta_1 = \{\theta_1\}$ and $\Theta_2 = \{\theta_2, \theta'_2\}$. Let p be a common prior over Ω , with $p(\omega) = \bar{p}$ and $p(\omega') = 1 - \bar{p}$. Suppose agent 2's type is θ_2^L when the state is ω and θ_2^H when the state is ω' , and this is known to both agents. Then the beliefs for agent 1 are

$$\begin{aligned}\pi_1(\{\omega\} \times \{\theta_2\} \mid \theta_1) &= \bar{p}, \\ \pi_1(\{\omega\} \times \{\theta'_2\} \mid \theta_1) &= 0,\end{aligned}$$

Likewise, the beliefs for agent 2 are

$$\begin{aligned}\pi_2(\{\omega\} \times \{\theta_1\} \mid \theta_2) &= 1, \\ \pi_2(\{\omega'\} \times \{\theta_1\} \mid \theta_2) &= 0, \\ \pi_2(\{\omega\} \times \{\theta_1\} \mid \theta'_2) &= 0, \text{ and} \\ \pi_2(\{\omega'\} \times \{\theta_1\} \mid \theta'_2) &= 1.\end{aligned}$$

As we've seen, Harsanyi type spaces form the foundation of Bayesian games. Just as subjective expected utility theory provides a coherent account of decision-making under uncertainty for isolated agents, the theory of Bayesian games is a coherent account of strategic interactions under uncertainty. But there are interesting questions to ask about Harsanyi type spaces that should shape how we approach modelling. How general are these spaces? Do they limit what kinds of belief hierarchy agents can have? Given predicted behaviour in a Bayesian game depends on the type space we choose, which predictions are consistent with some type space? And, since unions of type spaces are type spaces, is there some "largest" type space, that contains all the others?

Mertens & Zamir (1985) answer some of these questions, but Brandenburger & Dekel (1993) generalize the setting slightly and make the same points in a more accessible (but nevertheless technical) way. They construct a *universal type space* that contains every type space. The construction is as follows. As with the type space, we interpret Ω as a space of fundamentals, common to all players, that captures the underlying realm of uncertainty. This could be a space of payoffs, for example.

Definition 50 (Space of belief hierarchies). Let Ω denote a complete, separable metric space, and fix the number of players n . Define the spaces:

$$\begin{aligned}X_0 &:= \Omega, \\ X_m &:= X_{m-1} \times (\Delta(X_{m-1}))^{n-1},\end{aligned}$$

for each $m \in \mathbb{N}$. Define $\Theta_0 := \times_{k=1}^{\infty} \Delta(X_k)$. For each player i , a *belief hierarchy* is a vector $\theta_i = (\mu_1^i, \mu_2^i, \dots) \in \Theta_0$, and we call Θ_0 the *space of belief hierarchies*.

Each $\Delta(X_{k-1})$ is the space of k th order beliefs. For example, $\Delta(X_0)$ is the space of beliefs over opponents' fundamentals – first order beliefs – $\Delta(X_1)$ is the space of beliefs about opponents' first order beliefs, and so on. A belief hierarchy for each player i is simply the full hierarchy of i 's beliefs, an infinite-dimensional object. Brandenburger & Dekel (1993) define everything for the two-player case, but it extends easily to n players (as in the version we have here).

The idea here is that we are trying to create a type space where each agent's hierarchy of beliefs is that agent's type.³⁶ But we are not there yet. Players know their own belief hierarchy, but do not know the belief hierarchies of their opponents. Without further assumptions, the space we just defined faces an obvious problem – shouldn't players also have beliefs over their opponents' belief hierarchies? In general, taking $\Omega = \Theta_0$ generates a new space of hierarchies of belief, and for each new such space we face this same problem. Thus the model is not closed. To close the model, we need assumptions that ensure that each player i 's belief about their opponents' belief hierarchies is pinned down. Once we have pinned down these beliefs, we will have created a *universal type space*, and the hierarchies of beliefs will be the types of the agents.

Note that each $\theta \in \Theta_0$ is a sequence of probability measures μ_k . Denote the marginal distribution of μ_k on X_{k-2} by $\text{marg}_{X_{k-2}}(\mu_k)$.

Definition 51 (Coherency). A type $\theta \in \Theta_0$ is called *coherent* if for every $k \geq 2$, $\text{marg}_{X_{k-2}}(\mu_k) = \mu_{k-1}$.

We denote the set of all coherent belief hierarchies by Θ_1 .

Coherency simply requires that a player's higher order beliefs do not contradict their lower order beliefs (or vice versa). This pins down player i 's belief about her opponents' types.

Definition 52 (Homeomorphism). A function $f : X \rightarrow Y$ between two topological spaces X and Y is called a *homeomorphism* if

- (i) f is a bijection from X onto Y ,
- (ii) f is continuous, and
- (iii) f is an open mapping (i.e. the inverse f^{-1} is continuous).

The existence of a homeomorphism f from the set of coherent belief hierarchies onto the set of distributions on $X_0 \times \Theta_0$ tells us that the marginal probability $f(\mu_1, \mu_2, \dots)$ assigns to some event in X_{k-1} is the same as the probability μ_k assigns to that event.

Proposition 42. *There is a homeomorphism $f : \Theta_1 \rightarrow \Delta(X_0 \times \Theta_0)$.*

Proof. We require a technical lemma:

³⁶Brandenburger & Dekel (1993) jump straight to calling the elements of Θ_0 *types*. But Θ_0 is not a Harsanyi type space, and we will reserve “type” to mean only elements of a Harsanyi type space.

Lemma 11. Let $\{Z_n\}$ be a collection of Polish spaces,³⁷ let

$$M := \{(\mu_1, \mu_2, \dots) \mid \mu_k \in \Delta(Z_1 \times \dots \times Z_{k-1}) \text{ for all } k \geq 1\},$$

and let

$$M_0 := \{\mu \in M \mid \text{marg}_{Z_0 \times \dots \times Z_{k-2}}(\mu_k) = \mu_{k-2} \text{ for all } k \geq 2\}.$$

Then there is a homeomorphism $f : M_0 \rightarrow \Delta(\times_{k=0}^{\infty} Z_k)$.

Proof. Take any $\mu \in M_0$. By Kolmogorov's existence theorem, there is a unique measure $\hat{\mu} \in \Delta(\times_{k=0}^{\infty} Z_k)$ s.t. $\text{marg}_{Z_0 \times \dots \times Z_{k-2}}(\hat{\mu}) = \mu_k$ for all $k \geq 1$. Define f to be the mapping from $\mu = (\mu_1, \mu_2, \dots)$ into $\hat{\mu}$.

We claim that f is a homeomorphism. Suppose $\mu, \mu' \in M_0$ and $\mu \neq \mu'$. Then there is some k s.t. $\mu_k \neq \mu'_k$. If $f(\mu) = f(\mu')$ then we must have that $\mu_k = \text{marg}_{Z_0 \times \dots \times Z_{k-2}}(f(\mu)) = \text{marg}_{Z_0 \times \dots \times Z_{k-2}}(f(\mu')) = \mu'_k$, yielding a contradiction. Hence f is one-to-one. For any $\hat{\mu} \in \Delta(\times_{k=0}^{\infty} Z_k)$, $f(\text{marg}_{Z_0}(\hat{\mu}), \text{marg}_{Z_0 \times Z_1}(\hat{\mu}), \dots) = \hat{\mu}$, and so f is onto. Thus f is a bijection. Since the maps $\hat{\mu} \mapsto \text{marg}_{Z_0 \times \dots \times Z_{k-2}}(\hat{\mu})$ are all continuous, so is f^{-1} . Finally, take any sequence $\{\mu^r\} \subset M_0$ s.t. μ_k^r weakly converges to μ_k for all k with $\mu \in M_0$. Define $\hat{\mu}^r = f(\mu^r)$ for each r and $\hat{\mu} = f(\mu)$. Then $\hat{\mu}^r$ converges weakly to $\hat{\mu}$ since the cylinder sets form a convergence-determining class – that is, if $\hat{\mu}^r(C) \rightarrow \hat{\mu}(C)$ for all cylinder sets C such that $\hat{\mu}(\partial C) = 0$ then $\hat{\mu}^r$ converges weakly to $\hat{\mu}$. The values of $\hat{\mu}^r$ and $\hat{\mu}$ are respectively pinned down by the $\hat{\mu}_k^r$ s and $\hat{\mu}_k$ s. It follows that f is continuous. \square

Take $Z_0 = X_0$ and $Z_k = \Delta(X_{k-1})$ for $k \geq 1$ in the lemma. Then $Z_0 \times \dots \times Z_k = X_k$ and $\times_{k=0}^{\infty} Z_k = X_0 \times \Theta_0$. Any at most countable Cartesian product of a Polish space is a Polish space, so if Ω is a Polish space then so is X_0 . Moreover, if $Z_0 = X_0$ is a Polish space then so is $Z_1 = \Delta(X_0)$, and by extension, all of the Z_k s. By definition, the set of coherent belief hierarchies Θ_1 is simply M_0 in the lemma, so the lemma shows there is a homeomorphism $f : \Theta_1 \rightarrow \Delta(X_0 \times \Theta_0)$. \square

This tells us that coherence ensures that i 's belief about her opponents' belief hierarchies is pinned down by her own belief hierarchy. However, it does not follow that i 's higher order beliefs about her opponents' beliefs about her own type are pinned down by i 's type. For example, i might believe that some of her opponents' beliefs could be incoherent. To pin down players' higher order beliefs about belief hierarchies, we need to assume that it is common knowledge that players' beliefs are coherent.

Definition 53 (Universal type space). Consider n players, a space of fundamentals S and let $X_0 = \times_{k=1}^{n-1} S$. Let Θ_1 be a coherent type space with associated homeomorphism f . Define

$$\begin{aligned} \Theta_k &:= \{\theta \in \Theta_1 \mid f(\theta)(X_0 \times \Theta_{k-1}) = 1\} \quad \text{for } k \geq 2, \\ \Theta &:= \bigcap_{k=1}^{\infty} \Theta_k. \end{aligned}$$

Then we call Θ the *universal type space*.

³⁷A Polish space is a separable, completely metrizable topological space.

Because the model is now closed and beliefs are pinned down, the universal type space Θ is indeed a Harsanyi type space:

Proposition 43. *There is a homeomorphism $g : \Theta \rightarrow \Delta(X_0 \times \Theta)$.*

Proof. We have $\Theta = \{\mu \in \Theta_1 \mid f(\theta)(X_0 \times \Theta) = 1\}$. Hence $f(\Theta) = \{\mu \in \Delta(X_0 \times \Theta_0) \mid \mu(X_0 \times \Theta) = 1\}$, given f is onto. Since $f(\Theta)$ is homeomorphic to Θ and $\{\mu \in \Delta(X_0 \times \Theta_0) \mid \mu(X_0 \times \Theta) = 1\}$ is homeomorphic to $\Delta(X_0 \times \Theta)$, we have that Θ is homeomorphic to $\Delta(X_0 \times \Theta)$. \square

Hence all of a player i 's beliefs over their opponents types are pinned down in Θ , by the same argument we made for first-order beliefs being pinned down because of the existence of the homeomorphism f . As we might expect, the interpretation of Θ is that it is the set of types that are consistent with common knowledge of coherency.

This carries some lessons about what the limits of the type space approach are. Clearly, we have had to impose some restrictions on what kinds of higher order belief agents can hold. Coherency is quite a reasonable restriction – it rules out that agents hold contradictory higher order beliefs. While it is natural to think that agents might hold internally contradictory beliefs in the real world, it is obvious that to model this we would need to move away from the standard framework of Bayesian games. Common knowledge of coherency is a stronger restriction, but again no less unreasonable than, say, common knowledge of rationality.

In what sense does the universal type space “contain” all the other type spaces? Clearly, type spaces need not be subsets of the universal type space Θ . The types in the universal type space are hierarchies of beliefs, whereas a typical type space that we might work with in practice is closer to Example 30.

This said, because for any type space, the hierarchy of beliefs for each agent is pinned down by their type, we can “naturally” embed any type space $(\Omega, (\Theta_i, \pi_i)_{i=1}^n)$ in the universal type space. Denote player i 's first order belief by $\mu_i^1 : \Theta_i \rightarrow X_1$, where $X_1 = \Omega \times \Delta(\Omega)$, as in Definition 50. For each measurable subset $E \subset \Omega$, i 's first order belief if type θ_i is

$$\mu_i^1(E \mid \theta_i) = \pi_i(E \times \Theta_{-i} \mid \theta_i),$$

i.e. it is i 's belief integrating across all of the opponent's types. Now given embeddings up to the $(k-1)$ th order μ^{k-1} , define the projection

$$\gamma_i^k(\omega, \theta_{-i}) = (\omega, \mu^1(\theta_{-i}), \dots, \mu^{k-1}(\theta_{-i})) \in X_{k-1}.$$

Now a k th order belief is

The construction of the universal type space above relied heavily on the product topology, which we needed to apply Kolmogorov's extension theorem. Heifetz & Samet (1999) prove we can construct a universal type space, in the sense that there is an embedding of each type space in it, without making any topological assumptions at all.

Since we are relying on topological assumptions, however, it is worth asking whether the topological assumptions Brandenburger & Dekel (1993) made make sense. Dekel, Fudenberg & Morris (2006) argue they do not. In particular, the set of rationalizable

actions does not satisfy a lower hemicontinuity property in the product topology. This is obvious via a game introduced by Rubinstein (1989).

Example 31 (Rubinstein’s electronic mail game). Consider two generals (General 1 and General 2) commanding armies on either side of a valley, deciding whether to attack an enemy fort below. There is uncertainty over the state of the world, $\Omega = \{\omega_1, \omega_2\}$, with ω_1 occurring with probability $p \in (0, \frac{1}{2})$ and ω_2 occurring with probability $1 - p$. In the “good” state ω_1 , the enemy fort is not well-defended, so if both armies attack then the attack will be successful, whereas in the “bad” state ω_2 , the enemy fort is well-defended and an attack by both armies will result in large losses. A single army is by itself insufficient to take the fort so will be devastated if it attacks alone. Each general i chooses whether to retreat (R_i) or attack (A_i). We assume the payoffs in each state are:

| | State ω_1 | | | State ω_2 | |
|-------|------------------|---------|-------|------------------|---------|
| | R_2 | A_2 | | R_2 | A_2 |
| R_1 | 0, 0 | 0, $-L$ | R_1 | M, M | $M, -L$ |
| A_1 | $-L, 0$ | M, M | A_1 | $-L, M$ | 0, 0 |

where $L > M > 0$. Suppose the state is common knowledge. In the “bad” state ω_2 , retreating is always the strict best response for both generals, regardless of the other general’s decision. In the “good” state ω_1 , attacking is a best response for general i if general $j \neq i$ attacks, whereas retreating is a best response for i if j retreats. Hence there are two pure strategy equilibria: one in which $\sigma_i(\omega) = R_i$ for both i and one in which

$$\sigma_i(\omega) = \begin{cases} A_i & \text{if } \omega = \omega_1, \\ R_i & \text{if } \omega = \omega_2. \end{cases}$$

However, the state is not common knowledge. General 1 has a good view of the fort and so knows the state, but General 2’s side of the valley is shrouded in fog so she does not know whether the enemy fort is weakly-defended or well-defended. The generals communicate via ~~electronic mail~~ sending messengers on horseback to the other army.³⁸ In the bad state ω_2 , the enemy has patrols all across the valley and there is no hope that a messenger can reach General 2, so General 1 will not send a messenger. In the good state ω_1 , General 1 sends a messenger to General 2, but there is a probability $\epsilon > 0$ that the messenger is intercepted by an enemy patrol and so never reaches General 2. If General 2 receives the messenger, he knows the state is ω_1 and the fort is weakly defended, and he sends the messenger back to General 1 to let her know he knows the state. Again, the messenger is intercepted with probability ϵ , and if General 1 receives the messenger back, she knows that General 2 knows the state is ω_1 , and she sends the

³⁸Another popular version of the electronic mail game imagines two investors communicating about whether to invest in a project, with one having private information about the state. Rubinstein (1989) assumes there is some email system set up so that there is an automatic confirmation email sent when an email is received (with a small probability of error), hence the “electronic mail game” name. The email analogy probably makes more sense in that context, but the generals version is closer to the original formulation of coordinated attack problems in the distributed systems literature.

messenger back to General 2. This back and forth process continues indefinitely until the messenger is intercepted, and if that happens, communication ends.

Formally, each general i 's type is an integer $\theta_i \in \mathbb{N}$, with the interpretation that the type is the number of times i has received the messenger. The joint distribution of types and the state is partially summarized as follows:

| ω | θ_1 | θ_2 | Probability |
|------------|------------|------------|-----------------------------|
| ω_2 | 0 | 0 | $1 - p$ |
| ω_1 | 0 | 0 | $p\epsilon$ |
| ω_1 | 0 | 1 | $p(1 - \epsilon)\epsilon$ |
| ω_1 | 1 | 1 | $p(1 - \epsilon)^2\epsilon$ |
| ω_1 | 1 | 2 | $p(1 - \epsilon)^3\epsilon$ |
| ω_1 | 2 | 2 | $p(1 - \epsilon)^4\epsilon$ |
| ω_1 | 2 | 3 | $p(1 - \epsilon)^5\epsilon$ |
| \vdots | \vdots | \vdots | \vdots |

More completely, the probability that the state is ω_1 and the type profile is (k, k) for $k \geq 1$ is $p(1 - \epsilon)^{2k}\epsilon$, while the probability that the state is ω_1 and the type profile is $(k, k + 1)$ for $k \geq 0$ is $p(1 - \epsilon)^{2k+1}\epsilon$.

If general i is type $\theta_i > 0$ then $\omega_1 \in K_i^{\theta_i}(\{\omega_1\})$.³⁹ As $\theta_i \rightarrow \infty$, we have that beliefs converge pointwise to common knowledge, since $\omega_1 \in \bigcap_{k=1}^{\infty} K^k(\{\omega_1\})$. As we saw above, if $\{\omega_1\}$ is common knowledge in ω_1 then A_i is rationalizable for each general i . Yet for any type $\theta_i \in \mathbb{N}$, R_i is the unique rationalizable action: no matter how “close” the generals get to common knowledge, they will always strictly prefer to retreat.

To see this, we argue by induction. If the state is ω_2 , then General 1 knows the state is ω_2 and R_1 is strictly dominant. Now, if General 2's type is $\theta_2 = 0$, then either the state is the bad state ω_2 , which occurs with probability $1 - p > \frac{1}{2}$, or the state is the good state ω_1 but the messenger was intercepted, which occurs with probability $p\epsilon$. In ex ante terms, General 2's payoffs are

$$u_2(A_2, \sigma_1, \theta_2) \leq p\epsilon M - (1 - p)L < (1 - p)M = u_2(R_2, \sigma_1, \theta_2),$$

where the first inequality comes from noting that General 2's expected payoff from attacking is highest if General 1 always attacks if the state is ω_1 . We see the unique best response for General 2 if he has not received a message is R_2 . Given attacking is never a best response for type $\theta_2 = 1$, attacking is clearly not rationalizable for General 1 of type $\theta_1 = 0$, since General 2 is sure to retreat. This completes the base case.

Now, fix $k > 1$ and suppose $\sigma_2(\theta_2) = R_2$ for all $\theta_2 \leq k - 1$. If the type of General 1 is $\theta_1 = k$, then General 2 has type $\theta_2 \in \{k - 1, k\}$. In ex ante units, General 1's expected payoffs are

$$u_1(A_1, \sigma_2, \theta_1) \leq [-L + M(1 - \epsilon)]p(1 - \epsilon)^{2k-1}\epsilon < 0 = u_1(R_1, \sigma_2, \theta_1),$$

³⁹Recall Definition 9.

and so $\sigma_1(\theta_1) = R_1$ is the unique best response for General 1. A symmetric argument shows $\sigma_2(\theta_2) = R_2$ is the unique best response for General 2 if $\theta_2 = k$.

Hence it follows that attacking is not rationalizable for any type.

In Rubinstein’s electronic mail game, the rationalizable action correspondences $F_i : \Theta_i \rightrightarrows A_i$ are not lower hemicontinuous under the product topology. This suggests the product topology is probably not the natural topology to use to formalize the universal type space, because we should expect that “nearby” types have similar behaviour under reasonable solution concepts.⁴⁰ There is a literature that uses “finer” topologies than the product topology – see Dekel, Fudenberg & Morris (2006) or Chen et al. (2010, 2017).

3.6 Global games

In the previous section, we highlighted the role of higher order beliefs. In general, having to deal with the full hierarchies of belief makes analysis generally intractable – the type of a player in the universal type space is, after all, an infinite-dimensional object. It is very rare that you find games where the full hierarchy of belief is explicitly taken seriously. It is thus desirable to find settings that are tractable yet sufficiently rich to capture the idea that higher order beliefs matter. Global games, introduced in Carlsson & Van Damme (1993) provide such a setting. A global game is an incomplete information game where the payoff structure is fixed by a randomly drawn state, and players privately observe a signal about the state. A key insight in such games is that higher order beliefs do not need to be very sophisticated. Most of this section follows the (excellent) review of global games in Morris & Shin (2002).

3.6.1 Symmetric binary action global games

Suppose there is a continuum of players and each player chooses an action $a \in \{0, 1\}$. We assume all players have identical payoff function $u : \{0, 1\} \times [0, 1] \times \mathbb{R} \rightarrow \mathbb{R}$, where $u(a, x, \nu)$ denotes a player’s expected payoff if she chooses action $a \in \{0, 1\}$, if a proportion $x \in [0, 1]$ of her opponents choose action 1, and if she observes private signal $\nu \in \mathbb{R}$. Since only the difference in payoffs between $a = 1$ and $a = 0$ matters for a player’s best response, it is more convenient to parameterize the utility function using the function $\pi : [0, 1] \times \mathbb{R} \rightarrow \mathbb{R}$ defined by $\pi(x, \nu) := u(1, x, \nu) - u(0, x, \nu)$.

We first assume that players have a uniform prior about the initial state. Morris & Shin (2002) call these kinds of belief *Laplacian*, after Laplace’s (1824) *principle of insufficient reason* – the notion that agents’ beliefs should assign a uniform probability to unknown events.

Call an action a a *Laplacian action* if it is a best response to a uniform prior over opponents’ actions. That is, action 1 is the Laplacian action at ν if $\int_0^1 \pi(x, \nu) dx > 0$ and action 2 is the Laplacian action at ν if $\int_0^1 \pi(x, \nu) dx < 0$.

⁴⁰We could of course question whether “almost common knowledge”, as Rubinstein (1989) calls it, is indeed close to common knowledge.

Suppose a state $\theta \in \mathbb{R}$ is distributed according to the improper uniform density on \mathbb{R} , and each player i privately observes a signal $\nu_i = \theta + \alpha \epsilon_i$ where $\alpha > 0$ and ϵ_i is noise distributed with conditional density g with support on \mathbb{R} .⁴¹ On observing signal ν_i , player i places posterior density $f(\theta \mid \nu_i) = \frac{1}{\alpha} g((\nu_i - \theta)/\alpha)$ on θ .

Assume payoffs satisfy the following assumptions:

Assumption.

- (G1) *Action monotonicity.* $\pi(\cdot, \theta)$ is nonincreasing for all $\theta \in \mathbb{R}$.
- (G2) *State monotonicity.* $\pi(x, \cdot)$ is nonincreasing for all $x \in [0, 1]$.
- (G3) *Strict Laplacian state monotonicity.* There is a unique θ^* such that $\int_0^1 \pi(x, \theta^*) dx = 0$.
- (G4) *Limit dominance.* There are $\underline{\theta} \in \mathbb{R}$ and $\bar{\theta} \in \mathbb{R}$ such that
 - (i) $\pi(x, \nu) < 0$ for all $x \in [0, 1]$ and $\nu \leq \underline{\theta}$, and
 - (ii) $\pi(x, \nu) > 0$ for all $x \in [0, 1]$ and $\nu \geq \bar{\theta}$.
- (G5) *Continuity.* The mapping $(\nu, h) \mapsto \int_0^1 h(x) \pi(x, \nu) dx$, where h is a density, is continuous in ν and (with respect to the weak topology) in density h .

Action monotonicity (G1) ensures that the game is one of strategic complements. Note that the continuity assumption (G5) allows $\pi(\cdot, \nu)$ to be discontinuous in x (on a set of measure zero).

Let $G(\alpha)$ be the game of incomplete information just described, under the assumptions (G1)-(G5). For each player i , a strategy in this game is a mapping $\sigma_i : \mathbb{R} \rightarrow \{0, 1\}$, prescribing action $\sigma_i(\nu_i)$ for each signal $\nu_i \in \mathbb{R}$.

Proposition 44. *Consider the game $G(\alpha)$, and let θ^* solve $\int_0^1 \pi(x, \theta^*) dx = 0$. Then for each player i , there is an essentially unique strategy σ_i^* surviving iterated strict dominance. Moreover, $\sigma_i^*(\nu) = 0$ for all $\nu < \theta^*$ and $\sigma_i^*(\nu) = 1$ for all $\nu > \theta^*$.*

Proof. Let $\hat{\pi}(\nu_i, k)$ denote the expected payoff to player i if observing signal ν_i and knowing that all her opponents j will choose action 0 if they observe $\nu_j < k$. That is,

$$\hat{\pi}(\nu_i, k) = \int_{\mathbb{R}} \frac{1}{\alpha} g\left(\frac{\nu - \theta}{\alpha}\right) \pi\left(1 - G\left(\frac{k - \theta}{\alpha}\right), \nu\right) d\theta.$$

Now, $\pi(\nu, k)$ is continuous in ν, k , increasing in ν , decreasing in k , and $\hat{\pi}(\nu_i, k, \sigma) \leq 0$ if $x \leq 0$.

⁴¹An *improper density* f is such that $\int_{\mathbb{R}} f(\theta) d\theta = \infty$. The improper uniform density is $f(\theta) = 1$ for all $\theta \in \mathbb{R}$. Obviously, this is not a density, since it does not integrate to 1. However, as long as the marginal $m(\nu_i) := \int_{\mathbb{R}} f(\nu_i \mid \theta) f(\theta) d\theta$ is well-defined for all $\nu_i \in \mathbb{R}$ then the posterior $f(\theta \mid \nu_i)$ is well-defined.

Lemma 12. *A strategy σ survives m rounds of iterated strict dominance iff*

$$\sigma(\nu) = \begin{cases} 0 & \text{if } x < \underline{\xi}_m, \\ 1 & \text{if } x > \bar{\xi}_m, \end{cases}$$

where $\underline{\xi}_0 = -\infty$, $\bar{\xi}_0 = +\infty$, and

$$\begin{aligned} \underline{\xi}_{m+1} &= \min\{\nu \mid \hat{\pi}(\nu, \underline{\xi}_m) = 0\}, \\ \bar{\xi}_{m+1} &= \max\{\nu \mid \hat{\pi}(\nu, \bar{\xi}_m) = 0\}, \end{aligned}$$

for all $m \geq 1$.

Proof. Suppose this holds for rounds up to and including m . By **(G1)**, i.e. strategic complementarities, if action 1 was a best response to a strategy that survived m rounds, then it must also be a best response to the threshold strategy with threshold $\underline{\xi}_m$, and $\underline{\xi}_{m+1}$ gives the lowest signal for which this occurs. Likewise, if action 0 was a best response to a strategy surviving m rounds, then it must be a best response to the threshold strategy with threshold $\bar{\xi}_m$ and $\bar{\xi}_{m+1}$ gives the highest signal for which this occurs. \square

Now $\{\underline{\xi}_m\}$ is monotonically increasing and $\{\bar{\xi}_m\}$ is monotonically decreasing, and so $\underline{\xi} := \lim_{m \rightarrow \infty} \underline{\xi}_m$ and $\bar{\xi} := \lim_{m \rightarrow \infty} \bar{\xi}_m$ both exist. From continuity of $\hat{\pi}$, we have that $\hat{\pi}(\underline{\xi}, \underline{\xi}) = 0 = \hat{\pi}(\bar{\xi}, \bar{\xi})$.

It remains to show that θ^* is the unique solution to $\hat{\pi}(y, y) = 0$ and hence $\theta^* = \underline{\xi} = \bar{\xi}$. Let $\hat{\psi}(x, \nu, k)$ denote the probability that a player, on observing signal ν , assigns to a proportion less than x of the other players observing a signal greater than threshold k . If the state is θ , the proportion of players observing a signal exceeding k is $1 - G((k - \theta)/\alpha)$, which is less than x provided $\theta \leq k - \alpha G^{-1}(1 - x)$. Whence

$$\begin{aligned} \hat{\psi}(x, \nu, k) &= \int_{-\infty}^{k - \alpha G^{-1}(1 - x)} \frac{1}{\alpha} g\left(\frac{\nu - \theta}{\alpha}\right) d\theta \\ &= \int_{\frac{\nu - k}{\alpha} + G^{-1}(1 - x)}^{\infty} g(z) dz \\ &= 1 - G\left(\frac{\nu - k}{\alpha} + G^{-1}(1 - x)\right). \end{aligned}$$

If $\nu = k$ then $\hat{\psi}(x, \nu, k) = x$, the identity function, and so it is the cdf of the uniform density. Hence

$$\hat{\pi}(\nu, \nu) = \int_0^1 \pi(x, \nu) dx.$$

By **(G3)**, $\hat{\pi}(\nu, \nu) = 0$ gives $\nu = \theta^*$. \square

4 Mechanism design

In a typical Bayesian game, the structure of the game is fixed. Mechanism design is concerned with the study of how to optimally design a game to maximize some objective. For example, we might be concerned with how we can structure an auction to maximize the seller's revenue, how we should aggregate individual preferences to make collective decisions, as with voting mechanisms, or how to structure markets to achieve some social objective. The object of study, rather than games, is *mechanisms* or *game forms*.

Consider a finite set \mathcal{I} of n agents, and a set of potential decisions X . Agents i have private information summarized by a type $\theta_i \in \Theta_i$. As usual, a type profile $\theta = (\theta_1, \dots, \theta_n)$ lies in the product set $\Theta = \times_{i=1}^n \Theta_i$. Given agent i , we denote by $\Theta_{-i} = \times_{j \neq i} \Theta_j$ the set of type profiles of the other agents.

Agents' preferences are represented by a utility function $u_i : X \times \Theta \rightarrow \mathbb{R}$. This is a very general form – in particular, it might be that agents have interdependent values, so agent i 's benefit $u_i(x, \theta)$ from a decision x depends on the types of other agents. In the special case that $u_i(x, \theta_i, \theta_{-i}) = u_i(x, \theta_i, \theta'_{-i})$ for all $x \in X$, $\theta_i \in \Theta_i$ and $\theta_{-i}, \theta'_{-i} \in \Theta_{-i}$, we say that agents have *private values*, and we can instead summarize agent i 's preferences by the utility function $u_i : X \times \Theta_i \rightarrow \mathbb{R}$. Private values imply that the benefit an agent receives from a decision depends only on their own type.

Example 32. There are many situations where mechanism design is applicable. Here are some examples:

- (a) *Public decisions about projects.* Suppose a community \mathcal{I} decides whether to pursue some project, such as building a bridge, road, public hospital, defensive wall around a city, new clubhouse, etc. The cost of the project is c , and each individual i is asked to contribute c_i towards the cost, with $\sum_{i \in \mathcal{I}} c_i = c$. The decision is whether to pursue the project ($x = 1$) or not ($x = 0$), so the set of decisions is $X = \{0, 1\}$. Suppose each individual's type is the benefit they derive if the project goes ahead, and they receive payoff 0 if the project does not go ahead. Then the utility function of individual i is $u_i(x, \theta_i) = (\theta_i - c_i)x$.
- (b) *Public goods provision.* Suppose a society \mathcal{I} faces the problem of determining the production $x \in \mathbb{R}^n$ of a public good, with production split between the agents. Each agent i has preferences $u_i(x, \theta_i) = \theta_i \left(\sum_{j \in \mathcal{I}} x_j \right)^{1/\sigma} - c_i x_i$ ($\sigma > 0$), where $c_i > 0$ is the cost to agent i of producing x_i units of the public good, and type θ_i determining how much i values the public good.
- (c) *Elections.* Suppose a constituency \mathcal{I} of voters has to elect a representative from m alternative candidates $X = \{1, \dots, m\}$, and the decision to be made is which candidate x will be elected. Each voter has a preference ordering over the candidates determined by her type θ_i and receives payoff $u_i(x, \theta_i)$ if candidate x is elected.
- (d) *Allocating private goods.* Suppose an indivisible good is to be allocated to one of n agents. The set of decisions is $X = \{x \in \{0, 1\}^n \mid \sum_{i=1}^n x_i = 1\}$. Each agent i has a

value θ_i for the object, and so receives benefit $u_i(x, \theta_i) = x_i \theta_i$. This has motivated a lot of the mechanism design literature, particularly in relation to auction design.

We call the triple $(\mathcal{I}, X, (\Theta_i, u_i)_{i \in \mathcal{I}})$ a *mechanism design problem*.

Definition 54 (Social choice function). Given mechanism design problem $(\mathcal{I}, X, (\Theta_i, v_i)_{i \in \mathcal{I}})$, a *social choice function* is a function $f : \Theta \rightarrow X$.

Definition 55 (Mechanism).

- (a) *Mechanism*. Consider a mechanism design problem $(\mathcal{I}, X, (\Theta_i, v_i)_{i \in \mathcal{I}})$. For each $i \in \mathcal{I}$, let M_i denote a *message space* and define $M := \times_{i \in \mathcal{I}} M_i$. An element $m_i \in M_i$ is a *message*, and $m \in M$ is a *message profile*. Let $g : M \rightarrow X$ be an *outcome function*, mapping message profiles into a decision-transfer pair. Then we call the pair (M, g) a *mechanism* or *game form*.⁴²

Given a mechanism (M, g) , a (pure) *strategy* σ_i for an agent i is a mapping $\sigma_i : \Theta_i \rightarrow M_i$. We denote i 's strategy set by Σ_{Θ_i} .

- (b) *Direct mechanism*. A mechanism (M, g) is a *direct mechanism* if $M_i = \Theta_i$ for each agent $i \in \mathcal{I}$ and $g = f$ where f is the social choice function.

4.1 Incentive compatibility and the revelation principle

The space of mechanisms is very large, because in general, the possible message spaces can be arbitrary. Fortunately, it is often possible to restrict to a much smaller space of mechanisms – direct mechanisms – without loss of generality, thanks to the *revelation principle*.

There are three dominant approaches to mechanism design. The strongest – dominant strategy incentive compatibility – concerns implementation when it is a weakly dominant strategy to truthfully report one's type. A weaker notion is ex-post incentive compatibility, and weaker still is Bayesian incentive compatibility. In most settings of interest in mechanism design, we can apply revelation principles to restrict our setting to implementation in direct mechanisms.

4.1.1 Dominant strategy incentive compatibility

Given a mechanism (M, g) , we call a strategy $\sigma_i \in \Sigma_{\Theta_i}$ of player i a (weakly) *dominant strategy* at $\theta_i \in \Theta_i$ if

$$u_i(g(\sigma_i(\theta_i), m_{-i}), \theta) \geq u_i(g(m'_i, m_{-i}), \theta),$$

for all $m'_i \in M_i$ and $m_{-i} \in M_{-i}$. Note this differs from the standard notion of a weakly dominant strategy in the game theory literature, since we do not require that the inequality holds strictly for some profile of opponents' messages m_{-i} .

⁴²The distinction between a (normal-form) *game* (S, u) (where $S = \times_{i \in \mathcal{I}} S_i$ is the product space of strategies and $u = (u_i)_{i \in \mathcal{I}}$) and a *game form* (M, g) is that the former associates payoffs to actions, whereas the latter associates outcomes to actions.

Definition 56 (Dominant strategy implementation). We say a mechanism (M, g) *implements* a social choice function f in dominant strategies if there exists a strategy profile $\sigma = (\sigma_1, \dots, \sigma_n)$ such that

- (i) σ is a dominant strategy equilibrium of (M, g) , and
- (ii) $g_i(\sigma(\theta)) = f_i(\sigma(\theta))$ for every agent $i \in \mathcal{I}$ and every type profile $\theta \in \Theta$.

We call the social choice function f *dominant strategy implementable* if there exists some mechanism (M, g) that implements f .

Definition 57 (Dominant strategy incentive compatibility). We say a direct mechanism (Θ, f) is *dominant strategy incentive compatible* (DSIC) or *strategy-proof* if for every agent $i \in \mathcal{I}$, every $\theta_{-i}, \theta'_{-i} \in \Theta_{-i}$ and every $\theta_i, \theta'_i \in \Theta_i$, we have that

$$u_i(f(\theta_i, \theta'_{-i}), \theta) \geq u_i(f(\theta'_i, \theta'_{-i}), \theta).$$

That is, dominant strategy incentive compatibility implies truth-telling about one's type is a dominant strategy, regardless of whether other agents report truthfully. In the case of private values, strategy-proofness reduces to the requirement that for every i , every $\theta_{-i} \in \Theta_{-i}$, and every $\theta_i, \theta'_i \in \Theta_i$, we have that

$$u_i(f(\theta_i, \theta_{-i}), \theta_i) \geq u_i(f(\theta'_i, \theta_{-i}), \theta_i).$$

Proposition 45 (Revelation principle in dominant strategies; Myerson, 1979). *If a social choice function f is implementable in dominant strategies by some mechanism (M, g) , then the direct mechanism (Θ, f) is dominant strategy incentive compatible.*

Proof. Suppose the mechanism (M, g) implements f in dominant strategies, and each agent i has a dominant strategy σ_i under (M, g) . Fix agent i and type profiles $\theta, \theta' \in \Theta$. Suppose the profile of reported types of other agents in the direct mechanism is θ'_{-i} , and let $\sigma_i(\theta_i) = m_i$, $\sigma_i(\theta'_i) = m'_i$, and $\sigma_{-i}(\theta'_{-i}) = m'_{-i}$. Since f is implemented in dominant strategies by (M, g) , it must be the case that $u_i(g(m_i, m'_{-i}), \theta) \geq u_i(g(m'_i, m'_{-i}), \theta)$. Now, $f_i(\theta_i, \theta'_{-i}) = g_i(m_i, m'_{-i})$ and $f_i(\theta'_i, \theta'_{-i}) = g_i(m'_i, m'_{-i})$, so it follows that $u_i(f(\theta_i, \theta'_{-i}), \theta) \geq u_i(f(\theta'_i, \theta'_{-i}), \theta)$. Hence (Θ, f) is DSIC. \square

4.1.2 Ex post incentive compatibility

Dominant strategy incentive compatibility is particularly strong, since it requires that truthful reporting of one's type is always a best response to any strategy that an agents' opponents might play. Of course, this carries certain benefits – in particular, DSIC mechanisms are particularly robust to changing assumptions about the common knowledge structure of the setting. A weaker notion that maintains robustness is *ex post incentive compatibility*.

Definition 58 (Ex post implementation). We say a mechanism (M, g) *implements* a social choice function f in ex post equilibrium strategies if there exists a strategy profile $\sigma = (\sigma_1, \dots, \sigma_n)$ such that

- (i) σ is an ex post equilibrium of (M, g) , and
- (ii) $g_i(\sigma(\theta)) = f_i(\sigma(\theta))$ for every agent $i \in \mathcal{I}$ and every type profile $\theta \in \Theta$.

We say social choice function f is *ex post implementable* if there is some mechanism (M, g) that implements f in ex post equilibrium strategies.

Definition 59 (Ex post incentive compatibility). We say a direct mechanism (Θ, f) is *ex post incentive compatible* (EPIC) if for every agent $i \in \mathcal{I}$, every $\theta_{-i} \in \Theta_{-i}$ and every $\theta_i, \theta'_i \in \Theta_i$, we have that

$$u_i(f(\theta_i, \theta_{-i}), \theta) \geq u_i(f(\theta'_i, \theta_{-i}), \theta).$$

That is, ex post incentive compatibility implies an agent never regrets truth-telling about one's type *provided others tell the truth about their own types*. In the case of private values, EPIC reduces to the requirement that for every i , every $\theta_{-i} \in \Theta_{-i}$ and every $\theta_i, \theta'_i \in \Theta_i$, we have that

$$u_i(f(\theta_i, \theta_{-i}), \theta_i) \geq u_i(f(\theta'_i, \theta_{-i}), \theta_i),$$

and so under private values, EPIC is identical to DSIC.

Proposition 46 (Revelation principle in ex post equilibrium strategies). *If a social choice function f is implementable in ex post equilibrium strategies by some mechanism (M, g) , then the direct mechanism (Θ, f) is ex post incentive compatible.*

Proof. The proof is analogous to that of Proposition 45. Suppose the mechanism (M, g) implements f in ex post equilibrium strategies, and each agent i has a dominant strategy σ_i under (M, g) . Fix agent i and types $\theta_i, \theta'_i \in \Theta_i$. Suppose the profile of reported types of other agents in the direct mechanism is θ_{-i} , and let $\sigma_i(\theta_i) = m_i$, $\sigma_i(\theta'_i) = m'_i$, and $\sigma_{-i}(\theta_{-i}) = m_{-i}$. Since f is implemented in ex post equilibrium strategies by (M, g) , it must be the case that $u_i(g(m_i, m_{-i}), \theta) \geq u_i(g(m'_i, m_{-i}), \theta)$. Now, $f_i(\theta_i, \theta_{-i}) = g_i(m_i, m_{-i})$ and $f_i(\theta'_i, \theta_{-i}) = g_i(m'_i, m_{-i})$, so it follows that $u_i(f(\theta_i, \theta_{-i}), \theta) \geq u_i(f(\theta'_i, \theta_{-i}), \theta)$. Hence (Θ, f) is EPIC. \square

4.1.3 Bayesian incentive compatibility

Much weaker than dominant strategy and ex post incentive compatibility, Bayesian incentive compatibility requires only that truthful reporting is a Bayes Nash equilibrium. Much of classical mechanism design focussed on the Bayesian implementation. However, Bayesian implementation implies strong common knowledge assumptions – in particular, we assume common knowledge of the distribution of types (and consequently entire belief hierarchies).

Definition 60 (Bayesian implementation). We say a mechanism (M, g) *implements* a social choice function f in Bayes Nash equilibrium strategies if there exists a strategy profile $\sigma = (\sigma_1, \dots, \sigma_n)$ such that

- (i) σ is a Bayes Nash equilibrium of (M, g) , and
- (ii) $g_i(\sigma(\theta)) = f_i(\sigma(\theta))$ for every agent $i \in \mathcal{I}$ and every type profile $\theta \in \Theta$.

Definition 61 (Bayesian incentive compatibility). We say a direct mechanism (Θ, f) is *Bayesian incentive compatible* (BIC) if for every agent $i \in \mathcal{I}$ and every $\theta_i \in \Theta_i$, we have that

$$\mathbb{E}[u_i(f(\theta_i, \theta_{-i}), \theta) \mid \theta_i] \geq \mathbb{E}[u_i(f(\theta'_i, \theta_{-i}), \theta) \mid \theta_i]$$

for all $\theta'_i \in \Theta_i$. More explicitly, given distribution F over Θ of types,

$$\int_{\Theta_{-i}} u_i(f(\theta_i, \theta_{-i}), \theta_i, \theta_{-i}) F_i(d\theta_{-i} \mid \theta_i) \geq \int_{\Theta_{-i}} u_i(f(\theta'_i, \theta_{-i}), \theta_i, \theta_{-i}) F_i(d\theta_{-i} \mid \theta_i)$$

for all $\theta'_i \in \Theta_i$, where $F_i(\theta_{-i} \mid \theta_i)$ is the conditional distribution of the type profile θ_{-i} of agents other than i given i has type θ_i .

That is, Bayesian incentive compatibility implies truth telling about one's own types is optimal given others tell the truth about their own types *given the distribution of others' types*.

Proposition 47 (Revelation principle in Bayes Nash equilibrium strategies). *If a social choice function f is implementable in Bayes Nash equilibrium strategies by some mechanism (M, g) , then the direct mechanism (Θ, f) is Bayesian incentive compatible.*

Proof. Let F denote the distribution of types. Suppose the mechanism (M, g) implements f in Bayes Nash equilibrium strategies. Let (m_1, \dots, m_n) be a profile of Bayes Nash equilibrium strategies implementing f . Fix agent i and types $\theta_i, \theta'_i \in \Theta_i$. Implementability and Bayes Nash equilibrium implies that

$$\begin{aligned} \int_{\Theta_{-i}} u_i(f(\theta_i, \theta_{-i}), \theta_i, \theta_{-i}) F(d\theta_{-i} \mid \theta_i) &= \int_{\Theta_{-i}} u_i(g(m_i(\theta_i), m_{-i}(\theta_{-i})), \theta_i, \theta_{-i}) F(d\theta_{-i} \mid \theta_i) \\ &\geq \int_{\Theta_{-i}} u_i(g(m_i(\theta'_i), m_{-i}(\theta_{-i})), \theta_i, \theta_{-i}) F(d\theta_{-i} \mid \theta_i) \\ &= \int_{\Theta_{-i}} u_i(f(\theta_i, \theta_{-i}), \theta_i, \theta_{-i}) F(d\theta_{-i} \mid \theta_i) \end{aligned}$$

Hence (Θ, f) is BIC. □

4.1.4 Limits to the revelation principle

The revelation principle is a useful tool but should not be taken for granted. In particular, while the revelation principle is a necessary condition for implementability, it is not sufficient. It merely tells us that if a social choice function f is implementable, then the direct mechanism (Θ, f) has an equilibrium in which all agents report their type truthfully. But (with the exception of dominant strategies, which are necessarily unique), this does not preclude the direct mechanism having other, unwanted equilibria in which

agents do not report truthfully. Consider the following example from Moore & Repullo (1988). There are $n \geq 2$ agents, and suppose $u_i(f(\theta), \theta_i) \geq 0$ for each agent i and each type profile $\theta \in \Theta$. Suppose agents all know each others types and each agent announces a full type profile θ . If each agent announces the same type profile, $f(\theta)$ is implemented, whereas if agents do not all announce the same type profile, a bad outcome x is implemented which yields $u_i(x, \theta_i) = -1$ for each agent i . If θ^* is the true type profile, then it is an equilibrium for each agent to announce $\theta = \theta^*$. However, for any type profile θ , it is an equilibrium for all agents announce to announce θ , regardless of whether $\theta = \theta^*$.

4.2 Quasilinear preferences, transfers and private values

Typically, imposing structure on the mechanism design setting is important – indeed, there are general impossibility results that preclude us from making interesting positive conclusions without further structure. Here we will consider a setting where agents can make transfers. There are some important assumptions that discipline the setting:

- (i) *Private values.* Each agent $i \in \mathcal{I}$ receives utility from a decision $x \in X$ that depends only on their own type. That is $u_i(x, \theta_i, \theta_{-i}) = u_i(x, \theta_i, \theta'_{-i})$ for all $x \in X$, all $\theta_i \in \Theta_i$ and all $\theta_{-i}, \theta'_{-i} \in \Theta_{-i}$. Hence we can without loss take u_i to be a function $u_i : X \times \Theta_i \rightarrow \mathbb{R}$ (as opposed to $u_i : X \times \Theta \rightarrow \mathbb{R}$).
- (ii) *Quasilinear preferences.* We assume decisions consist of two components, an alternative $a \in A$ and a vector of transfers $t \in \mathbb{R}^n$. Each agent i of type $\theta_i \in \Theta_i$ has quasilinear preferences – that is, there exists a valuation function $v_i : A \times \Theta_i \rightarrow \mathbb{R}$ such that for $a, a' \in A$, $t_i, t'_i \in \mathbb{R}$, and each $\theta_i \in \Theta_i$, we have that $(a, t_i) \succ_{\theta_i} (a', t_i)$ iff $v_i(a, \theta_i) - t_i \geq v_i(a', \theta_i) - t'_i$.

We will call a mechanism design problem $(\mathcal{I}, A \times \mathbb{R}^n, (\Theta_i, u_i)_{i \in \mathcal{I}})$ that satisfies these assumptions a *QPV mechanism design problem*. We might also less formally refer to a *QPV setting* to mean situations where these assumptions apply.

The private values assumption implies that an agent's payoff does not depend on the types of the other agents (though it does not rule out that types may be correlated). This is reasonable in many settings involving allocations, public goods provision, certain matching settings, and so on: the value an object to a bidder typically does not depend on the value of the object to other bidders, and an agent's enjoyment of a public good does not depend on how much others value it.⁴³

⁴³Of course, we can contrive settings where the value to a bidder may depend on others' values. Consider the following setting. Alice and Bob are wine connoisseurs and neighbours. Alice's cute pet dog, Frodo, has taken a disliking to Bob and always barks at him whenever they encounter each other. Bob resents this, and the situation escalated until Alice and Bob had a serious falling out. Alice is a forgiving person, and does not hold a grudge, but Bob is petty and vindictive. Unfortunately, the two both frequent the same wine auctions. Alice's value for a bottle of wine does not depend on how much Bob values it. However, Bob, being petty and vindictive, derives satisfaction from depriving Alice of bottles of wine she values. His payoff from being allocated a particular bottle of wine in the auction is thus increasing in how much Alice values it.

Quasilinear preferences is a standard assumption in settings involving transfers. It is worth noting quasilinearity implies that, if we take the standard interpretation that a transfer is a payment, then agents are risk neutral with respect to wealth (their preferences over the alternatives A , on the other hand, are unrestricted).

In QPV settings, it is useful to disaggregate social choice functions $f : \Theta \rightarrow A \times \mathbb{R}^n$ into two components:

- (a) *Decision rule.* We call the function $f_a : \Theta \rightarrow A$, associating type profiles to alternatives, a *decision rule* (or in contexts where we are allocating objects, an *allocation rule*).
- (b) *Transfer rule.* We call the function $f_t : \Theta \rightarrow \mathbb{R}^n$, associating type profiles to a vector of transfers, a *transfer rule* (or *payment rule*). The transfer to/from agent i if the type profile is θ is $f_{ti}(\theta)$.

This decomposes $f = (f_a, f_t)$. A positive transfer $f_{ti}(\theta) > 0$ implies that agent i is making payments on net, whereas a negative transfer $f_{ti}(\theta) < 0$ implies agent i is receiving payments on net. We call the quantity $\sum_{i \in \mathcal{I}} f_{ti}(\theta)$ *total revenue*, since it is the value of total net transfers from agents to the designer (negative total revenue corresponds to a situation where the designer subsidizes the mechanism from some outside source.)

In many settings, it is natural to impose constraints on transfers. For example, typically the mechanism designer wants to avoid subsidising the system using revenue outside the mechanism. Some plausible constraints are detailed below. We do not impose these in all settings, but there are some settings where they are natural.

Definition 62. Let $f_t : \Theta \rightarrow \mathbb{R}^n$ be a transfer rule.

- (a) *Feasibility.* We call f_t *feasible* if $\sum_{i \in \mathcal{I}} f_{ti}(\theta) \geq 0$ for all $\theta \in \Theta$, that is, if total revenue is positive.
- (b) *No subsidy condition.* We say f_t satisfies the *no subsidy condition* if $f_{ti}(\theta) \geq 0$ for all $\theta \in \Theta$.
- (c) *Budget balance.* We say f_t is *ex post budget balanced* if

$$\sum_{i \in \mathcal{I}} f_{ti}(\theta) = 0$$

for all $\theta \in \Theta$. This requires that total revenue is zero for every type profile θ .

If the type profiles are distributed with distribution F , we say f_t is *ex ante budget balanced* if

$$\sum_{i \in \mathcal{I}} \int_{\Theta} f_{ti}(\theta) F(d\theta) = 0.$$

This requires only that total revenue is zero in expectation.

For now, let us focus on dominant strategy implementation. If a direct mechanism $f = (f_a, f_t)$ is strategyproof, then we can maintain strategyproofness with any other transfer rule that does not alter how each agent i 's own type θ_i determines i 's net transfer:

Proposition 48. *Suppose $f = (f_a, f_t)$ is strategyproof, and let $q : \Theta \rightarrow \mathbb{R}^n$ be defined by*

$$q_i(\theta) = f_{ti}(\theta_i, \theta_{-i}) + h_i(\theta_{-i})$$

for all $\theta \in \Theta$ and all $i \in \mathcal{I}$, where $h_i : \Theta_{-i} \rightarrow \mathbb{R}$ is an arbitrary real function on Θ_{-i} . Then the mechanism (f_a, q) is also strategyproof.

Proof. Fix agent i and type profile $\theta_{-i} \in \Theta_{-i}$. Now, take any $\theta_i, \theta'_i \in \Theta_i$. We have

$$\begin{aligned} v_i(f_a(\theta_i, \theta_{-i}), \theta_i) - q_i(\theta_i, \theta_{-i}) &= v_i(f_a(\theta_i, \theta_{-i}), \theta_i) - t_i(\theta_i, \theta_{-i}) - h_i(\theta_{-i}) \\ &\geq v_i(f_a(\theta_i, \theta_{-i}), \theta_i) - t_i(\theta'_i, \theta_{-i}) - h_i(\theta_{-i}) \\ &= v_i(f_a(\theta_i, \theta_{-i}), \theta_i) - q_i(\theta'_i, \theta_{-i}). \end{aligned}$$

□

Second, if $f = (f_a, f_t)$ is a strategyproof social choice function then the transfers depend only on the alternative selected by f_a (or in an allocation problem, the transfers depend only on the allocation). This is the *taxation principle*.

Proposition 49 (Taxation principle). *Suppose direct mechanism $f = (f_a, f_t)$ is strategyproof. Then for every agent $i \in \mathcal{I}$, every type profile $\theta_{-i} \in \Theta_{-i}$, and every pair of types $\theta_i, \theta'_i \in \Theta_i$,*

$$f_a(\theta_i, \theta_{-i}) = f_a(\theta'_i, \theta_{-i}) \quad \text{implies} \quad f_t(\theta_i, \theta_{-i}) = f_t(\theta'_i, \theta_{-i}).$$

Proof. Suppose f is strategyproof and fix agent i and type profile $\theta_{-i} \in \Theta_{-i}$. Take $\theta_i, \theta'_i \in \Theta_i$ s.t. $f_a(\theta_i, \theta_{-i}) = f_a(\theta'_i, \theta_{-i}) = a$. Now, strategyproofness implies

$$\begin{aligned} v_i(a, \theta_i) - f_{ti}(\theta_i, \theta_{-i}) &\geq v_i(a, \theta_i) - f_{ti}(\theta'_i, \theta_{-i}), \\ v_i(a, \theta'_i) - f_{ti}(\theta'_i, \theta_{-i}) &\geq v_i(a, \theta'_i) - f_{ti}(\theta_i, \theta_{-i}). \end{aligned}$$

Hence $f_{ti}(\theta_i, \theta_{-i}) = f_{ti}(\theta'_i, \theta_{-i})$.

□

4.2.1 Groves mechanisms

As ever, one of our main concerns in a QPV mechanism design problem is whether we can achieve efficient outcomes.

Definition 63 (Efficiency).

- (a) *Efficient decision rule.* A decision rule f_a is called *ex post efficient* (or *utilitarian*) if

$$f_a(\theta) \in \arg \max_{a \in A} \sum_{i \in \mathcal{I}} v_i(a, \theta_i)$$

for all type profiles $\theta \in \Theta$.

- (b) *Pareto optimal mechanism.* We call (direct) mechanism $f = (f_a, f_t)$ *Pareto optimal* if for each type profile $\theta \in \Theta$, there is no alternative $a \neq f_a(\theta)$ and no set of transfers $t \in \mathbb{R}^n$ such that

$$v_i(a, \theta_i) - t_i \geq v_i(f_a(\theta), \theta_i) - f_{ti}(\theta) \quad \text{for each } i \in \mathcal{I}, \text{ and}$$

$$\sum_{i \in \mathcal{I}} t_i \geq \sum_{i \in \mathcal{I}} f_{ti}(\theta),$$

with at least one of the above inequalities holding strictly.

An (ex post) efficient decision rule is one that maximizes total welfare, and hence why it is utilitarian. Pareto optimality may seem weaker, but because of the possibility of transfers it turns out the two are closely related. Note our definition of Pareto optimality compares a given mechanism only with those mechanisms that generate at least as much total revenue. There are a few motivations for this. First, if we only care about the welfare of the agents, then without this restriction, we can always increase the transfers from the designer to the agents and increase everyone's welfare, so there is no Pareto optimal mechanism. Second, if we interpret the designer as an agent who cares about total revenue (such as a seller in an auction), then any mechanism that yields lower total revenue leaves the designer worse off.

Theorem 17. *A (direct) mechanism $f = (f_a, f_t)$ is Pareto optimal iff f_a is an efficient decision rule.*

Proof. Towards contradiction, suppose f is Pareto optimal but f_a is not efficient. Then there is a type profile θ at which $f_a(\theta) \notin \arg \max_{a \in A} \sum_{i \in \mathcal{I}} v_i(a, \theta_i)$. Fix this θ , and let $a \in \arg \max_{a \in A} \sum_{i \in \mathcal{I}} v_i(a, \theta_i)$. Then

$$\sum_{i \in \mathcal{I}} v_i(f_a(\theta), \theta_i) < \sum_{i \in \mathcal{I}} v_i(a, \theta_i).$$

Take

$$\delta = \frac{1}{n} \sum_{i \in \mathcal{I}} (v_i(a, \theta_i) - v_i(f_a(\theta), \theta_i)) > 0.$$

Now, consider alternative transfers

$$t_i = v(a, \theta_i) - \delta - (v_i(f_a(\theta), \theta_i) - f_{ti}(\theta)).$$

Then $\sum_{i \in \mathcal{I}} t_i = \sum_{i \in \mathcal{I}} f_{ti}(\theta)$ but $v_i(a, \theta_i) - t_i = v_i(f_a(\theta), \theta_i) - f_{ti}(\theta) + \delta > v_i(f_a(\theta), \theta_i) - f_{ti}(\theta)$, and so f cannot be Pareto optimal.

Conversely, suppose f_a is an efficient decision rule. Let $t : \Theta \rightarrow \mathbb{R}^n$ be an arbitrary transfer rule. Towards contradiction, suppose (f_a, t) is not Pareto optimal. Then there is some type profile $\theta \in \Theta$ s.t. we have $f_a(\theta) \in \arg \max_{a \in A} \sum_{i \in \mathcal{I}} v_i(a, \theta_i)$ and there exists alternative $a \in A$ and transfers $q \in \mathbb{R}^n$ s.t. $\sum_{i \in \mathcal{I}} q_i \geq \sum_{i \in \mathcal{I}} t_i$ and $v_i(a, \theta_i) - q_i \geq v_i(f_a(\theta), \theta_i) - t_i$ for all $i \in \mathcal{I}$ with strict inequality for some i . Summing across agents gives

$$\sum_{i \in \mathcal{I}} v(a, \theta_i) - \sum_{i \in \mathcal{I}} q_i > \sum_{i \in \mathcal{I}} v(f_a(\theta), \theta_i) - \sum_{i \in \mathcal{I}} f_{ti}(\theta),$$

and so, since $\sum_{i \in \mathcal{I}} q_i \geq \sum_{i \in \mathcal{I}} t_i$, it follows that $\sum_{i \in \mathcal{I}} v_i(a, \theta_i) > \sum_{i \in \mathcal{I}} v_i(f_a(\theta), \theta_i)$, contradicting that f_a is an efficient decision rule. \square

Intuitively, from any given outcome, transfers allow any Kaldor-Hicks improvement to be realized by compensating those who would otherwise be left worse off. Hence the only outcomes that are Pareto optimal in this setting involve Kaldor-Hicks optimal decisions, implying an ex post efficient decision rule.⁴⁴

All this talk of efficiency begs the question: can we always implement an efficient decision rule in this setting? The answer is yes: the class of direct mechanisms known as *Groves mechanisms*, due to Groves (1973, 1977), implement the efficient decision rule in dominant strategies.

Definition 64 (Groves mechanism). Given a QPV mechanism design problem $(\mathcal{I}, A \times \mathbb{R}^n, (\theta_i, u_i)_{i \in \mathcal{I}})$, let $f_a : \Theta \rightarrow A$ be an efficient decision rule.

A direct mechanism $f = (f_a, f_t)$ is a *Groves mechanism* if the transfer function f_t satisfies

$$f_{ti}(\theta) = h_i(\theta_{-i}) - \sum_{j \neq i} v_j(f_a(\theta_i, \theta_{-i}), \theta_j)$$

for every type profile $\theta \in \Theta$ and every agent $i \in \mathcal{I}$, where $h_i : \Theta_{-i} \rightarrow \mathbb{R}$ is some function for each i .

The transfer function in a Groves mechanism has two components for each agent. The first, h_i , is some function of the types of other agents only. This function might be specific to agent i , but it is independent of the type θ_i that agent i chooses to report. The second component, $\sum_{j \neq i} v_j(f_a(\theta), \theta_j)$, sums the valuations of the other agents.

In the transfer rules

$$f_{ti}(\theta) = h_i(\theta_{-i}) - \sum_{j \neq i} v_j(f_a(\theta_i, \theta_{-i}), \theta_j),$$

Groves mechanisms impose no restrictions on $h_i(\theta_{-i})$: this is an arbitrary function for each i , and so there are uncountably many mechanisms within the class of Groves mechanisms. However, a particularly noteworthy special case within this class is the *Vickrey-Clarke-Groves mechanism* (also called the *pivotal mechanism*), owing its name to Vickrey (1961), Clarke (1971), and Groves (1973).

Definition 65 (Vickrey-Clarke-Groves mechanism). A Groves mechanism $f = (f_a, f_t)$ is the *Vickrey-Clarke-Groves (VCG) mechanism* if for each $i \in \mathcal{I}$ and type profile $\theta_{-i} \in \Theta_{-i}$,

$$h_i(\theta_{-i}) = \max_{a \in A} \sum_{j \neq i} v_j(a, \theta_j).$$

That is,

$$f_{ti}(\theta) = \max_{a \in A} \sum_{j \neq i} (v_j(a, \theta_j) - v_j(f_a(\theta), \theta_j)) \quad \text{for all } \theta \in \Theta.$$

⁴⁴ An outcome x is a *Kaldor-Hicks improvement* relevant to an outcome x' if with appropriate compensatory transfers, x Pareto dominates x' (the transfers need not be realized). If there is no Kaldor-Hicks improvement on an outcome, then we say it is *Kaldor-Hicks optimal*.

The transfer rule in the VCG mechanism is sometimes called the *Clarke pivot rule*. This rule has a nice intuitive interpretation: the transfer made by agent i in the Vickrey-Clarke-Groves mechanism amounts to the net externality i imposes on the other agents, i.e. the difference between the total value would receive from the efficient decision excluding i and the total value from the efficient decision including i . Moreover, the payoff to each agent in the VCG mechanism under truthful reporting is equal to their marginal contribution. That is, if $\theta \in \Theta$ is the type profile of agents, agent i 's payoff is the difference between total welfare with i present and total welfare were i not to participate (or exist):

$$\begin{aligned} v_i(f_a(\theta), \theta_i) - f_{ti}(\theta) &= v_i(f_a(\theta), \theta_i) - \max_{a \in A} \sum_{j \neq i} (v_j(a, \theta_j) - v_j(f_a(\theta), \theta_j)) \\ &= \max_{a \in A} \sum_{j \in \mathcal{I}} v_j(a, \theta_j) - \max_{b \in A} \sum_{j \neq i} v_j(b, \theta_j). \end{aligned}$$

Since agent i 's report of their own type does not directly affect the transfer they make/receive in a Groves mechanism, and the Pareto optimal outcome is implemented, there is no incentive for agents to misreport their types. Thus Groves mechanisms are dominant strategy incentive compatible:

Theorem 18 (Groves, 1977). *Every Groves mechanism is strategyproof.*

Proof. Let $f = (f_a, f_t)$ be a Groves mechanism. Fix an agent $i \in \mathcal{I}$, a type profile $\theta_{-i} \in \Theta_{-i}$, and any two types $\theta_i, \theta'_i \in \Theta_i$. Let $a = f_a(\theta_i, \theta_{-i})$ and $b = f_a(\theta'_i, \theta_{-i})$. Now,

$$\begin{aligned} v_i(a, \theta_i) - f_{ti}(\theta_i, \theta_{-i}) &= v_i(a, \theta_i) - \left[h_i(\theta_{-i}) - \sum_{j \neq i} v_j(a, \theta_j) \right] \\ &= \sum_{j \in \mathcal{I}} v_j(a, \theta_j) - h_i(\theta_{-i}) \\ &\geq \sum_{j \in \mathcal{I}} v_j(b, \theta_j) - h_i(\theta_{-i}) \\ &= v_i(b, \theta_i) - \left[h_i(\theta_{-i}) - \sum_{j \neq i} v_j(b, \theta_j) \right] \\ &= v_i(b, \theta_i) - f_{ti}(\theta'_i, \theta_{-i}), \end{aligned}$$

the inequality holding by efficiency of f_a . □

Green & Laffont (1977) prove a partial converse. We follow Holmström's (1979) generalization. First we need to set up some restrictions on preferences.

Definition 66 (Smooth path-connectedness). A set X is *smoothly path-connected* if for any points $x, x' \in X$, there exists a differentiable function $\tau : [0, 1] \rightarrow X$ such that $\tau(0) = x$ and $\tau(1) = x'$.⁴⁵

⁴⁵In general, a function $\tau : [0, 1] \rightarrow X$ s.t. $\tau(0) = x$ and $\tau(1) = x'$ is called a *path* between x and x' . If it is differentiable, then it is a *smooth path*.

We first prove a lemma that relies on i 's type space being $\Theta_i = [0, 1]$. We can then use this to extend to the case where Θ_i is smoothly path-connected.

Lemma 13 (Holmström, 1979). *Fix agent $i \in \mathcal{I}$ and suppose $\Theta_i = [0, 1]$. Let $f = (f_a, f_t)$ be a strategyproof direct mechanism such that f_a is an efficient decision rule. Define $V_i(\theta) = u_i(f(\theta), \theta_i) = v_i(f_a(\theta), \theta_i) - f_{ti}(\theta)$ for all $\theta \in \Theta$ and each agent i . Moreover, suppose $v_i(a, \cdot)$ is continuously differentiable for all $a \in A$.*

Then

- (i) *for every $\theta_{-i} \in \Theta_{-i}$ and $\theta_i \in \Theta_i$,*

$$V_i(\theta_i, \theta_{-i}) = V_i(0, \theta_{-i}) + \int_0^{\theta_i} \frac{\partial v_i(f_a(z, \theta_{-i}), z)}{\partial z} dz;$$

- (ii) *if $V_i(\cdot, \theta_{-i})$ is differentiable at θ_i then*

$$\frac{\partial V_i(\theta_i, \theta_{-i})}{\partial \theta_i} = \frac{\partial v_i(f_a(\theta_i, \theta_{-i}), \theta_i)}{\partial \theta_i};$$

- (iii) *the transfer rule f_t satisfies*

$$f_{ti}(\theta) = -V_i(0, \theta_{-i}) + v_i(f_a(\theta), \theta_i) - \int_0^{\theta_i} \frac{\partial v_i(f_a(z, \theta_{-i}), z)}{\partial z} dz$$

for all type profiles $\theta \in \Theta$.

Proof. Since f_a is an efficient decision rule, (i) follows directly from the Milgrom-Segal envelope theorem (Theorem 66), noting that for all $a \in A$, that $v_i(a, \cdot)$ is absolutely continuous is implied by continuous differentiability and the fact that $[0, 1]$ is compact. (ii) follows immediately from Theorem 65(iii). The third statement follows by substituting for $V_i(\theta_i, \theta_{-i})$ in (i). \square

Theorem 19 (Green-Laffont-Holmström). *For each agent $i \in \mathcal{I}$, let Θ_i be smoothly path-connected and let $v_i(a, \cdot)$ be continuously differentiable for all $a \in A$. Then any strategyproof direct mechanism implementing an efficient decision rule is a Groves mechanism.*

Proof. Let $f = (f_a, f_t)$ be a strategyproof direct mechanism such that f_a is an efficient decision rule, and let $g = (f_a, g_t)$ be the Vickrey-Clarke-Groves mechanism. Fix agent i , type profile $\theta_{-i} \in \Theta_{-i}$, and any two types $\theta_i, \theta'_i \in \Theta_i$. Let $\tau : [0, 1] \rightarrow \Theta_i$ be a smooth path between θ_i and θ'_i , i.e. a differentiable function with $\tau(0) = \theta_i$ and $\tau(1) = \theta'_i$.

Define $\hat{f}_a : [0, 1] \times \Theta_{-i} \rightarrow A$ by $\hat{f}_a(z, \theta_{-i}) = f_a(\tau(z), \theta_{-i})$ and define $\hat{g}_t : [0, 1] \times \Theta_{-i} \rightarrow \mathbb{R}^n$ by $\hat{g}_t(z, \theta_{-i}) = g_t(\tau(z), \theta_{-i})$. Let $\hat{V}_i(z, \theta_{-i}) = u_i(g(\tau(z), \theta_{-i}), \theta_i)$ for $z \in [0, 1]$, so $\hat{V}_i(0, \theta_{-i})$ is the payoff to i of type θ_i under the VCG mechanism and $\hat{V}_i(1, \theta_{-i})$ is the payoff to i of type θ'_i under the VCG mechanism. By Lemma 13,

$$\hat{g}_{ti}(1, \theta_{-i}) = -\hat{V}_i(0, \theta_{-i}) + v_i(\hat{f}_a(1, \theta_{-i}), \theta_i) - \int_0^1 \frac{\partial v_i(\hat{f}_a(z, \theta_{-i}), z)}{\partial z} dz.$$

Similarly, define $\hat{f}_t : [0, 1] \times \Theta_{-i} \rightarrow \mathbb{R}^n$ by $\hat{f}_t(z, \theta_{-i}) = f_t(\tau(z), \theta_{-i})$. Define $V_i(z, \theta_{-i}) = u_i(f(\tau(z), \theta_{-i}), \theta_i)$ for $z \in [0, 1]$, and define $h_i(\theta_{-i}) = \hat{V}_i(0, \theta_{-i}) - V_i(0, \theta_{-i})$. Again applying Lemma 13, we have

$$\begin{aligned} \hat{f}_{ti}(1, \theta_{-i}) &= -V_i(0, \theta_{-i}) + v_i(\hat{f}_a(1, \theta_{-i}), \theta_i) - \int_0^1 \frac{\partial v_i(f_a(z, \theta_i), z)}{\partial z} dz \\ &= h_i(\theta_{-i}) + \hat{g}_{ti}(1, \theta_{-i}). \end{aligned}$$

Since θ_i, θ'_i are arbitrary, this applies for all such type pairs. Hence any strategyproof mechanism f implementing an efficient allocation has a transfer rule f_t that differs from the VCG transfer rule only in some function of other agents' types for each agent i , and thus is a Groves mechanism. \square

Example 33 (VCG combinatorial auction). A canonical application of the VCG mechanism is in allocating indivisible goods to agents, in which case we call the VCG mechanism a *VCG auction*. Suppose there are three bidders, and two objects, a and b . Bidders value each object alone but also value owning both objects. Suppose the valuations are as follows:

| | \emptyset | $\{a\}$ | $\{b\}$ | $\{a, b\}$ |
|-------|-------------|---------|---------|------------|
| v_1 | 0 | 7 | 5 | 10 |
| v_2 | 0 | 10 | 3 | 12 |
| v_3 | 0 | 8 | 0 | 8 |

The efficient allocation assigns object a to Bidder 2 and object b to Bidder 1, generating total value 15. Bidder 3 is not allocated anything, and her absence would not change the efficient allocation, so her net transfer under the VCG mechanism is zero. If Bidder 2 were absent, then the efficient allocation would assign object a to Bidder 3 and b to Bidder 1, generating total value 13. The externality imposed by Bidder 2 is thus $15 - 13 = 2$, so Bidder 2 pays 2. If Bidder 1 were absent, the efficient allocation would assign both objects to Bidder 1, generating total value 12. Hence Bidder 1 pays $15 - 12 = 3$.

Corollary 8. *The VCG auction is the unique strategyproof direct mechanism implementing an efficient allocation rule such that bidders who are not assigned an object pay zero.*

Proof. By the Green-Laffont-Holmström theorem (Theorem 19), any such mechanism must be a Groves mechanism. Fix type profile $\theta \in \Theta$ and note that if i is not allocated the object then $f_{ti}(\theta) = \sum_{j \neq i} (v_j(f_a(\theta), \theta_j) - v_j(f_a(\theta), \theta_j)) = 0$, since the absence of i would not alter the efficient allocation. Any other Groves mechanism transfer function differs from VCG for some i in some function (not everywhere zero) of others' types θ_{-i} , and hence has $f_{ti}(\theta) \neq 0$ for some type profile in which i is not assigned the object. \square

In QPV settings, Groves mechanisms are a very theoretically attractive class of mechanisms: they allow the designer to implement an efficient decision rule, and because

they are strategyproof, they are ‘robust’ in the sense that they are not sensitive to changes in the distribution of agents’ types. However, there are notable drawbacks. One set of drawbacks are the “algorithmic” kind that concern computer scientists. In combinatorial allocation problems such as Example 33, there are two main algorithmic issues, stemming from the fact that as the number of objects to be assigned grows, the number of alternatives grows much faster:

- (i) *Computational complexity.* Computing the efficient allocation is an NP-complete problem. In large scale settings (such as the problem of allocating radio spectrum nationally), computing the exact efficient allocation is too computationally intensive to be practical.
- (ii) *Communication complexity.* Computing the efficient allocation and transfers is also communicationally intensive – when the number of alternatives is large, agents need to provide a lot of information to give a complete valuation schedule. For example, in a combinatorial auction with k objects, the set of possible allocations each agent might receive is the power set, so agents need to provide valuations for 2^k different alternatives. Again, this may not be practical. Agents may also have privacy concerns about giving full descriptions of their preferences (for example, if this information reveals commercially sensitive information).

Blumrosen & Nisan (2007) provide a nice overview of these issues.⁴⁶ A growing literature in algorithmic mechanism design aims to address these issues by finding approximations of Groves mechanisms that are computationally and communicatively simpler. These necessarily come at the cost of strategyproofness/efficiency.

A second set of drawbacks are those of more economic interest. While Groves mechanisms always implement the efficient decision rule, it does not follow that they are always efficient in a more general sense. In particular, suppose we impose that transfer functions must be feasible (that is, $\sum_{i \in \mathcal{I}} f_{ti}(\theta) \geq 0$ for all $\theta \in \Theta$), and we consider any positive revenue generated by the mechanism wasted. Then any feasible mechanism that is not budget-balanced will not always be efficient, even if it implements the efficient decision rule, since it will generate positive revenue for some type profiles θ : remitting this revenue to agents would improve total welfare. Unfortunately, Groves mechanisms suffer this issue: if the type spaces are sufficiently rich, there is no ex post budget balanced Groves mechanism.

5 Social choice theory

Aggregating individual preferences has long concerned economists, mathematicians and political scientists – two oft-mentioned names in the field, the Marquis de Condorcet and Jean-Charles de Borda were concerned with this all the way back in the 18th century

⁴⁶See chapter 11 in Nisan, Roughgarden, Tardos & Vazirani (eds.), *Algorithmic Game Theory*.

(the eponymous Borda count for which the latter is known is actually far older).⁴⁷ The importance of social choice theory in political economy is obvious, but it is applicable to almost all settings involving group decisionmaking. Indeed, Ken Arrow's original interest in social choice theory came not from thinking about elections but from thinking about how shareholders ought to reach a decision about production plans, given the shareholders might differ in their preferences or beliefs.

We consider a general setting as follows. Let \mathcal{I} be a set of n agents, and let $A = \{a_1, \dots, a_k\}$ denote a set of k alternatives. Each agent i has a *preference relation* \succsim_i over the set of alternatives A , that is, a binary relation on A . For any $a, b \in A$, $a \succsim_i b$ has the interpretation that i weakly prefers alternative a to alternative b . Given a preference relation \succsim_i , we define the strict preference relation \succ_i by $a \succ_i b$ if $a \succsim_i b$ but $b \not\succsim_i a$. This has the interpretation that i strictly prefers a to b . We define the preference relation \sim_i by $a \sim_i b$ if $a \succsim_i b$ and $b \succsim_i a$. This has the interpretation that i is indifferent between a and b .

We call \succsim_i an *ordering* if it is complete, reflexive and transitive (i.e. a transitive preorder). We denote the set of orderings over A by \tilde{R} , and the set of preference relations over A by \tilde{R} . We call $(\succsim_1, \dots, \succsim_n) \in R^n$ a *preference profile*. We call an ordering \succsim_i *linear* if it is antisymmetric. We call \succsim_i a *quasi-ordering* if the corresponding strict preference relation \succ_i is transitive.

A *social choice rule* is a mapping $f : R^n \rightarrow \tilde{R}$, and an *Arrovian social welfare function* is a mapping $f : R^n \rightarrow R$. Hence an Arrovian social welfare function is a social choice rule that always returns an ordering for any preference profile over A . The following gives some well-known social choice rules.

Example 34 (Social choice rules).

- (a) *Scoring rules.* Suppose \succsim_i is a linear ordering for each agent i . A *scoring vector* for the set of alternatives $A = \{a_1, \dots, a_k\}$ is a vector $s = (s_1, \dots, s_k) \in \mathbb{R}^k$ such that $s_1 \geq s_2 \geq \dots \geq s_k \geq 0$. For each agent $i \in \mathcal{I}$, for each ordering $\succsim_i \in R$, and for each alternative $a \in A$, define the *rank* of a in \succsim_i by

$$r(a, \succsim_i) = 1 + |\{b \in A \mid b \succ_i a\}|.$$

We define the *score* of rank $r(a, \succsim_i)$ by $s_{r(a, \succsim_i)}$, and for each profile $(\succsim_i) \in R^n$, we define the score of alternative a by

$$s(a, \succsim_1, \dots, \succsim_n) = \sum_{i \in \mathcal{I}} s_{r(a, \succsim_i)}.$$

A social choice rule $f : R^n \rightarrow \tilde{R}$ is a *scoring rule* if for all $a, b \in R$ and all $(\succsim_i)_{i \in \mathcal{I}} \in R^n$, the preference relation $\succsim := f(\succsim_1, \dots, \succsim_n)$ satisfies $a \succ b$ iff $s(a, \succsim_1, \dots, \succsim_n) \geq s(b, \succsim_1, \dots, \succsim_n)$.

⁴⁷The Borda count is used to elect an abbe in Ramon Llull's novel *Blanquerna*, written 1283-1285, and Nicholas of Cusa argued for using the Borda count to elect the Holy Roman Emperor in 1433. Borda proposed it as a way to elect members to *l'Académie des Sciences* in 1784. You don't need to be original to claim an idea.

- (i) *Plurality rule.* A scoring rule f is a *plurality rule* if the scoring vector is $s = (1, 0, 0, \dots, 0)$.
- (ii) *Anti-plurality rule.* A scoring rule f is an *anti-plurality rule* if the scoring vector is $s = (1, 1, \dots, 1, 0)$.
- (iii) *Borda rule.* A scoring rule f is a *Borda rule* if the scoring vector is $s = (k-1, k-2, \dots, 1, 0)$.

Scoring rules always define an ordering, and hence are Arrovian social welfare functions.

- (b) *Majority rules.* A social choice rule $f : R^n \rightarrow \tilde{R}$ is a *majority rule* if for all $a, b \in R$ and all $(\succsim_i)_{i \in \mathcal{I}} \in R^n$, the preference relation $\succsim := f(\succsim_1, \dots, \succsim_n)$ satisfies $a \succsim b$ iff

$$|\{i \in \mathcal{I} \mid a \succsim_i b\}| \geq |\{i \in \mathcal{I} \mid b \succsim_i a\}|.$$

Majority rules do not necessarily define an ordering, as we prove below.

- (c) *Oligarchic rules.* Given the set of agents \mathcal{I} , let some nonempty subset $G \subseteq \mathcal{I}$ be an *oligarchy*. A social choice rule $f : R^n \rightarrow \tilde{R}$ is an *oligarchic* social choice rule if for all $a, b \in A$ and all $(\succsim_i) \in R^n$, the preference relation $\succsim := f(\succsim_1, \dots, \succsim_n)$ satisfies $a \succsim b$ iff $a \succsim_i b$ for some $i \in G$.

If G is a singleton, i.e. $G = \{i\}$ for some $i \in \mathcal{I}$, then we call an oligarchic social choice rule *monarchic*.

Proposition 50 (Condorcet's paradox). *If f is a majority rule and $|\mathcal{I}| \geq 3$, then there exists a preference profile (\succsim_i) for which f does not define a quasi-ordering, and thus does not define an ordering.*

Proof. Take $\mathcal{I} = \{1, 2, 3\}$ and $A = \{a_1, a_2, a_3\}$. Consider the following strict preference relations for the agents:

$$\begin{array}{ccc} \succsim_1 & \succsim_2 & \succsim_3 \\ a_1 & a_3 & a_2 \\ a_2 & a_1 & a_3 \\ a_3 & a_2 & a_1 \end{array}$$

Each agent's preferences are a linear ordering. Let $\succsim := f(\succsim_1, \succsim_2, \succsim_3)$. Define $m(a, b) := |\{i \in \mathcal{I} \mid a \succsim_i b\}|$. We have $m(a_1, a_2) = m(a_2, a_3) = m(a_3, a_1) = 2$, while $m(a_2, a_1) = m(a_3, a_2) = m(a_1, a_3) = 1$. Hence we have $a_1 \succ a_2$, $a_2 \succ a_3$ and $a_3 \succ a_1$, which violates transitivity, so \succ is not a quasi-ordering (so also not an ordering).

To generalize to $|\mathcal{I}| > 3$, take all other agents to be indifferent over a_1, a_2, a_3 . \square

In the proof of Proposition 50, we showed we can get a situation where, given three alternatives a, b, c , majorities strictly prefer a over b , b over c and c over a . In such a situation, $\{a, b, c\}$ is called a *Condorcet cycle*.

Proposition 51. *Suppose f is an oligarchic social choice rule. Then for any preference profile (\succsim_i) , $f(\succsim_1, \dots, \succsim_n)$ is a quasi-ordering. Moreover, if $|A| \geq 3$, then $f(\succsim_1, \dots, \succsim_n)$ is an ordering for all preference profiles $(\succsim_i) \in R^n$ iff f is monarchic.*

Proof. If $|A| < 3$ then the proposition holds trivially. Hence suppose $|A| \geq 3$ and consider $a_1, a_2, a_3 \in A$. Suppose G is the oligarchy corresponding to f . Fix a preference profile $(\succsim_i) \in R^n$, and define $\succsim := f(\succsim_1, \dots, \succsim_n)$. Let $a_1 \succ a_2$ and $a_2 \succ a_3$. Then it must be that $a_1 \succ_i a_2$ and $a_2 \succ_i a_3$ for all $i \in G$, and since \succ_i is transitive for all i , we have that $a_1 \succ_i a_3$ for all $i \in G$, which gives $a_1 \succ a_3$, and so \succ is transitive. Thus \succsim is a quasi-ordering.

Suppose $G = \{i\}$ for some $i \in \mathcal{I}$. Then $a \succ b$ iff $a \succ_i b$ and $a \sim b$ iff $a \sim_i b$, so $\succsim = \succsim_i$, and the rhs is an ordering. Suppose instead that $|G| \geq 2$ and consider $a_1, a_2, a_3 \in A$. Consider the following strict preference relations for $G = \{1, 2\}$:

$$\begin{array}{cc} \succ_1 & \succ_2 \\ a_1 & a_3 \\ a_2 & a_1 \\ a_3 & a_2 \end{array}$$

We have $a_1 \succ a_2$, $a_2 \sim a_3$ and $a_3 \sim a_1$. If transitivity were to hold, we would require that $a_1 \sim a_2$, yielding a contradiction. Hence \succsim is not an ordering. This generalizes easily to $|G| > 3$. \square

5.1 Arrow's impossibility theorem

Arrow (1950, 1963) famously proved a general impossibility result that preferences cannot generally be aggregated in a way consistent with a collection of attractive (but not necessarily uncontroversial) axioms:

Axioms (Arrow axioms). Let f be an Arrowian social welfare function on R^n . Given any preference profile (\succsim_i) , take $\succsim := f(\succsim_1, \dots, \succsim_n)$.

- (P1) *Weak Pareto efficiency.* We say f satisfies *weak Pareto efficiency* (or *unanimity*) if for all $(\succsim_i) \in R^n$ and all $a, b \in A$, $a \succ_i b$ for all i implies $a \succ b$.
- (P2) *Independence of irrelevant alternatives.* Given preference profiles $(\succsim_i), (\succsim'_i) \in R^n$, we say the two profiles agree on $\{a, b\} \subseteq A$ if for all $i \in \mathcal{I}$ we have (i) $a \succ_i b$ iff $a \succ'_i b$ and (ii) $a \sim_i b$ iff $a \sim'_i b$. We say f satisfies *independence of irrelevant alternatives* if whenever (\succsim_i) and (\succsim'_i) agree on $\{a, b\}$, we have that $a \succsim b$ iff $a \succsim' b$, where $\succsim := f(\succsim_1, \dots, \succsim_n)$ and $\succsim' := f(\succsim'_1, \dots, \succsim'_n)$.
- (P3) *Non-dictatorship.* We say an agent $i \in \mathcal{I}$ is a *dictator* if for all $a, b \in A$ and for all preference profiles $(\succsim_i) \in R^n$, we have that if $a \succ_i b$ then $a \succ b$. We call f *nondictatorial* if there does not exist a dictator $i \in \mathcal{I}$.

A fourth axiom – *unrestricted domain*, also called *universality* – is often given in the context of Arrow's theorem. This states that for any preference profile $(\succsim_i) \in R^n$,

f should yield a unique ordering. In our setting, this axiom is encoded already in the definition of an Arrovian social welfare function.

Some comments are in order. First, note that if i is a dictator, it does not necessarily follow that $\succsim = \succsim_i$. For example, if $a \sim_i b$, then we could have that $a \succ b$. If \succsim_i is antisymmetric, then this possibility cannot arise, since $a \sim_i b$ implies $a = b$. Second, independence of irrelevant alternatives implies that given two candidates, agents' preferences over third candidates should not alter the ranking of the two candidates under f .

Theorem 20 (Arrow's impossibility theorem). *If $|A| \geq 3$, then there is no Arrovian social welfare function satisfying (P1)-(P3).*

Proof. Let f be an Arrovian social welfare function satisfying (P1) and (P2). Call a nonempty group of agents $G \subseteq \mathcal{I}$ *decisive* for alternatives $a, b \in A$ if for all profiles $(\succsim_i) \in R^n$, we have that $a \succ_i b$ for all $i \in G$ implies $a \succ b$, where $\succ := f(\succsim_1, \dots, \succsim_n)$. Call G *almost decisive* for $a, b \in A$ if for all profiles $(\succsim_i) \in R^n$, we have that if $a \succ_i b$ for all $i \in G$ and $b \succ_j a$ for all $j \notin G$, then $a \succ b$. Obviously, if G is decisive for a, b then it is also almost decisive for a, b .

Lemma 14. *For all nonempty subsets $G \subseteq \mathcal{I}$ and all alternatives $a, b, c, d \in A$, if G is almost decisive for a, b then it is decisive for c, d .*

Proof. There are seven cases to consider:

- (i) Suppose $a \neq b \neq c \neq d$. Take $(\succsim'_i) \in R^n$ s.t. $c \succ'_i d$ for all $i \in G$ and $(\succsim_i) \in R^n$ s.t. $c \succ_i a \succ_i b \succ_i d$ for all $i \in G$. Impose for $j \notin G$ that $c \succ_j a$, $b \succ_j d$, $b \succ_j a$ and \succsim_j and \succsim'_j agree on $\{c, d\}$. Then if G is almost decisive for a, b , it follows that $a \succ b$. Now (P1) implies $c \succ a$ and $b \succ d$. Transitivity gives $c \succ d$. Now (\succsim_i) and (\succsim'_i) agree on $\{c, d\}$, so (P2) implies $c \succ' d$. Thus G is decisive for c, d .
- (ii) Suppose $c \neq a \neq b$ and $b = d$. Take $(\succsim'_i), (\succsim_i) \in R^n$ so that $c \succ'_i b$ and $c \succ_i a \succ_i b$, for all $i \in G$. Impose for $j \notin G$ that $c \succ_j a$, $b \succ_j a$ and \succsim_j and \succsim'_j agree on $\{c, b\}$. Then if G is almost decisive for a, b , it follows that $a \succ b$. Now (P1) implies $c \succ a$, transitivity gives $c \succ b$, and (P2) implies $c \succ' b$. Thus G is decisive for c, b .
- (iii) Suppose $d \neq a \neq b$ and $c = b$. Take $(\succsim'_i), (\succsim_i) \in R^n$ so that $a \succ'_i d$ and $a \succ_i b \succ_i d$, for all $i \in G$. Impose for $j \notin G$ that $b \succ_j d$, $b \succ_j a$ and \succsim_j and \succsim'_j agree on $\{a, d\}$. Then if G is almost decisive for a, b , it follows that $a \succ b$. Now (P1) implies $b \succ d$, transitivity gives $a \succ d$, and (P2) implies $a \succ' d$. Thus G is decisive for a, d .
- (iv) Suppose $d \neq a \neq b$ and $c = a$. (iii) implies that if G is almost decisive for a, b then G is decisive for a, d and thus almost decisive for a, d . From (ii) it follows that G is decisive for b, d .
- (v) Suppose $c \neq a \neq b$ and $d = a$. (ii) implies that if G is almost decisive for a, b then it is decisive for c, b so almost decisive for c, b . From (iii) it follows that G is decisive for c, a .

- (vi) Suppose $c = a$ and $d = b$. Consider $d \neq a \neq b$. (iii) implies that if G is almost decisive for a, b then G is decisive for a, d and thus almost decisive for a, d . By (iii) again, it follows that G is decisive for a, b .
- (vii) Suppose $c = b$ and $d = a$. Consider $d \neq a \neq b$. By (v), if G is almost decisive for a, b then G is decisive for d, a and thus almost decisive for d, a . By (ii), it follows that G is decisive for b, a .

□

Lemma 15. *Suppose a nonempty subset $G \subseteq \mathcal{I}$ is decisive. If $|G| \geq 2$, then there is some nonempty proper subset of G that is also decisive.*

Proof. Let G be the disjoint union of two nonempty sets $G_1, G_2 \subset \mathcal{I}$. Let $a, b, c \in A$ and consider a preference profile $(\succsim_i) \in R^n$ such that $a \succ_i b \succ_i c$ for all $i \in G_1$, $c \succ_j a \succ_j b$ for all $j \in G_2$, and $b \succ_k c \succ_k a$ for all $k \in \mathcal{I} - G$. Since G is decisive and $a \succ_i b$ for all $i \in G$, we have that $a \succ b$. Now there are two cases to consider:

- (i) Suppose $a \succ c$. Note $a \succ_i c$ for all $i \in G_1$ but $c \succ_j a$ for all $j \in \mathcal{I} - G_1$. Hence G_1 is almost decisive for a, c , and so by Lemma 14, G_1 is decisive.
- (ii) Suppose $c \succsim a$. Since $a \succ b$, we have $c \succ b$ by transitivity. Note $c \succ_i b$ for all $i \in G_2$ and $b \succ_j c$ for all $j \in \mathcal{I} - G_2$. Hence G_2 is almost decisive for c, b , and so by Lemma 14, G_2 is decisive.

□

Now by (P1), \mathcal{I} is decisive. Iteratively applying Lemma 15 yields some agent $i \in \mathcal{I}$ such that $\{i\}$ is decisive, and thus i is a dictator, contradicting (P3). □

Arrow's impossibility theorem is sometimes wrongly interpreted as implying ranked voting systems are inherently flawed. It does not say this – it simply says that no ranked voting system for deciding between three or more alternatives satisfies certain attractive properties *all of the time*.

Wilson (1972) proves a slightly stronger version of Arrow's theorem. We need some new properties:

Axioms. Let f be an Arrovian social welfare function on R^n .

(P4) *Non-imposition.* We say f satisfies *non-imposition* if for all $a, b \in A$, there exists some profile $(\succsim_i) \in R^n$ such that $a \succsim b$ for $\succsim := f(\succsim_1, \dots, \succsim_n)$.

Definition 67. Let f be an Arrovian social welfare function on R^n . Given any preference profile (\succsim_i) , take $\succsim := f(\succsim_1, \dots, \succsim_n)$.

- (a) *Anti-dictatorship.* We say an agent $i \in \mathcal{I}$ is an *anti-dictator* if for all $a, b \in A$ and for all preference profiles $(\succsim_i) \in R^n$, we have that $a \succ_i b$ implies $b \succ a$. We say f is *anti-dictatorial* if there exists an anti-dictator.

(b) *Dictatorship*. We call f *dictatorial* if there exists a dictator.

(c) *Nullity*. We say f is *null* if $a \sim b$ for all $a, b \in A$ and all $(\succsim_i) \in R$.

Theorem 21 (Wilson's theorem; Wilson, 1972; Malawski & Zhou, 1994). *If $|A| \geq 3$ then any Arrovian social welfare function satisfying (P2) and (P4) must be anti-dictatorial, dictatorial or null.*

Proof. Suppose f is an Arrovian social welfare function satisfying (P2) and (P4). Say f is *Pareto consistent* for alternatives $a, b \in A$ if $a \succ_i b$ for all $i \in \mathcal{I}$ implies $a \succ b$. Say f is *Pareto anti-consistent* for $a, b \in A$ if $a \succ_i b$ for all $i \in \mathcal{I}$ implies $b \succ a$.

Lemma 16. *For all $a, b, c, d \in A$, if f is Pareto consistent for a, b then f is Pareto consistent for c, d .*

Proof. As in Lemma 14, we have many cases. We prove only the case where $c = a$ and $a \neq b \neq d$. Fix any profile (\succsim_i) s.t. $a \succ_i d$ for all $i \in \mathcal{I}$. By (P4), there exists a profile (\succsim'_i) s.t. $b \succ'_i d$. Define a profile $(\tilde{\succsim}_i)$ so $a \tilde{\succsim}_i b$ and $a \tilde{\succsim}_i d$ for all i , and $(\tilde{\succsim}_i)$ and (\succsim'_i) agree on b, d . Since f is Pareto consistent for a, b , $a \tilde{\succ} b$. Note that by (P2), $b \tilde{\succ} d$. Since $\tilde{\succsim}$ is transitive, $a \tilde{\succ} d$, and so (P2) implies $a \succ d$. \square

Lemma 17. *For all $a, b, c, d \in A$, if f is Pareto anti-consistent for a, b then f is Pareto anti-consistent for c, d .*

Proof. Again, there are many cases to check, and we prove only the case where $c = a$ and $a \neq b \neq d$. Fix any profile (\succsim_i) s.t. $a \succ_i d$ for all $i \in \mathcal{I}$. By (P4), there exists a profile (\succsim'_i) s.t. $d \succ'_i b$. Define a profile $(\tilde{\succsim}_i)$ so $a \tilde{\succsim}_i b$ and $a \tilde{\succsim}_i d$ for all i , and $(\tilde{\succsim}_i)$ and (\succsim'_i) agree on b, d . Since f is Pareto anti-consistent for a, b , $b \tilde{\succ} a$. By (P2), $d \tilde{\succ} b$. Since $\tilde{\succsim}$ is transitive, $d \tilde{\succ} a$, and hence (P2) implies $d \succ a$. \square

Lemma 18. *At least one of the following holds:*

- (i) f is null,
- (ii) there is some pair of alternatives a, b for which f is Pareto consistent, or
- (iii) there is some pair of alternatives a, b for which f is Pareto anti-consistent.

Proof. Towards contradiction, suppose none of (i), (ii) or (iii) hold. Since f is not null, there is some pair of alternatives a, b and some profile (\succsim_i) s.t. $a \succ b$. Choose $c \neq a, b$. Now choose a profile (\succsim'_i) s.t. $a \succ'_i c$ and $b \succ'_i c$ for all $i \in \mathcal{I}$ and (\succsim_i) and (\succsim'_i) agree on a, b . Since f is neither Pareto consistent nor Pareto anti-consistent for a, c , we have $a \sim' c$, and since f is neither Pareto consistent nor Pareto anti-consistent for b, c , we have $b \sim' c$. But since \succsim' is transitive, it follows that $a \sim' b$. By (P2), $a \sim b$, which contradicts $a \succ b$. \square

Now suppose f is not null. Then Lemma 18 implies f is either Pareto consistent or Pareto anti-consistent for some pair of alternatives a, b . If the former case holds, then weak Pareto efficiency (P1) holds and so Arrow's Theorem (Theorem 20) implies there is a dictator. If the latter case holds, a slight modification to Arrow's proof shows there is an anti-dictator. \square

Wilson's theorem relaxed the weak Pareto efficiency axiom – an anti-dictatorial social welfare function is obviously not Pareto efficient.

5.1.1 How much of a problem is Arrow's theorem?

Is Arrow's theorem a big problem for social choice theory, for democracy? The answer is “not really”. Firstly, if there are only two alternatives, majority rule has nice properties.

Axioms. Fix a set \mathcal{I} of n agents and a set of alternatives A such that $|A| = 2$. Let f be an Arrovian social welfare function over A .

- (P5) *Anonymity.* We say f satisfies *anonymity* if for any permutation $\sigma \in \Pi(\mathcal{I})$ and any preference profile $(\succsim_i)_{i=1}^n \in R^n$, we have $f(\succsim_1, \dots, \succsim_n) = f(\succsim_{\sigma(1)}, \dots, \succsim_{\sigma(n)})$.
- (P6) *Neutrality.* Given any preference relation \succsim , define by \succsim_- the preference relation such that $x(-\succsim)y$ iff $y \succsim x$. We say f satisfies *neutrality* if for any preference profile $(\succsim_i)_{i=1}^n \in R^n$, $f(-\succsim_1, \dots, -\succsim_n) = -f(\succsim_1, \dots, \succsim_n)$.
- (P7) *Strict monotonicity.* For any preference profile $(\succsim_i)_{i=1}^n \in R^n$, any pair of distinct alternatives $a, b \in A$ and any agent $i \in \mathcal{I}$, define $\tilde{\succsim}_i$ so that $a \tilde{\succsim}_i b$ if $b \succ_i a$ and $a \tilde{\succsim}_i b$ if $a \sim_i b$, and let $\tilde{\succsim}_j = \succsim_j$ for all $j \neq i$. Then we say the preference profile $(\tilde{\succsim}_i)_{i=1}^n$ *monotonically dominates* $(\succsim_i)_{i=1}^n$ for the pair a, b . Let $\tilde{\succsim} := f(\tilde{\succsim}_1, \dots, \tilde{\succsim}_n)$. We say f satisfies *strict monotonicity* if $a \succ b$ implies $a \tilde{\succ} b$ for any preference profile $(\tilde{\succsim}_i)_{i=1}^n$ that monotonically dominates $(\succsim_i)_{i=1}^n$ for a, b .

Theorem 22 (May, 1952). *If $|A| = 2$ then an Arrovian social welfare function f satisfies axioms (P5)-(P7) iff it is a majority rule.*

Proof. First, suppose f is a majority rule, so $a \succ b$ iff $|\{i \in \mathcal{I} \mid a \succsim_i b\}| \geq |\{i \in \mathcal{I} \mid b \succsim_i a\}|$. Then that f satisfies (P5)-(P7) is obvious.

Conversely, suppose f satisfies (P5)-(P7). Towards contradiction, suppose f is not a majority rule, i.e. $|\{i \in \mathcal{I} \mid a \succsim_i b\}| > |\{i \in \mathcal{I} \mid b \succsim_i a\}|$ but $b \succ a$. Let $n_a = |\{i \in \mathcal{I} \mid a \succsim_i b\}|$ and $n_b = |\{i \in \mathcal{I} \mid b \succsim_i a\}|$. Then $n_b < n_a$. Define $(\tilde{\succsim}_i)_{i=1}^n$ as follows. First let $\tilde{\succsim}_i = \succsim_i$ for all i . Pick any agent i for whom $a \succsim_i b$, and redefine $\tilde{\succsim}_i$ so that $b \tilde{\succsim}_i a$. Repeat this process until $|\{i \in \mathcal{I} \mid b \tilde{\succsim}_i a\}| = n_a$. That is, $(\tilde{\succsim}_i)$ is a preference profile in which some of the agents who originally preferred a instead prefer b . Now, strict monotonicity (P7) implies that $b \tilde{\succ} a$. Anonymity (P5) implies that f is invariant to permuting the identities of the agents, and thus we have that $b \succ a$ whenever $|\{i \in \mathcal{I} \mid b \succsim_i a\}| = n_a$. Now neutrality (P6) implies that if $|\{i \in \mathcal{I} \mid a \succsim_i b\}| = n_a$ then $a \succ b$, yielding a contradiction. \square

This is a comforting result because in many settings, voting concerns binary decisions – whether to approve a motion at a shareholder or council meeting, voting in referendums, whether to find a defendant guilty or not guilty, and so on.⁴⁸ It is nice to know (but also not surprising) that majority voting is sensible in such settings.

Corollary 9. *If $|A| = 2$ then there exists an Arrovian social welfare function satisfying (P1)-(P3), namely, majority rule.*

Proof. Strict monotonicity (P7) implies weak Pareto efficiency (P1), and if $|A| = 2$ then independence of irrelevant alternatives (P2) holds vacuously for any Arrovian social welfare function. Now, majority rule is non-dictatorial, so satisfies (P3), and by May’s theorem (Theorem 22) satisfies (P7) and hence (P1). \square

Of course, in most settings we face more than two alternatives, and there is no attractive way of reducing larger sets of alternatives to binary alternatives. For example, in a political setting, we might imagine that we can comply with Arrow’s theorem if we have a two party system, so there are only two candidates in any election. In that case, May’s theorem tells us majority rule performs well. But this just defers the problem up the chain – the two parties need to select candidates from a wide pool. In the United States, this is typically via competitive *primaries* involving many candidates, so Arrow’s theorem applies at the primaries stage. The way in which we reduce the choice set can also affect the outcome: a designer can thus manipulate the outcome by choosing appropriate binary alternatives.

Example 35. Consider the setting in the proof of Proposition 50 where we have a Condorcet cycle. Recall the set of alternatives was $A = \{a_1, a_2, a_3\}$ and the preferences of the three agents were

| \succ_1 | \succ_2 | \succ_3 |
|-----------|-----------|-----------|
| a_1 | a_3 | a_2 |
| a_2 | a_1 | a_3 |
| a_3 | a_2 | a_1 |

Suppose a designer wishes to restrict the choice to a set $A' \subset A$ consisting of only two alternatives. Then the designer can manipulate the outcome of majority rule to achieve any alternative as the most preferred outcome. For example, selecting $A' = \{a_1, a_2\}$, we get that a_1 is the most preferred alternative since two agents strictly prefer a_1 to a_2 . Selecting $A' = \{a_2, a_3\}$, we get that a_2 is the most preferred alternative, and selecting $A' = \{a_1, a_3\}$, we have that a_3 is the most preferred alternative.

Supposing we have at least three alternatives, we can get around Arrow’s theorem only by relaxing some of the axioms or changing the setting. Weak Pareto efficiency

⁴⁸It is worth noting that in some jurisdictions, juries are not presented with binary verdicts. For example, the Scottish legal system, as of 2023, affords juries three verdicts: guilty, not guilty and *not proven*. A “not proven” verdict typically signals that the jury considers the defendant likely to be guilty but does not consider the evidence presented by the prosecutor strong enough that the jury considers it beyond reasonable doubt that the defendant is guilty. In this case, Arrow’s theorem applies.

is clearly a desirable property, as is non-dictatorship. This leaves two options: we can either drop independence of irrelevant alternatives, or move away from Arrovian social welfare functions on the set of all preference profiles R^n .

In the literature, considering narrower domains than R^n has proven popular. In many settings, we typically think of preferences as naturally quite structured. In politics, for example, most voters' preferences align along a small number of dimensions (*cleavages*, to use the political science terminology). This restricts the set of plausible preferences voters might have. For example, consider a set of policies that differ in the amount of redistribution they propose (suppose there are not meaningful differences in other dimensions). It is quite odd to imagine that there are many voters who would prefer both very high levels of redistribution and no redistribution over a moderate level of redistribution.⁴⁹ This motivates a restriction to *single-peaked* preferences in political contexts. In other settings where transfers can be made, it is common to restrict to quasilinear preferences, as discussed extensively in the previous chapter. Proposals to take account of intensity of preference, such as quadratic voting, generally fit within this paradigm.

The other main approach in the spirit of “dropping the Arrovian social welfare function” is to consider *social welfare lotteries* instead of social welfare functions, that is, to allow the “social” preference relation to be random. Unfortunately, Arrow's result is deep enough that it effectively extends to social welfare lotteries. Barberá & Sonnenschein (1978) and Pattanaik & Peleg (1986), among others, generalize Arrow's theorem to this setting. The best we can do in a social welfare lottery setting is some form of weighted random dictatorship (that is, dictatorship power is probabilistically divided among agents). Despite not escaping Arrow's result, random dictatorship is much more attractive than deterministic dictatorship. We can think of deterministic non-dictatorship as a certain kind of fairness or democratic criterion. But arguably, random dictatorship satisfies a notion of fairness – indeed random oligarchy is a common practice in the real world, considered fair. The practice of selecting groups of decisionmakers at random is *sortition*. In the modern era, sortition only applies to the selection of juries: in common law jurisdictions, eligible adults are called to jury service at random. In ancient Athens, most officeholders were selected by a lottery among those eligible who nominated themselves, and a similar system was in place in parts of 10th century Tamil Nadu and in mediaeval Florence. In fact, the original notion of democracy in ancient Athens referred to sortition, not election.

The other option – relaxing independence of irrelevant alternatives – is also a pretty attractive move. It's worth asking what independence of irrelevant alternatives buys us. Maskin on spoilers. Kreps

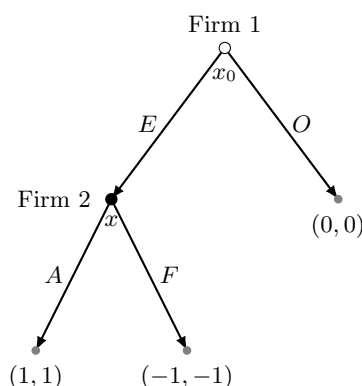
⁴⁹No doubt there are some people who would have such preferences. E.g. maybe you want to spark a communist revolution. There are three options – communism, liberal/social democracy, and reactionary conservatism. You most prefer communism, but see social democracy or liberalism as placating the proletariat and preventing revolution, so you would prefer reactionary conservatism over the middle option on accelerationist grounds. (Such a view would make you an opponent of Marx and Engels, who celebrate the destruction of feudalism by the bourgeoisie in *The Communist Manifesto*.)

6 Sequential games

We have touched on sequential games in some examples in the previous sections. We now discuss these in more detail and the solution concepts appropriate in these settings.

6.1 Sequential rationality

Example 36 (Market entry game III). Consider the following game of perfect information. Suppose Firm 1 is a potential entrant to a market and Firm 2 is an incumbent monopolist. Firm 1 chooses whether to enter (E) or stay out (O) and Firm 2 can choose to fight (F) or accommodate (A) if Firm 1 enters. This is a sequential game in that Firm 2 plays their strategy after Firm 1. The extensive form representation of this game is:



Alternatively, the normal form payoff matrix is:

| | A | F |
|-----|------|--------|
| E | 1, 1 | -1, -1 |
| O | 0, 0 | 0, 0 |

Clearly, there are two pure strategy Nash equilibria, (E, A) and (O, F) . Furthermore, there is a mixed strategy Nash equilibrium $((0, 1), (\frac{1}{2}, \frac{1}{2}))$ – that is, Firm 1 chooses not to enter with probability 1 and Firm 2 mixes $(\frac{1}{2}, \frac{1}{2})$ over accommodating and fighting.

However, the equilibrium in which Firm 2 plays F is unreasonable in the sense that once Firm 1 enters, Firm 2 is better off deviating from F to A . Provided Firm 2 is rational, its threat to play F is therefore *not credible*. Since Firm 2 cannot credibly commit to playing F if Firm 1 enters, and Firm 1 knows this, Firm 1 should enter.

Sequential rationality rules out these ‘unreasonable’ strategies.

Definition 68 (Sequential rationality). A strategy s_i (σ_i) for player i is called *sequentially rational* if at every information set at which i is to move, s_i (σ_i) maximizes i ’s expected payoff conditional on some belief of i over paths that lead to that information set.

This includes at those information sets precluded by i 's own strategy. The belief justifying some action a_i^* at an information set not reached must assign positive probability to paths that do not occur.

As in Example 36, Nash equilibrium is *ex ante* rational, but it is not necessarily optimal after each move by opponents.

Definition 69.

- (a) *On equilibrium path.* Given an equilibrium, we say that outcomes of the game lie *on the equilibrium path* if they occur with positive probability when players implement the equilibrium strategy.
- (b) *Off equilibrium path.* Given an equilibrium, we say that outcomes of the game lie *off the equilibrium path* if they occur with zero probability when players implement the equilibrium strategy.

6.1.1 Backward induction

In finite extensive form games of perfect information, we can use *backward induction* to determine a sequence of optimal actions. Suppose sequential rationality is common knowledge. Then players anticipate that other players will not take actions off the equilibrium path that are not optimal at the information set at which that action is taken. The backward induction algorithm works backward from terminal nodes to find an optimal set of actions.

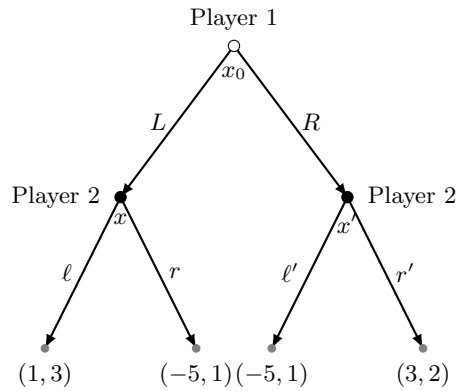
Definition 70 (Backward induction). The *backward induction algorithm* is as follows:

- (Step 1) Start at each node that is an immediate predecessor of a terminal node. Find the optimal action for the player who moves at this node, and modify this node to a terminal node with payoffs given by the optimal action.
- (Step k) Given the terminal nodes derived in the $(k - 1)$ th step, find each of these terminal nodes, find the optimal action for the player who moves at the predecessor of this node.

Continue until the initial node is reached. The list of optimal actions derived from the algorithm is called a *backward induction solution*.

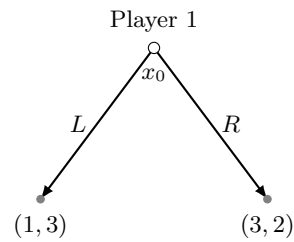
There is a unique backward induction solution if the payoffs across actions at each decision node (information set) differ. Otherwise, one needs to repeat the backward induction algorithm alternating the actions taken to derive all the backward induction solutions.

Example 37. Consider the extensive form game,



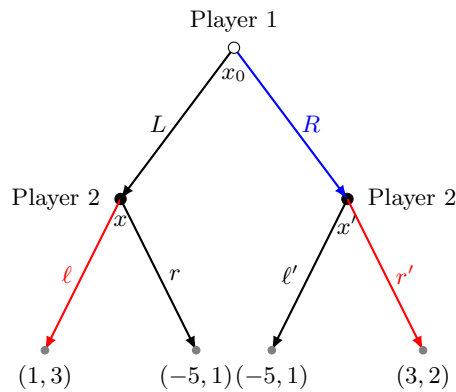
To implement the backward induction algorithm. We start at x and x' since these are the immediate predecessors to the terminal nodes. At x , Player 2's optimal action is ℓ , which yields a payoff of 3 instead of 1. At x' , Player 2's optimal action is r' , yielding payoff 2 instead of 1.

Replacing x and x' by terminal nodes whose payoffs are given by the optimal actions of Player 2 at x and x' , we have



At x_0 , Player 1's optimal action in this game is R , yielding payoff 3 instead of 1.

The backward induction solution is therefore $(R, (\ell, r'))$:



If Player 1 believes Player 2 is sequentially rational, he plays R , believing Player 2 will play ℓ at x and r' at x' .

Backward induction involves reasoning about opponents' future behaviour and beliefs, but not past choices – past choices are taken for granted.

Proposition 52. *Any backward induction solution of an extensive form game of perfect information is sequentially rational.*

Proof. In a game of perfect information, every decision node is an information set. From the first step of backward induction, the backward induction solution has player(s) at the nodes immediately preceding the terminal nodes taking actions that maximize the expected payoff. Given $k-1$ stages, at the k th stage, the backward induction solution has player(s) playing the optimal solution at the nodes immediately preceding the terminal nodes obtained from the $(k-1)$ th stage given subsequent play that would follow from each of those nodes under the backward induction solution, provided that node is reached. Hence these actions maximize the player's expected payoff under the backward induction strategy profile. Proof follows by induction. \square

Proposition 53. *Any backward induction solution of an extensive form game of perfect information is a Nash equilibrium.*

Proof. Any backward induction solution s^* results in an equilibrium path. At every node x on the equilibrium path, the player $\iota(x)$ plays the strategy maximizing $\iota(x)$'s payoff given that x is reached. Hence no player has a profitable deviation at any node on the equilibrium path. Thus s^* is a Nash equilibrium. \square

Indeed, the backward induction solution is a refinement of Nash equilibrium.

Theorem 23 (Zermelo). *Every finite extensive form game Γ of perfect information has a backward induction solution.*

Furthermore, if Γ has no ties – that is, if $u_i(z) \neq u_i(z')$ for all $z, z' \in X_\tau$ and all $i \in \mathcal{I}$ – then it has a unique backward induction solution.

Proof. Recall $\Gamma = (\mathcal{I}, T, P, \Phi, \mathcal{A}, (u_i)_{i \in \mathcal{I}}, \eta)$, with these objects defined as in Definition 2. In this case, Φ_i is the collection of all singleton subsets of X_i . Suppose there is no backward induction solution. Then for some $i \in \mathcal{I}$, there is some decision node $x \in X_i$ at which the backward induction algorithm prescribes no action $a_i \in A_i(x)$. However, every x precedes some terminal node z , and since T is a finite arborescence, the path from x_0 to z is unique and includes x . By structure of the backward induction algorithm, the optimal action at the node immediately preceding z will be determined in the first step. There are finitely many – say, k – nodes between x and z . Furthermore, at any node x' on the path between x and z , there is some (possibly non-unique) optimal strategy at each of these nodes given subsequent play, for there are finitely many strategies. Iterating from the first step, we see that node x will be reached by the algorithm at the $(k-1)$ th step. Hence the algorithm will choose some action $a \in A_i(x)$ at the $(k-1)$ th step given the actions chosen at steps $1, \dots, k-2$ for all nodes succeeding x on the path from x_0 to z . Since this holds for all nodes x , there is a backward induction solution.

Now suppose Γ has no ties. Then at every node x considered at the first step, there is a unique optimal action, since the set of actions at that node is finite and there are no $a, a' \in A_i(x)$ s.t. the payoff to player $i(x)$ from a is the same as the payoff from a' . Hence we can construct a strict order $a_1 > a_2 > \dots > a_M$ with $M = |A_i(x)|$ and $a > a'$ implies a has strictly higher payoff than a' : we thus have a unique action $a_1 \in A_i(x)$ that is optimal for $P(x)$. At every subsequent stage, the payoffs at each of the new terminal nodes is a subset of the payoffs of X_τ , so again we have no ties. Repeating the argument at each stage shows there is a unique backward induction solution. \square

Suppose Γ has infinite nodes. If the arborescence of Γ has an infinite path, then backward induction cannot be applied. However, if the arborescence does not have an infinite path, then we can apply backward induction, but a backward induction solution need not exist.

Example 38 (Ultimatum game). In the (two-player) ultimatum game, Player 1 proposes an offer $(x_1, v - x_1)$ dividing a surplus of total value $v > 0$ between Player 1 and Player 2. Given Player 1's offer, Player 2 can choose to accept the offer, in which case payoffs are $(x_1, v - x_1)$, or reject the offer, in which case payoffs are $(0, 0)$. We assume Player 1 cannot choose $x_1 > v$.

This is a situation where backward induction is possible in the presence of an (uncountably) infinite number of terminal nodes.

The set of feasible agreements is $Z = \{(z, v - z) \mid z \in [0, v]\}$. The (pure) strategy set for Player 1 is $S_1 = Z$, and for Player 2, the strategy set is

$$S_2 = \{s_2 : Z \rightarrow \{\text{Accept}, \text{Reject}\}\}.$$

There are infinitely many Nash equilibria in the ultimatum game. Indeed, any strategy profile $(s_1^*, s_2^*(s_1))$ such that, for any $z \in [0, v]$,

$$\begin{aligned} s_1^* &= (z, v - z), \\ s_2^*(x_1, v - x_1) &= \begin{cases} \text{Accept} & \text{if } x_1 \leq z, \\ \text{Reject} & \text{otherwise,} \end{cases} \end{aligned}$$

is a Nash equilibrium. For suppose Player 1 has alternative strategy $s'_1 = (z', v - z')$ with some $z' \neq z$. If $z' < z$, then Player 2 under s_2^* would accept and Player 1 receives payoff z' , yet Player 2 would also accept if Player 1 played s_1^* , in which case Player 1's payoff is $z > z'$. Hence $u_1(s_1^*, s_2^*(s_1)) \geq u_1(s'_1, s_2^*(s_1))$. If Player 1 chooses $z' > z$, then Player 2 rejects and so Player 1 receives payoff $u_1(s'_1, s_2^*(s_1)) = 0 < z = u_1(s_1^*, s_2^*(s_1))$. Turning to Player 2, under s_2^* Player 2 receives payoff $v - z$, which is clearly maximal given s_1^* .

The equilibrium outcome path in these equilibria is the offer $(z, v - z)$ followed by acceptance.

However, note that if $z < v$, then it is not sequentially rational for Player 2 to reject any offer $(x_1, v - x_1)$ with $z < x_1 < v$: given such an offer, Player 2 would, if

accepting, obtain a payoff $v - x_1 > 0$, and 0 if rejecting. In the Nash equilibrium, $x_1 > z$ is off the equilibrium path and so any response to $x_1 > z$ (including rejection) is *ex ante* optimal for Player 2. Since in the Nash equilibrium, Player 2 would reject in this off-equilibrium-path case, s_1^* is optimal for Player 1.

The only sequentially rational Nash equilibrium is

$$\begin{aligned} s_1^* &= (v, 0), \\ s_2^*(x_1, v - x_1) &= \text{Accept for all } x_1 \in Z. \end{aligned}$$

This is the unique backward induction solution. Player 2 is indifferent between accepting and rejecting $(v, 0)$. In general, indifferences can imply multiple backward induction solutions (per Theorem 23, ties are a necessary condition for multiple backward induction solutions.) However, in the ultimatum game, there is no backward induction solution with Player 2 rejecting $(v, 0)$: any other backward induction solution for Player 2 would entail strategy

$$s_2(x_1, v - x_1) = \text{Accept for all } x_1 \text{ s.t. } v - x_1 > 0.$$

Then Player 1 solves

$$\max_{x_1 \in [0, v]} x_1,$$

which is not well defined since for any $x_1 \in [0, v)$, there is an $x'_1 \in (x_1, v)$ yielding higher payoff to Player 1.

If the set of possible agreements were a discrete set,

$$Z = \{(z, v - z) \mid 0 \leq z \leq v, z = k\epsilon \text{ for } k \in \mathbb{N}\},$$

for some fixed $\epsilon > 0$, then other backward solutions may exist. For example, the strategy profile $(s_1^*, s_2^*(s_1))$ defined by

$$\begin{aligned} s_1^* &= (v - \epsilon, \epsilon), \\ s_2^*(s_1) &= \begin{cases} \text{Accept} & \text{if } x_1 \leq v - \epsilon, \\ \text{Reject} & \text{otherwise,} \end{cases} \end{aligned}$$

is a backward induction solution.

Example 39 (Stackelberg competition). Suppose there are two firms, $i = 1, 2$ which each produce $q_i \in [0, \infty)$ of a single good. Suppose the inverse demand function is

$$p(q_1, q_2) = \max\{a - b(q_1 + q_2), 0\},$$

for some $a, b > 0$, and that both firms have the same (linear) technology and so have the same marginal cost, $c < a$. We refer to Firm 1 as the *Stackelberg leader* and Firm 2 as the *follower*. The leader, Firm 1, chooses q_1 first. Firm 2 then chooses q_2 , knowing q_1 , and we assume Firm 2 is sequentially rational.

The strategy set for Firm 1 is $S_1 = [0, \infty)$ and the strategy set for Firm 2 is $S_2 = \{q_2 : S_1 \rightarrow [0, \infty)\}$.

Again, this is a case where there are infinite terminal nodes, but we can nevertheless apply backward induction. Fix any $q_1 \in [0, \frac{a-c}{b}]$. Firm 2's profit is

$$\pi_2(q_1, q_2) = (a - b(q_1 + q_2) - c)q_2.$$

We have first order condition,

$$\frac{\partial \pi_2(q_1, q_2)}{\partial q_2} = a - bq_1 - c - 2bq_2 = 0,$$

which yields solution

$$q_2^*(q_1) = \frac{a - bq_1 - c}{2b}.$$

Given the leader knows the follower is sequentially rational, the leader solves

$$\begin{aligned} \max_{q_1} \pi(q_1, q_2^*(q_1)) &= \max_{q_1} (a - b(q_1 + q_2^*(q_1)) - c)q_1 \\ &= \max_{q_1} \left(a - \frac{a + bq_1 - c}{2} - c \right) q_1. \end{aligned}$$

We have first order condition,

$$\frac{\partial \pi_1(q_1, q_2^*(q_1))}{\partial q_1} = \frac{1}{2}[a - 2bq_1 - c] = 0,$$

and so Firm 1's optimal strategy is $q_1^* = \frac{a-c}{2b}$.

The backward induction solution of the Stackelberg game is therefore

$$q_1^* = \frac{a-c}{2b}, \quad q_2^*(q_1) = \max \left\{ \frac{a - bq_1 - c}{2b}, 0 \right\}.$$

The backward induction outcome path is therefore $q_1^* = \frac{a-c}{2b}$, $q_2^* = \frac{a-c}{4b}$.

The equilibrium payoffs are

$$\begin{aligned} \pi_1^* &= \frac{(a-c)^2}{8b}, \\ \pi_2^* &= \frac{(a-c)^2}{16b}. \end{aligned}$$

Note the Cournot equilibrium payoff is $\pi_i^c = \frac{(a-c)^2}{9b}$. Thus we see that $\pi_1^* > \pi_1^c$ and $\pi_2^* < \pi_2^c$. Hence Firm 1 has a *first mover advantage*.

Note that while this is the unique backward induction solution, there are infinite (non-sequentially rational) Nash equilibria.

6.2 Subgame perfection

Note that while backward induction can be applied in (finite) extensive form games of perfect information, it cannot necessarily be applied in the presence of imperfect information, for a player's optimal decision depends on which node they are at in their information set, and they cannot distinguish between nodes within an information set. *Subgame perfect equilibrium* (introduced by Selten, 1965) is a generalization of the backward induction solution, applicable in games of imperfect information.

Definition 71 (Subgame perfect equilibrium).

(a) *Subgame*. Given an extensive form game

$$\Gamma = (\mathcal{I}, T, P, \Phi, \mathcal{A}, (u_i)_{i \in \mathcal{I}}, \eta),$$

a subgame

$$\Gamma' = (\mathcal{I}', T', P', \Phi', \mathcal{A}', (u'_i)_{i \in \mathcal{I}'}, \eta')$$

of Γ is a part of Γ such that:

- (i) There is a unique decision node x^* in Γ' such that the immediate predecessor of x^* is not in Γ' , and furthermore, $x^* \in \phi = \{x^*\}$ for some information set ϕ in Γ .
- (ii) A node x is in the subgame Γ' iff either $x = x^*$ or x is a successor of x^* .
- (iii) If $x \in \phi_i$ is in Γ' then every $x' \in \phi_i$ is in Γ' .

If Γ' is a subgame of Γ , we write $\Gamma' \subseteq \Gamma$. Note Γ is itself a subgame of Γ . If $\Gamma' \subseteq \Gamma$ and $\Gamma' \neq \Gamma$, then we call Γ' a *proper subgame* of Γ , and write $\Gamma' \subset \Gamma$.

- (b) *Restricted strategies*. For any (pure) strategy s_i of player i in the game Γ , the *restriction* of s_i to a subgame $\Gamma' \subset \Gamma$ is a mapping $s_{i|\Gamma'}$ such that $s_{i|\Gamma'}(\phi_i) = s_i(\phi_i)$ for every $\phi_i \in \Phi'_i$.

If $n = |\mathcal{I}|$, a restriction of a (pure) strategy profile s to Γ' is a profile $s_{|\Gamma'} = (s_{1|\Gamma'}, \dots, s_{n|\Gamma'})$.

Restrictions of mixed strategies are defined analogously.

- (c) *Subgame perfect equilibrium*. A (pure) strategy profile s^* is a *subgame perfect equilibrium* of Γ if $s_{|\Gamma'}^*$ is a Nash equilibrium of Γ' , for every subgame $\Gamma' \subseteq \Gamma$.

Subgame perfect equilibrium in mixed strategies is defined analogously.

If Γ' is a subgame of Γ , then it is a proper subgame iff the root node x'_0 of Γ' and the root node x_0 of Γ satisfy $x_0 \neq x'_0$.

Proposition 54. *In any game Γ , any $\Gamma' \subseteq \Gamma$ is an extensive form game in its own right.*

Proof. Let S be any set-valued object in the definition of Γ and let S' be the corresponding set-valued object in the definition of Γ' . Likewise, let f be any function-valued object in the definition of Γ , and let f' be the corresponding object in Γ' . Since Γ' is a part of Γ , $S' \subseteq S$ and f' is a restriction of f to a subdomain of f . Hence all the objects of Γ' are well-defined. By (i) in the definition, x^* is a root node, and since the precedence relation is preserved, (N', \succ, x^*) is thus a finite arborescence. It follows that Γ' satisfies the definition of an extensive form game. \square

Subgame perfect equilibria are sometimes called subgame perfect Nash equilibria, since quite trivially:

Proposition 55. *If s^* is subgame perfect equilibrium of Γ , then s^* is a Nash equilibrium of Γ .*

Proof. By definition, $s^*_{|\Gamma'}$ is a Nash equilibrium of Γ' for every $\Gamma' \subseteq \Gamma$. Since $\Gamma \subseteq \Gamma$, $s^* = s^*_{|\Gamma}$ is a Nash equilibrium of Γ . \square

Proposition 56. *Any subgame perfect equilibrium s^* of Γ is sequentially rational.*

Proof. Fix any subgame $\Gamma' \subseteq \Gamma$. Since $s^*_{|\Gamma'}$ is a Nash equilibrium of Γ' , we have that for every $i \in \mathcal{I}$ and at every information set $\phi_i \in \Phi'_i$ of Γ' , $s^*_{i|\Gamma'}(\phi_i) \in B_i(\phi_i)$, where $B_i(\cdot)$ is the best response correspondence. By definition, any strategy in the best response correspondence maximizes i 's payoff given belief that other players play $s^*_{-i|\Gamma}$. Since this holds for all subgames, it follows that under s^* , for all i , at every information set at which i moves, s^*_i maximizes i 's expected payoff (given the continuation play prescribed by s^* if that information set is reached). Hence s^* is sequentially rational. \square

Proposition 57. *A strategy profile s^* is a backward induction solution of a finite extensive form game of perfect information iff s^* is a subgame perfect equilibrium.*

Proof. Let Γ be a finite game of perfect information. Then every node x in Γ corresponds to a distinct information set $\{x\}$. Hence any part of Γ consisting of any decision node x and all its successors is a subgame of Γ . Consider any backward induction solution s^* of Γ . Since any subgame $\Gamma' \subseteq \Gamma$ is a game, the backward induction solution also identifies a backward induction solution s' for Γ' , and it is easy to see from the structure of Γ that s' and $s^*_{|\Gamma'}$ must coincide. By Proposition 53, $s^*_{|\Gamma'}$ is a Nash equilibrium. Since this holds for all $\Gamma' \subseteq \Gamma$, s^* is a subgame perfect equilibrium.

Conversely, suppose s^* is a subgame perfect equilibrium.

Define the *level* of a subgame of Γ as follows:

Definition 72 (Level). In a finite extensive form Γ , call a subgame $\Gamma' \subseteq \Gamma$ a *level 0* subgame if Γ' contains no proper subgame. Recursively, call Γ' a *level k* subgame if

- (i) There is a proper subgame Γ'' of Γ' such that Γ'' is level $k - 1$, and
- (ii) If Γ'' is a proper subgame of Γ' , then Γ'' is of level at most $k - 1$.

Since Γ is finite, it is of some level K . Let \mathcal{G}^k denote the set of all level k subgames of Γ , for $k = 0, 1, \dots, K$. First, consider \mathcal{G}^0 . Since Γ is a game of perfect information, any $\Gamma' \in \mathcal{G}^0$ consists of a single decision node, say x . Since s^* is a subgame perfect equilibrium, $s^*_{|\Gamma'}$ is a Nash equilibrium of Γ' and so prescribes that $\iota(x)$ play the action at x that yields the highest payoff. Since this is the optimal action for $\iota(x)$ at x , there is some backward induction solution s^b s.t. the restriction $s^b_{|\Gamma'}$ has $\iota(x)$ playing the same action, and so $s^*_{|\Gamma'} = s^b_{|\Gamma'}$. Since Γ' is an arbitrary member of \mathcal{G}^0 and the members of \mathcal{G}^0 are independent, it follows that there exists some backward induction solution s^b s.t. $s^b_{|\Gamma'} = s^*_{|\Gamma'}$ for all $\Gamma' \in \mathcal{G}^0$.

Fix any $k \geq 1$ and suppose we have that there exists some backward induction solution s^b s.t. $s_{|\Gamma'}^b = s_{|\Gamma'}^*$ for all $\Gamma' \in \bigcup_{j=0}^{k-1} \mathcal{G}^j$. Then consider any subgame $\Gamma' \in \mathcal{G}^k$. Γ' has some root node, say x , and all ‘downstream’ play is fixed by s^* . Hence $s_{|\Gamma'}^*$ has $\iota(x)$ playing the optimal action at x given the actions fixed by s^* at each of the successor nodes of x . This optimal action is equivalent to the optimal action of the modified game at which each decision node that is an immediate successor to x is replaced by a terminal node with payoffs fixed by the continuation play prescribed by s^* . Since by assumption, $s_{|\Gamma''}^b = s_{|\Gamma''}^*$ for all proper subgames Γ'' of Γ' , the same action at x as under s^* is also an optimal choice for the backward induction algorithm. Hence there is a backward induction solution s^b s.t. $s_{|\Gamma'}^b = s_{|\Gamma'}^*$. Proof follows by induction. \square

Corollary 10. *If Γ is a finite extensive form game of perfect information, then Γ has a subgame perfect equilibrium in pure strategies. Furthermore, if Γ has no ties, then Γ has a unique subgame perfect equilibrium (in pure strategies).*

Proof. Let Γ be a finite extensive form game of perfect information. By Theorem 23, Γ has a backward induction solution, which is unique if Γ has no ties. Proof follows by Proposition 57. \square

Proposition 58. *Any finite extensive form game has a (mixed strategy) subgame perfect equilibrium.*

Proof. Let Γ be any finite extensive form game. Since any extensive form game is equivalent to a normal form game, by Nash’s existence theorem (Theorem 8), every $\Gamma' \subseteq \Gamma$ has a mixed strategy Nash equilibrium. We need to show there is a mixed strategy Nash equilibrium σ^* of Γ s.t. $\sigma_{|\Gamma'}^*$ is a Nash equilibrium of Γ' for all $\Gamma' \subseteq \Gamma$.

Define the level of a subgame as in Definition 72. Since Γ is finite, Γ is level K for some integer K . Let \mathcal{G}^k denote the set of level k subgames of Γ .

We can apply a form of backward induction here. By the Nash existence theorem, each $\Gamma' \in \mathcal{G}^0$ has a Nash equilibrium. Since Γ' has no proper subgame, this is also a subgame perfect equilibrium. Very obviously:

Lemma 19. *σ^* is a subgame perfect equilibrium of Γ iff $\sigma_{|\Gamma'}^*$ is a subgame perfect equilibrium of Γ' for all subgames $\Gamma' \subseteq \Gamma$. The converse is equally trivial.*

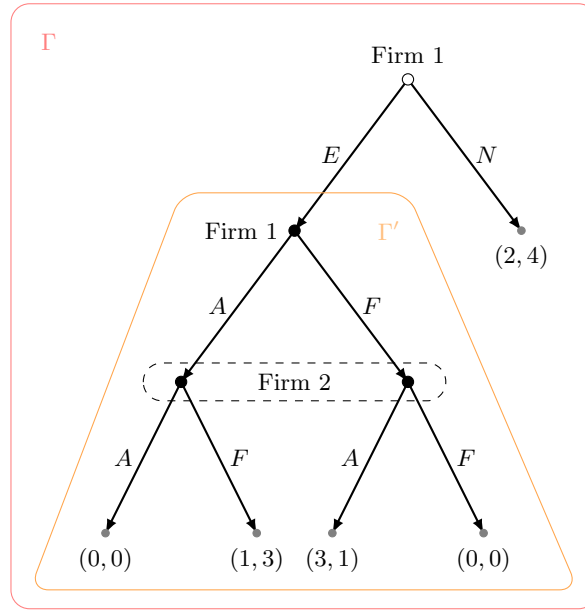
Proof. The profile σ^* is a subgame perfect equilibrium if $\sigma_{|\Gamma'}^*$ is a Nash equilibrium for every subgame Γ' . Any subgame Γ'' of Γ' is also a subgame of Γ , and hence $\sigma_{|\Gamma''}^*$ is a Nash equilibrium of Γ'' for all subgames $\Gamma'' \subseteq \Gamma'$. Thus $\sigma_{|\Gamma'}^*$ is a subgame perfect equilibrium of Γ' . \square

Now fix k and suppose every $\Gamma'' \in \mathcal{G}^{k-1}$ has a subgame perfect equilibrium $\sigma_{|\Gamma''}^*$. Fix some $\Gamma' \in \mathcal{G}^k$. For each $\Gamma'' \in \mathcal{G}^{k-1}$ s.t. $\Gamma'' \subseteq \Gamma'$, replace the root node x'' of Γ'' by a terminal node with payoffs fixed by the subgame perfect equilibrium $\sigma_{|\Gamma''}^*$. By Nash’s existence theorem, this modified game has a mixed strategy Nash equilibrium. Hence Γ' has a mixed strategy Nash equilibrium $\sigma_{|\Gamma'}^*$, s.t. $(\sigma_{|\Gamma'}^*)_{|\Gamma''} = \sigma_{|\Gamma''}^*$ for all $\Gamma'' \subseteq \Gamma'$. Thus

$\sigma_{|\Gamma'}^*$ is a subgame perfect equilibrium of Γ' . By induction, we can therefore construct a mixed strategy subgame perfect equilibrium σ^* of Γ . \square

Example 40 (Market entry game IV). Consider the following market entry game. Firm 1 is a potential entrant, who first decides whether to enter (E) or not enter (N). Conditional on entering, Firm 1 chooses whether to fight (F) or accommodate (A). The incumbent monopolist, Firm 2, simultaneously chooses whether to fight (F) or accommodate (A), not knowing Firm 1's action.

The set of pure strategies for Firm 1 are $\{EA, EF, NA, NF\}$ and the set of pure strategies for Firm 2 are $\{A, F\}$. The game has one proper subgame Γ' , circled in orange, and the entire game Γ is circled in pale red.

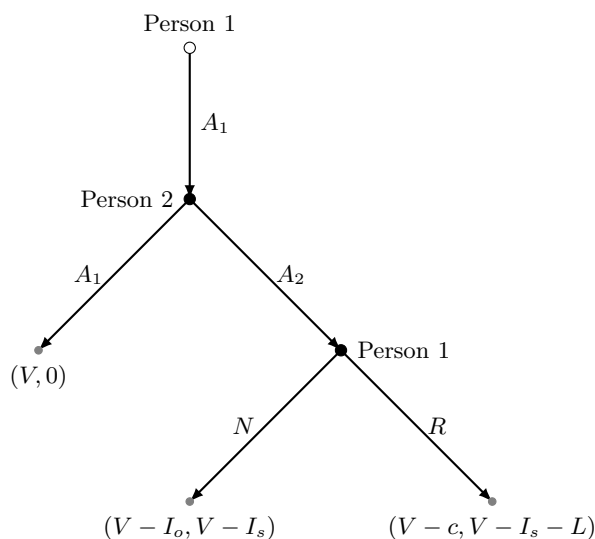


Consider subgame Γ' . There are two pure strategy Nash equilibria in this subgame, (A, F) and (F, A) . Furthermore, there is a mixed strategy Nash equilibrium $((\frac{1}{4}, \frac{3}{4}), (\frac{1}{4}, \frac{3}{4}))$.

Now considering Γ , we can identify those Nash equilibria whose restrictions are Nash equilibria in Γ' . First, if (A, F) is played in Γ' , then Firm 1 receives greater expected payoff from N than from E , so we have a subgame perfect equilibrium (NA, F) . Second, if (F, A) is played in Γ' , then Firm 1 receives greater expected payoff by playing E instead of N , so we have a subgame perfect equilibrium (EF, A) . Finally, if the mixed strategy equilibrium is played in Γ' , then the expected payoff to Firm 1 from E is $\frac{3}{4}$, versus 2 from N , so we have a subgame perfect equilibrium $((0, 0, \frac{1}{4}, \frac{3}{4}), (\frac{1}{4}, \frac{3}{4}))$.

Example 41 (Economics of identity: Akerlof & Kranton, 2000). Akerlof & Kranton (2000) are the first to formalize the notion that identity – an individual's sense of self – affects economic outcomes, via a simple sequential game. Suppose there are two possible activities – A_1 and A_2 – and agents earn payoff V if they engage in an activity that

accords with their tastes and zero otherwise. Suppose there are two social categories – Red and Green. Suppose everyone thinks of themselves and everyone else as Green, and Greens should engage in activity A_1 , or else they are not a “true Green” and lose the Green identity, which would result in a payoff loss of I_s . Reds meanwhile have a taste for activity A_2 . Moreover, suppose there are *identity externalities*: if i and j are paired and i engages in activity A_2 , this diminishes j ’s sense of Green identity and j incurs a loss I_o . After i commits to activity A_2 , j can respond (R) or not (N), and a response restores j ’s Green identity but incurs a cost c , while imposing a loss L on i . The following game tree illustrates an interaction between Person 1, who has a taste for activity A_1 and Person 2, who has a taste for activity A_2 :



The branch of the tree in which Person 1 chooses A_2 is suppressed since Person 1 never chooses A_2 in any subgame perfect equilibrium.

There are four possible subgame perfect equilibrium outcomes. First, if $I_s > V$ then Person 2 always chooses activity A_1 in equilibrium. Now suppose $c < I_o$. If $I_s < V < I_s + L$, then Person 1 deters Person 2 from engaging in A_2 . If instead $I_s + L < V$, then Person 1 does not deter Person 2 and Person 2 chooses A_2 . Finally, if $c > I_o$ and $I_s < V$ then Person 1 does not respond to A_2 and Person 2 engages in A_2 .

Kranton and Akerlof provide a number of examples of where this type of game arises in practice. Gender provides a good example. In stereotypically “male” lines of work such as coal handlers in power plants, men can be rather insecure – among coal handlers in power plants, for example, male coal handlers refused to train women and bullied them, seemingly because they felt threatened in their masculinity by women performing the same job.

6.3 One-shot deviation principle

The *one-shot deviation principle* is a particularly helpful dynamic programming result for finding subgame perfect equilibria in dynamic games (both sequential and repeated

games). First, note that given the strategies of other players $-i$, we have a maximization problem for player i . We can thus consider a single person decision tree. Recall that given a decision tree, a path $y = (z_0, \dots, z_n)$ is an ordered collection of nodes in the tree such that for each $z_j \in y$, z_{j+1} is an immediate successor to z_j . Given a node $x \in y$, define the *subpath* y_x to be the path with initial node x and all subsequent nodes identical to those successor nodes of x in y . We say that two paths y and y' *diverge at* x if they share the same nodes up to x but have distinct subpaths thereafter.

Let $\pi(y)$ be the *return* of path y , that is, the payoff to player i at the terminal node reached by y . We assume the following.

Assumption. Let π be a return function, mapping paths to payoffs. Assume

- (A1) *Consistency.* If $\pi(y) \geq \pi(y')$ and if y and y' diverge at x , then $\pi(y_x) \geq \pi(y'_x)$.
- (A2) *Continuity.* Fix any path y . For every $\epsilon > 0$, there exists an integer N such that if $n \geq N$ and if another path y' shares the first n nodes of y , then $|\pi(y) - \pi(y')| < \epsilon$.

Any finite game tree with payoffs at terminal nodes automatically satisfies both consistency and continuity assumptions. So do infinite problems with discounted additively separable payoffs: this covers infinitely repeated games of perfect information. Interpreting nodes as information sets and payoffs as expected payoffs, we can extend this to stochastic decision problems such as games of imperfect information.

Definition 73 (One-shot deviations).

- (a) Given a node x any strategy profile σ induces a probability distribution P over all paths y with initial node x . Define $\pi(\sigma, x) = \mathbb{E}_y[\pi(y) \mid \sigma]$.
- (b) *Optimal strategy.* A strategy σ is *optimal* if there is no strategy σ' and node x such that $\pi(\sigma', x) > \pi(\sigma, x)$.
- (c) *One-shot deviation.* Given a strategy σ and a node x , let σ_a denote the strategy obtained by substituting the deterministic choice $a \in A(x)$ in place of what was prescribed under σ , and leaving σ otherwise unchanged. We call σ_a a *one-shot deviation* from σ at x .
- (d) *Unimprovable strategy.* We call a strategy σ *unimprovable* if there is no one-shot deviation σ_a from σ at x such that

$$\pi(\sigma_a, x) > \pi(\sigma, x).$$

Theorem 24 (One-shot deviation principle). *Under assumptions (A1) and (A2), any unimprovable strategy must be optimal.*

Proof. Suppose σ is an unimprovable but non-optimal strategy. Then there exists a strategy σ' and a node x_0 s.t.

$$\pi(\sigma', x_0) > \pi(\sigma, x_0).$$

Since any best stochastic strategy is payoff-equivalent to some pure strategy, this is equivalent to assuming there is some path y with initial node x_0 s.t.

$$\pi(y) \geq \pi(\sigma, x_0) + 2\epsilon$$

for some $\epsilon > 0$. Under **(A2)**, there is some integer N s.t. if any path y' starting from x_0 shares the first $N + 1$ nodes with y , then

$$\pi(y') \geq \pi(y) - \epsilon.$$

Combining inequalities, we have that

$$\pi(y') \geq \pi(\sigma, x_0) + \epsilon.$$

Let x_0, \dots, x_N denote the first N ordered nodes of y . From the inequality, it follows that finitely many one-shot deviations at the nodes x_j , $j = 0, \dots, N$ are sufficient to generate a payoff improvement from σ .

Define a family of N different strategies α_j for $j = 1, \dots, N - 1$ s.t. α_j chooses x_{k+1} at node x_k for every $k \in \{0, \dots, j\}$, and coincides with σ elsewhere. Then by the argument of the previous paragraph,

$$\pi(\alpha_{N-1}, x_0) > \pi(\sigma, x_0).$$

Note α_{N-2} coincides with σ at node x_{N-1} and all subsequent nodes, while α_{N-1} is a one-shot deviation from σ at node x_{N-1} . Since σ is unimprovable by assumption,

$$\pi(\alpha_{N-2}, x_{N-1}) = \pi(\sigma, x_{N-1}) > \pi(\alpha_{N-1}, x_{N-1}).$$

Applying assumption **(A1)**, we have that

$$\pi(\alpha_{N-2}, x_0) \geq \pi(\alpha_{N-1}, x_0),$$

since α_{N-2} and α_{N-1} share the same nodes x_0, \dots, x_{N-2} along all paths they generate. But since we also had $\pi(\alpha_{N-1}, x_0) > \pi(\sigma, x_0)$, we have

$$\pi(\alpha_{N-2}, x_0) > \pi(\sigma, x_0).$$

Reapplying this argument iteratively for α_{N-3} , then α_{N-4} , and so on, we reach

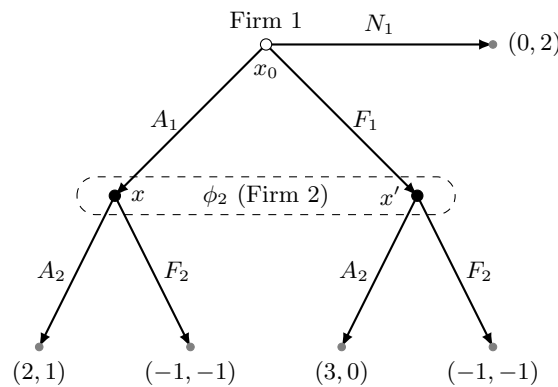
$$\pi(\alpha_0, x_0) > \pi(\sigma, x_0).$$

But α_0 is a one-shot deviation from σ : in particular, $\alpha_0 = \sigma_{x_1}$. Hence this inequality contradicts the unimprovability of σ . \square

6.4 Perfect Bayesian equilibrium

There are settings where subgame perfection lacks ‘bite’, most notably where there is no proper subgame: in this case, any Nash equilibrium is a subgame perfect equilibrium, yet we can see that in some situations, this will result in unreasonable subgame perfect equilibria.

Example 42 (Market entry game V). The following market entry game is very similar to Example 40. Firm 1 is a potential entrant, who decides either to not enter (N_1), to enter and fight (F_1) or to enter and accommodate (A_1). Firm 2 is an incumbent monopolist who simultaneously chooses to either fight (F_2) or accommodate (A_2). The game in extensive form representation is as follows:



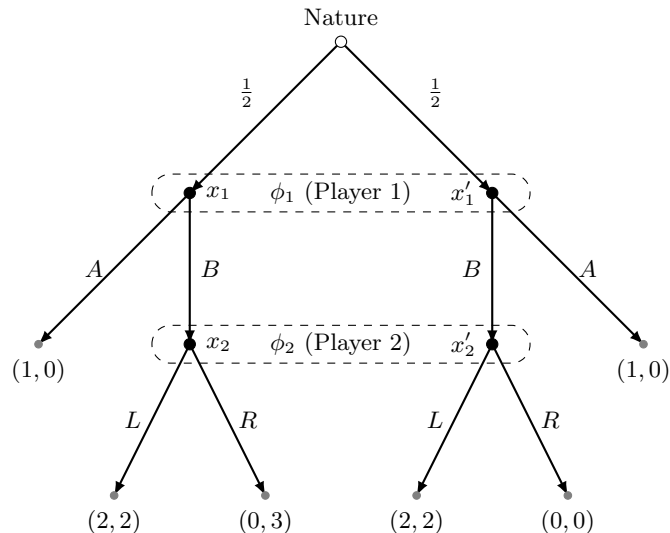
Unlike in Example 40, there is now only one subgame – the whole game itself. Hence any Nash equilibrium is a subgame perfect equilibrium. There are two pure strategy Nash equilibria – (F_1, A_2) and (N_1, F_2) . Furthermore, there is a mixed strategy Nash equilibrium. Suppose Player 2's strategy σ_2 mixes between A_2 with probability p and F_2 with probability $1 - p$. Then Player 1's payoffs are

$$\begin{aligned} u_1(A_1, \sigma_2) &= 3p - 1, \\ u_1(F_1, \sigma_2) &= 4p - 1, \\ u_1(N_1, \sigma_2) &= 0. \end{aligned}$$

Hence $u_1(N_1, \sigma_2) \geq u_1(F_1, \sigma_2) > u_1(A_1, \sigma_2)$ provided $p \leq \frac{1}{4}$, so we have a family of mixed strategy Nash equilibria $\{\sigma^* = ((0, 0, 1), (p, 1 - p)) : 0 \leq p \leq \frac{1}{4}\}$ where mixing is over $\{A_1, F_1, N_1\}$ and $\{A_2, F_2\}$ in that order respectively.

While all of these Nash equilibria are trivially also subgame perfect, the mixed strategy equilibria and the pure strategy equilibrium (N_1, F_2) are not sequentially rational, since at both nodes x and x' , Player 2 will obtain a higher payoff by playing A_2 for certain than by playing F_2 with any positive probability. Subgame perfection does not capture this, because for these equilibria, the information set ϕ_2 neither lies on the equilibrium path nor forms a part of a proper subgame off the equilibrium path.

Example 43. A second problem is that subgame perfect equilibrium can allow actions that are only possible if players have unreasonable beliefs. This arises in the following game:



This game has only one subgame. It has pure strategy subgame perfect equilibria (B, L) and (A, R) . The latter is intuitively unreasonable. R is optimal for Player 2 if he believes he is at x_2 with probability at least $\frac{2}{3}$ given he is at $\phi_2 = \{x_2, x'_2\}$. Yet this would be an unreasonable belief: if Player 1 chooses B , then Player 2 should place equal probability on being at x_2 and x'_2 .

Perfect Bayesian equilibrium addresses some of these issues by generalizing backward induction-like reasoning to games of imperfect information. In particular, perfect Bayesian equilibrium imposes sequential rationality at every information set. Recall that a behavioural strategy σ_i for player i is a mapping $\sigma_i : \phi_i \mapsto \Delta(A_i(\phi_i))$.

Definition 74 (Perfect Bayesian equilibrium).

- (a) *System of beliefs.* A *system of beliefs* $\mu = (\mu_i)_{i \in \mathcal{I}}$ is a profile of belief functions μ_i that each assign to each information set $\phi_i \in \Phi_i$ a probability distribution over the nodes of ϕ_i . We write $\mu_i(x \mid \phi_i)$ for i 's belief that i is at node $x \in \phi_i$ given that i is at information set ϕ_i .

Given a strategy profile σ and a system of beliefs μ , player i 's *continuation payoff* at information set ϕ_i is

$$u_i(\sigma \mid \phi_i, \mu) = \sum_{x \in \phi_i} u_i(\sigma \mid x) \mu_i(x),$$

where $u_i(\sigma \mid x)$ is the expected payoff at node x under the actions prescribed by strategy profile σ for all successor information sets of x .

(b) *Assessment.* An *assessment* (σ, μ) consists of a (behavioural) strategy profile σ and a system of beliefs μ .

(c) *Perfect Bayesian equilibrium.* An assessment (σ^*, μ^*) is a (weak) *perfect Bayesian equilibrium* if

(i) For all $i \in \mathcal{I}$, σ_i^* is sequentially rational, that is, for all information sets $\phi_i \in \Phi_i$,

$$u_i(\sigma_i^*, \sigma_{-i}^* \mid \phi_i, \mu_i) \geq u_i(\sigma_i, \sigma_i^* \mid \phi_i, \mu_i)$$

for all strategies σ_i of i .

(ii) Beliefs μ_i are updated by Bayes' rule whenever it applies – that is, given σ^* , for all $i \in \mathcal{I}$, all $\phi_i \in \Phi_i$ such that $\mathbb{P}\{\phi_i \mid \sigma^*\} > 0$, and all $x \in \phi_i$,

$$\mu_i(x) = \frac{\mathbb{P}\{x \mid \sigma^*\}}{\mathbb{P}\{\phi_i \mid \sigma^*\}} = \frac{\mathbb{P}\{x \mid \sigma^*\}}{\sum_{x' \in \phi_i} \mathbb{P}\{x' \mid \sigma^*\}}. \quad 50$$

(d) *Strong perfect Bayesian equilibrium.* A profile of assessments (σ, μ) is a *strong perfect Bayesian equilibrium* of a game Γ if $(\sigma|_{\Gamma'}, \mu|_{\Gamma'})$ is a perfect Bayesian equilibrium in every subgame $\Gamma' \subseteq \Gamma$.

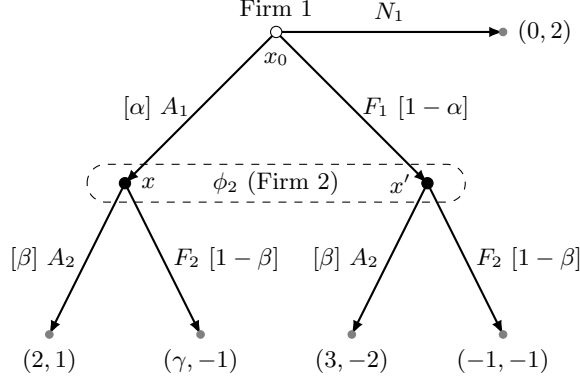
Perfect Bayesian equilibrium is now a bit dated as a solution concept. Sequential equilibrium, discussed in the next section, is preferred.

Example 42 (continued). A system of beliefs μ must assign to each information set a probability distribution over that set. We have information sets $\{x_0\}$ and $\phi_2 = \{x, x'\}$. Hence μ assigns probability $\mu_1(x_0) = 1$ and probabilities $\mu_2(x), \mu_2(x')$ such that $\mu_2(x) + \mu_2(x') = 1$. If Firm 2 maximizes its expected payoff at ϕ_2 given belief μ_2 . We see that A_2 is the optimal action for any belief Firm 2 may have. Hence (F_1, A_2) is the only perfect Bayesian equilibrium.

Example 43 (continued). Suppose $\sigma_1(B) > 0$. Then Bayesian updating implies $\mu_2(x_2) = \mu_2(x'_2) = \frac{1}{2}$. Under these beliefs, Player 2's optimal action at ϕ_2 is L . Player 1's optimal strategy at ϕ_1 is then B . The belief profile consistent with Bayes' rule given strategies (B, L) is $\mu = ((\mu_1(x_1), \mu_1(x'_1)), (\mu_2(x_2), \mu_2(x'_2))) = \left(\left(\frac{1}{2}, \frac{1}{2}\right), \left(\frac{1}{2}, \frac{1}{2}\right)\right)$. Hence the assessment $((B, L), \mu)$ is a perfect Bayesian equilibrium.

Example 44. Slightly modifying the payoffs in Example 42 to create a less computationally trivial example, consider the following game:

⁵⁰To see how this follows from Bayes' theorem, note that $\mu_i(x) = \mathbb{P}\{x \mid \phi_i, \sigma^*\} = \frac{\mathbb{P}\{\phi_i \mid x, \sigma^*\} \mathbb{P}\{x \mid \sigma^*\}}{\mathbb{P}\{\phi_i \mid \sigma^*\}}$, and $\mathbb{P}\{\phi_i \mid x, \sigma^*\} = 1$.



Assume that $\gamma > 0$, and let μ be a system of beliefs assigning $\mu_1(x_0) = 1$ and $\mu_2(x) + \mu_2(x') = 1$. At ϕ_2 , Firm 2's expected payoffs given these beliefs are

$$\begin{aligned} u_2(s_1, A_2 \mid \phi_2, \mu) &= \mu_2(x) - 2(1 - \mu_2(x)) = 3\mu_2(x) - 2, \\ u_2(s_1, F_2 \mid \phi_2, \mu) &= -1, \end{aligned}$$

where s_1 is the strategy of Firm 1 (irrelevant for the calculation of expected payoffs here since continuation play from ϕ_2 does not involve Firm 1 taking decisions.)

Given belief system μ , we see that Firm 2's optimal action is A_2 if $3\mu_2(x) - 2 \geq -1 \Rightarrow \mu_2(x) \geq \frac{1}{3}$ and Firm 2's optimal action is F_2 if $\mu_2(x) \leq \frac{1}{3}$, noting that if $\mu_2(x) = \frac{1}{3}$, then any arbitrary randomization over $\{A_2, F_2\}$ is optimal.

First, suppose $\mu_2(x) < \frac{1}{3}$ and $\sigma^* = (A_1, F_2)$. Then applying Bayes' rule, we have

$$\mu_2(x) = \frac{\mathbb{P}\{x \mid \sigma^*\}}{\mathbb{P}\{x \mid \phi_2\}} = \frac{1}{1} = 1 \not< \frac{1}{3},$$

so $\mu_2(x) < \frac{1}{3}$ is not consistent with Bayes' rule.

Second, suppose $\mu_2(x) > \frac{1}{3}$ and $\sigma^* = (F_1, A_2)$. Then applying Bayes' rule, we have

$$\mu_2(x) = \frac{\mathbb{P}\{x \mid \sigma^*\}}{\mathbb{P}\{x \mid \phi_2\}} = \frac{0}{1} = 0 \not> \frac{1}{3},$$

and thus $\mu_2(x) > \frac{1}{3}$ is not consistent with Bayes' rule.

This leaves $\mu_2(x) = \frac{1}{3}$. Suppose Firm 1's strategy σ_1 mixes over $\{A_1, F_1\}$, playing A_1 with probability $\alpha \in [0, 1]$ and F_1 with complementary probability. Suppose Firm 2's strategy σ_2 mixes over $\{A_2, F_2\}$, playing A_2 with probability $\beta \in [0, 1]$ and F_2 with complementary probability. Then

$$\begin{aligned} u_1(A_1, \sigma_2 \mid \{x_0\}, \mu) &= 2\beta + \gamma(1 - \beta) = (2 - \gamma)\beta + \gamma, \\ u_1(F_1, \sigma_2 \mid \{x_0\}, \mu) &= 3\beta - (1 - \beta) = 4\beta - 1. \end{aligned}$$

Hence Firm 1 mixes iff $4\beta - 1 = (2 - \gamma)\beta + \gamma$, that is, iff $\beta = \frac{1+\gamma}{2+\gamma}$.

Now, Firm 1 randomizes between A_1 and F_1 only since N_1 is strictly dominated by A_1 given $\gamma > 0$. It is optimal for Firm 2 to randomize over A_2 and F_2 only if $\alpha = \frac{1}{3}$. Finally, we check that $\mu_2(x) = \frac{1}{3}$, is consistent with Bayes' rule:

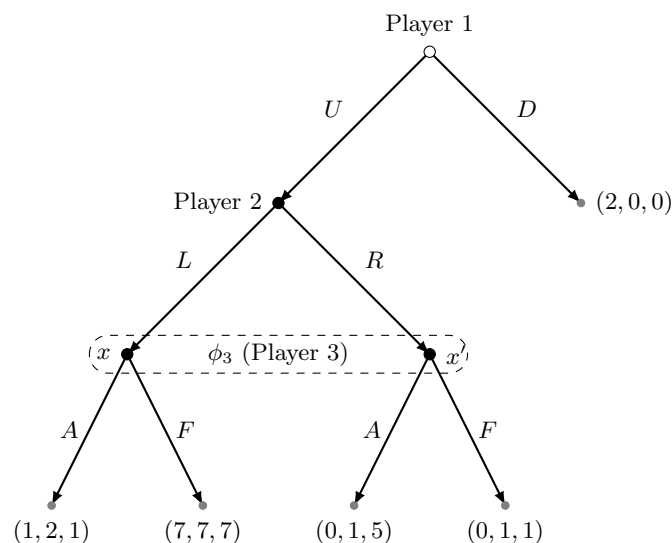
$$\mu_2(x) = \frac{\mathbb{P}\{x \mid \sigma\}}{\mathbb{P}\{\phi_2 \mid \sigma\}} = \frac{\alpha}{1} = \alpha = \frac{1}{3}.$$

Hence we have verified that the unique perfect Bayesian equilibrium is

$$\left((\sigma_1, \sigma_2) = \left(\left(\frac{1}{3}, \frac{2}{3}, 0 \right), \left(\frac{1+\gamma}{2+\gamma}, \frac{1}{2+\gamma} \right) \right), \mu = \left(1, \left(\frac{1}{3}, \frac{2}{3} \right) \right) \right).$$

A defect of (weak) perfect Bayesian equilibrium is that a perfect Bayesian equilibrium need not be subgame perfect. This can result in intuitively unreasonable perfect Bayesian equilibria. A straightforward refinement addressing this – strong perfect Bayesian equilibrium – imposes subgame perfection by requiring that the assessment induces a perfect Bayesian equilibrium in every subgame, thus restricting some beliefs off the equilibrium path.

Example 45. Consider the following game:



Suppose $\mu_3(x') = 1$. Then Player 3's optimal action at ϕ_3 is A , since this yields expected payoff 3 to Player 3, versus 2 if playing F . Now, for Player 2, L is a strictly dominant strategy. Thus if Player 1 chooses U , his expected payoff is $u_1(U, L, A) = 0 < u_1(D, L, A) = 2$, so D is optimal. Clearly, $\mu_3(x') = 1$ is unreasonable, since Player 2's strictly dominant strategy is L , so were ϕ_3 to be reached, Player 3 ought to reason that he is at x for certain. But since the information set ϕ_3 is never reached on the equilibrium path, this doesn't matter – so the assessment (σ, μ) where $\sigma = (D, L, A)$ and $\mu = (1, 1, (0, 1))$ is a perfect Bayesian equilibrium.

Strong perfect Bayesian equilibrium does not suffer this defect, since any strong perfect Bayesian equilibrium induces a perfect Bayesian equilibrium in all subgames. Considering the proper subgame here, we have that ϕ_3 is reached with positive probability. Given L is strictly dominant for Player 2, Player 3 must reason that $\mu_3(x) = 1$. Given this belief, the optimal action at ϕ_3 is F . Player 1's optimal action given continuation play is then U . We therefore have a strong perfect Bayesian equilibrium involving strategy profile (U, L, F) and belief system $(1, 1, (1, 0))$. This is of course also a perfect Bayesian equilibrium.

6.5 Sequential equilibrium

Even strong perfect Bayesian equilibrium allows for equilibria that can seem quite unreasonable. In Example 43, for example, there is only one subgame, so problems with unreasonable off-the-equilibrium-path beliefs are left unchecked:

Example 43 (continued). We previously showed (B, L) is a perfect Bayesian equilibrium strategy profile. However, note that there is another. Suppose $\sigma_1(B) = 0$. Then Bayes' rule does not apply, so μ_2 can be arbitrary. If $\mu_2(x_2) > \frac{2}{3}$, then it is sequentially rational for Player 2 to play R . Given Player 2 plays R at ϕ_2 , it is optimal for Player 1 to play A at ϕ_1 . Hence we also have a family of perfect Bayesian equilibria

$$\left\{ \left(\sigma^* = (A, R), \mu = \left(\left(\frac{1}{2}, \frac{1}{2} \right), (\alpha, 1 - \alpha) \right) \right) : \alpha \in \left(\frac{2}{3}, 1 \right] \right\}.$$

Since there is only one subgame, these are also strong perfect Bayesian equilibria.

Yet as we had originally noted, Player 2's beliefs here are unreasonable.

Kreps and Wilson (1982) introduce *sequential equilibrium*, which imposes strong restrictions on the beliefs that players can hold off the equilibrium path – stronger than strong perfect Bayesian equilibrium, which only restricts some beliefs in subgames off the equilibrium path, but slightly weaker than perfect equilibrium. The main innovation in sequential equilibrium is that it imposes that perfect Bayesian equilibrium holds even if players tremble. Sequential equilibrium has displaced perfect Bayesian equilibrium in the literature as the go-to solution concept in games of incomplete information.

Definition 75.

- (a) *Consistency*. An assessment (σ, μ) is called *consistent* if there exists some sequence of assessments (σ^n, μ^n) such that σ^n is totally mixed, μ^n is derived from σ^n by Bayes' rule, and $(\sigma, \mu) = \lim_{n \rightarrow \infty} (\sigma^n, \mu^n)$.
- (b) *Sequential rationality*.⁵¹ An assessment (σ, μ) is *sequentially rational* if, for every player $i \in \mathcal{I}$ and every information set $\phi_i \in \Phi_i$,

$$u_i(\sigma_i, \sigma_{-i} \mid \phi_i, \mu_i) \geq u_i(\sigma'_i, \sigma_{-i} \mid \phi_i, \mu_i)$$

⁵¹We defined sequential rationality in the context of perfect Bayesian equilibrium, but restate it here for convenience.

for all strategies σ'_i of i .

(c) *Sequential equilibrium*. An assessment (σ^*, μ^*) is a *sequential equilibrium* if

- (i) (σ^*, μ^*) is sequentially rational, and
- (ii) (σ^*, μ^*) is consistent.

To relate this to perfect Bayesian equilibrium, a sequential equilibrium is a PBE with an additional consistency requirement on off-path beliefs.

The requirement that beliefs are consistent is pretty strong. Namely, consistency requires that players have common beliefs following a deviation from equilibrium behaviour. This is sometimes criticised for being too strong – if something goes wrong, then should we expect that different players have the same conjectures about what happened?

Example 43 (continued). The unique sequential equilibrium in Example 43 is

$$(\sigma^*, \mu^*) = \left((B, L), ((1/2, 1/2), (1/2, 1/2)) \right).$$

To see this, note that for any strategy with $\sigma_1^n(A), \sigma_1^n(B) > 0$, we must have that $\mu_2^n(x) = \frac{1}{2}$. Yet then $\lim_{n \rightarrow \infty} \mu_2^n(x) = \frac{1}{2}$, so $\mu_2(x) = \mu_2(x') = \frac{1}{2}$ is the only possible beliefs of Player 2 in any consistent assessment. Choosing some sequence $\{\sigma_1^n, \sigma_2^n\}$ such that $\sigma_1^n(B) \rightarrow 1$ and $\sigma_2^n(L) \rightarrow 1$, we see that (σ^*, μ^*) is a consistent assessment and therefore a sequential equilibrium. Since none of the other perfect Bayesian equilibria involved $\mu_2 = (1/2, 1/2)$, they cannot be sequential equilibria.

As with the Nash equilibrium correspondence, we can ask whether in a sequence of games, limits of sequential equilibria are sequential equilibria of the limit game. For the family of extensive form games $\mathcal{G}_\Lambda = \{\Gamma(\lambda) \mid \lambda \in \Lambda\}$ parameterized by Λ , define the *sequential equilibrium correspondence* $\text{SE} : \Lambda \rightrightarrows \times_{i=1}^n \Delta(S_i)$ by $\text{SE}(\lambda) = \{\sigma \mid \sigma \text{ is a sequential equilibrium of } \Gamma\}$. As with the Nash equilibrium correspondence, we assume the payoff function of each player is continuous in both strategy profiles and parameters.

Proposition 59. *Let \mathcal{G}_Λ be a family of extensive form games parameterized by Λ . The sequential equilibrium correspondence of \mathcal{G}_Λ has a closed graph.*

Proof. Let $\{\Gamma^k\}$ be a sequence of games in \mathcal{G}_Λ s.t. $\Gamma^k \rightarrow \Gamma$, let $\{u^k\}$ be the corresponding sequence of payoff functions with $u^k \rightarrow u$, and let $\{(\sigma^k, \mu^k)\}$ be a corresponding sequence of sequential equilibria of the games Γ^k s.t. $(\sigma^k, \mu^k) \rightarrow (\sigma, \mu)$. Since the expected payoffs conditional on reaching any information set are continuous in the payoff functions and beliefs, we have that (σ, μ) is sequentially rational.

Now, since each (σ^k, μ^k) is consistent, there exists a sequence of totally mixed strategies $\{\sigma^{m,k}\}$ s.t. $\sigma^{m,k} \rightarrow \sigma^k$, and corresponding induced beliefs $\{\mu^{m,k}\}$ s.t. $\mu^{m,k} \rightarrow \mu^k$. For each k , there is m_k sufficiently large that $\sigma^{m,k}$ lies in a $(1/k)$ -ball about σ^k and $\mu^{m,k}$ lies in a $(1/k)$ -ball about μ^k for all $m \geq m_k$. Now since $\sigma^k \rightarrow \sigma$ and $\mu^k \rightarrow \mu$, it follows that $\sigma^{m_k,k} \rightarrow \sigma$ and $\mu^{m_k,k} \rightarrow \mu$, and thus (σ, μ) is consistent. \square

It turns out that sequential equilibria nest the perfect equilibria (recall that for extensive form games, a perfect equilibrium is a perfect equilibrium of the corresponding agent-normal form game):

Theorem 25. *Consider a finite extensive form game Γ . For every perfect equilibrium σ of Γ , there exists a system of beliefs μ such that (σ, μ) is a sequential equilibrium of Γ .*

Proof. Suppose σ is a perfect equilibrium of Γ . Then there is some sequence of totally mixed strategies $\{\sigma^m\}$ in the corresponding agent-normal form game s.t. $\sigma^m \rightarrow \sigma$, and σ_ϕ^m is a best response to $\sigma_{-\phi}^m$ for every information set ϕ and each $m \in \mathbb{N}$. For each m , the induced system of beliefs μ^m is well-defined, and there exists some subsequence of $\{\mu^{m_k}\} \subseteq \{\mu^m\}$ such that $\mu^{m_k} \rightarrow \mu$ for some system of beliefs μ . Thus (σ, μ) is consistent, and since σ_ϕ is a best response given $\sigma_{-\phi}^{m_k}, \mu^{m_k}(\cdot | \phi)$ for each m_k , we have that (σ, μ) is sequentially rational. \square

Theorem 26. *Every finite extensive form game Γ has a sequential equilibrium.*

Proof. By Theorem 25, for every perfect equilibrium σ of Γ , there exists beliefs μ s.t. (σ, μ) is a sequential equilibrium. Now by Theorem 10, there is some perfect equilibrium σ of Γ . \square

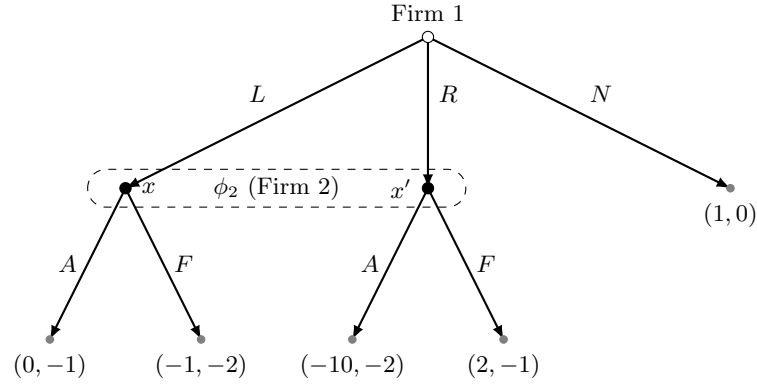
Sequential equilibrium may still fail to eliminate certain kinds of unreasonable equilibrium. Techniques such as forward induction (see sections 6.6 and 7.1.2) can be used to further refine the set of equilibria.

6.6 Forward induction

Backward induction captures the notion that whatever a player chooses, players making subsequent decisions will behave rationally. *Forward induction* extends this to previous decisions – it captures the notion that players assume that even if confronted with an unexpected event, their opponents chose rationally in the past and will continue to choose rationally in subsequent play. If a player finds herself off the equilibrium path, she assumes this is the result of opponents having maximized their utility, as long as such an assumption is reasonable.

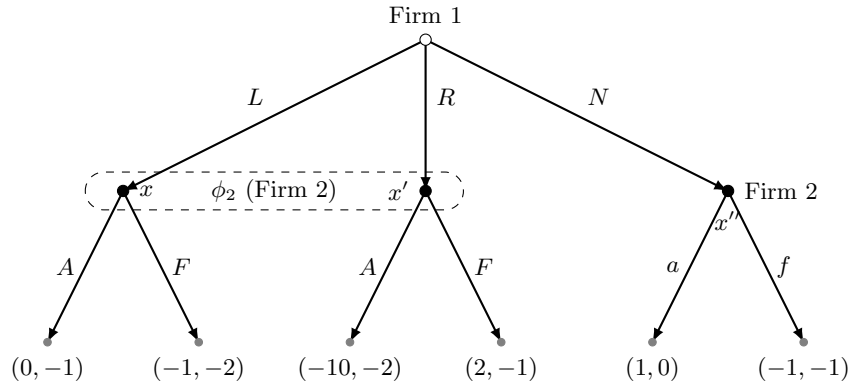
Example 46.

- (a) Here is yet another iteration of the market entry game:



There are two sets of pure strategy sequential equilibria: $((R, F), \mu_2(x') = 1)$ and $((N, A), \mu_2(x) > 1/2)$. We can apply forward induction by a dominance argument. If Firm 2 finds itself at ϕ_2 , then Firm 1 must have played L or R , yet L is strictly dominated by N . Since Firm 2 assumes Firm 1 must act rationally, Firm 2 deduces that Firm 1 played R , so $\mu_2(x) = 0$ and $\mu_2(x') = 1$. The family of equilibria involving strategy profile (N, A) thus fail forward induction. Given μ_2 , Firm 2's optimal strategy is F , and thus $((R, F), \mu(x') = 1)$ is consistent with forward induction.

(b) Suppose that if Firm 1 chooses N , Firm 2 faces a decision between a and f :



Again, we have two sets of pure strategy equilibria: $((R, Fa) : \mu_2(x') = 1)$ and $((N, Aa), \mu_2(x) > 1/2)$. In this case, N does not strictly dominate L . However, we can apply forward induction by appeal to equilibrium payoffs. Consider the family of equilibria involving the strategy profile (N, Aa) . Firm 1 achieves a payoff of 1 in equilibrium. If Firm 1 deviates to play L , then the maximum payoff Firm 1 can achieve is 0. Thus if Firm 1 is rational, it will not deviate to play L since it obtains a strictly worse payoff, regardless of Firm 2's action. On the other hand, if Firm 1 deviates to play R , the maximum payoff Firm 1 can achieve is 2. Thus, if Firm 2 finds itself at ϕ_2 , it can reason that Firm 1 chose R in the belief that

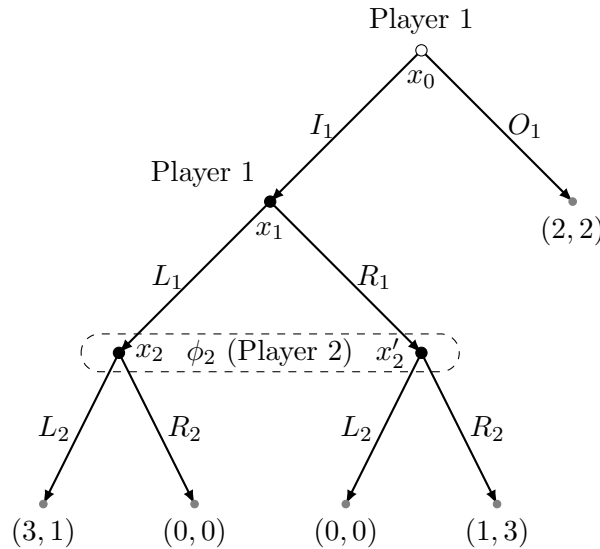
Firm 2 would choose F . Thus the only belief consistent with forward induction is $\mu_2(x) = 0, \mu_2(x') = 1$. This eliminates the family of equilibria involving (N, Aa) . Firm 2's optimal action at ϕ_2 given these beliefs is F . Hence the equilibrium $((R, Fa), \mu_2(x') = 1)$ is consistent with forward induction.

Like sequential rationality (at least, for a long time), forward induction reasoning in its entirety is somewhat abstract and not rigorously defined. Capturing the notion formally is difficult and there have been many attempts. Iterated weak dominance (see section 2.6) partially captures both sequential rationality and forward induction reasoning.

In games of perfect information, iterated weak dominance implies backward induction. At a penultimate node, any action that is not optimal is weakly dominated, so does not survive iterated weak dominance. Considering sets of immediate predecessor nodes in turn, we see that the set of strategy profiles that survive iterated weak dominance are precisely the set of backward induction solutions. In games of imperfect information, however, iterated weak dominance lacks the bite of, say, sequential equilibrium.

Now, iterated weak dominance also partially captures forward induction reasoning. To see this, consider the following extensive form version of Bach or Stravinsky (Example ??):

Example 47. Suppose two players are playing Bach or Stravinsky. Player 1 first has the option of whether to leave before the game is played, in which case the payoff is $(2, 2)$:



This game has a subgame perfect equilibrium $((O_1, R_1), R_2)$. Yet a forward induction argument implies that this equilibrium is unreasonable. For Player 1, I_1 is only optimal if she expects Player 2 to play L_2 , in which case she plays L_1 . Hence if information set ϕ_2 is reached, Player 2 can infer that Player 1 must have played L_1 , and so his best response is L_2 .

In reduced normal form, the game has the representation:

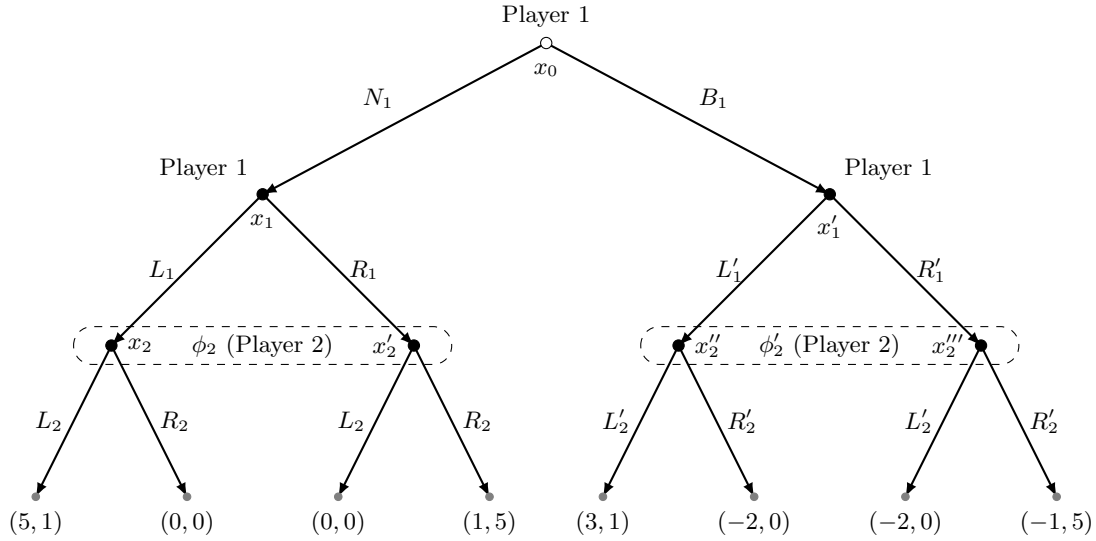
| | | |
|-----------|-------|-------|
| | L_2 | R_2 |
| O_1 | 2, 2 | 2, 2 |
| $I_1 L_1$ | 3, 1 | 0, 0 |
| $I_1 R_1$ | 0, 0 | 1, 3 |

Consider applying iterated weak dominance to this game. $I_1 R_1$ is strictly dominated by O_1 , and thus we eliminate it. No other strategies of Player 1 are weakly dominated at this stage. Having eliminated $I_1 R_1$, we see that R_2 is weakly dominated by L_2 , so we eliminate R_2 . Finally, of the remaining strategies, O_1 is strictly dominated, leaving only $(I_1 L_1, L_2)$. Iterated weak dominance in the reduced normal form game thus successfully captures our forward induction argument here.

There are many other concepts that formalize forward induction reasoning in a partial way. Stable equilibria, which we are about to discuss, capture a lot of the force of forward induction. In signalling games, the (possibly iterated) intuitive criterion and divine equilibrium also capture forward induction reasoning – see section 7.1.2.

Before we move on to stable equilibrium, it is worth mentioning a striking example of forward induction reasoning due to Ben-Porath & Dekel (1992):

Example 48 (Burning money). Consider the following game:



Two players are playing battle-of-the-sexes. Before the game, Player 1 has the choice to “burn” two units of utility. The subgame in which Player 1 has not burned utility (played N_1) is on the left, and the subgame in which Player 1 has burned two units of utility (played B_1) is on the right. There are four pure-strategy subgame perfect equilibria (which are also all sequential equilibria for some consistent system of beliefs): $s^1 = ((N_1, L_1, R'_1), (L_2, R'_2))$, $s^2 = ((N_1, L_1, L'_1), (L_2, L'_2))$, $s^3 = ((N_1, R_1, R'_1), (R_2, R'_2))$ and $s^4 = ((B_1, R_1, L'_1), (R_2, L'_2))$.

Let \succ_1 be a ranking of these equilibria by how favourable the outcome is for Player 1. Clearly, $s^1 \sim_1 s^2 \succ_1 s^4 \succ_1 s^3$. At first glance, Player 1 has an equilibrium selection problem. If she plays N_1 , then there is nothing in the notion of subgame perfection or sequential rationality to fix Player 2's beliefs at the information set ϕ_2 . If $\mu_2(x_2 | \phi_2) \geq \frac{5}{6}$, then L_2 is optimal and if $\mu_2(x_2 | \phi_2) \leq \frac{1}{6}$ then R_2 is optimal for Player 2.

A forward induction argument eliminates all but one of these equilibria, however, as we see if we apply iterated weak dominance. Note that any strategy of Player 1 that involves (B_1, R'_1) is strictly dominated – Player 1 will only burn if she intends to play L'_1 . Having eliminated these strategies, R'_2 is weakly dominated in the right subgame by L'_2 for Player 2. Thus Player 1 can guarantee herself a payoff of 3 by burning. Having eliminated these strategies, we see that any strategy involving (B_1, L'_1) strictly dominates any strategy involving (N_1, R_1) , and that (N_1, L_1, L'_1) strictly dominates (B_1, L_1, L'_1) . The only subgame perfect equilibrium surviving weak iterated dominance is thus $s^2 = ((N_1, L_1, L'_1), (L_2, L'_2))$.

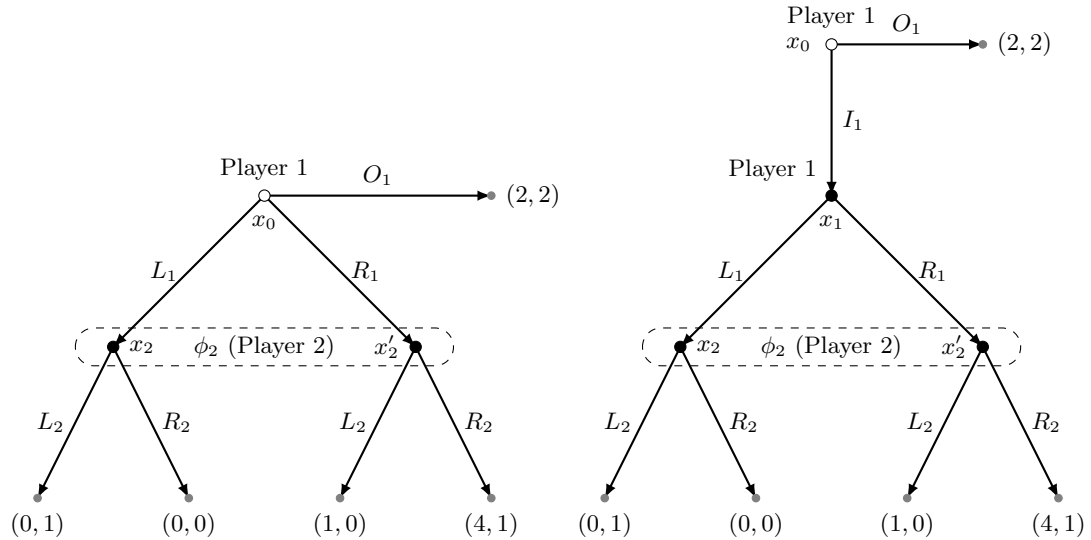
In the language of forward induction, at ϕ_2 , Player 2 can conclude that ϕ_2 is only reached because Player 1 has played (N_1, L_1) , since if Player 1 instead played (N_1, R_1) , she would have been better off burning and playing (B_1, L'_1) to achieve a payoff of 3. Thus Player 2 concludes that his optimal action at ϕ_2 is L_2 .

The interesting conclusion is that the *option* to burn cash allows Player 1 to select their most preferred equilibrium, even though Player 1 never exercises the option.

6.7 Stable equilibria

One principle we might hope an equilibrium concept satisfies is that irrelevant (“strategically neutral”) changes to the game tree do not affect the equilibrium outcome. Kohlberg & Mertens (1986) point out that sequential equilibrium fails this principle.

Example 49. Consider the following two games:



These games are identical, except that in the right game, Player 1's actions have been split in a way that seems irrelevant to play. Rather than one information set $\{x_0\}$, Player 1 now has two information sets $\{x_0\}$ and $\{x_1\}$, and chooses whether to play I_1 or O_1 at the first, and L_1 or R_1 at the second.

Intuitively, we would expect equilibrium outcomes to be the same in both games. Now, (O_1, L_2) can be supported as a sequential equilibrium strategy profile in the left game for beliefs $\mu_1 = (1, 0)$ over (L_2, R_2) and $\mu_2 = (1, 0, 0)$ over (O_1, L_1, R_1) . To see $((O_1, L_2), \mu)$ is consistent, consider a sequence $\epsilon_n \rightarrow 0$ and the totally mixed strategy profiles $\sigma^n = ((1 - \epsilon_n, 2\epsilon_n/3, \epsilon_n/3), (1 - \epsilon_n, \epsilon_n))$. Then the limit of the induced beliefs is $\lim_{n \rightarrow \infty} \mu^n = \mu$, and $\sigma^n \rightarrow \sigma$.

Yet there is no sequential equilibrium strategy profile involving O_1 in the right game. Indeed, the unique sequential equilibrium strategy profile in the right game is $((I_1, R_1), R_2)$. To see this, note L_1 is strictly dominated by R_1 for Player 1 (in both games). In the right game, therefore, R_1 is optimal at x_1 for any beliefs of Player 1, and thus in a consistent system of beliefs μ , $\mu_2(x'_2) = 1$, so Player 2's optimal action at ϕ_2 is R_2 .

Indeed, there are good reasons to think that in the left game, (O_1, L_2) is not a credible equilibrium strategy profile, on a forward induction argument: if ϕ_2 is reached, then because R_1 dominates L_1 for Player 1, Player 2 must conclude he is at x'_2 , so $\mu_2(x'_2) = 1$. Player 2's best response given this belief is R_2 . We see then that the reason that sequential equilibrium gives us different answers in the left and right games is because it does not capture forward induction arguments. Whereas the left game requires a forward induction argument to eliminate (O_1, L_2) , the structure of the right game is such that the forward induction argument is unnecessary.

Kohlberg & Mertens (1986) set out to develop an set-valued equilibrium concept that captures what they term *strategic stability*. A strategy profile σ is strategically stable if no player $i \in \mathcal{I}$ ever has an incentive to deviate from σ_i . Put another way, a strategy profile is strategically stable if it is *credible*, in the sense that at no point in a game, for any history, can any rational player make inferences that would lead them to prefer to do otherwise than σ . What does this require?

- *Sequential rationality*. As we discussed in e.g. Example 36, certain Nash equilibria are not credible because they rely on threats that, if it comes to it, a player would prefer not to follow through with. Sequential rationality rules this out. However, it is clearly not sufficient for strategic stability, as Example 49 illustrates.
- *Invariance to extensive form representation*. The addition or removal of irrelevant elements of the game tree should not affect the equilibrium outcome. Kohlberg and Mertens conceptualize this as follows. Recall that any extensive form game has a “reduced” normal form representation (section 1.3). Kohlberg and Mertens’ argue that the normal form representation of the game is all that should matter. This is a controversial take. Selten (1975) argues that a lot of information is potentially lost in the normal form representation of an extensive form game, and Kreps & Wilson (1982) argue analyses of normal form representations, in ignoring the role of beliefs

off-equilibrium-path, lack the power of extensive form analyses – for example, the reduced normal form shines no light on the question of what agents’ beliefs at the root node should be.

- *Weak dominance.* Kohlberg and Mertens argue that players never have a good reason to play a weakly dominated strategy, and strategic stability thus requires that no weakly dominated strategy is played in equilibrium. This could be extended to rule out equilibria that do not survive iterated weak dominance. But they don’t quite go this far when they define sets of stable equilibria, because then such sets are not guaranteed to exist. Thus they just require that strategic stability rules out that weakly dominated strategies are played in equilibrium.
- *Forward induction.* A credible Nash equilibrium should survive forward induction reasoning.

Definition 76 (Strategically stable equilibrium). Let G be a game. Let \mathcal{E} be the set of all sets E of Nash equilibria of G with the property that for every $\delta > 0$, there exists an ϵ -perturbation G_ϵ of G such that there is a Nash equilibrium σ^ϵ of G_ϵ within distance δ of some equilibrium $\sigma \in E$.⁵²

The *strategically stable set* S of Nash equilibria is the set

$$S = \bigcap_{E \in \mathcal{E}} E.$$

In the definition G can be a normal form or extensive form game. Kohlberg and Mertens show that:

- For any game, a strategically stable set exists.
- There exists a stable set that lies in a connected component of the set of Nash equilibria (and thus equilibrium outcomes are the same throughout this set; only off-equilibrium-path actions differ).
- Stable sets are invariant to irrelevant changes in the structure of the game tree.
- Any stable set contains a proper equilibrium (and thus a perfect equilibrium). However, if the game is extensive form, then the stable set is only guaranteed to contain a perfect equilibrium of the reduced normal form, not of the corresponding agent-normal form, and thus a stable set of an extensive form game does not necessarily contain a sequential equilibrium (Gul provides a counterexample.)
- No weakly dominated strategy is played in any strategy profile in any stable set.
- Any stable set of a game also contains the stable set of the game obtained by deleting a weakly dominated strategy (and thus contains the stable set of the game corresponding to those strategies that survive iterated weak dominance).

⁵²That is, the Hausdorff distance between σ^ϵ and E is less than δ . See Definition 123.

- Any stable set of the game also contains the stable set of the game obtained by deleting a strategy s_i of some player i that is never a best response to any of the strategy profiles s_{-i} of opponents that lie in the set. This captures forward induction.

Note that Kohlberg & Mertens' interpretation of strategic stability doesn't quite capture credibility (they need to sacrifice this to ensure that e.g. stable sets always exist), though strategically stable sets do contain all the credible equilibria.

6.8 Noncooperative theory of bargaining

The typical noncooperative approach to modelling bargaining envisages bargaining as a dynamic game, in which players make offers sequentially until an agreement is reached. The cooperative approach to bargaining is discussed in a later section.

A *bilateral bargaining* model involves two players dividing a surplus of size $v > 0$. In each period, one of the players is assigned to the position of *proposer*, and the other is assigned the position of *responder*. The proposer proposes an *offer* – a division of the surplus between the players – and the responder can accept or reject the offer. If accepting, the offer is implemented and payoffs are obtained. If rejecting, then (i) if the period is not a final period, then the period is ended and the next period begins; (ii) if the period is the final period, payoffs corresponding to failure to reach an agreement are realized. There is of course only a final period in finite-period bargaining games. Payoffs for each player i are typically discounted according to some discount factor $\delta_i \in (0, 1)$, capturing impatience – players would rather reach an agreement sooner rather than later.

Wlog, we can normalize the surplus v to be divided among the players to 1. The *set of feasible agreements* is

$$Z = \{(z, 1 - z) \mid z \in [0, 1]\},$$

where z is Player 1's share and $1 - z$ is Player 2's share.

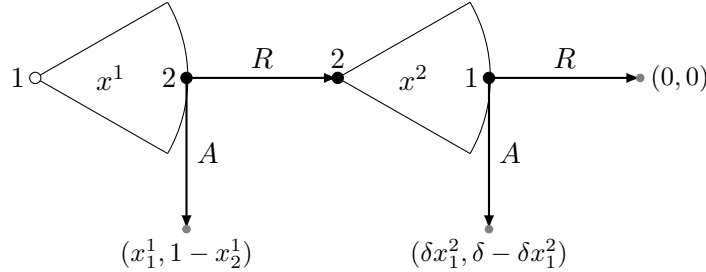
6.8.1 Finite period alternating bargaining

Ståhl (1972) analyses a finite period *alternating bargaining model*. This assigns one player (say, Player 1) as the proposer in odd periods and the other player (Player 2) in even periods. We can think of this as an *offer-counteroffer* situation – one player makes an offer, if the other player rejects they propose a counteroffer, and the process of haggling continues until an agreement is reached (or not). In period 1, Player 1 offers $x^1 = (x_1^1, 1 - x_1^1) \in Z$. For any offer $x^1 \in Z$, Player 2 accepts or rejects. If Player 2 rejects, then in period 2, Player 2 proposes an offer $x^2 = (x_1^2, 1 - x_1^2) \in Z$, which Player 1 accepts or rejects. In general, if Player 1 rejects, then in period 3, Player 1 makes an offer, and so on.

The simplest case is one period, in which case this is simply the ultimatum game. The two periods case is slightly more interesting:

Example 50 (Two-period alternating bargaining). Suppose there are two periods. If no offer is accepted by the end of the second period, the payoff to both players is 0. Assume

payoffs are discounted at a common discount rate $\delta \in (0, 1)$ so if the offer $(x_1^1, 1 - x_1^1)$ is accepted in period 1, payoff profile is $(x_1^1, 1 - x_1^1)$ and if the offer $(x_1^2, 1 - x_1^2)$ is accepted in period 2, the payoff profile is $(\delta x_1^2, \delta - \delta x_1^2)$. We can represent the game in extensive form as:



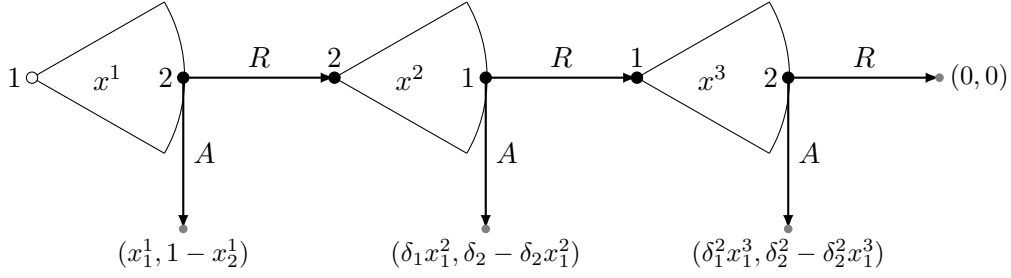
As in the ultimatum game, there are infinite Nash equilibria. However, there is a unique backward induction solution (and so a unique subgame perfect equilibrium):

$$\begin{aligned}
 s_1^*(\emptyset) &= (1 - \delta, \delta), \\
 s_1^*(x^1, R, x^2) &= \text{Accept for all } x^2 \in Z, \\
 s_2^*(x^1) &= \begin{cases} \text{Accept} & \text{if } x_1^1 \leq 1 - \delta, \\ \text{Reject} & \text{if } x_1^1 > 1 - \delta, \end{cases} \\
 s_2^*(x^1, R) &= (0, 1).
 \end{aligned}$$

The situation in the second period is simply that of the ultimatum game. As previously discussed, Player 2's optimal strategy in this case is to offer $(0, 1)$, and Player 1's optimal strategy is to accept. The payoff profile if the second round is reached and these strategies were implemented would be $(0, \delta)$. In period 1, given an offer $(x_1^1, 1 - x_1^1)$, Player 2 knows that if he rejects, he can obtain payoff δ . Hence it is optimal for him to reject any offer $x_1^1, 1 - x_1^1$ such that $1 - x_1^1 < \delta \Rightarrow x_1^1 > 1 - \delta$. Conversely, it is optimal for him to accept any offer such that $1 - x_1^1 \geq \delta \Rightarrow x_1^1 \leq \delta$. Given this, Player 1 can maximize her payoff by offering $(1 - \delta, \delta)$ in the first round.

The equilibrium outcome path is thus that Player 1 offers $(1 - \delta, \delta)$ in the first period and Player 2 accepts.

The more general setting involves $T = 2k + 1$ periods for some $k \in \mathbb{N}$, and unlike the example, we do not assume a common discount factor. Player 1 proposes in odd periods and Player 2 proposes in even periods.



We focus on the subgame perfect equilibrium.

The situation in the final period $T = 2k + 1$ is again equivalent to the ultimatum game, so Player 1 offers $x^{2k+1} = (1, 0)$ and Player 2 accepts any offer in Z .

In the penultimate period $t = 2k$, Player 1 knows she be assured of payoff δ_1 if she rejects, so she finds it optimal to accept any offer $x^{2k} = (x_1^{2k}, 1 - x_1^{2k})$ such that $x_1^{2k} \geq \delta_1$ and reject any offer such that $x_1^{2k} < \delta_1$. Hence Player 2's optimal offer is $x^{2k} = (\delta_1, 1 - \delta_1)$.

In period $t = 2k - 1$, Player 2 knows he can be assured of payoff $\delta_2(1 - \delta_1)$, so his optimal strategy is to accept any offer x^{2k-1} such that $1 - x_1^{2k-1} \geq \delta_2(1 - \delta_1) \Rightarrow x_1^1 \leq 1 - \delta_2(1 - \delta_1)$ and reject any offer x^{2k-1} such that $1 - x_1^{2k-1} < \delta_2(1 - \delta_1) \Rightarrow x_1^1 > 1 - \delta_2(1 - \delta_1)$. Knowing this, Player 1's optimal strategy is to offer $x^{2k-1} = (1 - \delta_2 + \delta_2\delta_1, \delta_2 - \delta_2\delta_1)$.

Recursively, we have that in every odd period $t = 2k - 2\ell + 1$, Player 1 offers

$$x^{2(k-\ell)+1} = \frac{1 - \delta_2 + \delta_1^\ell \delta_2^{\ell+1} (1 - \delta_1)}{1 - \delta_1 \delta_2},$$

and Player 2 will accept this offer. In particular, in period 1, Player 1 offers,

$$x_1^1 = \frac{1 - \delta_2 + \delta_1^k \delta_2^{k+1} (1 - \delta_1)}{1 - \delta_1 \delta_2}.$$

Hence the backward induction outcome path is that Player 1 offers x^1 as above and Player 2 will accept.

Note that as k increases, x_1^1 decreases. Indeed, as $k \rightarrow \infty$, $x_1^1 \rightarrow \frac{1-\delta_2}{1-\delta_2\delta_1}$, so $x^1 \rightarrow \left(\frac{1-\delta_2}{1-\delta_2}, \frac{\delta_2(1-\delta_1)}{1-\delta_1\delta_2}\right)$. If $\delta_1 = \delta_2 = \delta$, this simplifies to $x^1 \rightarrow \left(\frac{1}{1+\delta}, \frac{\delta}{1+\delta}\right)$. If additionally players become very patient, so $\delta \rightarrow 1$, we have $\lim_{\delta \rightarrow 1} \lim_{k \rightarrow \infty} x^1 = \left(\frac{1}{2}, \frac{1}{2}\right)$.

Example 51 (Ståhl model extension: costly proposing, no discounting). Suppose that players do not discount, but every time it is a player i 's turn to propose, they must pay a cost $c \in (0, 1)$. The game proceeds for T periods and if no agreement is reached in any period, the payoffs are $(0, 0)$. As before, Player 1 proposes in odd rounds and Player 2 proposes in even rounds. For any given T , this game has a unique subgame perfect equilibrium, which differs depending on whether T is odd or even.

First, suppose T is odd. Then in the final round, Player 1 proposes. Player 1's optimal proposal is clearly $(1, 0)$, which Player 2 will accept if sequentially rational given the alternative is that Player 2 receives 0. The payoffs in this period would be $(1 - c, 0)$. In period $T - 1$, Player 1 can obtain $1 - c$ if rejecting Player 2's offer, and therefore Player 2's optimal offer is $(1 - c, c)$, and accepting is optimal for Player 1. Player 2's payoff is $c - c = 0$. Iterating, we have that in every odd period, Player 1 proposes $(1, 0)$ and Player 2 accepts, and in every even period, Player 2 offers $(1 - c, c)$ and Player 1 accepts. The equilibrium path has an agreement in the first period, with equilibrium payoffs $(1 - c, 0)$.

If T is even, then in the final round, it is Player 2 who proposes. The argument is identical to the case where T is odd, except with the role of Players 1 and 2 reversed: we will have that Player 1 offers $(c, 1 - c)$ in every odd period, and Player 2 accepts, while in every even period, Player 2 offers $(0, 1)$ and Player 1 accepts. On the equilibrium path, an agreement is reached in the first period and the equilibrium payoffs are $(0, 1 - c)$. Unlike the standard Ståhl model, in this version of the game there is a *last-mover advantage*.

Now suppose $T \rightarrow \infty$. Then there is no unique subgame perfect equilibrium. Let Γ^T denote the game of T periods and let π_i^T denote the subgame perfect equilibrium payoff of player i in Γ^T . Consider the sequence $\{\pi_i^T\}$. For player 1,

$$\pi_1^T = \begin{cases} 1 - c & \text{if } T \text{ is odd,} \\ 0 & \text{if } T \text{ is even.} \end{cases}$$

Hence the subsequence $\{\pi_1^{2k+1}\}_{k \in \mathbb{N}}$ has $\pi_1^{2k+1} = 1 - c$ for all k , so $\lim_{k \rightarrow \infty} \pi_1^{2k+1} = 1 - c$. Conversely, the subsequence $\{\pi_1^{2k}\}_{k \in \mathbb{N}}$ has $\pi_1^{2k} = 0$ for all k so $\lim_{k \rightarrow \infty} \pi_1^{2k} = 0$. We therefore have that

$$\limsup_{T \rightarrow \infty} \pi_1^T = 1 - c > 0 = \liminf_{T \rightarrow \infty} \pi_1^T,$$

so the set of subgame perfect equilibrium payoffs for Player 1 in the limit as $T \rightarrow \infty$ is not a singleton. The same can be derived for Player 2.

6.8.2 Infinite period alternating bargaining

Rubinstein (1982) analyses an infinite period alternating bargaining model. The structure of the model is identical to Ståhl's: two players bargain to divide a surplus of normalized value 1, with Player 1 proposing in odd periods and Player 2 proposing in even periods. We let Player 1 have discount rate δ_1 and Player 2 have discount rate δ_2 .

Proposition 60. *There is a unique subgame perfect equilibrium in the Rubinstein sequential bargaining game. Furthermore, the subgame perfect equilibrium is as follows: Whenever Player 1 proposes, she offers a division $(x, 1 - x)$ with $x = \frac{1 - \delta_2}{1 - \delta_1 \delta_2}$. Player 2 accepts any division giving her at least $1 - x$. Whenever Player 2 proposes, she offers a division $(y, 1 - y)$ with $y = \frac{\delta_1(1 - \delta_2)}{1 - \delta_1 \delta_2}$. Player 1 accepts any division giving her at least y . Thus on the equilibrium path, bargaining ends in the first round with division $(x, 1 - x)$.*

Proof. First, we confirm the proposed strategies constitute a subgame perfect equilibrium. By the one-shot deviation principle (Theorem 24), it is sufficient to check one-shot deviations. Note $y = \delta_1 x$. First consider any period where Player 1 makes an offer. Player 1 has no profitable one-shot deviation: suppose she offers $(x', 1 - x')$; if $x' < x$, Player 2 accepts and Player 1 receives a strictly lower payoff $x' < x$; if $x' > x$, Player 2 rejects and Player 1 receives, in present value terms, $\delta_1 y = \delta_1^2 x < x$. Likewise, if Player 2 rejects Player 1's offer, Player 2 receives, in present value terms $\delta_2(1 - y) = \delta_2(1 - \delta_1 x) = \delta_2 - \delta_1 \delta_2 x = 1 - x$ [to see this final equality, note it rearranges to $x = \frac{1 - \delta_2}{1 - \delta_1 \delta_2}$ which is definitional.] The argument for periods where Player 2 makes an offer is similar.

Next, we show the equilibrium is unique. Let Π_1 be the set of subgame perfect equilibrium payoffs for Player 1. Let $\underline{v}_1 = \inf \Pi_1$ and $\bar{v}_1 = \sup \Pi_1$. Consider any period where Player 2 makes an offer. Player 1 will accept any offer greater than $\delta_1 \bar{v}_1$ and reject any offer less than $\delta_1 \underline{v}_1$. Hence Player 2 can secure at least $1 - \delta_1 \bar{v}_1$ by proposing division $(\delta_1 \bar{v}_1, 1 - \delta_1 \bar{v}_1)$, and can secure at most $1 - \delta_1 \underline{v}_1$ by proposing $(\delta_1 \underline{v}_1, 1 - \delta_1 \underline{v}_1)$.

Now consider a period where Player 1 makes an offer. For Player 2 to accept Player 1's offer, Player 1 must offer at least $\delta_2(1 - \delta_1 \bar{v}_1)$, and hence

$$\bar{v}_1 \leq 1 - \delta_2(1 - \delta_1 \bar{v}_1),$$

implying

$$\bar{v}_1 \leq \frac{1 - \delta_2}{1 - \delta_1 \delta_2}.$$

Likewise, Player 2 will certainly accept if offered more than $\delta_2(1 - \delta_1 \underline{v}_1)$. Hence

$$\underline{v}_1 \geq 1 - \delta_2(1 - \delta_1 \underline{v}_1),$$

implying

$$\underline{v}_1 \geq \frac{1 - \delta_2}{1 - \delta_1 \delta_2}.$$

Combining inequalities, we have

$$\bar{v}_1 \leq \frac{1 - \delta_2}{1 - \delta_1 \delta_2} \leq \underline{v}_1.$$

Since by definition, $\underline{v}_1 \leq \bar{v}_1$, it follows that in any subgame perfect equilibrium, Player 1 receives $v_1 = \underline{v}_1 = \bar{v}_1 = \frac{1 - \delta_2}{1 - \delta_1 \delta_2}$. Making a similar argument for Player 2 completes the proof. \square

As with the Ståhl model, relatively more patient players in the Rubinstein model receive higher payoffs – Player 1's payoff is increasing in δ_1 and decreasing in δ_2 , whereas Player 2's payoff is increasing in δ_2 and decreasing in δ_1 . If there is a common discount factor $\delta = \delta_1 = \delta_2$, then Player 1's offer in the first period is $\left(\frac{1}{1 + \delta}, \frac{\delta}{1 + \delta}\right)$. This illustrates Player 1's *first-mover advantage* – since the players are impatient, Player 2 is willing to accept a smaller slice of the pie now than a slightly larger slice in the next period – and likewise, Player 1 would be willing to accept a slightly smaller slice were we to reach a period where Player 2 was proposing.

Example 52 (Rubinstein extension: outside options). Suppose now that when considering a proposal, the responding player i has an option to either continue to the next period or to end the game immediately in which case Player 1 and Player 2 receive payoffs $c_1, c_2 \in [0, 1]$ respectively, with $c_1 + c_2 \leq 1$. We assume a common discount rate δ . The game is otherwise identical to the standard Rubinstein model

Proposition 61. *In the modified Rubinstein bargaining game with outside options and a common discount rate δ , there is a unique subgame perfect equilibrium. Furthermore, the subgame perfect equilibrium is as follows: Whenever Player 1 proposes, she offers a division $(x, 1 - x)$ with*

$$x = \begin{cases} \frac{1}{1+\delta} & \text{if } c_i \leq \frac{\delta}{1+\delta} \text{ for all } i = 1, 2, \\ 1 - \delta(1 - c_1) & \text{if } c_1 > \frac{\delta}{1+\delta} \text{ and } c_2 \leq \delta(1 - c_1), \\ 1 - c_2 & \text{otherwise,} \end{cases}$$

in every odd period. Player 2 accepts any division giving her at least $1 - x$. Whenever Player 2 proposes, she offers a division $(y, 1 - y)$ with

$$y = \begin{cases} \frac{\delta}{1+\delta} & \text{if } c_i \leq \frac{\delta}{1+\delta} \text{ for all } i = 1, 2, \\ \delta(1 - c_2) & \text{if } c_2 > \frac{\delta}{1+\delta} \text{ and } c_1 \leq \delta(1 - c_1), \\ c_1 & \text{otherwise,} \end{cases}$$

and Player 1 accepts any division giving her at least y .

Proof. First we verify the proposed strategies constitute a subgame perfect equilibrium. Consider any period where Player 1 proposes. There are three cases to check:

- (i) $c_1, c_2 \leq \frac{\delta}{1+\delta}$. Suppose Player 1 offers $(x', 1 - x')$. If $x' < x$, then Player 2 accepts and Player 1's payoff will clearly be lower. If $x' > x$, then Player 2 rejects. Player 1 accepts $y = \delta x$ in the next period which in present value terms is worth $\delta y = \delta^2 x < x$. Hence Player 1 has no profitable deviation. Suppose Player 2 rejects $(x, 1 - x)$. If she opts to end the game, she receives $c_2 \leq \frac{\delta}{1+\delta} = 1 - x$, so this is not a profitable deviation. If she advances the game to the next period, she receives, in present value terms, a payoff of $\delta(1 - y) = \delta(1 - \delta x) = \delta - \delta^2 x = 1 - x$ [to verify the final equality holds, note that rearranging it yields $x = \frac{1-\delta}{1-\delta^2} = \frac{1-\delta}{(1-\delta)(1+\delta)} = \frac{1}{1+\delta}$]. This is no greater than accepting, so is not a profitable deviation.
- (ii) $c_1 > \frac{\delta}{1+\delta}, c_2 \leq \delta(1 - c_1)$. Note that $c_2 \leq \delta(1 - c_1)$ implies $c_2 < \frac{\delta}{1+\delta}$. Suppose Player 1 offers $(x', 1 - x')$. Again, if $x' < x$ the Player 2 accepts and Player 1 receives strictly lower payoff. If $x' > x$, Player 2 rejects and continues to the next period, so Player 1 receives, in present value terms, $\delta^2 x < x$. Hence Player 1 has no profitable deviation. Suppose Player 2 rejects $(x, 1 - x)$. If Player 2 ends the game, she receives $c_2 \leq \delta(1 - c_1) = 1 - x$, so this is not a profitable deviation. If Player 2 continues the game, she receives in present value terms $\delta(1 - c_1)$, which is precisely what she receives if accepting the offer, so this is not a profitable deviation.

- (iii) $c_1 > \frac{\delta}{1+\delta}$, $c_2 > \delta(1 - c_1)$. Suppose Player 1 offers $(x', 1 - x')$. As before, if $x' < x$ then Player 2 accepts and Player 1 is strictly worse off compared to the division $(x, 1 - x)$. If $x' > x$, then Player 2 rejects and ends the game to receive outside payoff c_2 . In this case, Player 1 receives $c_1 \leq 1 - c_2$, so this is not a profitable deviation. Suppose Player 2 rejects $(x, 1 - x)$. Since Player 2 would receive c_2 if accepting, ending the game immediately is not a profitable deviation. If advancing to the next period, Player 2 receives, in net present value terms, $\delta(1 - c_1) < c_2$ or $\delta(1 - \delta(1 - c_2)) < c_2$, where the inequality follows from $c_2 > \frac{\delta}{1+\delta}$. Hence this is not a profitable deviation for Player 2.

The situation in even periods is symmetric.

Next, we show the equilibrium is unique. Let Π_i be the set of subgame perfect equilibrium payoffs for Player i , and define $\underline{v}_i = \inf \Pi_i$ and $\bar{v}_i = \sup \Pi_i$. Fixing $i = 1, 2$ consider any period where i is proposing. Player $j \neq i$ will certainly accept any offer that gives her at least $\max\{\delta\bar{v}_j, c_j\}$, and hence player i can obtain at least $\underline{v}_i \geq 1 - \max\{\delta\bar{v}_j, c_j\}$. Likewise, for player j to accept, she must receive at least $\max\{\delta\underline{v}_i, c_j\}$, so player i receives at most $\bar{v}_i \leq 1 - \max\{\delta\underline{v}_j, c_j\}$. Hence

$$\bar{v}_i - \underline{v}_i \leq \max\{\delta\bar{v}_j, c_j\} - \max\{\delta\underline{v}_j, c_j\}.$$

If $c_j \geq \delta\bar{v}_j \geq \delta\underline{v}_j$, then we have $\bar{v}_i - \underline{v}_i \leq 0$, and since $\bar{v}_i \geq \underline{v}_i$, it follows that $\bar{v}_i = \underline{v}_i$. If $c_j \leq \delta\underline{v}_j \leq \delta\bar{v}_j$, then we have $\bar{v}_i - \underline{v}_i \leq \delta(\underline{v}_j - \bar{v}_j) \leq 0$. Since $\bar{v}_i - \underline{v}_i \geq 0$, it follows that $\bar{v}_i = \underline{v}_i$. Finally, if $\delta\underline{v}_j \leq c_j \leq \delta\bar{v}_j$, then we have $\bar{v}_i - \underline{v}_i \leq c_j - \delta\bar{v}_j \leq 0$, and since $\bar{v}_i - \underline{v}_i \geq 0$, it again follows that $\bar{v}_i = \underline{v}_i$. Hence Π_i contains only one value for both i . This completes the proof. \square

7 Communication games

Most economic settings (and non-economic settings in life) involve communication. There is often a tension between the incentives of agents and the truthful revealing of private information, and so credible communication is important. High productivity workers will want to communicate this to firms in order to get hired or secure high pay. But it is not just enough to tell an employer you will make a great employee, because bad employees will also want to claim they are great employees.

Credibly signalling information therefore often requires costly signals.

7.1 Signalling games

Signalling games, introduced by Spence (1974), are incomplete information dynamic games that capture the idea that signals can sometimes be used to credibly communicate one's type. There are two periods and two players – a sender and a receiver. The structure of a signalling game proceeds as follows:

- Stage 0: Nature chooses a type $\theta \in \Theta$ of Player 1 from a distribution p with support Θ .

- Stage 1: Player 1 (the *sender*) observes θ and chooses a *message* $m \in M$ (or *signal*).
- Stage 2: Player 2 (the *receiver*) observes m and chooses action $a \in A$.
- Payoffs: at the end of stage 2, payoffs $u_1(m, a, \theta)$ and $u_2(m, a, \theta)$ are realized.
- Strategies: A pure strategy of the sender is a mapping $s_1 : \Theta \rightarrow M$, and a pure strategy of the receiver is a mapping $s_2 : M \rightarrow A$.

Signalling games arise in many economic settings. Some notable examples:

Example 53.

- (a) *Job market signalling.* This is Spence's original example. Player 1 is a prospective worker and Player 2 is a perfectly competitive labour market. The prospective worker's type θ is her ability, and she can choose to acquire a level of education e , the cost of which is decreasing in ability. The firm observes the worker's education level and makes a wage offer based on her expected ability conditional on her education level. The worker hopes to signal her ability via her education level.
- (b) *Initial public offerings.* Player 1 is a founder of a private company and Player 2 is a set of potential investors. The founder's type θ is the future profitability of the company, and she chooses what fraction of the company to float publicly and the price at which these shares will be offered. The potential investors choose whether to accept or reject the founder's offer in response. The founder hopes to signal high future expected profitability via offer price and the size of the stake to be sold.
- (c) *Monetary policy.* Player 1 is a central bank, and Player 2 is the set of firms. The central bank's type is its policy preferences over unemployment and inflation. In the first period, the central bank sets an inflation level $m \in M$. Firms form expectations a about inflation in the second period given inflation in the first period. The central bank wishes to signal its policy preferences so that period 2 inflation is low.
- (d) *Pretrial negotiations.* Player 1 is a defendant in a civil case, and Player 2 is the plaintiff. The defendant's type is the strength of his defence. The defendant makes a settlement offer $m \in M$, and the plaintiff responds by accepting or rejecting. If the plaintiff rejects, the parties go to trial. The defendant hopes to signal he has a strong case.
- (e) *Signalling preferences.* In the US academic job market (at least in economics), most candidates are often less willing to take jobs at liberal arts colleges than they are at universities, but candidates have idiosyncratic preferences and so some prefer liberal arts colleges. A similar situation is faced by universities in less attractive locations, universities outside the US, and so on. The AEA lets job

markets candidates send two signals to universities and colleges indicating they are willing to work there. There is an opportunity cost to using one of these signals, because then the signal cannot be used elsewhere. Player 1 is a job market candidate and Player 2 is a liberal arts college. The type of Player 1 is their willingness to take a liberal arts college AP position.

A similar situation often arises in dating apps. In heterosexual online dating, there is a lot of congestion – men swipe right on women a lot without necessarily being particularly interested. Many dating apps give users a limited number of signals they can send to signal genuine interest.

Definition 77 (Sequential equilibrium in signalling games). In a (finite) signalling game, a sequential equilibrium is an assessment (s, μ) such that

- (i) Player 1's strategy $s_1(\theta)$ is optimal given Player 2's strategy $s_2(m)$, that is, $s_1(\theta)$ solves

$$\max_{m \in M} u_1(m, s_2(m), \theta) \quad \text{for all } \theta \in \Theta.$$

- (ii) Player 2's beliefs are compatible with Bayes' rule, that is, if $\mathbb{P}\{s_1(\theta) = m\} > 0$ then

$$\mu_2(\theta | m) = \frac{\mathbb{P}\{s_1(\theta) = m\}p(\theta)}{\sum_{\theta' \in \Theta} \mathbb{P}\{s_1(\theta') = m\}p(\theta')}.$$

- (iii) Player 2's strategy is optimal given μ_2 and m , that is, $s_2(m)$ solves

$$\max_{a \in A} \sum_{\theta \in \Theta} u_2(m, a, \theta) \mu_2(\theta | m) \quad \text{for all } m \in M.$$

The definition can be easily extended if the type space or action spaces are infinite.⁵³

It is helpful to partition the sequential equilibria in signalling games into three categories:

- (a) *Separating equilibria*. Player 1 will play different actions in equilibrium depending on her type, and so Player 2 learns Player 1's type perfectly.
- (b) *Pooling equilibrium*. Player 1 will play the same action in equilibrium regardless of her type, so no information is transmitted to Player 2.

⁵³In particular, let $f_{s_1(\theta)}$ be the pdf of the distribution over M induced by $s_1(\theta)$ and let p be the pdf of the distribution from which θ is drawn. (ii) becomes: If $f_{s_1(\theta)}(m) > 0$ then

$$\mu_2(\theta | m) = \frac{f_{s_1(\theta)}(m)p(\theta)}{\int_{\Theta} f_{s_1(\theta')}(m)p(\theta') d\theta'},$$

where μ_2 now represents a density. (iii) becomes: $s_2(m)$ solves

$$\max_{a \in A} \int_{\Theta} u_2(m, a, \theta) \mu_2(\theta | m) d\theta.$$

- (c) *Semi-separating equilibrium.* Some actions are chosen by several types of Player 1, and others are chosen by a single type. Player 2 thus learns Player 1's type imperfectly.

7.1.1 Job market signalling

Spence's original model concerned job market signalling. Player 1 is a *worker* whose ability (i.e. productivity) is given by $\theta \in \{\theta_L, \theta_H\}$, where $\theta_H > \theta_L > 0$. The worker knows her own ability, and the labour market assigns probability λ that she has productivity θ_H . The worker chooses an education level $e \in E$. To acquire e costs $c(e, \theta)$.

Assumption. $\frac{\partial c(e, \theta)}{\partial e} > 0$ and $\frac{\partial c(e, \theta)}{\partial e \partial \theta} < 0$.

This assumption captures the idea that (i) on the margin, it costs more to acquire more education since this takes more time and money, and (ii) higher ability find education less costly on the margin – i.e. they learn more efficiently and so can complete education with less effort or time.

Once the worker chooses her education level, the perfectly competitive firm make wage offers. For simplicity, assume the reservation wage of workers is 0, regardless of ability.

Player 2 is a competitive labour market – for example, a set of two or more identical firms in Bertrand competition. Let $\mu(e)$ denote the market's belief that the worker with education level e is high productivity (type θ_H). Since each firm in the market is perfectly competitive, it offers a wage equal to expected productivity,

$$w(e) = \mu(e)\theta_H + (1 - \mu(e))\theta_L.$$

This generates a wage schedule $w : E \rightarrow \mathbb{R}_+$.

Consider the problem facing the worker given the market's beliefs $\mu(e)$ and the wage schedule. The worker with ability θ chooses e to solve

$$\max_e w(e) - c(e, \theta).$$

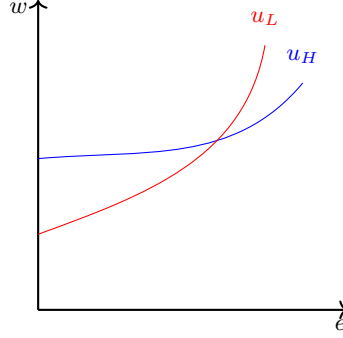
Consider the worker's indifference curves in (e, w) -space. Implicitly differentiating $u_1(w, e, \theta) = \bar{u}$ gives

$$\left. \frac{dw}{de} \right|_{u=\bar{u}} = \frac{\partial c(e, \theta)}{\partial e} > 0,$$

so the indifference curves are upward-sloping. Furthermore,

$$\left. \frac{d}{d\theta} \frac{dw}{de} \right|_{u=\bar{u}} = \frac{\partial c(e, \theta)}{\partial e \partial \theta} < 0.$$

Thus indifference curves are flatter for higher ability workers.



An example where $u_H(e, w, \theta_H) = u_L(e, w, \theta_L) = \bar{u}$. The two indifference curves cross only once. Note this implies that a worker of type θ_H will choose a weakly higher education level than a worker of type θ_L .

To make this argument more formally, consider the worker's utility function $u_1(e, w(e), \theta) = w(e) - c(e, \theta)$. Note this is a function in two arguments, (e, θ) . Because $\frac{\partial c(e, \theta)}{\partial e \partial \theta} < 0$ and $\frac{\partial w(e)}{\partial e} = 0$, we have that $\frac{\partial u_1}{\partial e \partial \theta} > 0$ and thus u_1 is supermodular in (e, θ) . Applying Topkis' monotonicity theorem (Theorem 69), we immediately have that the optimal education level for a type- θ_H worker is weakly greater than for a type- θ_L worker.

On the equilibrium path, the wage schedule $w(e)$ (or equivalently, the market's belief $\mu(e)$) is pinned down by the worker's choice. For those levels of education e that are not chosen in equilibrium, however, $w(e)$ can be anything in the interval (θ_L, θ_H) , since perfect Bayesian equilibrium does not restrict beliefs off the equilibrium path. This allows for many possible equilibria.

Example 54 (Equilibria in the Spence job market signalling game).

- (a) *Separating equilibria.* First, consider separating equilibria, i.e. equilibria for which $e(\theta_L) \neq e(\theta_H)$.

Lemma 20. *In a separating equilibrium, $w(e(\theta)) = \theta$ for $\theta \in \{\theta_L, \theta_H\}$.*

Proof. Since $e(\theta_L) \neq e(\theta_H)$, if μ is to be consistent with Bayes' rule, we must have

$$\begin{aligned}\mu(e(\theta_L)) &= \frac{\mathbb{P}\{e(\theta_H) = \theta_L\}\lambda}{\mathbb{P}\{e(\theta_H) = \theta_L\}\lambda + \mathbb{P}\{e(\theta_L) = \theta_L\}(1 - \lambda)} = \frac{0}{1 - \lambda} = 0, \\ \mu(e(\theta_H)) &= \frac{\mathbb{P}\{e(\theta_H) = \theta_L\}\lambda}{\mathbb{P}\{e(\theta_H) = \theta_H\}\lambda + \mathbb{P}\{e(\theta_L) = \theta_H\}(1 - \lambda)} = \frac{\lambda}{\lambda} = 1.\end{aligned}$$

Given these beliefs, we have wages $w(e(\theta_H)) = \theta_H$ and $w(e(\theta_L)) = \theta_L$. □

Lemma 21. *In a separating equilibrium, $e(\theta_L) = 0$.*

Proof. Suppose otherwise. By the previous lemma, $w(e(\theta_L)) = \theta_L$. Suppose a type- θ_L worker instead chooses $e = 0$. Now, $\mu(0) \geq 0$ so $w(0) \geq \theta_L$, and $c(0, \theta_L) < c(e(\theta_L), \theta_L)$. Hence $u_1(0, w(\cdot), \theta) > u_1(e(\theta_L), w(\cdot), \theta)$, so for any $e(\theta_L) > 0$, $e = 0$ is a profitable deviation. □

Finally, we need to find $e(\theta_H)$ s.t. high ability workers prefer to select $e(\theta_H)$ over $e = 0$ and earn the higher wage $w = \theta_H$, while lower ability workers prefer to select $e(\theta_L) = 0$.

Lemma 22. *For any $e(\theta_H)$ that satisfies*

$$\begin{aligned}\theta_H - c(e(\theta_H), \theta_H) &\geq \theta_L - c(0, \theta_H), \\ \theta_L - c(0, \theta_L) &\geq \theta_H - c(e(\theta_H), \theta_L),\end{aligned}$$

there is a separating equilibrium in which type- θ_H workers earn θ_H .

Proof. If the two conditions hold, then clearly no type of worker benefits from choosing the signal of the other type. The only other thing to check is that no type of worker benefits by choosing some $e \notin \{e(\theta_L), e(\theta_H)\}$. This requires setting $\mu(e)$ sufficiently low for all $e \notin \{e(\theta_L), e(\theta_H)\}$, so that

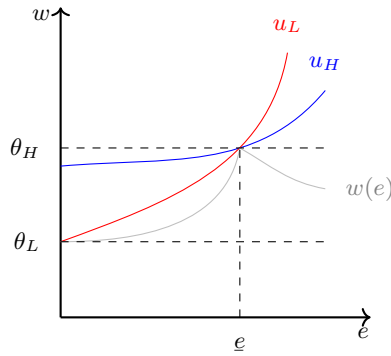
$$\begin{aligned}w(e) - c(e, \theta_H) &\leq \theta_H - c(e(\theta_H), \theta_H), \\ w(e) - c(e, \theta_L) &\leq \theta_L - c(0, \theta_L).\end{aligned}$$

This can be assured trivially by setting $\mu(e) = 0$ for all $e < e(\theta_H)$. Then $w(e) = \theta_L$ for all $e < \theta_H$. \square

We can find bounds on the range of values of $e(\theta_H)$ that are possible in a separating equilibrium. By the single crossing property, if the second condition in Lemma 22 is satisfied with equality then the first condition holds with strict inequality. Thus define \underline{e} so that

$$\theta_L - c(0, \theta_L) = \theta_H - c(\underline{e}, \theta_H).$$

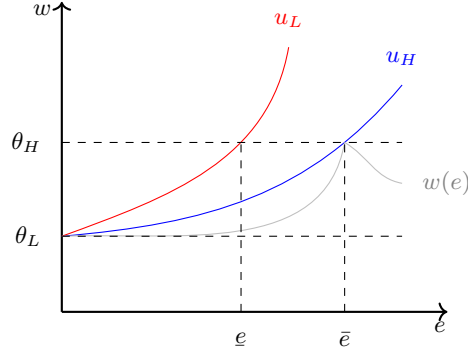
We have a separating equilibrium with $e(\theta_H) = \underline{e}$:



Likewise, the single crossing property implies that if the first condition in Lemma 22 holds with equality, then the second condition holds with strict inequality. Thus define \bar{e} so that

$$\theta_H - c(\bar{e}, \theta_H) = \theta_L - c(0, \theta_H).$$

We have a separating equilibrium with $e(\theta_H) = \bar{e}$:



Clearly, \underline{e} is the minimum value of $e(\theta_H)$ that can be supported in a separating equilibrium, and \bar{e} is the maximum. Any value in $[\underline{e}, \bar{e}]$ is a value of $e(\theta_H)$ in some separating equilibrium.

The single crossing property is key to separating equilibrium – that allows high ability workers to acquire education that is relatively low cost for them but prohibitively costly for lower ability workers. The differential costs allow separation.

- (b) *Pooling equilibria.* In a pooling equilibrium, every worker chooses the same education level $e^P = e(\theta_L) = e(\theta_H)$. Since the market's beliefs in equilibrium must be consistent with Bayes' rule, we must have $\mu(e^P) = \lambda$. Therefore the equilibrium wage is

$$w(e^P) = \lambda\theta_H + (1 - \lambda)\theta_L =: \mathbb{E}\theta.$$

Define \hat{e} to satisfy

$$\mathbb{E}\theta - c(\hat{e}, \theta_L) = \theta_L = c(0, \theta_L),$$

i.e. \hat{e} is such that a lower-ability worker is indifferent between acquiring \hat{e} and no education at all.

Proposition 62. *For any $e^P \in [0, \hat{e}]$, there is a pooling equilibrium in which all types of worker choose e^P with probability 1.*

Proof. Fix $e^P \in [0, \hat{e}]$. Suppose $\mu(e^P) = \lambda$ and $\mu(e) = 0$ for all $e \neq e^P$, with corresponding wages $w(e) = \theta_L$. Then $w(e) < w(e^P)$. Hence no worker benefits by deviating to $e > e^P$, and by definition of \hat{e} , type- θ_L workers prefer e^P to any $e < e^P$. By the single crossing property, type- θ_H workers must also prefer e^P to any $e < e^P$. \square

An interesting consequence of Proposition 62 is that there can be inefficiency even in a pooling equilibrium. If $e^P > 0$, then all workers are needlessly acquiring some costly education.

- (c) *Hybrid equilibria.* In a hybrid equilibrium, one or both types of worker randomize. For example, suppose $e(\theta_L) = 0$ but that type θ_H chooses $e(\theta_H) = 0$ with probability q and $e(\theta_H) = \bar{e}$ for some $\bar{e} > 0$ with probability $1 - q$. The market knows

on observing education level \bar{e} that the worker is type θ_H , so $\mu(\bar{e}) = 1$ and the market pays wage $w(\bar{e}) = \theta_H$. By Bayes' rule, the market's belief that a worker with education level 0 is type θ_H is $\mu(0) = \frac{\lambda q}{\lambda q + 1 - \lambda}$, and so the market pays wage $w(0) = \frac{\lambda q}{\lambda q + 1 - \lambda} \theta_H + \frac{1 - \lambda}{\lambda q + 1 - \lambda} \theta_L$. Now we need that type θ_H is indifferent between education levels 0 and \bar{e} , or else θ_H will not randomize in equilibrium. Thus we require that \bar{e} solves $w(0) = w(\bar{e}) - c(\bar{e}, \theta_H)$. Finally, let $\mu(e) = 0$ for all $e \notin \{0, \bar{e}\}$. Then this constitutes a hybrid equilibrium.

All three types of equilibrium exist in the Spence model, because sequential equilibrium imposes no restrictions on off-equilibrium-path beliefs in this setting.

In the standard job market signalling model we have discussed, education has no effect on productivity. In the separating equilibrium, education reveals information, and so correlates with wages. In the pooling equilibrium, education reveals no information, and yet workers may still end up acquiring some positive level of education.

Any equilibrium in which some type of worker acquires any positive level of education is inefficient. If firms in the labour market were able to observe types directly, agents would not need to burden the cost of acquiring an education. Again, this reflects the (unrealistic) assumption that education does not improve productivity.

7.1.2 Forward induction in signalling games

“Despite the name we have given it, the Intuitive Criterion is not completely intuitive” – In-Koo Cho and David M. Kreps, 1987.

While there are generally infinitely many equilibria in the signalling game, some of these equilibria are arguably less appealing than others. Cho & Kreps (1987) introduce the *intuitive criterion*, a criterion for refinement of equilibrium in signalling games to eliminate these arguably less appealing equilibria. This is an application of forward induction reasoning.

Definition 78 (Intuitive criterion).

- (a) *Best responses.* Let $BR(T, m)$ denote the set of best responses of Player 2 if Player 1 has chosen m and Player 2's beliefs have support in $T \subseteq \Theta$. That is,

$$BR(T, m) = \bigcup_{\mu \in \Delta(T)} \arg \max_{a \in A} \sum_{\theta \in T} u_2(m, a, \theta) \mu(\theta).$$

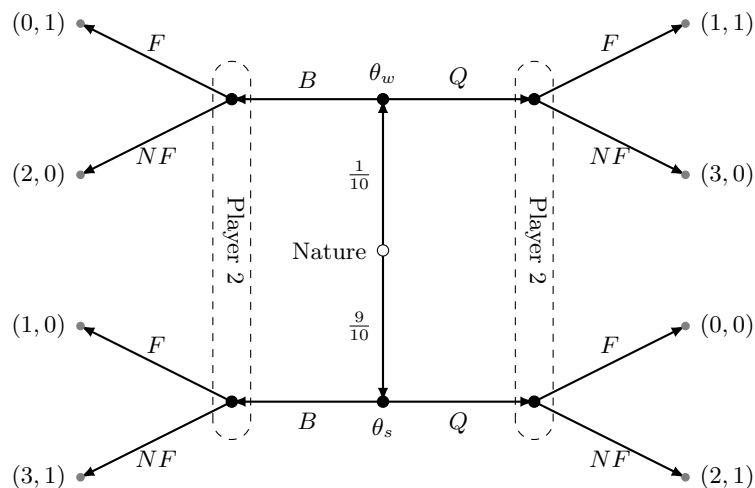
- (b) *Intuitive criterion.* A sequential equilibrium (or PBE) (s^*, μ) of a signalling game G is said to *fail the intuitive criterion* if there exists $m \in M$, $\theta' \in \Theta$ and $J \subseteq \Theta$ such that

$$\begin{aligned} u_1(s^*, \theta) &> \max_{a_2 \in BR(\Theta, a_1)} u_1(m, a, \theta) && \text{for all } \theta \in J, \text{ and} \\ u_1(s^*, \theta) &< \min_{a \in BR(\Theta - J, m)} u_1(m, a, \theta'). \end{aligned}$$

The first condition implies types in J would not play m since they would do strictly worse if they do so, even if they can persuade Player 2 that they are of a particular type. The second condition implies type θ' performs strictly better by playing m than the equilibrium provided that type θ' player can convince Player 2 that her type is not in J .

The key idea of the intuitive criterion is that if a certain action is dominated in equilibrium for some subset J of types, then Player 2 should reason that if Player 1 takes that action, then Player 1 cannot be of a type in J . We can see this is a forward induction argument – Player 2 is assuming Player 1 has acted rationally, and making inferences about Player 1's type on that basis.

Example 55 (Beer-quiche game; Cho & Kreps, 1987). In the beer-quiche game Player 1 is wimpy (θ_w) with probability $\frac{1}{10}$ and surly (θ_s) with probability $\frac{9}{10}$. Player 2 is a bully who would prefer to fight (F) the wimpy type and prefers to not fight (NF) the surly type. Player 1 orders breakfast – either beer B or quiche Q – and Player 2 then decides whether to fight.⁵⁴ The game, with payoffs, is depicted below.



In the beer-quiche game, there are no separating equilibria:

- Suppose $\sigma_1(\theta_w) = Q$ and $\sigma_1(\theta_s) = B$. The only consistent beliefs of Player 2 are then

$$\begin{aligned}\mu_2(\theta_w | Q) &= 1, & \mu_2(\theta_s | Q) &= 0, \\ \mu_2(\theta_w | B) &= 0, & \mu_2(\theta_s | B) &= 1.\end{aligned}$$

Player 2's best responses are thus $\sigma_2(B) = NF$ and $\sigma_2(Q) = F$. But then $u_1(B, \sigma_2(B), \theta_w) = 2 > 1 = u_1(Q, \sigma_2(Q), \theta_w)$, so this cannot be an equilibrium.

- Suppose $\sigma_1(\theta_w) = B$ and $\sigma_1(\theta_s) = Q$. The only consistent beliefs of Player 2 are now

$$\begin{aligned}\mu_2(\theta_w | Q) &= 0, & \mu_2(\theta_s | Q) &= 1, \\ \mu_2(\theta_w | B) &= 1, & \mu_2(\theta_s | B) &= 0.\end{aligned}$$

⁵⁴This story seems very typical of Kreps' style.

Player 2's best responses are thus $\sigma_2(B) = F$ and $\sigma_2(Q) = NF$. But then $u_1(Q, \sigma_2(Q), \theta_w) = 3 > 0 = u_1(B, \sigma_2(B), \theta_w)$, so this cannot be an equilibrium.

However, there are two classes of sequential equilibrium, both of which are pooling equilibria:

- *Pooling on quiche.* Suppose $\sigma_1(\theta_w) = \sigma_1(\theta_s) = Q$. Then any consistent beliefs of Player 2 must involve

$$\mu_2(\theta_w | Q) = \frac{1}{10} \quad \text{and} \quad \mu_2(\theta_s | Q) = \frac{9}{10}.$$

Thus Player 2 has a best response $\sigma_1(Q) = NF$, since this has expected payoff $u_2(NF, Q | \mu_2) = \frac{9}{10} > \frac{1}{10} = u_2(F, Q | \mu_2)$.

However, off-equilibrium-path beliefs are not pinned down by consistency here. We have $u_2(NF, B | \mu_2) = \mu_2(\theta_s | B)$ and $u_2(F, B | \mu_2) = \mu_2(\theta_1 | B)$. Hence if $\mu_2(\theta_w | B) \geq \mu_2(\theta_s | B)$, i.e. $\mu_2(\theta_w | B) \geq \frac{1}{2}$, then Player 2's best response is $\sigma_2(B) = F$. In this case,

$$\begin{aligned} u_1(Q, \sigma_2(Q), \theta_w) &= 3 > 0 = u_1(B, \sigma_2(B), \theta_w), \\ u_1(Q, \sigma_2(Q), \theta_s) &= 2 > 1 = u_1(B, \sigma_2(B), \theta_s), \end{aligned}$$

and thus if $\sigma = (Q, (F, NF))$, where $\sigma_2 = (\sigma_2(B), \sigma_2(F))$, and $\mu_2 = ((p, 1 - p), (1/10, 9/10))$ with $p \in [1/2, 1]$, then (σ, μ_2) is a sequential equilibrium.

- *Pooling on beer.* Suppose $\sigma_1(\theta_w) = \sigma_1(\theta_s) = B$. Now any consistent beliefs of Player 2 must involve

$$\mu_2(\theta_w | B) = \frac{1}{10} \quad \text{and} \quad \mu_2(\theta_s | B) = \frac{9}{10}.$$

Player 2's best response is then $\sigma_1(B) = F$, since this has expected payoff $u_2(NF, B | \mu_2) = \frac{9}{10} > \frac{1}{10} = u_2(F, B | \mu_2)$.

Again, off-equilibrium-path beliefs are not pinned down by consistency here, and we have $u_2(NF, Q | \mu_2) = \mu_2(\theta_s | Q)$ and $u_2(F, Q | \mu_2) = \mu_2(\theta_1 | Q)$. So if $\mu_2(\theta_w | Q) \geq \mu_2(\theta_s | Q)$, i.e. $\mu_2(\theta_w | Q) \geq \frac{1}{2}$, then Player 2's best response is $\sigma_2(Q) = F$. Thus,

$$\begin{aligned} u_1(B, \sigma_2(Q), \theta_w) &= 2 > 1 = u_1(Q, \sigma_2(Q), \theta_w), \\ u_1(B, \sigma_2(Q), \theta_s) &= 3 > 1 = u_1(Q, \sigma_2(Q), \theta_s), \end{aligned}$$

and so if $\sigma = (B, (NF, F))$ and $\mu_2 = ((1/10, 9/10), (p, 1 - p))$ with $p \in [1/2, 1]$, then (σ, μ_2) is a sequential equilibrium.

The pooling equilibrium on quiche seems unreasonable on the forward induction argument captured by the intuitive criterion. It is never a best response for type θ_w to choose B : θ_w cannot conceivably want to deviate to Q , since he receives payoff 3 from the quiche equilibrium. Type θ_s can conceivably choose B , however, if he believes Player 2 will not fight if he does so. Hence Player 2, on observing B , should conclude that Player 1's type is θ_s . Thus the pooling equilibrium on quiche fails the intuitive criterion.

Proposition 63. *In the job market signalling game with two types, the only equilibrium outcome surviving the intuitive criterion is the separating equilibrium in which $e(\theta_L) = 0$ and $e(\theta_H) = \underline{e}$.*

Proof. First, we claim any pooling equilibrium fails the intuitive criterion. Consider a pooling equilibrium with equilibrium education level e^P . We require an education level e satisfying

$$\begin{aligned}\mathbb{E}\theta - c(e^P, \theta_L) &> \theta_H - c(e, \theta_L), \\ \mathbb{E}\theta - c(e^P, \theta_H) &< \theta_H - c(e, \theta_H).\end{aligned}$$

Fix $\bar{u} = u_1(e^P, w = \mathbb{E}\theta, \theta_L)$, i.e. \bar{u} is the utility level a type- θ_L attains in the pooling equilibrium. Let \hat{e} be s.t. $u_1(\hat{e}, w = \theta_H, \theta_L) = \bar{u}$, i.e. so that a type- θ_L player is indifferent between the equilibrium and $(\hat{e}, w = \theta_H)$. Since $c(\hat{e}, \theta_H) < c(\hat{e}, \theta_L)$, a type- θ_H player will strictly prefer $(\hat{e}, w = \theta_H)$ to the equilibrium. Hence choosing some $e = \hat{e} + \epsilon$ for some sufficiently small $\epsilon > 0$ achieves our two conditions.

Next, we claim all separating equilibria with $e(\theta_H) > \underline{e}$ fail the intuitive criterion. Consider any $e' \in (\underline{e}, e(\theta_H))$. We have $\theta_L - c(0, \theta_L) < \theta_H - c(e', \theta_L)$ (c.f. Lemma 22 and the single-crossing property) but $\theta_H - c(e', \theta_H) > \theta_H - c(e(\theta_H), \theta_H)$, so the equilibrium fails the intuitive criterion.

Thus, the only equilibrium surviving is the separating equilibrium with $e(\theta_H) = \underline{e}$. \square

In the Spence model with n types $\theta_1 < \dots < \theta_n$, the *Riley outcome* is the separating equilibrium in which each type θ_k chooses the best education level for themselves assuming they will be paid according to their type, subject to the constraint that types $\theta_j < \theta_k$ do not prefer to acquire that level of education and pretend to be type θ_k . This is the most efficient of the separating equilibria.

The intuitive criterion is very effective in the Spence model with two types – it selects the Riley equilibrium. Once we enlarge the number of possible types of worker, however, the intuitive criterion no longer isolates the Riley outcome.

The intuitive criterion can be applied iteratively (imaginatively, the *iterated intuitive criterion*). Suppose we have a proposed equilibrium (s, μ) . One can apply the intuitive criterion to eliminate type-message pairs (θ, m) (messages m that type θ cannot conceivably ever want to send). Having done so, one can then eliminate actions for Player 2 that are not best responses to some belief about the type-message pairs that have not yet been eliminated. We can then carry out subsequent rounds of eliminating type-message pairs and then actions of Player 2, until we arrive at a set of type-message pairs and actions that survive this iterative procedure.

Banks & Sobel (1987) introduce an alternative to the intuitive criterion: *divinity*.

Definition 79 (Divinity). Consider a signalling game.

- (a) *Mixed best responses.* Let $\text{MBR}(T, m)$ denote the set of mixed strategy best responses of Player 2 if Player 1 has chosen m and Player 2's beliefs have support in

$T \subseteq \Theta$. That is,

$$\text{MBR}(T, m) = \bigcup_{\mu \in \Delta(T)} \arg \max_{\sigma_2 \in \Delta(A)} \sum_{\theta \in T} u_2(m, \sigma_2, \theta) \mu(\theta).$$

- (b) *Divinity criteria.* Let $T(m) \subseteq \Theta$ denote the set of types θ that might have sent message m . Fix a sequential equilibrium s^* . For each type $\theta \in \Theta$ and message $m \in M$, define

$$D_\theta(m, \theta) = \left\{ \sigma_2 \in \text{MBR}(T(m), m) \mid u_1(s^*, \theta) < \sum_{a \in A} u_1(m, a, \theta) \sigma_2(a) \right\},$$

$$D_\theta^0(m, \theta) = \left\{ \sigma_2 \in \text{MBR}(T(m), m) \mid u_1(s^*, \theta) = \sum_{a \in A} u_1(m, a, \theta) \sigma_2(a) \right\}.$$

Define the following criteria:

- (D1) Type-message pair (θ, m) is said to *fail criterion (D1)* if there exists a type θ' such that $D_\theta(m) \cup D_\theta^0(m) \subseteq D_{\theta'}(m)$.

In this case, we say that type θ' is *infinitely more likely* to choose the out-of-equilibrium message m than type θ .

- (D2) Type-message pair (θ, m) is said to *fail criterion (D2)* if $D_\theta(m) \cup D_\theta^0(m) \subseteq \bigcup_{\theta' \neq \theta} D_{\theta'}(m)$.

- (c) *Universal divinity.* We say a sequential equilibrium (s^*, μ) is a *universally divine equilibrium* if it survives iterated application of (D2).

Banks & Sobel (1987) also define a weaker notion of *divine equilibrium*, but it is difficult to map into the above framework. See their paper and Cho & Kreps (1987) for more details.

7.2 Cheap talk

Sometimes, effective costly signals may not be available to players, but players may be able to transmit costless signals. These signals are *cheap talk* – signals that have no direct effect on payoffs. The seminal paper on this is Crawford & Sobel (1982), though Green & Stokey (2007) were also considering similar ideas around the same time.⁵⁵

The structure of a cheap talk model is very similar to the signalling model:

- Stage 0: Nature chooses a type $\theta \in \Theta$ of Sender (Player 1) from a distribution p with support Θ .
- Stage 1: Sender observes θ and chooses a *message* $m \in M$ (or more generally, a signal, which can be randomly distributed).

⁵⁵A not so fun fact: Green and Stokey's paper took 29 years to be published!

- Stage 2: Receiver (Player 2) observes m and chooses $a \in A$.
- Payoffs: at the end of stage 2, payoffs $u_S(a, \theta)$ for Sender and $u_R(a, \theta)$ for Receiver are realized.
- Strategies: A pure strategy of the sender is a mapping $s_S : \Theta \rightarrow M$, and a pure strategy of the receiver is a mapping $s_R : M \rightarrow A$.

As before, the solution concept of interest here is perfect Bayesian equilibrium or sequential equilibrium.

Whether cheap talk can be effective depends on the setting. If the interests of the sender and receiver are sufficiently opposed or if different types of sender have sufficiently similar preferences to each other, then

Example 56.

- (a) *Different types of senders have same preferences over actions.* Consider a job market signalling setting with two types, but suppose instead of acquiring costly education, the worker can only make a costless announcement about their type, possibly not truthfully. Suppose we have payoffs:

| | a_L | a_M | a_H |
|--------------|-------|-------|-------|
| $\theta = L$ | 1, 1 | 2, 0 | 3, -2 |
| $\theta = H$ | 1, -2 | 2, 0 | 3, 1 |

There is no separating perfect Bayesian equilibrium in this game. Suppose otherwise, i.e. that $s_1(\theta = L) = m$ and $s_1(\theta = H) = m'$ with $m \neq m'$. Bayes' rule implies $\mu_R(L | m) = 1$ and thus $a(m) = a_L$ and $\mu_R(H | m) = 0$ so $a(m) = a_H$. Yet then Sender, if type L , has a profitable deviation by playing m' , since $u_S(a(m), L) = 1 < 3 = u_S(a(m'), L)$. Thus this is not a perfect Bayesian equilibrium.

- (b) *Different types of senders have completely opposed preferences over actions.* Consider a game where Sender wants to accept a job if he is of type A and not if he is of type B . Receiver, conversely, wants to hire (H) Sender if Sender is of type B and does not want to hire (N) if Sender is of type A . We have payoffs:

| | H | N |
|--------------|-------|------|
| $\theta = A$ | 2, -2 | 0, 0 |
| $\theta = B$ | -2, 2 | 0, 0 |

There is no separating perfect Bayesian equilibrium in this game. Suppose $s_1(A) = m$ and $s_1(B) = m'$ with $m \neq m'$. By Bayes' rule, we have that $\mu_R(A | m) = 1$ so the optimal action is $a(m) = N$, and $\mu_R(A | m') = 0$ so the optimal action is $a(m') = H$. But both types of Player 1 gain from switching messages, since

$$\begin{aligned} u_S(a(m), A) &= 0 < 2 = u_S(a(m'), A), & \text{and} \\ u_S(a(m), B) &= 0 > -2 = u_S(a(m'), B). \end{aligned}$$

Hence this cannot be a perfect Bayesian equilibrium.

At the other extreme, cheap talk is most plausibly effective in games where the sender and receiver's interests are perfectly aligned, such as coordination games. If you tell someone where you plan to meet them, they are typically better off going to that place rather than someplace else!

Example 57. Consider a version of the meeting game in Example 17 where players go to either Grand Central Station (G) or the Empire State Building (E). Suppose Player 1 is already at one location and cannot move from there – this is Player 1's type θ . Sender can text Receiver his location. Receiver reads the text and then chooses which location to go to. The payoffs are:

| | G | E |
|-----|------|------|
| G | 1, 1 | 0, 0 |
| E | 0, 0 | 1, 1 |

This game has a separating equilibrium:

$$\begin{aligned} s_1(G) &= m, & s_1(E) &= m', \\ \mu_2(\theta = G \mid m) &= 1, & \mu_2(\theta = G \mid m') &= 0, \\ a_2(m) &= G & a_2(m') &= E. \end{aligned}$$

Note that there is also a *babbling equilibrium*, where Sender announces the same message regardless of his type and Receiver does not change her beliefs based on the message, or where Sender randomizes the message independent of type.

Crawford & Sobel (1982) provide a general analysis. Rather than send a message, they have the sender send a signal $\pi : \Theta \rightarrow \Delta(X)$, where X is a signal realization space. A signal realization can be interpreted as a noisy estimate of the sender's type.

For each $i = S, R$ and for each type $\theta \in \Theta$, we assume $u_i(a, \theta) = 0$ for some action $a \in A$. We assume $\frac{\partial^2 u_i(a, \theta)}{\partial a^2} < 0$ so for each θ , there is a unique optimal action a for Sender and Receiver (not necessarily coinciding for the two players). We assume $\frac{\partial^a u_i(a, \theta)}{\partial \theta \partial a} > 0$, i.e. u_i has strictly increasing differences in (a, θ) .

7.3 Bayesian persuasion

A lot of resources (a quarter of GDP, according to McCloskey & Klammer (1995)) go into trying to persuade people to do things they might not otherwise do. Advertisements try to persuade us to buy or vote for things, prosecutors try to persuade the court to convict, lobbyists try to persuade legislators and regulators to let their clients dump toxic waste in the water supply, etc. Attempts at persuasion usually involve structuring arguments and selectively presenting information to best make a case.

When can you persuade someone else to change their action, even when they know you are trying to manipulate them? If you are rational Bayesians (and why wouldn't you be), then Gentzkow and Kamenica (2011) supply an answer.

There are two players, Sender and Receiver. Receiver has a continuous utility function $u : A \times \Omega \rightarrow \mathbb{R}$, where A is a set of actions Receiver can take, and Ω is a set of

states of the world. Given belief μ , we denote Receiver's expected utility from taking action $a \in A$ by $u(a \mid \mu) := \mathbb{E}_\mu u(a, \omega)$. Sender also has a continuous utility function $v : A \times \Omega \rightarrow \mathbb{R}$, depending on Receiver's action and the state of the world. Both Sender and Receiver have a common prior $\mu_0 \in \Delta(\Omega)^\circ$ (i.e. the interior of the set of probability distributions on Ω). Moreover, A is compact and Ω is finite.

Sender chooses a signal π , which consists of a family of distributions $\{\pi(\cdot \mid \omega)\}_{\omega \in \Omega}$ with realizations in a finite space X (chosen by Sender). Receiver observes the signal π and signal realization $x \in X$, derives a posterior μ via Bayes' rule, and takes an action $a \in a^*(\mu_x) := \arg \max_{a \in A} u(a \mid \mu_x)$. Assume $|A| \geq 2$ and for each $a \in A$, there is some μ for which $a \in a^*(\mu)$. If $|a^*(\mu_x)| > 1$, i.e. if Receiver has more than one best response, then assume Receiver chooses the best response that maximizes Sender's expected utility $v(a \mid \mu_x)$. Denote Receiver's optimal action given belief μ by $\hat{a}(\mu)$ and call $\hat{a}(\mu_0)$ the *default action*. Sender knows Receiver's behaviour and chooses a signal π to maximize Sender's expected utility. The solution concept is thus *Sender-preferred subgame perfect equilibrium*.

For a given signal π , each realization x induces a posterior belief $\mu_x \in \Delta(\Omega)$ given by Bayes' rule:

$$\mu_x(\omega) = \frac{\pi(x \mid \omega) \mu_0(\omega)}{\sum_{\omega' \in \Omega} \pi(x \mid \omega') \mu_0(\omega')} \quad \text{for all } x \in X \text{ and } \omega \in \Omega.$$

Thus each signal induces a distribution $\tau \in \Delta(\Delta(X))$ over posterior beliefs. The distribution τ is induced by signal π if $\text{supp } \tau = \{\mu_x \mid x \in X\}$ and

$$\tau(\mu) = \sum_{x: \mu_x = \mu} \sum_{\omega' \in \Omega} \pi(x \mid \omega') \mu_0(\omega') \quad \text{for all } \mu.$$

Definition 80.

- (a) *Bayes plausibility*. A distribution τ of posteriors is called *Bayes plausible* if

$$\sum_{\mu \in \text{supp}(\tau)} \mu \tau(\mu) = \mu_0,$$

that is, if the expected posterior probability under τ is equal to the prior.

- (b) *Value of a signal*. The *value* of a signal π is defined by $v^*(\pi) := \mathbb{E}_\pi v(\hat{a}(\mu_x), \omega)$, the expected value of $v(a, \omega)$ under signal π given Receiver takes the optimal action given the beliefs induced by the signal realization x of π .

If both Sender and Receiver share belief μ , we denote the expected utility for Sender, given Receiver plays optimally, by $\hat{v}(\mu) := \mathbb{E}_\mu v(\hat{a}(\mu), \omega)$. The *gain* from a signal π is $v^*(\pi) - \hat{v}(\mu_0)$. If there is some signal π for which $v^*(\pi) - \hat{v}(\mu_0) > 0$, then we say that Sender *benefits from persuasion*.

- (c) *Straightforwardness*. We call a signal π *straightforward* if $X \subseteq A$ and $\hat{a}(\mu_x) = x$ for all $x \in X$. That is, Receiver's optimal action is equal to the signal realization.

Only Bayes plausible distributions are consistent with Bayesian rationality. Conversely, it turns out that for any Bayes plausible distribution, there is some signal that induces it. Moreover, a revelation principle-type result holds – we can wlog consider only straightforward signals.

Proposition 64. *The following are equivalent:*

- (i) *There exists a signal with value $v^*(\pi) = \bar{v}$;*
- (ii) *There exists a straightforward signal with value $v^*(\pi) = \bar{v}$;*
- (iii) *There is a Bayes plausible distribution of posteriors τ such that $\mathbb{E}_\tau \hat{v}(\mu) = \bar{v}$.*

Proof. (ii) implies (i) and (iii) by definition. Let π be a signal with $v^*(\pi) = \bar{v}$. Define $X^a = \{x \mid \hat{a}(\mu_x) = a\}$ for each $a \in A$. Now consider signal π' such that $X' = A$ and $\pi'(a \mid \omega) = \sum_{x \in X^a} \pi(y \mid \omega)$. Since a is an optimal response for Receiver to each $x \in X^a$, it is also an optimal response to the realization a of π' , and so the distribution of Receiver's actions under π' and under π are the same. Thus (i) implies (ii).

Now, fix \bar{v} and a Bayes plausible distribution τ such that $\mathbb{E}_\tau \hat{v}(\mu) = \bar{v}$. Given Ω is finite, Carathéodory's theorem implies the existence of a Bayes plausible τ^* with $\text{supp } \tau^*$ finite s.t. $\mathbb{E}_{\tau^*} \hat{v}(\mu) = \bar{v}$. Taking X s.t. $\text{supp}(\tau^*) = \{\mu_x \mid x \in X\}$ and letting $\pi(x \mid \omega) = \frac{\mu_x(\omega)\tau^*(\mu_x)}{\mu_0(\omega)}$, we have that π is a signal inducing τ^* . Hence (iii) implies (i). \square

Corollary 11. *Sender benefits from persuasion iff there is a Bayes plausible distribution τ such that*

$$\mathbb{E}_\tau \hat{v}(\mu) > \hat{v}(\mu_0),$$

and the value of an optimal signal is

$$\max_{\tau \in \Delta(\Delta(X))} \mathbb{E}_\tau \hat{v}(\mu) \quad \text{subject to} \quad \sum_{\mu \in \text{supp } \tau} \mu \tau(\mu) = \mu_0.$$

Proof. Follows directly from the equivalence of (i) and (iii) in Proposition 64. \square

Thus we only need to look at the Bayes plausible distributions of posteriors to find the optimal signal for sender. By restricting to sender-preferred subgame perfect equilibria, Sender's expected utility \hat{v} is upper semicontinuous, which is enough to guarantee an optimal signal exists.

Definition 81 (Concave closure). Given Sender's expected utility \hat{v} , we define the *concave closure* of \hat{v} by

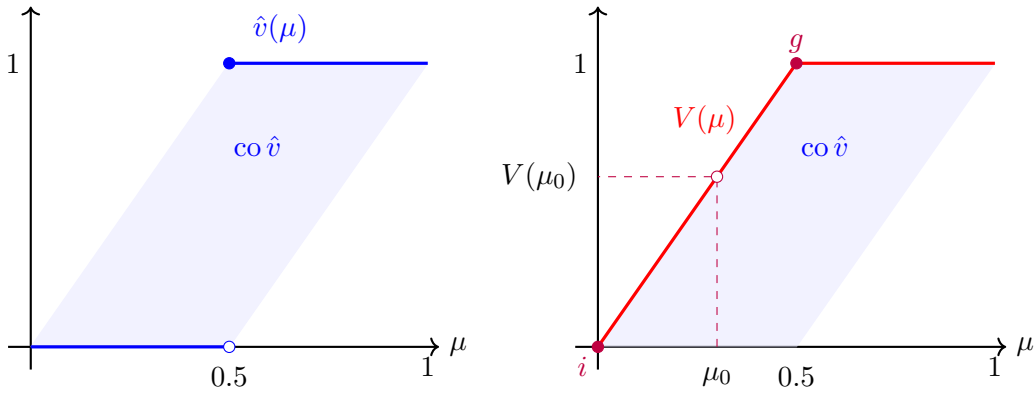
$$V(\mu) := \sup\{z \mid (\mu, z) \in \text{co}(\hat{v})\},$$

where $\text{co}(\hat{v})$ is the convex hull of the graph of \hat{v} .

Corollary 12. *The value of an optimal signal is $V(\mu_0)$.*

Proof. If $(\mu', z) \in \text{co } \hat{v}$ then there is a distribution of posteriors τ s.t. $\mathbb{E}_\tau \mu = \mu'$ and $\mathbb{E}_\tau \hat{v}(\mu) = z$. By Proposition 64, it follows that $\text{co } \hat{v}$ is the set of (μ_0, z) s.t. there exists a signal with value z if the prior is μ_0 . Thus $V(\mu)$ is the maximal payoff that Sender can achieve with any signal under the prior μ_0 . \square

Example 58 (Judge and prosecutor). The motivating example in Kamenica & Gentzkow (2009) involves a judge (Receiver) who chooses whether to acquit or convict a defendant. There are two states – either the defendant is innocent (I) or guilty (G) – and the judge receives utility 1 from choosing the state-appropriate action and 0 otherwise (i.e. the judge wants to convict if guilty and acquit if innocent). Sender is a prosecutor who receives payoff 1 if the judge convicts and 0 if the judge acquits, regardless of state. Both the judge and prosecutor have a prior $\mu_0(G) = 0.3$.



In the figure, μ is the probability that the defendant is guilty. The expected utility for the prosecutor, \hat{v} , takes value 0 for $\mu < \frac{1}{2}$ and 1 if $\mu \geq \frac{1}{2}$, since the judge will convict iff her posterior belief that the defendant is guilty (weakly) exceeds $\frac{1}{2}$. The optimal signal for the prosecutor to send involves two signal realizations, i or g . If the judge observes i , her posterior belief is $\mu = 0$ giving $\hat{v}(\mu) = 0$ and if she observes g , her posterior belief is $\mu = \frac{1}{2}$, giving $\hat{v}(\mu) = 1$. The optimal signal induces distribution $\tau = (0.4, 0.6)$ over $(0, 0.5)$, which is Bayes plausible because $\mu_0 = 0.3 = 0 \times 0.4 + 0.6 \times 0.5$. The concave closure of \hat{v} is given in the figure, and the value of the optimal signal is $V(\mu_0) > \hat{v}(\mu_0)$, so the prosecutor benefits from persuasion. Indeed, the prosecutor is more likely than not to persuade the judge to convict, even though both start with a prior that the defendant is innocent!

When exactly will the Sender benefit from persuasion? Clearly, if \hat{v} is concave, then Corollary 12 tells us that the Sender will not benefit from persuasion, because the concave closure of concave \hat{v} is simply \hat{v} . Conversely, if \hat{v} is strictly convex, then $V(\mu)$ always lies above $\hat{v}(\mu)$ in the interior of the space of priors, so the Sender always benefits from persuasion as long as the prior is non-degenerate. However, in Example 58, \hat{v} is not concave or strictly convex, and the prosecutor benefits from persuasion if $\mu_0 < 0.5$ but does not benefit from persuasion if $\mu_0 \geq 0.5$. Under the prior that $\mu_0 = 0.3$, the prosecutor benefits from persuasion because the judge's default action (to acquit) is not

the prosecutor's preferred action (to convict) and because the judge's action is constant in a neighbourhood around the prior.

We say there is *information Sender would share* if there is some prior μ such that $\hat{v}(\mu) > \mathbb{E}_\mu v(\hat{a}(\mu), \omega)$. That is, if Sender has private information that causes her to believe μ , she would prefer to share this information with Receiver. We say Receiver's *preference is discrete* at μ if there is an $\epsilon > 0$ such that for every action $a \neq \hat{a}(\mu)$, we have $\mathbb{E}_\mu u(\hat{a}(\mu), \omega) > \mathbb{E}_\mu u(a, \omega) + \epsilon$. That is, Receiver's expected payoff from her preferred action is bounded away from her expected payoff from any other action. For a finite action space A , this holds if Receiver is not indifferent between some pair of actions. Thus Receiver's preference is almost always discrete at the prior if A is finite.

Proposition 65. *If there is no information that Sender would share, Sender does not benefit from persuasion. If there is information that Sender would share and, under prior μ_0 , Receiver's preference is discrete at μ_0 , then Sender benefits from persuasion.*

Proof. Suppose there is no information Sender would share. Then $\hat{v}(\mu) \leq \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu(\omega)$ for all μ . Given a signal π induces some posterior τ , the value of signal π is $\sum_{x \in X} \tau_x \hat{v}(\mu_x) \leq \sum_{x \in X} \tau_x \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu_x(\omega) = \hat{v}(\mu_0)$, and so Sender does not benefit from persuasion.

Suppose there is information Sender would share. Since u is continuous in ω , it follows that $\sum_{\omega \in \Omega} u(\hat{a}(\mu_0), \omega) \mu(\omega)$ is continuous in μ . Supposing Receiver's preference is discrete at μ_0 , there is some $\delta > 0$ s.t. for all μ in a δ -ball B_δ about μ_0 , $\hat{a}(\mu) = \hat{a}(\mu_0)$. Since there is information Sender would share, there is some μ_h s.t. $\hat{v}(\mu_h) > \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu_h(\omega)$. Fix any ray R through μ_h and μ_0 . Since μ_0 does not lie on the boundary of $\Delta(\Omega)$, there is some $\mu_\ell \in R$ s.t. $\mu_\ell \in B_\delta$ and $\mu_0 = \lambda \mu_\ell + (1 - \lambda) \mu_h$ for some $\lambda \in (0, 1)$. Consider Bayes-plausible posterior distribution $\tau(\mu_\ell) = \lambda$ and $\tau(\mu_h) = 1 - \lambda$. Given $\hat{a}(\mu_0) = \hat{a}(\mu_\ell)$, we have $\hat{v}(\mu_\ell) = \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu_\ell(\omega)$. Thus $\mathbb{E}_\tau \hat{v}(\mu) = \lambda \hat{v}(\mu_\ell) + (1 - \lambda) \hat{v}(\mu_h) > \lambda \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu_\ell(\omega) + (1 - \lambda) \sum_{\omega \in \Omega} v(\hat{a}(\mu_0), \omega) \mu_h(\omega) = \hat{v}(\mu_0)$, so Sender benefits from persuasion.

Finally, we claim if A is finite than Receiver's preference at a belief μ is generically discrete. If Receiver's preference is not discrete at μ , then there must be some action $a \neq \hat{a}(\mu)$ s.t. $\sum_\omega u(\hat{a}(\mu), \omega) \mu(\omega) = \sum_\omega u(a, \omega) \mu(\omega)$. To see this, note that if there is no such action, then defining

$$\epsilon := \frac{1}{2} \min_{a \neq \hat{a}(\mu)} \left\{ \sum_\omega u(\hat{a}(\mu), \omega) \mu(\omega) - \sum_\omega u(a, \omega) \mu(\omega) \right\},$$

(which is well-defined since A is finite), we have that $\sum_\omega u(\hat{a}(\mu), \omega) \mu(\omega) > \sum_\omega u(a, \omega) \mu(\omega) + \epsilon$ for all $a \neq \hat{a}(\mu)$, but this implies Receiver's preference is discrete at μ .

Now given there are only finitely many pairs a, a' and the set of all such pairs has measure zero, we need only show $\{\mu \mid \sum_\omega u(a, \omega) \mu(\omega) = \sum_\omega u(a', \omega) \mu(\omega)\}$ has measure zero. Fix any distinct a, a' . Index states by i . Let $\beta = (\beta_1, \dots, \beta_n)$ and $\mu = (\mu(\omega_1), \dots, \mu(\omega_n))$. We need only show $\{\mu \mid \beta' \mu = 0\}$ is measure zero. Since for any action a , there is some μ s.t. $a^*(\mu) = \{a\}$, there is some ω s.t. $u(a, \omega) \neq u(a', \omega)$. Hence there is at least one $\beta_i \neq 0$, so β is a linear transformation of rank 1, and the kernel of β is a vector space of dimension $n - 1$. The claim that $\{\mu \mid \beta' \mu = 0\}$ is measure zero follows. \square

8 Repeated games

A repeated game is a (typically normal form) game that is played repeatedly among the same set of players, with discounted payoffs aggregated over time. For the moment, we consider only repeated games with *perfect monitoring*.

Repeated games are very well-studied, and they matter because there are many real-world strategic interactions where agents interact repeatedly. A classic application is industrial organization, where repeated games give us a guide to how collusion can arise between firms, sometimes even in the absence of an agreement between them (this is *tacit collusion*, and is very difficult for competition regulators to deal with, given the absence of evidence of collusion). Repeated interactions obviously extend much further than industry however – people interact repeatedly all the time, as partners, friends, family members, work colleagues, neighbours, and so on. Some (naïve) theories of altruism use the analysis of repeated games to explain why humans are kind to each other and tend to cooperate, even when it might be in their immediate interests not to do so. I say these theories are naïve because there is strong empirical evidence of *strong reciprocity* – altruism when there is no benefit in doing so, such as under anonymous conditions or directed towards strangers with whom you almost certainly won't interact again.⁵⁶

Definition 82 (Repeated games).

- (a) *Stage game*. In the context of a repeated game, a *stage game* G is a normal form game $G = (\mathcal{I}, (A_i, u_i)_{i \in \mathcal{I}})$, or less commonly a finite extensive form game Γ . We will typically assume a normal form game. Moreover, we assume the payoffs of the stage game are bounded.
- (b) *Perfect monitoring*. We say there is *perfect monitoring* if in every period t , every player perfectly recall the full history of actions by all players in all periods $0 \leq t' < t$.
- (c) *History*. In a repeated game, a *history* at time t is a profile $h^t = (a^0, a^1, \dots, a^{t-1})$, where each $a^\tau \in A_1 \times \dots \times A_n$ is the profile of actions played in period τ . The set of histories at time t is defined as

$$\mathcal{H}^t = \{(a^0, \dots, a^{t-1}) \mid a^\tau \in A_1 \times \dots \times A_n \text{ for all } \tau = 0, \dots, t-1\}.$$

- (d) *Repeated game*. Fix $T \in \mathbb{N} \cup \{+\infty\}$. A $(T+1)$ -period repeated game G^T , with corresponding to stage game $G = (\mathcal{I}, (A_i, u_i)_{i \in \mathcal{I}})$, is a game $(\mathcal{I}, \mathcal{H}, (\Sigma_i, \succsim_i)_{i \in \mathcal{I}})$, where
 - (i) $\mathcal{H} = \bigcup_{t=0}^{T+1} \mathcal{H}^t$ is the set of all histories of the game. We call \mathcal{H}^{T+1} the set of *terminal histories* of the game;

⁵⁶Bowles, Gintis, Boyd and others, develop an interesting game theoretic approach to understanding why humans are altruistic, via evolutionary models. A good summary is Bowles & Gintis (2011), *A Cooperative Species: Human Reciprocity and its Evolution*.

- (ii) $\Sigma_i = \{\sigma_i \mid \sigma_i : \mathcal{H} \rightarrow \Delta(A_i)\}$, that is, Σ_i is the set of all mappings from the set of histories into i 's stage game mixed strategy set. A strategy $\sigma_i \in \Sigma_i$ of the repeated game specifies for each t and each history $h^t \in \mathcal{H}^t$ a strategy $\sigma_i^t(h^t)$ in the t th stage game G ;
- (iii) \succsim_i is a preference relation over $\Sigma = \times_{j \in \mathcal{I}} \Sigma_j$, the set of all strategy profiles σ .

If T is finite, we call the game a *finitely repeated game*. If T is infinite, then we call the game an *infinitely repeated game*.

For convenience, we will denote by $\sigma(t)$ the action profile played at time t on the path induced by the strategy profile σ .

- (e) *Preference relations in repeated games.* There are a number of different ways to conceptualize the preference relations \succsim_i in repeated games.

- (i) *Discounted payoff criterion.* Fix some $\delta \in (0, 1)$ and $T \in \{1, \dots, \infty\}$. Define $U_i(\sigma) = \sum_{t=0}^T \delta^t u_i(\sigma(t))$. We say \succsim_i satisfies the *discounted payoff criterion* with discount rate δ if $\sigma \succsim_i \sigma'$ iff $U_i(\sigma) \geq U_i(\sigma')$. If $T < \infty$ then we can also take $\delta = 1$.
- (ii) *Limit-of-means criterion.* Define $U_i(\sigma) = \liminf_{T \rightarrow \infty} \left[\frac{1}{T} \sum_{t=0}^T u_i(\sigma(t)) \right]$. Then \succsim_i satisfies the *limit-of-means criterion* if $\sigma \succsim_i \sigma'$ iff $U_i(\sigma) \geq U_i(\sigma')$.
- (iii) *Overtaking criterion.* We say \succsim_i satisfies the *overtaking criterion* if $\sigma \succsim_i \sigma'$ iff

$$\liminf_{T \rightarrow \infty} \sum_{t=0}^T (u_i(\sigma(t)) - u_i(\sigma'(t))) \geq 0.$$

This is often written in terms of the induced strict preference relation \succ_i , for which we have $\sigma \succ_i \sigma'$ iff

$$\liminf_{T \rightarrow \infty} \sum_{t=0}^T (u_i(\sigma(t)) - u_i(\sigma'(t))) > 0.$$

Given a stage game G , we refer to the repeated game $G^T(\delta)$ or $G^\infty(\delta)$ as a *supergame*.

Note we defined repeated games so that the stage game is always a mixed strategy game. We assume that A_i is finite throughout. Note that we have defined the strategies σ_i of the repeated game to be behavioural strategies. By Theorem 1, this is wlog – these strategies are each equivalent to some mixed strategy and vice versa.

It is worth discussing the payoff criteria in more detail. We are usually interested in repeated games where the preferences \succsim_i for each player can be characterized by one of the payoff criteria given in Definition 82, though this is obviously quite restrictive (we require that payoffs are additively separable, and so on.) In finitely repeated games, we assume the discounted payoff criterion applies, either with some discount rate $\delta < 1$ or without discounting (i.e. $\delta = 1$). We call infinitely repeated games for which the

discounted payoff criterion applies with $\delta \in (0, 1)$ *discounted infinitely repeated games*. If u_i is bounded, $\sum_{t=0}^{\infty} \delta^t u_i(\sigma(t))$ converges and so discounted payoff criterion preferences are well-defined. It is conventional to normalize $U_i(\sigma) = \sum_{t=0}^{\infty} \delta^t u_i(\sigma(t))$ in discounted infinitely repeated games by defining $U_i(\sigma) = (1 - \delta) \sum_{t=0}^{\infty} \delta^t u_i(\sigma(t))$. This scales the supgame payoffs so that they are comparable with the stage game payoffs.

If $\delta = 1$, we cannot guarantee that $\sum_{t=0}^{\infty} u_i(\sigma(t))$ converges.⁵⁷ In this case, we typically rely on either the limit-of-means criterion or overtaking criterion, and the games in which these apply are called *undiscounted infinitely repeated games*. We defer discussion of these two criteria later.

We say that the preferences \succsim_i admit payoff representation $U_i : \Sigma \rightarrow \mathbb{R}$ if $\sigma \succsim_i \sigma'$ iff $U_i(\sigma) \geq U_i(\sigma')$. When preferences admit a payoff representation U_i and we are concerned with subgame perfect equilibria or their refinements, it is often very convenient to work with continuation payoffs:

Definition 83 (Continuation payoff). In a repeated game for which preferences \succsim_i for each player i admit a payoff representation U_i , we define the *continuation payoff* $U_i(\sigma \mid h^t)$ for i of strategy profile σ at history h^t as the payoff to i on the path induced by σ in the subgame induced by history h^t .

In English, the history h^t induces a subgame, a new repeated game with initial history h^t . The continuation payoff is simply the payoff if the strategies prescribed by the strategy profile σ are followed by all players in this subgame.

Example 59. In a finitely repeated game under the discounted payoff criterion, the (normalized) continuation payoff for player i given history h^t is given by

$$U_i(\sigma \mid h^t) = \mathbb{E}_{\sigma} \left[\frac{1}{T-t} \sum_{\tau=t}^T \delta^{\tau-t} u_i(\sigma^{\tau}(h^{\tau})) \mid h^t \right].$$

In a discounted infinitely repeated game under the discounted payoff criterion, the (normalized) continuation payoff for player i given history h^t is given by

$$U_i(\sigma \mid h^t) = \mathbb{E}_{\sigma} \left[(1 - \delta) \sum_{\tau=t}^{\infty} \delta^{\tau-t} u_i(\sigma^{\tau}(h^{\tau})) \mid h^t \right].$$

In the definitions of repeated games, we assumed a common discount rate δ . This can be relaxed by replacing δ with a vector of individual discount rates $(\delta_i)_{i \in \mathcal{I}}$.

As with extensive form sequential games, there are generally many Nash equilibria in both finitely and infinitely repeated games, many of which are not sequentially rational. Indeed, in infinitely repeated games we can show there are infinitely many Nash equilibria. Subgame perfection is usually the suitable refinement, though even then, the number of subgame perfect equilibria is potentially very large.

⁵⁷For example, suppose $u_i(\sigma(t)) = -1$ if t is odd and $u_i(\sigma(t)) = 1$ if t is even.

8.1 Finitely repeated games

We will always assume preferences in a finitely repeated game satisfy the discounted payoff criterion. As with finite extensive form games more generally, finitely repeated games benefit from being solvable by backward induction.

First, note that even in finitely repeated games, there may be many subgame perfect equilibria.

Example 60. Consider the following stage game G :

| | A_2 | B_2 | C_2 |
|-------|--------------|---------------------|---------------------|
| A_1 | 4, 4 | 0, 0 | 0, 5 |
| B_1 | 0, 0 | 1 , 1 | 0, 0 |
| C_2 | 5 , 0 | 0, 0 | 3 , 3 |

This game has three Nash equilibria – two pure strategy equilibria (B_1, B_2) and (C_1, C_2) and a mixed strategy equilibrium $((0, \frac{3}{4}, \frac{1}{4}), (0, \frac{3}{4}, \frac{1}{4}))$.⁵⁸

Now suppose the stage game is repeated twice and that both players discount the future at discount rate δ . Player i 's payoff given strategy profile σ is therefore

$$U_i(\sigma) = u_i(\sigma_1(1), \sigma_2(1)) + \delta u_i(\sigma_1(2), \sigma_2(2)).$$

There are a large number of subgame perfect equilibria here. First, any σ^* such that $\sigma^*(h)$ is one of the three Nash equilibria for each history h is a subgame perfect equilibrium in this game. Furthermore, by making the choice of which Nash equilibrium to play in the second period contingent on first period play, we can construct a subgame perfect equilibrium where play in the first period differs from Nash play. Consider the following strategy profile:

1. Play (A, A) in the first period;
2. If first period play is (A, A) , play (C, C) in the second period. Else play (B, B) .

To confirm this is a subgame perfect equilibrium, consider one-shot deviations. In the second period, (C, C) and (B, B) are Nash equilibria so there is no gain to deviating in either case. In period 1, if player i follows the strategy they receive $4 + 3\delta$. The best deviation in period 1 is C , which yields 5, but this is followed in period 2 by (B, B) , and thus the payoff to deviating is $5 + \delta$. We thus have that the strategy above is a subgame perfect equilibrium provided that $4 + 3\delta \geq 5 + \delta$, i.e. $\delta \geq \frac{1}{2}$.

Example 61 (Repeated prisoners' dilemma). In the previous example, there were multiple Nash equilibria in the stage game. If there is only a single Nash equilibrium in the

⁵⁸ A_1 and A_2 are never-best responses so we can exclude mixing over them. Suppose under σ_2 , Player 2 plays B_2 with probability p and C_2 with probability $1 - p$. Then for Player 1 to mix, we need $u_1(B_1, \sigma_2) = p = 3(1 - p) = u_1(C_1, \sigma_2)$, implying $p = \frac{3}{4}$. The game is symmetric so Player 1 mixes with the same probability in equilibrium.

stage game, then there is a unique subgame perfect equilibrium in the finitely repeated game, which involves that Nash equilibrium being played in every period.

Suppose the stage game G is the prisoner's dilemma:

| | | |
|-------|-------|-------|
| | C_2 | D_2 |
| C_1 | 1, 1 | -1, 2 |
| D_1 | 2, -1 | 0, 0 |

Consider the repeated game $G^T(\delta)$, i.e. the prisoner's dilemma repeated T times with common discount rate δ . We proceed by backward induction. In the final period T , a Nash equilibrium must be played since both players choose optimal strategies, and the only Nash equilibrium is (D_1, D_2) . In period $T - 1$, both players know that in period T , (D_1, D_2) will be played. Hence D_i is the strictly dominant strategy for each player i in period $T - 1$ and so (D_1, D_2) is played in equilibrium. Recursively, we have that D_i is the strictly dominant strategy for i in every period $t = 0, 1, \dots, T$ and thus (D_1, D_2) is played in each period t in equilibrium. This is therefore the unique subgame perfect equilibrium.

Since this holds for any $T \in \mathbb{N}$, we have that as $T \rightarrow \infty$, the unique subgame perfect equilibrium is (D_1, D_2) every period. This is a subgame perfect equilibrium of the limit game $G^\infty(\delta)$, but importantly, it is not unique.⁵⁹

As we might expect, this result holds generally:

Proposition 66. *Suppose G is an n -player stage game possessing a unique Nash equilibrium α^* . Then the strategy profile $\sigma^* = (\sigma_1^*, \dots, \sigma_n^*)$ where $\sigma_i^*(h) = \alpha_i^*$ for each player $i \in \mathcal{I}$ is the unique subgame perfect equilibrium of the $(T + 1)$ -period repeated game $G^T(\delta)$.*

Proof. The proof is almost identical to the example immediately above, applying backward induction. Since α^* is a unique Nash equilibrium, it must be played in the T th subgame under any subgame perfect equilibrium. Suppose σ^* is played in all periods $k + 1, \dots, T$. Applying the backward induction algorithm in the k th game, we can substitute $u_i(\alpha_i, \alpha_{-i}) + \sum_{j=1}^T \delta^j u_i(\sigma^*)$ for $u_i(\alpha_i, \alpha_{-i})$. Since α^* is the unique Nash equilibrium,

$$u_i(\alpha_i^*, \alpha_{-i}^*) \geq u_i(s_i, \alpha_{-i}^*)$$

for all $s_i \in S_i$ and all i , and so

$$u_i(\alpha_i^*, \alpha_{-i}^*) + \sum_{j=1}^T \delta^j u_i(\alpha^*) \geq u_i(s_i, \alpha_{-i}^*) + \sum_{j=1}^T \delta^j u_i(\alpha^*),$$

for all $s_i \in S_i$ and all i so α^* is optimal in period k , regardless of the history of play. Furthermore, no other strategy profile α is optimal in period k since α^* is the unique Nash equilibrium of the stage game. Proof follows by induction. \square

⁵⁹This comes down to the fact that the subgame perfect equilibrium correspondence is not lower hemicontinuous.

In practice, the stark prediction that if a stage game has a dominant strategy equilibrium, the strictly dominant action will be played at every stage in a finitely repeated game does not hold – evidence from numerous lab experiments suggests that players tend to cooperate in the early stages of play and only begin to defect as the final period comes close. Kreps, Milgrom, Roberts & Wilson (1982) explain such behaviour via a game of incomplete information:

Example 62 (Finitely repeated prisoner’s dilemma with type uncertainty). Consider the repeated prisoner’s dilemma as in Example 61, with slightly more general payoffs:

| | | |
|-------|--------|--------|
| | C_2 | D_2 |
| C_1 | 1, 1 | b, a |
| D_1 | a, b | 0, 0 |

where $a > 1$, $b < 0$ and $a + b < 2$, and the game is repeated for T periods. Suppose payoffs are not discounted. As before, the unique subgame perfect equilibrium path has both players play D_i in every period.

Suppose there is some small positive probability $p > 0$ that Player 2 (column) is not a rational type and instead has only the tit-for-tat strategy available to him, i.e. only the strategy of playing C_2 in the first period and mimicking Player 1’s previous action thereafter. With probability $1 - p$, Player 2 is a rational type that has C_2, D_2 available each period and maximizes his payoff. Kreps et al prove that in any sequential equilibrium of this game, there is an upper bound on the number of rounds that either of the players i play D_i .

Suppose before some period t , it becomes common knowledge that Player 1 is rational. Then in any sequential equilibrium, both players i play D_i in periods $t, t + 1, \dots, T$, by a standard backward induction argument. Now suppose Player 2 plays D_2 in period t . Then Player 1 plays D_1 in period $t + 1$ surely – tit-for-tat calls for the irrational type to play D_1 , and so if Player 1 plays C_1 , it becomes common knowledge that he is a rational type, and so the continuation payoff to Player 1 would be 0, whereas the continuation payoff from playing D_1 , and D_1 strictly dominates C_1 in the stage game.

Consider any period t history h^t for which Player 2 believes that Player 1 is the irrational type with probability q . First, if Player 1 played C_1 in the previous period under h^t , the expected payoff to Player 1 for the remainder of the game must be at least $q(T - t) + b$. For suppose she instead plays C_1 until Player 2 next plays D_2 , and playing D_1 thereafter. Against the irrational type, such a strategy yields payoff $T - t$ and against the rational type, this yields payoff no worse than b , so the expected payoff for Player 1 from such a strategy is $q(T - t) + (1 - q)b \geq q(T - t) + b$.

Second, if Player 1 played D_2 in the previous period under h^t , the expected payoff for Player 1 for the remainder of the game must be at least her minimum expected payoff from playing C_1 which is $q(T - t - 1) + 2b$, since Player 2 plays D_2 for sure in the next period and thus Player 1’s belief next period is q .

Third, if rational Player 2’s expected payoff under h^t is weakly greater than $q(T - t - 1) + 3b - a$, then Player 1 does no worse if Player 2 plays tit-for-tat than if rational Player 2

plays his equilibrium strategy, by induction. Hence by playing tit-for-tat, rational Player 2's payoff is within $b - a$ of Player 1's payoff.

Now if $t \leq T - \frac{2a-4b+2q}{q}$, then Player 2 plays tit-for-tat surely, and so Player 1's belief remains q in each such period. We need only consider the case where Player 1 previously cooperated at every period. Then if Player 2 plays D_2 , it becomes common knowledge he is rational and so the total payoff from D_2 is a . By playing C_2 , Player 2 receives at least b immediately and continuation payoff at least $q(T - t - 2) + 3b - a$, so C_2 is strictly better than D_2 for the rational type if $\frac{2a-4b+2q}{q} \leq T - t$.

Since Player 2 plays tit-for-tat for all but the final $\lfloor \frac{2a-4b+2p}{p} \rfloor$ periods, then if Player 1 plays C_1 until Player 2 plays D_2 then Player 1 receives at least $T - \frac{2a-4b+2p}{p}$, whereas if Player 1 plays D_1 at some earlier date t and returns to playing C_1 only m periods later, she receives a in period t , b in period $t + m$ and 0 in the interim, sacrificing $1 + m - a + b$ (or in per period terms, $1 + (1 - a - b)/m$) compared to if playing C_1 each period. Hence each time Player 1 plays D_1 in the tit-for-tat phase, it costs her at least $\min\{2 - a - b, 1\}$, so if Player 1 plays D_1 k times before the final $\lfloor \frac{2a-4b+2p}{p} \rfloor$ periods, her payoff is less than $T - k \min\{2 - a - b, 1\}$. Combining bounds gives us $k \leq \frac{2a-4b+2p}{p \min\{2-a-b, 1\}}$, and since playing D_1 causes Player 2 to play D_2 under tit-for-tat, there are at most $2k$ periods involving playing D_i . Thus in any sequential equilibrium, the total number of periods involving some player playing D_i is bounded from above by $\frac{2a-4b+2p}{p} \left(1 + \frac{2}{\min\{2-a-b, 1\}}\right)$.

Restricting attention to the Pareto-undominated sequential equilibria, on the equilibrium path of such equilibria neither player i plays D_i in all but the final $\lfloor 1 + \frac{2a-4b+2p}{p} \rfloor$ periods. This follows since if in sequential equilibrium some player plays D_i in an earlier period, then there is a Pareto dominating sequential equilibrium that involves playing (C_1, C_2) on the equilibrium path in that period instead.

8.2 Discounted infinitely repeated games

In discounted repeated games, the discounted payoff criterion always applies. Note that in a discounted infinitely repeated game, δ has two possible interpretations:

- *Time preference.* As in a finitely repeated game, δ can represent time preference, measuring how patient players are. If there is an interest rate r , for example, then the discount rate is $\delta = \frac{1}{1+r}$.
- *Uncertainty over end date.* An alternative interpretation is that δ is the probability that the interaction will be repeated in the subsequent period. Here, the game ends almost surely in finite time, but the end date is stochastic. $1 - \delta$ could represent the probability that an agent dies before the next period or that a firm becomes exogenously bankrupt before the next period, for example.

The one-shot deviation principle is always applicable to discounted repeated games:

Proposition 67. *Consider any discounted infinitely repeated game $G^\infty(\delta)$, with corresponding stage game G and $\delta \in (0, 1)$. Then the one-shot deviation principle applies, that is, any unimprovable strategy is optimal.*

Proof. By Theorem 24, we need only show that the consistency and continuity assumptions (A1) and (A2) introduced in section 6.3 hold. Fix player i and strategy profile σ_{-i} of i 's opponents. Define the return function $\pi(\sigma_i, h) = u_i(\sigma_i, \sigma_{-i} \mid h)$, where $u_i(\sigma \mid h)$ denotes the continuation payoff under σ given history h . For any terminal history y for which there is some $\sigma_i \in \Sigma_i$ such that y is the outcome path of (σ_i, σ_{-i}) , let σ_i^y denote any of the strategies σ_i that ensures y is the equilibrium path. Let y_h denote the subpath of y initiating at history h , and let $\sigma_i^{y|h}$ denote any of the strategies σ_i that ensure y_h is the equilibrium path in the subgame induced by history h . Then we can define $\pi(y) = \pi(\sigma_i^y)$ and $\pi(y_h) = \pi(\sigma_i^{y|h}, h)$.

Now, suppose y, y' are terminal histories diverging at a history h and with $\pi(y) \geq \pi(y')$. Suppose divergence occurs at time t . Then since the first t terms in the sum of discounted payoffs are the same for both y and y' , it must be that $\pi(y_h) \geq \pi(y'_h)$. Hence (A1) is satisfied. Next, fix any terminal history y (wlog) and $\epsilon > 0$. Let $B = \max\{u_i(s) - u_i(s') \mid s, s' \in S\}$, i.e. B is the maximal difference in payoffs in the stage game across pairs of stage game action profiles. Now for any terminal history y' that is identical to y for the first T periods, for sufficiently large T we have $|\pi(y) - \pi(y')| \leq \sum_{t=T}^{\infty} \delta^t B = \frac{B}{1-\delta} \delta^T < \epsilon$, so (A2) is satisfied. \square

Example 61 (continued). Continuing the repeated prisoners' dilemma example, suppose now that the game is repeated for infinite periods with discounting. Recall we have stage game G with payoff matrix:

| | | |
|-------|-------|-------|
| | C_2 | D_2 |
| C_1 | 1, 1 | -1, 2 |
| D_1 | 2, -1 | 0, 0 |

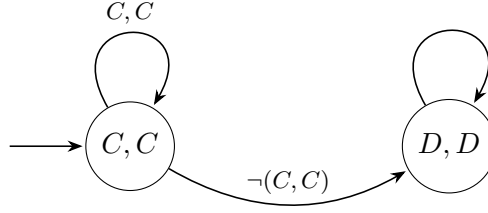
As in the finite case, playing (D_1, D_2) in every period constitutes a subgame perfect equilibrium. However, if $\delta \geq \frac{1}{2}$, there is a subgame perfect equilibrium in which (C_1, C_2) is played every period.

Consider the following symmetric strategies, known as *grim trigger* strategies, since they involve playing the minimax strategies if any player deviates:

1. Play C_i in the first period and every subsequent period provided no player j ever plays D_j ;
2. In each period, if any player j has ever played D_j , play D_i .

These strategies involve *Nash reversion* – the punishment phase involves playing a stage game Nash equilibrium. In the prisoner's dilemma, the Nash equilibrium happens to be the minimax

It can be convenient to visualize these strategies by means of an automaton:



Call this strategy s .

To see this is a subgame perfect equilibrium, consider any period t and suppose that D_j has already been played by some player. Then there is one one-shot deviation for player i (play C_i) which yields payoff

$$\pi(s_{C_i}, t) = -1 > 0 = \pi(s, t),$$

so the only one-shot deviation is not profitable. Now suppose neither player j has played D_j . Then the one shot deviation for player i is to play D_i , which gives return $\pi(s_{D_i}, t) = 2(1 - \delta)$, against $\pi(s, t) = (1 - \delta) \sum_{k=0}^{\infty} \delta^k = 1$. Rearranging, we thus have that $\pi(s, t) \geq \pi(s_{D_i}, t)$ provided $\delta \geq \frac{1}{2}$.

8.3 Undiscounted infinitely repeated games

In undiscounted infinitely repeated games, we typically apply either the limit-of-means criterion or the overtaking criterion. These are not generally equivalent. For example, suppose we have two strategy profiles σ and σ' such that the stage game payoffs for player i are $(u_i(\sigma(t)))_{t=0}^{\infty} = (0, 0, 0, 0, \dots)$ and $(u_i(\sigma'(t)))_{t=0}^{\infty} = (-1, 3, 0, 0, \dots)$. Under the limit-of-means criterion, player i is indifferent between σ and σ' , since

$$\liminf_{T \rightarrow \infty} \frac{1}{T+1} \sum_{t=0}^T u_i(\sigma(t)) = 0 = \liminf_{T \rightarrow \infty} \frac{1}{T+1} \sum_{t=0}^T u_i(\sigma'(t)).$$

However, under the overtaking criterion, player i strictly prefers σ' to σ , since

$$\liminf_{T \rightarrow \infty} \sum_{t=0}^{\infty} (u_i(\sigma'(t)) - u_i(\sigma(t))) = 2 > 0.$$

This points to a philosophical difference. Since players do not discount (i.e. they are perfectly patient), they do not care about the timing of payoffs. However, the notion of patience that is captured by the limit-of-means criterion implies that players also do not care about their payoffs in any finite number of periods – any temporary increase or decrease in payoffs relative to the long-run average evaporates when we take the limit of the mean. A reasonable complaint is that this is unintuitive – we might reckon that perfectly patient players will prefer payoff stream $(1, 0, 0, \dots)$ to $(0, 0, 0, \dots)$. The overtaking criterion captures this intuition. However, it does so at a cost – we cannot represent preferences \succsim_i that satisfy the overtaking criterion by means of a utility functional U_i .

Example 61 (continued). Once again consider the infinitely repeated prisoners' dilemma, except now without discounting.

Under both the limit-of-means criterion and overtaking criterion, the grim trigger strategy profile $\sigma = (\sigma_1, \sigma_2)$ we constructed in section 8.2 is a subgame perfect equilibrium. Fix player i .

First, consider the limit-of-means criterion. Suppose i chooses any other strategy σ'_i that first differs on the path induced by (σ_i, σ_{-i}) in period τ (i.e. prescribes D_i on path in τ). Then

$$U_i(\sigma'_i, \sigma_{-i}) \leq \liminf_{T \rightarrow \infty} \frac{1}{T+1}(\tau+2) = 0 < 1 = U_i(\sigma_i, \sigma_{-i}),$$

so σ'_i cannot be optimal. Now suppose i chooses any strategy σ'_i that agrees with σ_i on path but differs from σ_i after some history h for which σ prescribes the punishment phase. Then the continuation payoffs satisfy $U_i(\sigma'_i, \sigma_{-i} \mid h) \leq 0 = U_i(\sigma_i, \sigma_{-i} \mid h)$, and so σ_i is a best response in any subgame. Hence σ is a subgame perfect equilibrium.

Now consider the overtaking criterion. Suppose i chooses any strategy $\sigma'_i \neq \sigma_i$ that first prescribes D_i on the path induced by σ in period τ . Then on path,

$$\liminf_{T \rightarrow \infty} \sum_{t=0}^T (u_i((\sigma_i, \sigma_{-i})(t)) - u_i((\sigma'_i, \sigma_{-i})(t))) = \liminf_{T \rightarrow \infty} \sum_{t=\tau}^T (-1) = -\infty < 0,$$

and so σ'_i is not optimal for i . Now suppose i chooses any strategy σ'_i that agrees with σ_i on path but differs from σ_i after some history h for which σ prescribes the punishment phase. Then in the subgame induced by h , we have that

$$\liminf_{T \rightarrow \infty} \sum_{t=0}^T (u_i((\sigma_i, \sigma_{-i})(t) \mid h) - u_i((\sigma'_i, \sigma_{-i})(t) \mid h)) \leq \liminf_{T \rightarrow \infty} \sum_{t=0}^T 0 = 0,$$

and so σ_i is a best response for i in this subgame. Hence σ is a subgame perfect equilibrium.

8.4 The folk theorems

There are many folk theorems, most of which say something like the following:

If players are sufficiently patient in an infinitely repeated game, then any feasible, strictly individually rational payoff profile in the stage game can be supported as an equilibrium average payoff profile in the repeated game.

Not all of them say this – some apply to finitely repeated games, some apply to infinitely repeated games without discounting, and some restrict to feasible payoffs that are strictly greater than some Nash equilibrium of the stage game. The name derives from the fact that the original folk theorem was part of the “folk wisdom” of game theory, not being attributed to any particular authors.

We defined feasible and individually rational payoff profiles in Definition 33. Since it is very important here, we will restate it:

Definition 33 (Feasibility and individual rationality).

- (a) *Feasible payoffs*.⁶⁰ In a finite game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a payoff profile $v = (v_1, \dots, v_n)$ is *feasible* if there is some probability distribution $p \in \Delta(S)$ such that

$$v_i = \sum_{s \in S} u_i(s) p(s) \quad \text{for all } i \in \mathcal{I}.$$

- (b) *Individually rational payoffs*. In a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, a payoff profile $v = (v_1, \dots, v_n)$ is *individually rational* if for each $i \in \mathcal{I}$,

$$v_i \geq \min_{\sigma_{-i} \in \Delta_{-i}(S_{-i})} \max_{\sigma_i \in \Delta(S_i)} u_i(\sigma_i, \sigma_{-i}) = v_i.$$

A payoff profile v is *strictly individually rational* if this holds with strict inequality for all i .

The game G in the definition refers to a stage game in the current context, and v to an average payoff profile. We will use $\underline{v} = (\underline{v}_1, \dots, \underline{v}_n)$ to denote the profile of minimax payoffs throughout this section.

The interpretation of the minimax payoff \underline{v}_i of a player i is that it is i 's reservation utility. This is only true under the assumption of perfect monitoring:

Proposition 68. *Consider a repeated game of perfect monitoring with stage game G and set of players \mathcal{I} . For any player $i \in \mathcal{I}$, if α^* is a Nash equilibrium α^* of the stage game G then $u_i(\alpha^*) \geq \underline{v}_i$, and if σ^* is a Nash equilibrium of the repeated game then $U_i(\sigma^*) \geq \underline{v}_i$.*

Proof. First consider α^* and fix player $i \in \mathcal{I}$. Let $\hat{\alpha} \in \times_{i=1}^n \Delta(S_i)$ be the strategy profile that solves $\min_{\alpha_{-i}} \max_{\alpha_i} u_i(\alpha_i, \alpha_{-i})$. By definition, $\hat{\alpha}_i$ is a best response to $\hat{\alpha}_{-i}$. Now,

$$u_i(\alpha_i^*, \alpha_{-i}^*) \geq u_i(\hat{\alpha}_i, \alpha_{-i}^*) \geq u_i(\hat{\alpha}_i, \hat{\alpha}_{-i}) = \underline{v}_i.$$

Next consider σ^* . Construct a strategy $\hat{\alpha}_i$ for player i as follows. For each history $h^t \in \mathcal{H}$, suppose player i myopically chooses $\hat{\alpha}_i(h^t) \in \arg \max_{a_i \in A_i} u_i(a_i, \sigma_{-i}^*(h^t))$.⁶¹ Since the information available at t is identical across players,⁶² i 's opponents can randomize only independently, and so \underline{v}_i is the minimum payoff that i 's opponents can enforce on i in any period, and thus $u_i(\hat{\alpha}_i(h^t), \sigma_{-i}^*(h^t)) \geq \underline{v}_i$ for all histories $h^t \in \mathcal{H}_t$ and every period t . Finally, since σ_i^* is a best response to σ_{-i}^* , we have $U_i(\sigma_i^*, \sigma_{-i}^*) \geq U_i(\hat{\alpha}_i, \sigma_{-i}^*) \geq \underline{v}_i$. \square

One problem in dealing with the folk theorems is that the set of feasible payoffs is not necessarily convex if δ is sufficiently small. For δ is sufficiently close to 1, Sorin (1986)

⁶⁰More generally, a payoff profile v is *feasible* if there is some probability distribution $p \in \Delta(S)$ such that $v_i = \int_S u_i dp$.

⁶¹This is not necessarily optimal (for example, consider the grim trigger strategy in Example 61).

⁶²This is important because the minimax payoff assumes independent strategies for the opponents, not correlated strategies.

and Fudenberg & Maskin (1991) show that the set of feasible payoffs is convex. They show this by showing that any convex combination of pure strategy payoff profiles in the stage game can be supported as average payoffs by some time-varying deterministic strategy.

To avoid the issue of dealing with time-varying deterministic strategies, it is convenient to instead assume the existence of a public randomization device. A *public randomization device* is a probability space $([0, 1], \mathcal{B}, p)$, where \mathcal{B} is the Borel σ -algebra on $[0, 1]$. We assume the public randomization device is common knowledge and that p is a common prior. At the start of each period t , a signal $\omega^t \in \Omega$ is observed by all players, and players can condition strategies on this signal. This allows us to “convexify” the set of feasible payoff profiles. To see this, suppose that v and v' are two pure strategy payoff profiles in the stage game, and let $F(x) = p(\omega \leq x)$. Now any convex combination $\lambda v + (1 - \lambda)v'$ can be supported as an average payoff profile in the repeated game if a strategy profile supporting v is played whenever $\omega \leq F^{-1}(\lambda)$ is observed and a strategy profile supporting v' is played otherwise.

A *history* in the repeated game is now a profile

$$h^t = ((a^0, \omega^0), (a^1, \omega^1), \dots, (a^{t-1}, \omega^{t-1}), \omega^t),$$

where each $a^\tau \in A_1 \times \dots \times A_n$ is the profile of actions played in period τ and $\omega^\tau \in [0, 1]$ is the signal observed in period τ . We now define the set of histories at time t as

$$\mathcal{H}^t = \left\{ ((a^0, \omega^0), \dots, (a^{t-1}, \omega^{t-1}), \omega^t) \mid a^\tau \in \prod_{i=1}^n A_i \text{ for all } \tau = 0, \dots, t-1 \right. \\ \left. \text{and } \omega^\tau \in [0, 1] \text{ for all } \tau = 0, \dots, t \right\}.$$

The definition of the repeated game is otherwise unchanged. For sufficiently low discount factor δ , the set of feasible payoffs will differ between the repeated game without a public randomization and the repeated game with the device. However, Fudenberg & Maskin (1991) show that the assumption there exists a public randomization device is innocuous for δ sufficiently close to 1.

8.4.1 In discounted infinitely repeated games

The most commonly studied class of folk theorems are for infinitely repeated games with $\delta \in (0, 1)$. First, any feasible, strictly individually rational payoff vector can be supported as a Nash equilibrium (despite the name, Nash did not come up with this one):

Theorem 27 (Nash folk theorem). *Consider a finite stage game G . If v is a feasible and strictly individually rational payoff profile of G , then there exists a $\underline{\delta} < 1$ such that for any $\delta \in [\underline{\delta}, 1)$, there is a subgame perfect equilibrium σ^* of $G^\infty(\delta)$ inducing payoff profile v .*

The next folk theorems concern subgame perfect equilibria. One of the simpler of these folk theorems (at least in terms of proof) is down to Friedman:

Theorem 28 (Friedman, 1971). *Consider a finite stage game G . If e is a payoff profile of a Nash equilibrium of G , and if v is a feasible payoff profile such that $v_i > e_i$ for every player $i \in \mathcal{I}$, then there exists a $\underline{\delta} < 1$ such that for any $\delta \in [\underline{\delta}, 1)$, there is a subgame perfect equilibrium σ^* of the infinitely repeated game $G^\infty(\delta)$ inducing payoff profile v .*

Proof. In general, the proof requires detailing public randomizations. To simplify matters, suppose there is some action profile a s.t. $u_i(a) = v_i$ for all players $i \in \mathcal{I}$. Let α^* be the Nash equilibrium of the stage game that yields payoff profile e . Consider the following grim trigger strategy σ_i . for each player $i \in \mathcal{I}$:

- (I) Play a_i in the first period and every subsequent period provided no player j ever plays some action $a'_j \neq a_j$.
- (II) In each period, if any player j has ever played some action $a'_j \neq a_j$, play α_i^* .

We claim that σ_i is a best response to σ_{-i} for sufficiently large $\delta < 1$. Let σ'_i be a one-shot deviation from σ_i in some period t , and let h^t be the history in which σ has been played in all periods prior to t . Then i 's continuation payoff is at most $U_i(\sigma'_i, \sigma_{-i} \mid h^t) = (1 - \delta) \max_{a_i \in A_i} u_i(a_i, \sigma_{-i}) + \delta e_i$, whereas her continuation payoff from sticking to σ_i is $U_i(\sigma_i, \sigma_{-i} \mid h^t) = v_i$. Defining $\underline{\delta} = \frac{\max_{a_i \in A_i} u_i(a_i, \sigma_{-i}) - v_i}{\max_{a_i \in A_i} u_i(a_i, \sigma_{-i}) - e_i} < 1$ and applying the one-shot deviation principle, we have that σ is optimal for each player on the equilibrium path, and thus is a Nash equilibrium. Off-the-equilibrium path, the Nash equilibrium α^* is played in every stage game, and thus it is never optimal for any player to deviate. Thus σ is a subgame perfect equilibrium.

In the case that there is no action profile a generating payoffs v_i for every player i , we need to randomize action profiles on the equilibrium path via the public randomization device such that v_i is the expected payoff to i in each period. The strategies and the application of the one-shot deviation principle are otherwise unchanged. \square

Fudenberg & Maskin (1986) prove the following folk theorem, which has become standard:

Theorem 29 (Folk theorem). *Suppose the set of feasible and strictly individually rational payoff profiles of stage game G has dimension $|\mathcal{I}|$. Then for any feasible and strictly individually rational payoff profile v of G , there exists a $\underline{\delta} < 1$ such that for any $\delta \in [\underline{\delta}, 1)$, there is a subgame perfect equilibrium σ^* of the infinitely repeated game $G^\infty(\delta)$ inducing payoff profile v .*

Proof. Let V denote the set of feasible and strictly individually rational payoff profiles and let $n = |\mathcal{I}|$. For simplicity, suppose there is some action profile \bar{a} such that $u_i(\bar{a}) = \underline{v}_i$ for all $i \in \mathcal{I}$. For any $v \in V$, let σ be a strategy profile such that $\pi_i(\sigma) = v_i$ for every player $i \in \mathcal{I}$. Let v' be another payoff profile in the interior of V , such that $v_i > v'_i$ for each player i . Since v' is in the interior of V and V has dimension $n = |\mathcal{I}|$, there is some $\epsilon > 0$ such that

$$(v'_1 + \epsilon, \dots, v'_{j-1} + \epsilon, v'_j, v'_{j+1} + \epsilon, \dots, v'_n + \epsilon) \in V$$

for each player j . Let $l^j = (l_1^j, \dots, l_n^j)$ denote a strategy profile that realizes these payoffs. Normalize j 's minmax payoff to $\underline{v}_j = 0$, let $m^j = (m_1^j, \dots, m_n^j)$ be the strategy profile that minmaxes player j and let $w_i^j = u_i(m^j)$, i 's payoff when minmaxing j .

Again for simplicity, assume that for each $i \in \mathcal{I}$, there exists a pure action profile $a(i)$ so that $u_j(a) = l_j^i$ for all $j \in \mathcal{I}$. This is not necessary to prove the theorem, it just removes the need to spend a long time on the details of public randomizations. Choose an integer T with $\max_a u_i(a) + Tv_i < \min_a u_i(a) + Tv'_i$.

Consider the following multiple phase strategy for player i :

(I) Play σ_i in each period t (generating stage payoffs v) provided σ is played in period $t - 1$. If a single player j deviates from (I) then go to (II $_j$).

(II $_j$) play m_i^j for T periods, then go to (III $_j$).

(III $_j$) play l_i^j thereafter.

If a player k deviates in phase (II $_j$) or (III $_j$), go to (II $_k$). If more than one player deviates in any phase, remain in that phase.

Suppose a player i deviates in phase (I) and conforms thereafter. She receives \bar{v}_i at most, zero for the subsequent T periods and v'_i thereafter, i.e. her gain from deviating here is at most

$$\bar{v}_i + \frac{\delta^{T+1}}{1-\delta} v'_i.$$

Not deviating instead yields $\frac{v_i}{1-\delta}$. Hence the gain from deviating is

$$\bar{v}_i - \frac{1-\delta^{T+1}}{1-\delta} v'_i$$

at most. Note $\frac{1-\delta^{T+1}}{1-\delta} \rightarrow T+1$ as $\delta \rightarrow 1$. By the condition that $\frac{\bar{v}_i}{v'_i} < T_i + 1$, the gain from deviating will therefore be negative for any $\delta > \underline{\delta}$ for some lower bound $\underline{\delta} < 1$.

If i deviates while being punished himself in phase (II $_i$), then i obtains 0 at most in this period, and lengthens punishment, postponing positive payoff v'_i (and hence receiving a strictly lower discounted payoff). If i deviates when $j \neq i$ is being punished in phase (II $_j$), then he receives at most

$$\bar{v}_i + \frac{\delta^{T+1}}{1-\delta} v'_i,$$

as for deviation in phase (I). However, by not deviating, he receives at least

$$w_i^j \frac{1-\delta^{T'}}{1-\delta} + \frac{\delta^{T'+1}}{1-\delta} (v'_i + \epsilon),$$

where $1 < T' < T$. The gain to deviating is at most

$$\bar{v}_i + \frac{1-\delta^{T'+1}}{1-\delta} (v'_i - w_i^j) - \frac{\delta^{T'+1}}{1-\delta} \epsilon - \delta^{T'} w_i^j.$$

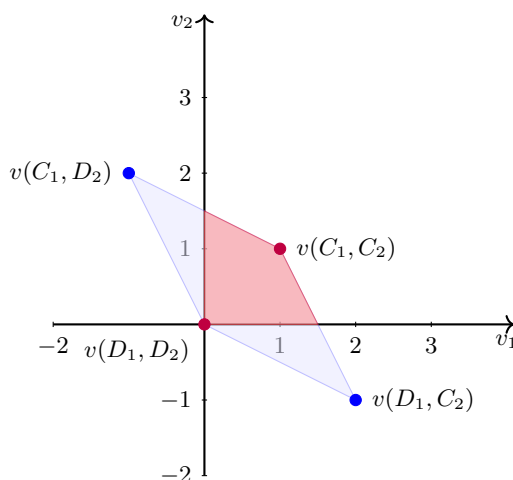
As $\delta \rightarrow 1$, the second term stays finite but the third converges to $-\infty$ and so the gain to deviating stays negative. Hence there is some $\delta_i < 1$ such that i will not deviate for any $\delta > \delta_i$. The proof ruling out deviations in phases (III_{*i*}) and (III_{*j*}) are very similar to that of phase (I).

If instead, the minimax strategies are mixed, we need to alter the proof so that in phase (II_{*j*}), player i is indifferent between all of the possible length- T realizations of sequences of actions that are prescribed by the strategy profile for minimaxing j . This can be achieved by tailoring the reward ϵ to each possible such sequence, making i indifferent by promising a greater future payoff in phase (III_{*j*}) for those sequences that yield him a lower payoff in phase (II_{*j*}). See Fudenberg & Maskin (1986) for details. \square

The proof hinges on a punishment scheme whereby a player i 's opponents are rewarded for minimaxing player i if he deviates: in phase (III_{*i*}), they all receive strictly greater payoffs than they would if they had returned to (I).

Maskin & Fudenberg (1986) assumption of full dimensionality, i.e. that if V is the set of feasible payoff profiles, then $|V| = |\mathcal{I}| = n$. The motivation for this condition is that the strategies constructed in the proof of the folk theorem require that the other players can be rewarded for punishing a player who deviates. We thus want to avoid a situation where rewarding one of the players implementing the punishment also rewards the deviating player. Full dimensionality is sufficient to avoid this, but not necessary. Abreu, Dutta & Smith (1994) show that a necessary and sufficient condition is that for no pair of players $i, j \in \mathcal{I}$, is u_i an affine transformation of u_j – that is, there does not exist $\alpha_{ij} > 0$ and $\beta_{ij} \in \mathbb{R}$ such that $u_i = \alpha_{ij}u_j + \beta_{ij}$.

Example 61 (continued). In the infinitely repeated prisoners' dilemma, the set of payoff profiles can be depicted as follows:

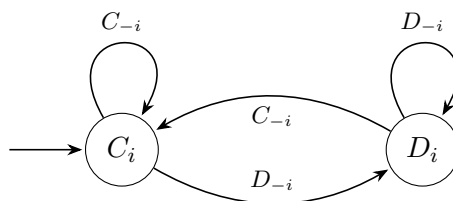


The set of feasible payoff profiles is shaded blue and the set of feasible and individually rational payoff profiles is shaded red. By the folk theorem, any payoff profile in the red set can be sustained as a profile of average payoffs for some subgame perfect equilibrium of the infinitely repeated prisoners' dilemma provided players are sufficiently patient.

We previously considered a subgame perfect equilibrium consisting of grim trigger strategies such that on the equilibrium path, (C_1, C_2) was played in every period. Consider instead the following tit-for-tat strategy σ_i for each player i :

1. Play C_i in the first period. In any subsequent period, play C_i if player $-i$ played C_{-i} in the previous period.
2. Play D_i if player $-i$ played D_{-i} in the previous period.

As an automaton, we can visualize this strategy as follows:



The nodes in this case are player i 's own actions, rather than the action profiles of both players. These tit-for-tat strategies do not constitute a Nash equilibrium: for any $\delta < 1$, we have $U_i(\sigma_i, \sigma_{-i} \mid C_i, C_{-i}) = 1 - \delta < (1 - \delta)(2 - \delta) + \delta^2 = U_i(D_i, \sigma_{-i} \mid C_i, C_{-i})$, so a deviation is profitable. Now suppose that instead, for each player i , σ_i calls for i to play D_i for 3 periods, and only return to playing C_i if player $-i$ plays C_i for 3 consecutive periods. Then this constitutes a subgame perfect equilibrium for sufficiently high δ . Now we have that $U_i(\sigma_i, \sigma_{-i} \mid C_i, C_{-i}) = 1 - \delta \geq (1 - \delta)(2 - \delta - \delta^2 - \delta^3) + \delta^4 = U_i(D_i, \sigma_{-i} \mid C_i, C_{-i})$ for all $\delta \geq \underline{\delta}$ where $\underline{\delta}$ solves $0 = 1 - \underline{\delta} - \underline{\delta}^2 - \underline{\delta}^3 + (1 - \underline{\delta})\underline{\delta}^4$, which gives us $\underline{\delta} \approx 0.780$. The other states are straightforward to check.

9 Cooperative game theory

9.1 Cooperative games

In non-cooperative game theory, solutions can always be envisaged as a self-enforcing agreement, in the sense that no rational player will choose to deviate. In cooperative game theory, solutions can be envisaged as an agreement which is enforced externally, such as via a system of legal institutions that enforce contracts. The interesting question is which groups will come together to form an agreement and what their payoffs will be – a more high-level approach, since how these coalitions come about is not important.

Definition 84 (Cooperative game). A *cooperative game* or *coalitional game* is a tuple $(N, (A_S)_{S \in 2^N}, (u_i)_{i \in N})$, where:

- (i) N is a nonempty finite set of players. Each $S \subseteq N$ is a *coalition*, and we call N the *grand coalition*.
- (ii) A_S is a nonempty set of actions available to coalition S .

- (iii) An *outcome* $(\mathcal{P}, (a_S)_{S \in \mathcal{P}})$ consists of a partition \mathcal{P} of N and a list of actions $a_S \in A_S$ for each coalition $S \in \mathcal{P}$. Let X denote the set of possible outcomes.
- (iv) $u_i : X \rightarrow \mathbb{R}$ is a payoff function for player i , mapping outcomes to payoffs.

Some examples:

Example 63.

- (a) *Marriage market.* Suppose D is a group of doctors and H is a group of hospital vacancies. Doctors have preferences over which vacancy they would like to fill, if any, and each hospital hiring team (which is different for each vacancy) has preferences over which doctors they wish to hire. A *matching* is a partition of $D \cup H$ into doctor-vacancy pairs and single doctors/vacancies. Each doctor and hiring team cares only about their own match.
- (b) *Control of a resource.* Suppose there are three shepherds who have access to a well, which they need to supply water to their sheep. Unfortunately, a large rock is blocking access to the well, and it is too large for any one of the shepherds to move. However, if any group of two or more of the shepherds work together, they can move the rock and claim control of the well, deciding how to divide access to the water within the group. Each shepherd only cares about what share of the water they receive, and they prefer more water to less. An action for any coalition S is an allocation of water among the members of the coalition.

Definition 85. A cooperative game is *cohesive* if for every outcome $(\mathcal{P}, (a_S)_{S \in \mathcal{P}})$, there exists some outcome $(\{N\}, a_N)$ generated by the grand coalition such that $u_i(\{N\}, a_N) \geq u_i(\mathcal{P}, (a_S)_{S \in \mathcal{P}})$ for every player $i \in N$.

In other words, a cooperative game is cohesive if for any outcome, the grand coalition can always produce an outcome that every player weakly prefers to it.

9.2 Cooperative theory of bargaining

Bargaining nicely illustrates the difference between cooperative and noncooperative approaches. In the noncooperative approach, we specified the structure of bargaining, and this could take many forms. The structure of bargaining heavily influenced equilibrium outcomes. In the ultimatum game, the proposer gets away with awarding herself the entire surplus. Where there is multiperiod bargaining and alternating proposers, outcomes depended on whether the bargaining process continued indefinitely or for finitely many rounds, and on the patience of the players. The cooperative approach to bargaining abstracts away from this.

Definition 86 (Bargaining problem). An n -player bargaining problem is a pair $\mathcal{B} = (U, d)$, consisting of

- (i) a *feasible set* $U \subseteq \mathbb{R}^n$, also called a *utility possibility set*, and

(ii) a *disagreement point* $d \in U$.

We call any point $u \in U$ an *agreement*.

We assume the set U is nonempty, convex and compact and that d lies in the interior of U . This implies there is some $u \in U$ such that $u \gg d$ (i.e. $u_i > d_i$ for all i), and thus the game is cohesive.

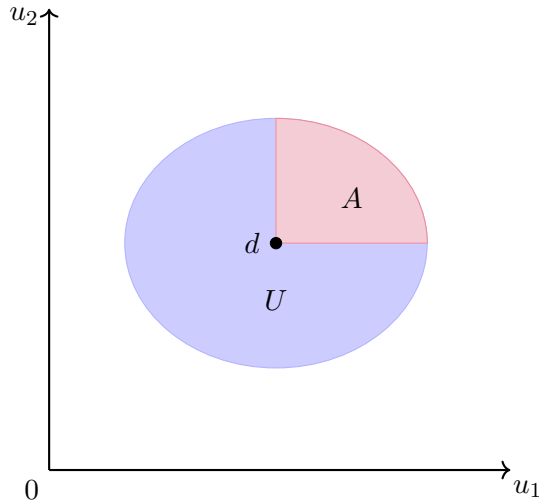
The disagreement point is interpretable as the utility agents would receive if bargaining fails, so it is equivalent to a reservation utility.

Hereafter, we consider bargaining problems with two agents, though everything can be generalized to $n \geq 2$ agents.

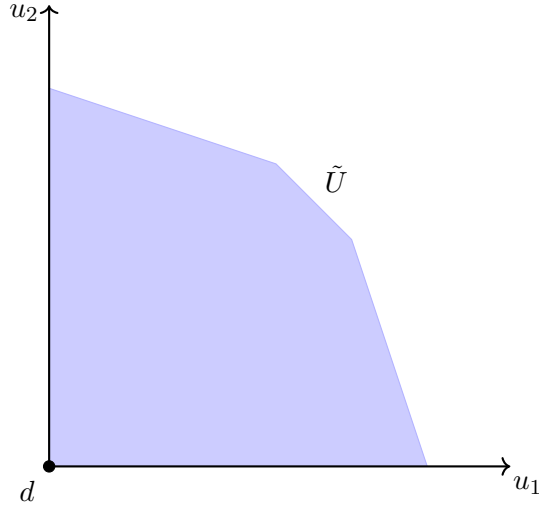
Definition 87 (Set of possible agreements). In a bargaining problem $\mathcal{B} = (U, d)$, the *set of possible agreements* A is the set of individually rational agreements, that is,

$$A = \{u \in U : u \gg d\}.$$

Example 64 (A two-agent bargaining problem). The feasible set U is given in blue, and the set of possible agreements $A \subseteq U$ is given in red.



Intuitively, if agents are rational then no bargaining process can result in an agreement outside the set of possible agreements, since then some agent would receive a higher payoff from the disagreement point. Thus we can restrict attention to the set of possible agreements and renormalize the disagreement point to the origin:



We can always translate the feasible set U so that the disagreement point is at the origin, i.e. by considering $\tilde{U} = \{(u_1 - d_1, u_2 - d_2) : (u_1, u_2) \in U\}$. For the solution concepts we consider, there is no loss of generality to such translations.

Definition 88 (Bargaining solution). A *bargaining solution* f is a function mapping any bargaining problem $\mathcal{B} = (U, d)$ to a feasible outcome $f(\mathcal{B}) \in U$.

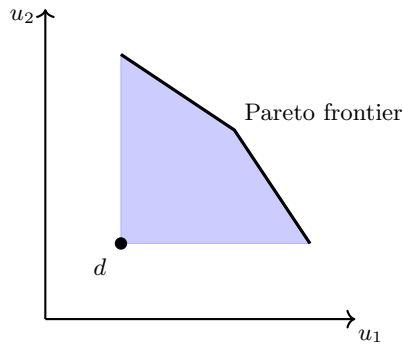
A bargaining solution is chosen to satisfy a set of desirable axioms.

9.2.1 Nash bargaining solution

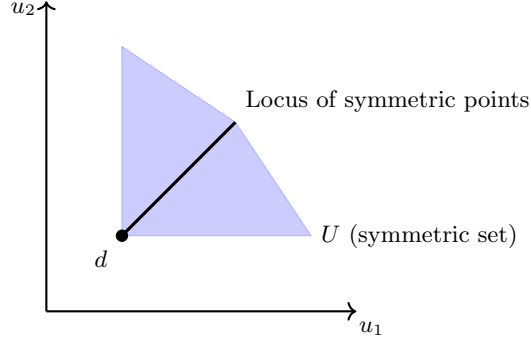
Introduced by Nash (1950), the Nash bargaining solution satisfies the following four axioms:

Axioms.

(B1) *Weak Pareto optimality*. If $u^* = f(U, d)$ then there is no point $u' = (u'_1, u'_2) \in U$ such that $u'_i > u_i^*$ for each $i = 1, 2$.

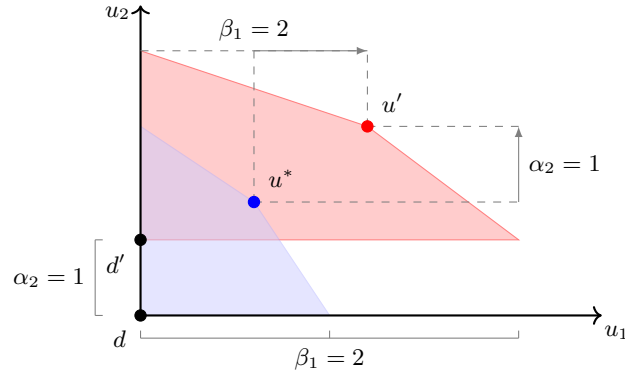


- (B2) *Symmetry.* If (U, d) is a symmetric problem [i.e. $d_1 = d_2$ and if $(u, u') \in U$ then $(u', u) \in U$] then $u^* = f(U, d)$ is such that $u_1^* = u_2^*$.

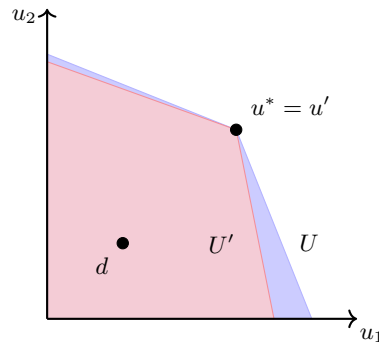


- (B3) *Invariance to equivalent payoff representations.* For any $U' = \{\beta' u + \alpha : u \in U\}$ where $\alpha, \beta \in \mathbb{R}^2$ with $\beta \gg 0$ and $d' = \beta' d + \alpha$,

$$u^* = f(U, d) \quad \text{iff} \quad \beta' u^* + \alpha = f(U', d').$$



- (B4) *Independence of irrelevant alternatives.* If $U' \subseteq U$ and $f(U, d) \in U'$ then $f(U', d) = f(U, d)$.



Of these axioms, the independence of irrelevant alternatives (A4) is the most controversial. For example, an asymmetric change in the feasible set intuitively might change the relative bargaining power of the two players and thus the outcome, yet under the independence of irrelevant alternatives axiom it does not (as in the above figure). Kalai & Smorodinsky (1975) make this criticism.

Definition 89 (Nash bargaining solution). Given a bargaining problem $\mathcal{B} = (U, d)$, the Nash bargaining solution $f^N(\mathcal{B})$ is the point u^* that solves

$$\max_{u_1 \geq d_1, u_2 \geq d_2} (u_1 - d_1)(u_2 - d_2) \quad \text{s.t.} \quad (u_1, u_2) \in U, u_1 \geq d_1, u_2 \geq d_2.$$

Theorem 30 (Nash, 1950). *Every bargaining problem $\mathcal{B} = (U, d)$, where U is a compact convex set containing d , has a unique Nash bargaining solution $f^N(\mathcal{B})$, and the Nash bargaining solution is the unique bargaining solution satisfying axioms (B1)-(B4).*

Proof. First, we prove f^N is well-defined and unique. Fix any bargaining problem (U, d) . First, since $g(u_1, u_2) = (u_1 - d_1)(u_2 - d_2)$ is a continuous function on a compact set $V = \{U \in u : u \geq d\}$, g attains a maximum on V by the extreme value theorem, so a Nash solution exists. Suppose there exist points $u, u' \neq u'$ in U s.t. u and u' are both Nash bargaining solutions. Then

$$(u_1 - d_1)(u_2 - d_2) = (u'_1 - d_1)(u'_2 - d_2),$$

and we must thus have that if $u_1 - u'_1 < 0$ then $u_2 - u'_2 > 0$, or if $u_1 - u'_1 > 0$ then $u_2 - u'_2 < 0$. Since $u \neq u'$, one of these two sets of inequalities must hold. Because U is convex, $u'' = \frac{1}{2}(u + u')$ lies in U . Now, let $x = 2(u''_1 - d_1)(u''_2 - d_2) - (u_1 - d_1)(u_2 - d_2) - (u'_1 - d_1)(u'_2 - d_2)$. Then

$$\begin{aligned} x &= \frac{1}{2}(u_1 + u'_1)(u_2 + u'_2) - u_1 u_2 - u'_1 u'_2 \\ &= \frac{1}{2}[u'_1 u_2 + u_1 u'_2 - u_1 u_2 - u'_1 u'_2] \\ &= \frac{1}{2}(u_1 - u'_1)(u'_2 - u_2) > 0. \end{aligned}$$

Hence u'' has a strictly greater Nash product than u and u' , so u and u' do not maximize the Nash product, yielding a contradiction. Thus the Nash bargaining solution is unique.

Now to prove f^N satisfies the four axioms:

- (B1) Let $g(u_1, u_2) = (u_1 - d_1)(u_2 - d_2)$. On the set $\{u \in U : u \gg d\}$, this is strictly increasing in u_1 and u_2 . Consider any non-Pareto optimal point $u \gg d$ in U . There is some u' where $u'_i \geq u_i$ with strict inequality for at least one $i = 1, 2$. We have $g(u') > g(u)$, and hence u does not maximize g . Thus f^N is Pareto optimal.
- (B2) Suppose $u = f^N(U, d)$ is asymmetric given a symmetric problem (U, d) . Let $\bar{d} = d_1 = d_2$. Then either $u_1 > u_2$ or $u_2 > u_1$. Define $u' = (u'_1, u'_2)$ with $u'_1 = u'_2 =$

$\frac{1}{2}(u_1 + u_2)$. Since U is convex, $u' \in U$. We have

$$\begin{aligned}(u_1 - \bar{d})(u_2 - \bar{d}) &= u_1 u_2 - \bar{d} u_1 - \bar{d} u_2 + \bar{d}^2 \\ &< \frac{1}{4}(u_1^2 + 2u_1 u_2 + u_2^2) - u_1 \bar{d} + u_2 \bar{d} + \bar{d}^2 \\ &= (u'_1 - \bar{d})(u'_2 - \bar{d}),\end{aligned}$$

where the inequality follows since

$$u_1^2 + 2u_1 u_2 + u_2^2 - 4u_1 u_2 = u_1^2 + u_2^2 - 2u_1 u_2 = (u_1 - u_2)^2 > 0.$$

(B3) Fix a bargaining problem (U, d) . For any $\alpha, \beta \in \mathbb{R}^2$ such that $\beta \gg 0$ consider the transformed problem (U', d') with $U' = \{(\beta_1 u_1 + \alpha_1, \beta_2 u_2 + \alpha_2) : (u_1, u_2) \in U\}$ and $d' = (\beta_1 d_1 + \alpha_1, \beta_2 d_2 + \alpha_2)$. The Nash solution for (U', d') solves

$$\begin{aligned}&\max_{u' \in U' : u' \gg d'} (u'_1 - d'_1)(u'_2 - d'_2) \\ &= \max_{u \in U : u \gg d} (\beta_1 u_1 + \alpha_1 - \beta_1 d_1 - \alpha_1)(\beta_2 u_2 + \alpha_2 - \beta_2 d_2 - \alpha_2) \\ &= \max_{u \in U : u \gg d} \beta_1 \beta_2 (u_1 - d_1)(u_2 - d_2) \\ &= \max_{u \in U : u \gg d} (u_1 - d_1)(u_2 - d_2),\end{aligned}$$

where the final line follows because $\beta_1 \beta_2 > 0$.

(B4) Suppose u^* is the Nash solution for bargaining problem (U, d) . Then

$$(u_1^* - d_1)(u_2^* - d_2) \geq (u_1 - d_1)(u_2 - d_2) \quad \text{for all } u \in U.$$

This inequality thus holds for all $u \in U' \subseteq U$, so if $u^* \in U'$ it must maximize the Nash product within U' .

Finally, suppose f is a bargaining solution satisfying **(B1)**-**(B4)**. Consider any bargaining problem (U, d) and let $u^* = f^N(U, d)$. Define

$$U' = \{\beta' u + \alpha \mid u \in U, \beta' u + \alpha = (1/2, 1/2)' \text{ and } \beta' d + \alpha = (0, 0)'\}.$$

That is, U' results from applying an affine transformation to U such that u^* is mapped to $(1/2, 1/2)$ and d is mapped to $(0, 0)$. Since f and f^N satisfy the invariance axiom **(B3)**, $f^N(U', 0) = (1/2, 1/2)$ and so we need only show $f(U', 0) = (1/2, 1/2)$.

First, note there is no $u \gg 0$ in U' with $u_1 + u_2 > 1$. Suppose otherwise, and define $t(\lambda) = [1 - \lambda](1/2, 1/2) + \lambda(u_1, u_2)$ for $\lambda \in (0, 1)$. By convexity of U' , $t(\lambda) \in U'$ for all $\lambda \in (0, 1)$. We have Nash product

$$t_1(\lambda)t_2(\lambda) = (1 - \lambda)^2/4 + \lambda(1 - \lambda)[u_1 + u_2]/2 + \lambda^2 u_1 u_2 > (1 - \lambda)^2/4 + \lambda(1 - \lambda) + \epsilon$$

where $\epsilon \in (0, u_1 u_2]$. This exceeds $1/4$, the Nash product of $f^N(U', 0)$ for any $\lambda \in (0, 1)$ s.t.

$$(1 - \lambda)^2 + 2\lambda(1 - \lambda) + \epsilon > 1,$$

i.e. for $\lambda \in (0, \sqrt{\epsilon})$. But this implies $f^N(U', 0)$ does not maximize the Nash product on U , yielding a contradiction. Thus U' is bounded.

Since U' is bounded and convex, we can find a rectangle $U'' \supset U'$ s.t. U'' is symmetric about the line $\{u \mid u_1 = u_2\}$ and $(1/2, 1/2)$ lies on the boundary of U'' .⁶³ By the weak Pareto and symmetry axioms (B1) and (B2), $f(U'', 0) = (1/2, 1/2)$. Since $U' \subseteq U''$, by independence of irrelevant alternatives (B4), we must have that $f(U', 0) = (1/2, 1/2)$. Thus $f(U', 0) = f^N(U', 0)$. By the invariance axiom (B3), it follows that $f(U, d) = f^N(U, d)$. \square

The Nash bargaining solution and theorem easily extends for all finite n -player bargaining problems, not just the 2 player setting here. It is straightforward to generalize the Nash axioms to n players, and the Nash bargaining solution f^N is then the point u^* solving

$$\max_{u \in U \cap \{\bar{u} \mid \bar{u} \gg d\}} \prod_{i=1}^n (u_i - d_i).$$

Example 65 (Risk aversion and the Nash bargaining solution). Consider the constant relative risk aversion preferences $u(x) = x^\rho$ where $\rho \in (0, 1]$.

The Nash bargaining solution generally favours less risk averse players. For example, suppose player 1 is risk neutral ($\rho = 1$) and player 2 is strictly risk averse ($\rho < 1$). That is, $u_1(x_1) = x_1$ and $u_2(x_2) = x_2^\rho$. As usual, there is a surplus of 1 to be divided between the two players. Let the disagreement point be $d = (0, 0)$. The Nash bargaining solution f^N solves

$$\max_{x_1} u_1(x_1)u_2(1 - x_1) = \max_{x_1} x_1(1 - x_1)^\rho,$$

or equivalently,

$$\max_{x_1} \log x_1 + \rho \log(1 - x_1).$$

We have first order condition,

$$\frac{1}{x_1} - \frac{\rho}{1 - x_1} = 0,$$

which yields solution $x_1^* = \frac{1}{1+\rho}$, so the risk neutral player secures a larger share of the surplus.

We previously discussed non-cooperative models of bargaining at length – models in which a lot of care was taken to explicitly describe a bargaining process. A natural question to ask is whether these models and the Nash bargaining solution coincide. Binmore, Rubinstein & Wolinsky (1985) show that the Nash bargaining solution can be obtained as the limit of the subgame perfect equilibrium of the Rubinstein alternating bargaining model as the rate of time preference δ tends to 1.

⁶³The rectangle has an edge passing through $(1/2, 1/2)$ of slope -1 .

Proposition 69. *In the Rubinstein alternating bargaining model, suppose $\delta := \delta_1 = \delta_2$, and let $x^*(\delta)$ and $y^*(\delta)$ denote the unique subgame perfect equilibrium offers of players 1 and 2 respectively. Let $(U, 0)$ be the bargaining problem corresponding to this game. Then $\lim_{\delta \rightarrow 1} x^*(\delta) = \lim_{\delta \rightarrow 1} y^*(\delta) = f^N(U, 0)$.*

Proof. The bargaining problem $(U, 0)$ is symmetric, and therefore the Nash bargaining solution is $(1/2, 1/2)$ by the symmetry and Pareto axioms.

Write $x^*(\delta) = (x(\delta), 1 - x(\delta))$ and $y^*(\delta) = (y(\delta), 1 - y(\delta))$. By Proposition 60 where $x(\delta) = \frac{1-\delta}{1-\delta^2}$ and $y(\delta) = \frac{\delta(1-\delta)}{1-\delta^2}$. Applying l'Hôpital's rule, we have that $\lim_{\delta \rightarrow 1} x(\delta) = \frac{1}{2} = \lim_{\delta \rightarrow 1} y(\delta)$, and so $\lim_{\delta \rightarrow 1} x^*(\delta) = (1/2, 1/2) = \lim_{\delta \rightarrow 1} y^*(\delta)$. \square

9.2.2 Raiffa-Kalai-Smorodinsky bargaining solution

We could take issue with the independence of irrelevant alternatives axiom, because an increase in one's feasible set could arguably improve one's bargaining position, even if the Nash solution does not change. This leads Kalai & Smorodinsky (1975) to suggest an alternative to the Nash bargaining solution, which was first introduced by Raiffa (1953). The Raiffa-Kalai-Smorodinsky solution drops the independence of irrelevant alternatives axiom in favour of *individual monotonicity*.

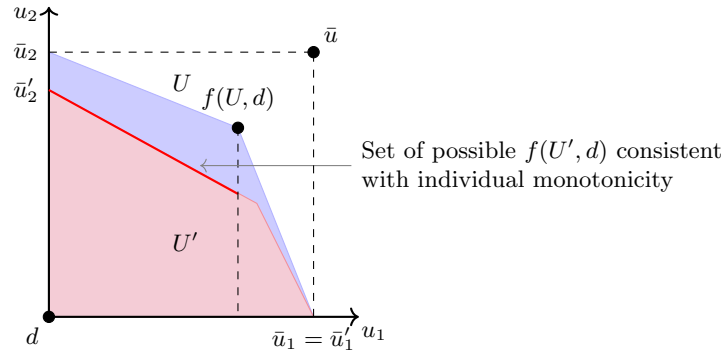
Given a bargaining problem $\mathcal{B} = (U, d)$, define the *utopia point* of \mathcal{B} by

$$\begin{aligned} \bar{u}(U, d) &:= (\max\{u_1 : u \in U, u_1 \geq d_1\}, \max\{u_2 : u \in U, u_2 \geq d_2\}) \\ &= (\bar{u}_1, \bar{u}_2). \end{aligned}$$

Now consider the following axiom:

Axiom.

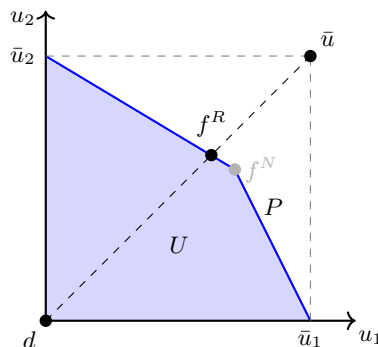
(B5) Individual monotonicity. If $U' \subseteq U$, $\bar{u}_i(U', d) = \bar{u}_i(U, d)$ for some $i \in \{1, 2\}$, then $f_j(U', d) \leq f_j(U, d)$ for $j \neq i$.



Definition 90 (Raiffa-Kalai-Smorodinsky solution). The *Raiffa-Kalai-Smorodinsky solution* $f^R(U, d)$ is given by the intersection of the straight line between d and \bar{u} with the

weak Pareto optimal boundary of U . That is, if $P \subseteq U$ is the Pareto frontier of U , then $f^R(U, d)$ is the point $u^* \in P$ such that

$$\frac{u_1^* - d_1}{u_2^* - d_2} = \frac{\bar{u}_1 - d_1}{\bar{u}_2 - d_2}.$$



An example of the Raiffa-Kalai-Smorodinsky bargaining solution f^R . In this bargaining problem, f^R prescribes a different allocation to the Nash bargaining solution f^N .

Theorem 31 (Kalai & Smorodinsky, 1975). *The Raiffa-Kalai-Smorodinsky bargaining solution is the unique bargaining solution satisfying axioms (B1)-(B3) and (B5).*

9.3 Transferable utility games

We say that utility is *transferable* if a player can losslessly transfer part of her utility to any other player. In the context of a cooperative game, this implies that the payoff for a coalition can be summarized by a single number, and this payoff can then be distributed in some way to coalition members. In a characteristic function game, there are no externalities – players only care about the actions of the coalition to which they belong, and not to the actions of other coalitions.

Definition 91 (TU games). A *characteristic function game* or *TU-game* is a pair (N, v) where $N = \{1, \dots, n\}$ is a set of n players and $v : 2^N \rightarrow \mathbb{R}$ is a function assigning to each subset $S \subseteq N$ a real number $v(S)$ such that $v(\emptyset) = 0$.

A subset $S \subseteq N$ is called a *coalition* and the set N is called a *grand coalition*. The function v is called the *characteristic function* and the number $v(S)$ is called the *worth* of S .

An *allocation* is a vector $x \in \mathbb{R}^n$.

We use \mathcal{G}^N to denote the set of all TU-games with set of players N .

Given a coalition S and an allocation $x = (x_1, \dots, x_n)$, we write $x(S) := \sum_{i \in S} x_i$ for the total payoff of coalition S under allocation x . To avoid needlessly complicating notation, we write $v(i)$ for $v(\{i\})$, the worth of a player i on her own, and we write $S \cup i$ for $S \cup \{i\}$.

The set of actions available to a coalition S in a game (N, v) is $A_S = \{x \in \mathbb{R}^n \mid x(N) \leq v(S)\}$. That is, coalitions' actions are allocations. The coalition generates worth $v(S)$, and thus S can only implement an allocation x if $x(N) = \sum_{i \in N} x_i \leq v(S)$.

In much of the literature, the actions of a coalition S are instead taken to be vectors $(x_i)_{i \in S}$ such that $\sum_{i \in S} x_i = v(S)$. Conceptually, the difference in approach is that we allow S to implement allocations that give something to players who are not in S . This is without consequence in the games we consider, because players only care about their own payoffs. When we consider coalitional deviations from the grand coalition, the players in coalition S will always most prefer some allocation that splits $v(S)$ only among themselves.

Definition 92 (Imputation). Given a TU-game (N, v) , a vector $x \in \mathbb{R}^n$ is called an *imputation* if

- (i) x is *individually rational*, that is, if $x_i \geq v(i)$ for all $i \in N$, and
- (ii) x is *efficient*, that is, $x(N) = v(N)$.

We denote the set of imputations of (N, v) by $I(v)$.

An imputation $x \in I(v)$ is an allocation of the worth of the grand coalition N such that each player is given at least as great a payoff as she could earn on her own.

Because $v : 2^N \rightarrow \mathbb{R}$ already encodes N (for v is defined on the power set of N), we often call v a game instead of (N, v) . We detail some properties of v :

Definition 93 (Properties of the characteristic function). Consider a characteristic function $v : 2^N \rightarrow \mathbb{R}$.

- (a) *Convexity*. We call a game v *convex* (or *supermodular*) if

$$v(S \cup T) + v(S \cap T) \geq v(S) + v(T) \quad \text{for all } S, T \subseteq N.$$

- (b) *Superadditivity*. We call a game v *superadditive* if

$$v(S \cup T) \geq v(S) + v(T) \quad \text{for all } S, T \subseteq N \text{ s.t. } S \cap T = \emptyset,$$

that is, for all disjoint coalitions S, T .

If $v(S \cup T) = v(S) + v(T)$ for all disjoint coalitions S, T , then v is called *additive*.

- (c) *Monotonicity*. We call a game v *monotone* if $S \subseteq T$ implies $v(S) \leq v(T)$.
- (d) *Essential game*. We call a game v *essential* if $v(N) \geq \sum_{i=1}^n v(i)$.

An essential game is one in which the worth of the grand coalition is weakly greater than the sum of the payoffs that individual players can earn on their own. If a game is not essential, then clearly there is no imputation of the game, since any efficient allocation x must involve $x_i < v(i)$ for some i .

Proposition 70. *Let (N, v) be a TU-game.*

- (i) *If v is convex then it is superadditive.*
- (ii) *If v is superadditive, it is essential.*
- (iii) *If v is superadditive or monotone, it is cohesive.*
- (iv) *If v is nonnegative and superadditive, then it is monotone.*
- (v) *v is convex iff for every player $i \in N$,*

$$v(S \cup i) - v(S) \leq v(T \cup i) - v(T)$$

for all coalitions $S \subseteq T \subseteq N - \{i\}$.

Proof. (i) If v is convex and S, T are disjoint coalitions, then $S \cap T = \emptyset$, and thus by convexity, $v(S \cup T) = v(S \cup T) + v(\emptyset) \geq v(S) + v(T)$.

(ii) Let $S_1 = N - \{1\}$ and let $S_k = S_{k-1} - \{k\}$ for $k = 2, \dots, n-1$. By superadditivity, we have $v(N) \geq v(S_1) + v(1)$ and for each $k = 1, \dots, n-1$, we have $v(S_k) \geq v(S_{k+1}) - v(k)$. Combining gives $v(N) \geq \sum_{i=1}^n v(i)$.

(iii) Monotonicity immediately implies cohesiveness. Suppose v is superadditive. Applying the definition repeatedly, we see that for any partition \mathcal{P} of N , we have $v(N) \geq \sum_{S \in \mathcal{P}} v(S)$. Hence for any allocation y such that $y(S) = v(S)$ for all $S \in \mathcal{P}$, there is some imputation x such that $x_i \geq y_i$ for all players i . Thus v is cohesive.

(iv) Suppose v is nonnegative and superadditive. Suppose $S \subseteq T$. Then $v(T) = v(S \cup (T - S)) \geq v(S) + v(T - S)$, and $v(T - S) \geq 0$, so $v(T) \geq v(S)$.

(v) Suppose v is convex. Then $v((S \cup i) \cup T) + v((S \cup i) \cap T) \geq v(S \cup i) + v(T)$, and since $S \subseteq T$ and $i \notin S \cup T$, $v((S \cup i) \cap T) = v(S)$ and $v((S \cup i) \cup T) = v(T \cup i)$. Hence $v(S \cup i) - v(S) \leq v(T \cup i) - v(T)$. Conversely, suppose $v(S \cup i) - v(S) \leq v(T \cup i) - v(T)$ for all $S \subseteq T \subseteq N - \{i\}$, for all players i . Fix $S_0 \subseteq T_0 \subseteq N$ and any $R = \{i_1, \dots, i_k\} \subseteq N - T_0$. Repeatedly applying the inequality gives

$$v(S_0 \cup \{i_1, \dots, i_j\}) - v(S_0 \cup \{i_1, \dots, i_{j-1}\}) \leq v(T_0 \cup \{i_1, \dots, i_j\}) - v(T_0 \cup \{i_1, \dots, i_{j-1}\})$$

for $j = 1, \dots, k$. Combining the inequalities gives $v(S_0 \cup R) - v(S_0) \leq v(T_0 \cup R) - v(T_0)$ for any $R \subseteq N - T_0$. Taking arbitrary S, T and setting $S_0 = S \cap T$, $T_0 = T$ and $R = S - T$ yields the convex inequality, and thus v is convex. □

In essential games, the set of imputations is a convex set that we can easily derive:

Proposition 71. Consider any essential game (N, v) . For each $i \in N$, define $f^i = (f_1^i, \dots, f_n^i)$ by

$$f_k^i = \begin{cases} v(k) & \text{if } k \neq i, \\ v(N) - \sum_{j \in N - \{i\}} v(j) & \text{if } k = i. \end{cases}$$

Then the set of imputations $I(v)$ is given by

$$I(v) = \text{co}(f^1, \dots, f^n),$$

that is, the convex hull of the points f^1, \dots, f^n .

Proof. Consider any f^i . For $k \neq i$, k 's payoff is $f_k^i = v(k)$ and i 's payoff is $v(N) - \sum_{j \neq i} v(j) \geq v(i)$, where the inequality holds since the game is essential. Hence f^i is individually rational. Furthermore, $f^i(N) = \sum_{k \neq i} v(k) + v(N) - \sum_{k \neq i} v(k) = v(N)$, so f^i is efficient. It follows that f^i is an imputation.

Next, consider any convex combination $x = \lambda_1 f^1 + \lambda_2 f^2 + \dots + \lambda_n f^n$ (with $\sum_{i \in N} \lambda_i = 1$, $\lambda_i \geq 0$). For each $k \in N$, since $f_k^i \geq v(k)$ for all $i \in N$, we have that $x_k \geq v(k)$. Thus x is individually rational. Furthermore, $x(N) = \sum_{i=1}^n \lambda_i f^i(N) = \sum_{i=1}^n \lambda_i v(N) = v(N) \sum_{i=1}^n \lambda_i = v(N)$, so x is efficient. Thus x is an imputation. Since $x \in \text{co}(f^1, \dots, f^n)$ iff x is a convex combination of f^1, \dots, f^n , it follows that $\text{co}(f^1, \dots, f^n) \subseteq I(v)$.

If $v(N) = \sum_k v(k)$, then trivially, the only imputation is x defined by $x_k = v(k)$ for all k , and we have $f^1 = \dots = f^n = x$. If $v(N) > \sum_k v(k)$, then $\{f^1, \dots, f^n\}$ is linearly independent and spans \mathbb{R}^n . Consider any $x \in I(v)$. We must have $x_k \geq v(k)$ and $\sum_i x_i = v(N)$. Let F be an $n \times n$ matrix with i th column f^i , and let x be a column vector and $\lambda = (\lambda_1, \dots, \lambda_n)$ be a column vector. Since the columns of F form a basis for \mathbb{R}^n , the system of equations $F\lambda = x$ has a unique solution $\bar{\lambda}$, and so $x = \sum_i \bar{\lambda}_i f^i$. If $\sum_i \bar{\lambda}_i \neq 1$, then $x(N) = \sum_{i=1}^n \bar{\lambda}_i f^i(N) = v(N) \sum_{i=1}^n \bar{\lambda}_i \neq v(N)$, which contradicts $x \in I(v)$. Likewise, if $\bar{\lambda}_i < 0$ for some i , then for player i ,

$$\begin{aligned} x_i &= -|\lambda_i| \left[v(N) - \sum_{j \neq i} v(j) \right] + \sum_{j \neq i} \bar{\lambda}_j v(j) \\ &= -|\lambda_i| \left[v(N) - \sum_{j \neq i} v(j) \right] + (1 + |\lambda_i|)v(i) \\ &< v(i), \end{aligned}$$

where the inequality follows from $v(N) - \sum_{j \neq i} v(j) > v(i)$. Thus x is not individually rational, yielding a contradiction. It follows that $I(v) \subseteq \text{co}(f^1, \dots, f^n)$. This completes the proof. \square

9.4 The core and stable sets

The core and stable set are two set-valued solution concepts. In cohesive games, the grand coalition always forms, and thus a natural question to ask is which actions

the grand coalition chooses. These solution concepts lack much interpretation in non-cohesive games.

9.4.1 The core

Definition 94 (Domination). Let (N, v) be a game, let $y, z \in I(v)$ and let $S \subseteq N$ be nonempty. Then y *dominates* z in coalition S if

- (i) $y_i > z_i$ for all $i \in S$, and
- (ii) $y(S) \leq v(S)$.

For $y, z \in I(v)$, we say that y *dominates* z if there exists a (nonempty) coalition $S \subseteq N$ such that y dominates z in S .

For each coalition S , we define the *set of imputations dominated in S* by

$$D(S) := \{z \in I(v) \mid \text{there exists } y \in I(v) \text{ s.t. } y \text{ dominates } z \text{ in } S\}.$$

We say that an imputation $x \in I(v)$ is *undominated* if $x \in I(v) - \bigcup_{S \subseteq N: S \neq \emptyset} D(S)$, that is, if x is not dominated in any coalition S .

Clearly $D(\{i\}) = \emptyset$ for any singleton $\{i\}$, since for any $y, z \in I(v)$, $y_i \geq v(i)$ and $z_i \geq v(i)$ by individual rationality. For y to dominate z in $\{i\}$, we have $y_i \leq v(i)$ so $y_i = v(i)$ and $y_i > z_i \geq v(i)$, which would be a contradiction.

Likewise, $D(N) = \emptyset$. Any $y, z \in I(v)$ has $y(N) = z(N) = v(N)$. Hence if $y_i > z_i$ for some $i \in N$ then there must be some $j \in N$ such that $y_j < z_j$, and so y cannot dominate z in N .

Domination gives rise to a set-valued solution concept, the *domination core*, which is simply the set of all undominated imputations. This relates closely to the *core*, the set of imputations that cannot be improved on by any coalition.

Definition 95 (Core).

- (a) *Domination core*. In a game (N, v) , the *domination core* or D-core is the set

$$DC(v) := I(v) - \bigcup_{S \subseteq N: S \neq \emptyset} D(S).$$

- (b) *Core*. In a game (N, v) , the *core* is the set

$$C(v) := \{x \in I(v) \mid x(S) \geq v(S) \text{ for all nonempty } S \subseteq N\}.$$

This definition of the core is specialized to TU-games. More generally, we say that a coalition S *blocks* an action a_N of the grand coalition if there is some action a_S of S that all members of S prefer to a_N . The core is the set of actions of the grand coalition that are not blocked by any coalition. We say an action in the core is *stable*, and *unstable* otherwise. In these terms, in a TU-game (N, v) , a coalition S blocks the allocation $x \in I(v)$ iff $x(S) < v(S)$.

In general, the core may be empty. If the core is empty, we might want to find an approximation:

Definition 96 (ϵ -core). In a game (N, v) , the ϵ -core is the set $\{x \in I(v) \mid x(S) \geq v(S) - \epsilon \text{ for all nonempty } S \subseteq N\}$.

The ϵ -core has the interpretation that it is the set of stable allocations if any coalition has to pay a cost ϵ in order to block an allocation.

Proposition 72. For any TU-game (N, v) , $C(v) \subseteq DC(v)$.

Proof. Consider any $x \in I(v)$ s.t. $x \notin DC(v)$. Then there is a $y \in I(v)$ and a coalition $S \neq \emptyset$ s.t. y dominates x in S . Thus $v(S) \geq y(S) > x(S)$, implying $x \notin C(v)$. \square

Example 66 (Game with empty core). Consider the shepherd's problem in Example 63(b). We can model the situation as a game (N, v) where $N = \{1, 2, 3\}$ and $v(S) = 1$ if $|S| \geq 2$ and $v(S) = 0$ otherwise.

For any imputation $x \in I(v)$, we have $x \geq 0$ and $x_1 + x_2 + x_3 = 1$. Suppose $x_i < \frac{1}{3}$. Since $\sum_j x_j = 1$, we cannot have that $x_j \geq \frac{2}{3}$ for both $j \neq i$. Fix $j \neq i$ with $x_j < \frac{2}{3}$ and let $S = \{i, j\}$. Consider an imputation y with $y_i = \frac{1}{3}$ and $y_j = \frac{2}{3}$. Then y dominates x in S . Thus, for x to be in the core, we need $x_i \geq \frac{1}{3}$ for all i , and thus the only imputation that can lie in the core is $x = (1/3, 1/3, 1/3)$. But this is dominated in $\{1, 2\}$ by $x = (1/2, 1/2, 0)$. Hence the core is empty.

Intuitively, if any two shepherds come to an agreement about how to split the water, cutting out the third shepherd, then the third shepherd can make a better offer to one of the two shepherds and the agreement breaks down. And if all three shepherds came to an agreement, then any two of the shepherds could increase their water rights by cutting out the other shepherd. The shepherds stay squabbling and the poor sheep stay thirsty.

Under a fairly mild condition (implied by superadditivity), the core and domination-core coincide:

Theorem 32. Consider any TU-game (N, v) . If

$$v(N) \geq v(S) + \sum_{i \notin S} v(i) \quad \text{for all } S \in 2^N - \{\emptyset\},$$

then $C(v) = DC(v)$.

Proof. Since $C(v) \subseteq DC(v)$ by Proposition 72, we need only show $DC(v) \subseteq C(v)$ if the condition holds.

Lemma 23. Suppose the condition of Theorem 32 holds. Then if $x \in I(v)$ and $x(S) < v(S)$ for some nonempty coalition S , then there is a $y \in I(v)$ such that y dominates x in S .

Proof. If $i \in S$, define $y_i := x_i + |S|^{-1}(v(S) - x(S))$. If $i \notin S$, define $y_i := v(i) + (v(N) - v(S) - \sum_{j \notin S} v(j))|N - S|^{-1}$. Then $y \in I(v)$, where $y_i \geq v(i)$ follows from the condition of the theorem. Clearly, y dominates x in S since $y_i \geq x_i$ for all $i \in S$. \square

Now, suppose $x \in DC(v)$. Then there is no $y \in I(v)$ s.t. y dominates x . From the lemma, it follows that $x(S) \geq v(S)$ for all $S \in 2^N - \{\emptyset\}$. Hence $x \in C(v)$. We thus prove that $DC(v) \subseteq C(v)$. \square

Note that the condition of this theorem is implied by superadditivity.

Proposition 73. *If (N, v) is a game with a nonempty core, then $v(N) \geq v(S) + \sum_{i \notin S} v(i)$ for all $S \in 2^N - \{\emptyset\}$, and thus $C(v) = DC(v)$.*

Proof. If $C(v) \neq \emptyset$, then there exists some imputation x s.t. $x(S) \geq v(S)$ for all nonempty coalitions S . Fix some (nonempty) coalition S and define $T := N - S$. Since x is an imputation, $x(N) = v(N)$ and $x_i \geq v(i)$ for all $i \in N$, so $x(T) \geq \sum_{i \in T} v(i)$. Since $x(S) \geq v(S)$, it follows that

$$v(N) = x(N) = x(S) + x(T) \geq v(S) + \sum_{i \in T} v(i).$$

\square

Definition 97 (Simple game).

- (a) *Simple game.* A game (N, v) is a *simple game* if $v(N) = 1$ and $v(S) \in \{0, 1\}$ for all $S \in 2^N - \{\emptyset\}$.

In a simple game, we call a coalition S a *winning coalition* if $v(S) = 1$ and a *losing coalition* if $v(S) = 0$. We say a winning coalition S is *minimal* if every nonempty proper subset $S' \subseteq S$ is losing.

- (b) *Dictator.* In a simple game, a player i is called a *dictator* if $v(S) = 1$ iff $i \in S$.
- (c) *Veto player.* In a simple game, a player i is called a *veto player* if i is a member of all winning coalitions. The *set of veto players* is given by

$$\text{veto}(v) = \bigcap \left\{ S \in 2^N \mid v(S) = 1 \right\}.$$

Simple games arise in many real-world contexts:

- *United Nations Security Council.* The Security Council has 15 members, 5 of which are permanent members who are veto players. Excluding abstentions for simplicity, the set of winning coalitions is the set of coalitions involving all 5 permanent members and at least 8 members in total.
- *Treaty or contract negotiations.* The agreement of a treaty between governments or a contract between two or more parties, requires unanimity among parties. In these settings, typically all parties to the treaty or contract are veto players.

- *Legislatures.* In most legislatures, bills typically pass if a majority of members vote in favour. Some player(s) (e.g. the Speaker, certain government members) may have veto power via agenda-setting powers, and in some legislatures, rules – such as those governing the filibuster in the US Senate (as of 2022) – may give conditional veto powers to all members.

Example 67.

- (a) *Dictator game.* Fix $i \in N$. The *dictator game* δ_i is the simple game with characteristic function

$$\delta_i(S) = \begin{cases} 1 & \text{if } i \in S, \\ 0 & \text{otherwise.} \end{cases}$$

The set of imputations is $I(\delta_i) = \{e^i\}$, where e^i is the matrix with i th entry 1 and all other entries 0. Furthermore, $\text{veto}(\delta_i) = \{i\}$ and $C(\delta_i) = DC(\delta_i) = \{e^i\}$.

- (b) *Majority game.* A n -player majority game has characteristic function

$$v(S) = \begin{cases} 1 & \text{if } |S| \geq \lceil n/2 \rceil, \\ 0 & \text{otherwise.} \end{cases}$$

With three players, $v(S) = 1$ iff $|S| \in \{2, 3\}$. We have

$$\text{veto}(v) = \{1, 2\} \cap \{2, 3\} \cap \{1, 3\} \cap \{1, 2, 3\} = \emptyset,$$

and

$$C(v) = DC(v) = \emptyset.$$

- (c) *T-unanimity game.* Let T be a nonempty coalition. The T -unanimity game is the simple game with characteristic function

$$u_T(S) = \begin{cases} 1 & \text{if } T \subseteq S, \\ 0 & \text{otherwise.} \end{cases}$$

We have $\text{veto}(u_T) = T$ and

$$C(u_T) = DC(u_T) = \text{co}\{e^i \mid i \in T\}.$$

A simple game has a nonempty core iff it has veto players:

Theorem 33. *Let (N, v) be a simple game. Then*

- $C(v) = \text{co}\{e^i \in \mathbb{R}^n \mid i \in \text{veto}(v)\};$
- If $\text{veto}(v) = \emptyset$ and $\{i \in N \mid v(i) = 1\} = \{k\}$ then $C(v) = \emptyset$ and $DC(v) = \{k\}$. Otherwise, $DC(v) = C(v)$.

Proof. Suppose $i \in \text{veto}(v)$. Let $S \in 2^N - \{\emptyset\}$. If $i \in S$, then $e^i(S) = 1 \geq v(S)$; else $e^i(S) = 0 = v(S)$. Clearly, $e^i(N) = 1 = v(N)$. Thus $e^i \in C(v)$. Since $C(v)$ is convex, it follows that $C(v) \supset \text{co}\{e^i \in \mathbb{R}^n \mid i \in \text{veto}(v)\}$.

Now, let $x \in C(v)$. To show $C(v) \subseteq \text{co}\{e^i \in \mathbb{R}^n \mid i \in \text{veto}(v)\}$, we need only prove that if $i \notin \text{veto}(v)$, then $x_i = 0$. Suppose otherwise, i.e. $x_i > 0$ for some $i \notin \text{veto}(v)$. Take S with $v(S) = 1$ and $i \notin S$ (were such an S not to exist, then i would be a veto player). Then $x(S) = x(N) - x(N - S) \leq 1 - x_i < 1$, contradicting $x \in C(v)$. This proves (i).

For (ii), if $\text{veto}(v) = \emptyset$ and $\{i \in N \mid v(i) = 1\} = \{k\}$ then $C(v) = \emptyset$ follows immediately from (i). Suppose instead that $\text{veto}(v)$ is nonempty. Then (i) implies that $C(v)$ is nonempty by (i). By Proposition 73, it follows that $C(v) = DC(v)$. \square

Lastly, convex games are guaranteed to have a nonempty core:

Theorem 34. *If (N, v) is a convex game, then it has a nonempty core.*

Proof. We claim the allocation x defined by $x_i = v(\{1, \dots, i\}) - v(\{1, \dots, i-1\})$ lies in the core. For every $i_1 < i_2 < \dots < i_j$, since v is convex we have that

$$\begin{aligned} \sum_{k=1}^j x_{i_k} &= \sum_{k=1}^j [v(\{i_1, \dots, i_k\}) - v(\{i_1, \dots, i_{k-1}\})] \\ &\geq \sum_{k=1}^j [v(\{i_1, \dots, i_{k-1}\} \cup i_k) - v(\{i_1, \dots, i_{k-1}\})] \\ &= v(\{i_1, \dots, i_j\}). \end{aligned}$$

\square

9.4.2 Stable sets

Another solution concept based on domination is the *stable set*, due to von-Neumann and Morgenstern (1944).

Definition 98 (Stable set). Let v be a game and let $A \subseteq I(v)$. The set A is called a *stable set* if

- (i) If $x, y \in A$ then x does not dominate y (*internal stability*);
- (ii) If $x \in I(v) - A$ then there exists a $y \in A$ that dominates x (*external stability*).

Example 68. Suppose v is a three person game, and the value of any multiperson coalition has value 1, but all single person coalition has value 0. Consider the set $A = \{x, y, z\}$ made up of three imputations: $x = (\frac{1}{2}, \frac{1}{2}, 0)$, $y = (\frac{1}{2}, 0, \frac{1}{2})$ and $z = (0, \frac{1}{2}, \frac{1}{2})$. We show this set is stable.

- (i) *Internal stability.* For each vector $a, b \in A$, $a_i > b_i$ for one i only. Hence a could only dominate b in a coalition consisting of one person. But $v(i) = 0$ for all i , and so $a(i) \not\geq v(i)$. Hence a does not dominate b in any coalition, so a does not dominate b .
- (ii) *External stability.* Consider any $h \in I(v)$ such that $h \notin A$. Consider any coalition $S = \{i, j\}$ of two people and let $\{k\}$ be the third person, excluded from S . Then $v(S) = 1$. Now there is $a \in A$ with $a_i = a_j = \frac{1}{2}$ and $a_k = 0$, so $a(S) = 1 = v(S)$. If $h_i < \frac{1}{2}$ and $h_j < \frac{1}{2}$ then we have that $a_i > h_i$ for all $i \in S$ and so a dominates h in S and thus a dominates h . Hence we need only consider the case where $h_i > \frac{1}{2}$ or $h_j > \frac{1}{2}$. Suppose wlog that $h_i > \frac{1}{2}$. Then $h_j < \frac{1}{2}$ and $h_k < \frac{1}{2}$, since $h(N) = 1$ and therefore $h(S) = h_i + h_j \leq 1$. Now consider the coalition $T = \{j, k\}$. There is a $b \in A$ s.t. $b_j = b_k = \frac{1}{2}$ and $b_i = 0$, so $b(T) = 1 = v(T)$. Since $h_j < \frac{1}{2}$, $b_j > h_j$ and since $h_k < \frac{1}{2}$, $b_k > h_k$. Thus b dominates h in T , so b dominates h .

We have thus shown that for any imputation $h \notin A$, one of the elements of A dominates h .

Hence A is stable. However, A is not a unique stable set. Indeed, if $c \in [0, \frac{1}{2})$ and $B = \{x \in I(v) \mid x_3 = c\}$ then B is stable (proof is a slight modification of the above).

In general, stable sets are non-unique. However, even if unique, a selection has to be made from the stable set – stability is a property of sets and not particular allocations. The core on the other hand is by definition a single set. Unlike stable sets, selection from the core can be plausibly motivated (e.g. the nucleolus).

The question of existence of stable sets is partially solved. First, note any essential simple game has a stable set, namely the set of imputations that assign nothing to players outside a minimal winning coalition:

Proposition 74. *Let v be a simple game and suppose S is a minimal winning coalition. Define $\Delta^S := \{x \in I(v) \mid x_i = 0 \text{ for all } i \notin S\}$. Then Δ^S is a stable set if it is nonempty.*

Proof. Since S is a minimal winning coalition, $v(S) = 1$ and $v(S') = 0$ for any proper subset $S' \subseteq S$. Suppose some $x \in \Delta^S$ dominates some $y \in \Delta^S$. Since $x_i = y_i = 0$ for all $i \notin S$, x can only dominate y in some coalition $T \subseteq S$. If T is a proper subset of S , then $v(T) = 0$. If $x(T) > 0$ then x cannot dominate y in T , and if $x(T) = 0$ then we must have $x_i = y_i = 0$ for all $i \in T$, and again x cannot dominate y . Finally, if $T = S$, then $x(T) = y(T) = v(T) = 1$. Since $\sum_{i \in S} x_i = \sum_{i \in S} y_i$, we must have that if $x_i > y_i$ for some $i \in S$ then $x_j < y_j$ for at least one $j \in S$. Hence x does not dominate y .

We have thus established internal stability. Now we turn to external stability. Consider any $y \notin \Delta^S$. By definition, $y(S) < 1$, and so $\epsilon := 1 - y(S) > 0$. Define, x by $x_i = 0$ for all $i \notin S$ and $x_i = y_i + \frac{1}{|S|}\epsilon$ for all $i \in S$. Then $x \in \Delta^S$, $x(S) = 1 = v(S)$ and $x_i > y_i$ for all $i \in S$, so x dominates y in S . Hence we have established external stability. \square

Example 69 (Stable sets in zero-one games). A *zero-one game* (N, v) is a game in which $v(i) = 0$ for all $i \in N$ and $v(N) = 1$. This is not necessarily a simple game, because in

games of three or more players, coalitions that are not the grand coalition but consist of more than one person can have worth not in $\{0, 1\}$.

Consider a three person zero-one game (N, v) where $N = \{1, 2, 3\}$, $v(i) = 0$ for all $i \in N$, $v(N) = 1$, and for every two-person coalition S , $v(S) = \alpha$ for some $\alpha \in [0, 1]$.

- (i) If $\alpha \geq \frac{2}{3}$, then the set

$$A = \left\{ (\lambda, \lambda, 1 - 2\lambda), (\lambda, 1 - 2\lambda, \lambda), (1 - 2\lambda, \lambda, \lambda) \mid \frac{\alpha}{2} \leq \lambda \leq \frac{1}{2} \right\}$$

is a stable set.

Each $x \in A$ has $x(N) = 1$. Consider any $x, y \in A$ and let λ_x and λ_y be the values of λ associated with each respectively. Since $x(i) = \lambda_x > 0 = v(i)$, x cannot dominate y in any singleton coalition. In the grand coalition, we also have that since $y(N) = x(N) = 1$ and $y_i < x_i$ for at least one $i \in N$, we must have $y_j > x_j$ for some $j \in N$, so x does not dominate y in N . Consider any two-person coalition $S = \{i, j\}$. We have several cases:

- (1) $x_i = \lambda_x$ and $y_i = \lambda_y$ (note the case where $x_j = \lambda_x$ and $y_j = \lambda_y$ is also covered because we can just switch i, j). Suppose x dominates y in S . We then require $x_i > y_i$ and thus $\lambda_x > \lambda_y \geq \frac{\alpha}{2}$. If $x_j = \lambda_x$ then we would have $x(S) > \alpha = v(S)$ and so x would not dominate y in S . Hence suppose $x_j = 1 - 2\lambda_x$. If $y_j = 1 - 2\lambda_y$ then because $\lambda_x > \lambda_y$, $y_j > x_j$ and so x does not dominate y in S . Hence $y_j = \lambda_y$. Since $\lambda_x > \frac{\alpha}{2}$, $1 - 2\lambda_x < 1 - \alpha \leq \frac{1}{3}$, and $\lambda_y \geq \frac{\alpha}{2} \geq \frac{1}{3}$. Hence $y_j > x_j$, yielding a contradiction.
- (2) $x_i = \lambda_x$, $x_j = 1 - 2\lambda_x$, $y_i = 1 - 2\lambda_y$, $y_j = \lambda_y$ (again, we can swap i, j labels to cover the converse case). Suppose x dominates y in S . Then $x_i = \lambda_x > 1 - 2\lambda_y = y_i$ and $x_j = 1 - 2\lambda_x > \lambda_y = y_j$. Rearranging these gives us $\lambda_x + 2\lambda_y > 1 > 2\lambda_x + \lambda_y$. Subtracting $\lambda_x + \lambda_y$ from both sides gives us $\lambda_y > \lambda_x \geq \frac{\alpha}{2}$ and so $2\lambda_x + \lambda_y > \frac{3}{2}\alpha \geq 1$, yielding a contradiction.

Thus A is internally stable. Next, we consider external stability. Consider any imputation $y \notin A$. Let $i = \arg \min_{k \in N} y_k$ and let $j = \arg \min_{k \in N - \{i\}} y_k$. Consider the two-person coalition $S = \{i, j\}$ and let $\{\ell\} = N - S$. Now, $y_i + y_j \leq \frac{2}{3}$, since if $y_i, y_j > \frac{2}{3}$ then $y_\ell < \frac{1}{3}$ and since $y_i < y_j < y_\ell$, it follows that $y_i, y_j < \frac{1}{3}$, yielding a contradiction. Hence $y(S) \leq \alpha$. Consider the $x \in A$ with $x_i = x_j = \frac{\alpha}{2}$. Then we have $x_i > y_i$, $x_j > y_j$ and so x dominates y . This completes the proof.

- (ii) If $\alpha < \frac{2}{3}$ then A is not stable. We show that external stability does not hold here (A remains internally stable). Consider $y = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. Note $y \in A$ iff $\alpha = \frac{2}{3}$, and by hypothesis $\alpha < \frac{2}{3}$. For any singleton $\{i\}$, we have $y_i = \frac{1}{3} > 0 = v(i)$, so there is no $x \in A$ dominating y in $\{i\}$. For the grand coalition we have $y(N) = 1 = v(N)$ but for any $x \in A$, $x_i < \frac{1}{3}$ for some i since $y \notin A$ and $\sum_{i=1}^3 x_i = 1$. Finally, in any two person coalition S , $y(S) = \frac{2}{3} > \alpha = v(S)$, so there is no $x \in A$ dominating y in S . It follows that y is undominated. Note if $\alpha > \frac{2}{3}$ then y would be dominated by $(\alpha/2, \alpha/2, 1 - \alpha)$.

(iii) If $\alpha \leq \frac{1}{2}$ then the core is the unique stable set.

Proposition 75. *Let (N, v) be a TU-game. Then*

- (i) *The D-core of v is a subset of any stable set;*
- (ii) *If A and B are stable and $A \neq B$ then $A \not\subseteq B$;*
- (iii) *If the D-core of v is a stable set, then it is the unique stable set of v .*

Proof. (i) Recall the D-core is the set of all undominated imputations. Suppose there is some undominated imputation $x \in I(v)$ and some stable set A of v s.t. $x \notin A$. Since x is not dominated in any coalition, there is no $y \in A$ s.t. y dominates x , contradicting the external stability of A .

(ii) Suppose $A \subseteq B$. Then there is some $x \in B$ s.t. $x \notin A$. For A to satisfy external stability, we must have that some $y \in A$ dominates x . Yet since $x, y \in B$, internal stability of B requires that y does not dominate x , yielding a contradiction.

(iii) Suppose $DC(v)$ and A are stable and $A \neq DC(v)$. Since $DC(v)$ consists of all undominated imputations, $A \supset DC(v)$, for were there some $x \in DC(v)$ s.t. $x \notin A$, then there would be no $y \in A$ s.t. y dominates x , which would violate external stability. But by (ii), if $A \supset DC(v)$ and A is stable then either $DC(v) = A$ or $DC(v)$ is not stable, both of which contradict the hypotheses. □

9.4.3 Balanced games

Definition 99 (Balanced games). Let $N = \{1, \dots, n\}$.

(a) *Characteristic vector.* Given a coalition $S \subseteq N$, the vector e^S defined by

$$e_i^S = \begin{cases} 1 & \text{if } i \in S, \\ 0 & \text{if } i \in N - S, \end{cases}$$

is called the *characteristic vector* for S .

(b) *Balanced map.* Call a map $\lambda : 2^N - \{\emptyset\} \rightarrow \mathbb{R}_+$ a *balanced map* if

$$\sum_{S \in 2^N - \{\emptyset\}} \lambda(S) e^S = e^N.$$

(c) *Balanced coalition.* A collection B of nonempty coalitions is called a *balanced coalition* if there exists a balanced map λ such that

$$B = \left\{ S \in 2^N - \{\emptyset\} \mid \lambda(S) > 0 \right\}.$$

- (d) *Balanced game.* A game (N, v) is called a *balanced game* if for each balanced map λ , we have

$$\sum_S \lambda(S) v(S) \leq v(N).$$

- (e) *Totally balanced game.* A game (N, v) is called *totally balanced* if for every coalition $S \subseteq N$, the restriction $v|_S$ of v to S is a balanced game.

To give balancedness an intuitive interpretation, suppose each player is endowed with one unit of time to distribute over the various coalitions of which she is a member. The distribution is balanced if it corresponds to a balanced map λ , where we interpret $\lambda(S)$ as the length of time for which the coalition S exists. Balancedness of the map λ requires that each player spends their entire time endowment of 1 over the coalitions. Under this interpretation, a balanced game is a game in which operating the grand coalition the entire time is at least as productive as any balanced distribution of the time endowment involving smaller coalitions, where we interpret the worth of a coalition as productivity. Intuitively, in any balanced game, it is natural to form the grand coalition. In fact, the Bondareva-Shapley theorem demonstrates that the core of a game is nonempty iff the game is balanced:

Theorem 35 (Bondareva-Shapley). *Let (N, v) be a TU-game. Then the following are equivalent:*

- (i) *The core $C(v)$ is nonempty;*
- (ii) *(N, v) is a balanced game.*

Proof. $C(v)$ is nonempty iff

$$v(N) = \min \left\{ \sum_{i=1}^n x_i \mid x \in \mathbb{R}^N, \ x(S) \geq v(S) \text{ for all } S \in 2^N - \{\emptyset\} \right\}.$$

By Theorem 54 (duality theorem), this equality holds iff

$$v(N) = \max \left\{ \sum_S \lambda(S) v(S) \mid \sum_S \lambda(S) e^S = e^N, \ \lambda \geq 0 \right\},$$

if we take the matrix with characteristic vectors e^S as columns for A , define $c = e^N$ and let b be the vector of coalitional worths. This holds iff $\sum_S \lambda(S) v(S) \leq v(N)$, and so (i) is equivalent to (ii). \square

9.4.4 Implementing the core

The core tells us which imputations – efficient allocations of the worth of the grand coalition – we can expect might be implemented in a cohesive TU-game. For any unstable

imputation – not in the core – some subset of players can form a coalition and obtain an allocation that each strictly prefers. As with Nash bargaining, we can map this into a non-cooperative setting – there must be some process by which such a blocking coalition can form over time. Perry & Reny (1994) introduce a non-cooperative dynamic game that implements any and all core allocations.

Let (N, v) be a superadditive TU-game with n players. Suppose time is continuous and that players do not discount.⁶⁴ At time t , any player i has several options:

- i can *remain quiet* (q).
- i can *make a proposal* (x, S) , which consists of a coalition S and an allocation x , which has $x(N) \leq v(S)$ (i.e. $\sum_{j \in N} x_j \leq v(S)$). If i makes a proposal at t , then this proposal becomes active, replacing any previous proposal. A proposal remains active until a new proposal is made.
- i can *accept* (a) the current proposal, if one is on the table. Note if i makes a proposal that does not necessarily mean she accepts it.
- i can *announce she is leaving* (ℓ), in which case she leaves and then consumes.

A proposal (x, S) becomes *binding* if all members of S accept it while it is active. If a proposal (x, S) becomes binding, we say S has *formed and accepted* (x, S) . Once this happens, any player in S can either leave and consume x_i , or remain at the bargaining table, in the hope of a more favourable proposal. If S has formed and accepted (x, S) and some player i in S leaves to consume x_i , then all players j in S are forced to leave and consume x_j . Any player i in S thus guarantees herself a payoff of at least x_i . If a player i never leaves to consume, she receives payoff $d_i \leq v(i)$.

A binding proposal is no longer active, and thus we can potentially have several binding proposal and an active proposal on the table together at any point in time. However, if making a proposal that includes some member of a binding proposal, we assume the proposal must include all members of that binding proposal, so that all members of the coalition in the binding proposal agree to annul it before it can be superceded. If all members of a binding proposal agree to an active proposal, then the binding proposal is annulled and removed from consideration. A player can only be associated with one binding proposal at a time.

For the game to be well-defined, we need to make a technical assumption. Assume that at every time t , for every history up to t , and for every vector of actions available to the players at t , players strategies are such that they cannot make a proposal, accept a proposal or leave (i.e. take a non-quiet action) just before or just after t . Formally, we assume there exists some $\epsilon > 0$ such that players cannot, according to their strategies, take a non-quiet action in the intervals $(t - \epsilon, t)$ and $(t, t + \epsilon)$. This ensures that for any history up to t , the players' strategies induce a unique continuation path and that payoffs are well-defined. Moreover, if players can take non-quiet actions instantaneously, then

⁶⁴Perry & Reny (1994) have a nice discussion of why they think continuous time is the natural setting here.

no player can convince them not to leave by placing an appropriately timed blocking proposal. Since ϵ is not bounded away from 0, players can react arbitrarily quickly.

A history h^t up to time t in this game is an n -tuple $h^t = (h_1^t, \dots, h_n^t)$, where each h_i^t is a function $h_i^t : [0, t) \rightarrow A_i$, where A_i is the space of player i 's actions, i.e. $A_i = \{q, a, \ell\} \cup P$ with $P := \{(x, S) \mid S \subseteq N, x_j \geq v(j) \text{ for all } j \in S, x(N) \leq v(S)\}$ being the set of possible proposals. Let:

- $H(t)$ denote the set of all histories up to time t with $t \geq 0$, where $H(0) = \emptyset$;
- $H = \bigcup_{t=0}^{\infty} H(t)$ denote the set of all histories;
- $p(h) \in P$ denote the current active proposal under history $h \in H(t)$;
- $\tau(h)$ denote the length of time that has elapsed since $p(h)$ was proposed under $h \in H(t)$;
- $N(h) \subseteq N$ denote the set of players who have not yet left under $h \in H(t)$;
- $A(h) \subseteq N(h)$ denote the set of players who have accepted $p(h)$ under $h \in H(t)$;
- $\Pi(h)$ denote the set of current binding proposals among the players in $N(h)$, and
- $\eta(h) = (p(h), \tau(h), N(h), A(h), \Pi(h))$ denotes the state under h .

A strategy $s_i : H \rightarrow \{q, a, \ell\} \cup P$ for a player i specifies an action $s_i(h)$ for each history $h \in H$, and we write $s(h) = (s_1(h), \dots, s_n(h))$. We let S_i denote the set of strategies of player i . The solution concept here is stationary subgame perfect equilibrium (SSPE):

Definition 100 (Stationary subgame perfect equilibrium). A strategy profile $s^* = (s_1^*, \dots, s_n^*)$ is a *stationary subgame perfect equilibrium* if

- (i) for every player $i \in N$ and every history $h \in H$, $u_i(s_i^*, s_{-i}^* \mid h) \geq u_i(s_i, s_{-i}^* \mid h)$ for all $s_i \in S_i$ (subgame perfection), and
- (ii) for all $h, h' \in H$, if $\eta(h) = \eta(h')$ then $s^*(h) = s^*(h')$, that is, players' strategies depend only on the state (stationarity).

Proposition 76. *Every stationary subgame perfect equilibrium allocation of the game is in the core of (N, v) .*

Proof. Suppose towards contradiction that x is an SSPE allocation but that x is not in the core $C(v)$. Suppose s is an SSPE that supports x as an outcome. We must have $x_i \geq v(i)$ for every player i . Since $x \notin C(v)$, there is some proposal (y, S) with $y_i > x_i$ for all $i \in S$. Wlog, let $S = \{1, \dots, k\}$.

At time t , consider any history $h \in H(t)$ s.t. $N(h) = N$, $\Pi(h) = \emptyset$, players $1, \dots, k-1$ have accepted (y, S) under h , and $s_i(h) = q$ for all $i \in N$. Now, under s , player k will accept (y, S) before a new proposal is made. To see this, suppose, under s , a proposal (z, T) is made before player k accepts (y, S) . The state then becomes $\bar{\eta} = ((z, T), 0, \emptyset, N, \emptyset)$. Since k could have accepted (y, S) and guaranteed herself a payoff

$y_k > x_k$ but chose to let it be replaced by (z, T) , the continuation for k must have payoff $w_k \geq y_k$. By stationarity, whenever $\eta(h) = \bar{\eta}$, the continuation payoff for k under s is $w_k \geq y_k > x_k$. But since s is an SSPE, this yields a contradiction, because player k can always ensure the state $\bar{\eta}$ is reached by proposing (z, T) herself near enough to $t = 0$. Thus under s , player k accepts (y, S) before any new proposal is made.

Now at time t , consider any history $h \in H(t)$ s.t. $N(h) = N$, $\Pi(h) = \emptyset$, players $1, \dots, k-2$ have accepted (y, S) under h , and $s_i(h) = q$ for all $i \in N$. Now, under s , both players $k-1$ and k will accept (y, S) before a new proposal is made. To see this, note that $k-1$ can guarantee himself payoff at least $y_{k-1} > x_{k-1}$ by accepting the proposal at time t , because by the previous paragraph, this will induce player k to accept (y, S) so the proposal becomes binding. Suppose instead that (y, S) is not accepted by $k-1$ and k in the continuation of play after h . If under s , no new proposal is made in the continuation, then player $k-1$ receives $d_{k-1} < x_{k-1} < y_{k-1}$, so this cannot be optimal. If under s , a new proposal (z, T) is made, then we arrive at a contradiction by a similar argument to the previous paragraph: player $k-1$ must receive a continuation payoff of at least y_{k-1} , and by stationarity, $k-1$ can ensure this by making the proposal himself close enough to $t = 0$.

Proceeding inductively, we can conclude that at any time t and any history $h \in H(t)$ with $N(h) = N$, $\Pi(h) = \emptyset$ player 1 accepting (y, S) under h , and $s_i(h) = q$ for all players i , the continuation has every player in the coalition S accepting the proposal (y, S) . But then x cannot be an SSPE outcome, because player 1 can propose (y, S) and accept it at a time sufficiently close to $t = 0$, ensuring herself a payoff $y_1 > x_1$. \square

Proposition 77. *If (N, v) is a totally balanced game, then any imputation in its core can be supported as a stationary subgame perfect equilibrium.*

Proof. See Perry & Reny (1994). Repeating the proof would take several pages. \square

9.4.5 Competitive equilibrium and the core

There is a close relationship between the core and general equilibrium theory. In particular, competitive equilibria lie in the core. This is easiest to see in the case of a pure-exchange economy, though the result extends to production economies provided we envisage blocking opportunities appropriately.

Definition 101 (Economy).

- (a) *Pure exchange economy.* A pure exchange economy is a tuple $\mathcal{E} = (H, (X_i, u_i, \omega_i)_{i \in H})$, where
 - (i) H is a finite set of n consumers;
 - (ii) $X_i \subseteq \mathbb{R}^k$ is a consumption set for consumer i , and k is the number of commodities;
 - (iii) $u_i : X_i \rightarrow \mathbb{R}$ is a utility function for consumer h ;
 - (iv) $\omega_i \in X_i$ is consumer i 's endowment vector.

(b) *Private ownership economy.* A private ownership economy is a tuple

$$\mathcal{P} = (H, F, (X_i, u_i, \omega_i, s_i)_{i \in H}, (Y_j)_{j \in F})$$

where

- (i) $H, X_i \subseteq \mathbb{R}^k$, $u_i : X_i \rightarrow \mathbb{R}$ and $\omega_i \in X_i$ are defined as above;
- (ii) F is a finite set of m firms;
- (iii) $Y_j \subseteq \mathbb{R}^k$ is a *production set* for firm j ;
- (iv) $s_i = (s_{ij})_{j \in F}$ is consumer i 's vector of *shareholdings*. We assume $s_{ij} \in \mathbb{R}_+$ and $\sum_{j=1}^m s_{ij} = 1$ for each firm $j \in F$.

For each consumer $i \in H$, we assume i 's consumption set is $X_i = \mathbb{R}_+^k$ and i 's utility function u_i is continuous. For each firm $j \in F$, we assume j 's production set Y_j is nonempty and closed. We define an economy's *total endowment* to be $\bar{\omega} = \sum_{i=1}^n \omega_i$.

In a pure-exchange economy \mathcal{E} , an allocation is a vector $x = (x_1, \dots, x_n) \in X_1 \times \dots \times X_n$. An allocation x is *feasible* if $\sum_{i=1}^n x_i = \bar{\omega}$.

In a private ownership economy \mathcal{P} , an allocation is a vector $(x, y) = (x_1, \dots, x_n, y_1, \dots, y_m) \in X_1 \times \dots \times X_n \times Y_1 \times \dots \times Y_m$, consisting of a *consumption allocation* x and a vector of *production plans* y . An allocation (x, y) is *feasible* if $\sum_{i=1}^n x_i = \bar{\omega} + \sum_{j=1}^m y_j$.

Definition 102 (Competitive equilibrium).

- (a) In a pure-exchange economy \mathcal{E} , a *competitive equilibrium* is a price vector $p \in \mathbb{R}^k$ and an allocation x^* such that

- (i) for each consumer $i \in H$, x_i^* is a solution to

$$\max_{x_i \in X_i} u_i(x_i) \quad \text{subject to} \quad p \cdot x \leq p \cdot \omega_i;$$

- (ii) markets clear, that is, $\sum_{i=1}^n x_i = \bar{\omega}$.

- (b) In a private ownership economy \mathcal{P} , a *competitive equilibrium* is a price vector $p \in \mathbb{R}^k$ a consumption allocation x^* , and a vector of production plans y^* such that

- (i) for each consumer $i \in H$, $x_i^* \in X_i$ is a solution to

$$\max_{x_i \in X_i} u_i(x_i) \quad \text{subject to} \quad p \cdot x \leq p \cdot \omega_i + \sum_{j=1}^m s_{ij} p \cdot y_j;$$

- (ii) for each firm $j \in F$, $y_j^* \in Y_j$ is a solution to

$$\max_{y_j \in Y_j} p \cdot y_j;$$

- (iii) markets clear, that is, $\sum_{i=1}^n x_i = \bar{\omega} + \sum_{j=1}^m y_j$.

This is also called *Walrasian equilibrium*. Existence conditions for Walrasian equilibrium are very well-known. A good overview is e.g. Chapter 14.4 of Kreps (2013) or Chapter 17.C in MWG (1995).

Focus for now on a pure-exchange economy \mathcal{E} . The pure-exchange economy has an interpretation as a cooperative game, where the set of actions A_S of each coalition $S \subseteq H$ is the set of all allocations of their total endowment $\sum_{i \in S} \omega_i$:

$$A_S = \left\{ x \in X_1 \times \cdots \times X_n \mid \sum_{i \in H} x_i \leq \sum_{i \in S} \omega_i \right\}.$$

The payoff to a consumer $i \in S$ from action x is $u_i(x_i)$. Note that this is not, in general, a TU-game.

Recall that the core is the set of allocations that cannot be blocked by any coalition. In the context of a pure-exchange economy, we say that a coalition S can block a feasible allocation x if there exists an allocation x' such that $\sum_{i=1}^n x'_i \leq \sum_{i \in S} \omega_i$ and $u_i(x'_i) > u_i(x_i)$ for every consumer $i \in S$. That is, S can block an allocation x if they can redistribute their own endowments among themselves in such a way that they are all better off than under x .

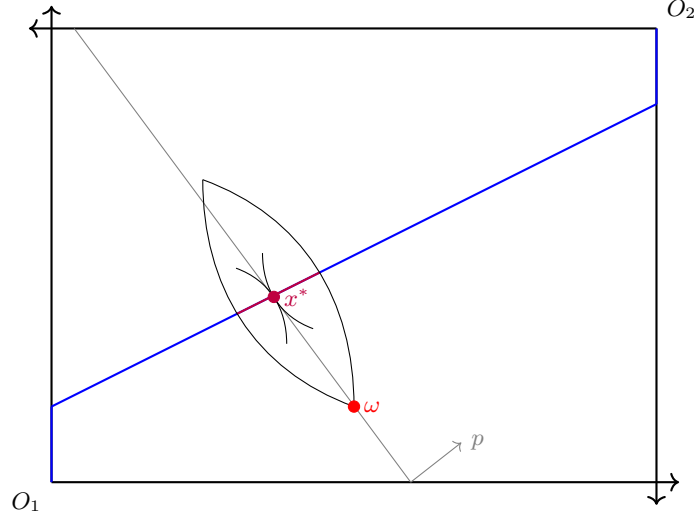
Theorem 36. *Any competitive equilibrium of a pure-exchange economy \mathcal{E} with locally insatiable consumers is in the core.*

Proof. Fix a pure-exchange economy \mathcal{E} , and let (p, x^*) be a competitive equilibrium of \mathcal{E} , where p is a price vector and x^* is an allocation. Suppose there is some coalition S that can block x^* . Then there exists an allocation x s.t. $\sum_{i=1}^n x_i \leq \sum_{i \in S} \omega_i$ and $u_i(x_i) > u_i(x_i^*)$ for all $i \in S$. But by definition of competitive equilibrium, this implies that x_i is unaffordable for each consumer $i \in S$, since $x_i^* = \max u_i(x_i)$ subject to the consumer's budget constraint. Hence $p \cdot x_i > p \cdot \omega_i$ for all $i \in S$. Thus $p \cdot \sum_{i \in H} x_i \geq p \cdot \sum_{i \in S} x_i > p \cdot \sum_{i \in S} \omega_i$, so $\sum_{i=1}^n x_i > \sum_{i \in S} \omega_i$, yielding a contradiction. \square

The converse is not true – allocations in the core are not necessarily competitive equilibria, as the following Edgeworth box economy illustrates.⁶⁵

Example 70 (The core in an Edgeworth box economy). Consider the following Edgeworth box economy.

⁶⁵However, because any allocation in the core is Pareto optimal, under the conditions of the second theorem of welfare economics, any core allocation can be supported as a quasi-equilibrium.



The initial endowment point is ω , and the competitive equilibrium allocation is x^* . The Pareto set is shown in blue, and the purple line, the segment of the Pareto set enclosed between the indifference curves of the two consumers at ω , is the core. This is of course identical to the contract curve. While the competitive equilibrium is a single point, the core is a continuum of points.

A similar result to Theorem 36 holds for private ownership economies, if we conceive of blocking in the right way. The problem is that it is not obvious how to allocate production to coalitions that are smaller than the grand coalition. For example:

- (a) *Peasantry.* We could imagine that coalitions have no control over production at all, and so can only redistribute their total endowment. The set of actions A_S of coalition S is thus the same as in the pure-exchange economy case, namely,

$$A_S = \left\{ x \in X_1 \times \cdots \times X_n \mid \sum_{i \in H} x_i \leq \sum_{i \in S} \omega_i \right\}.$$

- (b) *Revolution!* Another option is that any coalition seizes the means of production and controls the productive technology of all firms, so the set of actions available to coalition S is

$$A_S = \left\{ x \in X_1 \times \cdots \times X_n \mid \sum_{i \in H} x_i \leq \sum_{i \in S} \omega_i + \sum_{j=1}^m y_j \text{ for some } y \in Y_1 \times \cdots \times Y_m \right\}.$$

- (c) *Unanimous control.* We could also imagine that coalitions can only employ the productive capacity of firms they fully control. For a coalition S , let $F(S)$ denote the set of firms $j \in F$ for which $\sum_{i \in S} s_{ij} = 1$. The set of actions available to S is

$$A_S = \left\{ x \in X_1 \times \cdots \times X_n \mid \sum_{i \in H} x_i \leq \sum_{i \in S} \omega_i + \sum_{j \in F(S)} y_j \text{ for some } y \in Y_1 \times \cdots \times Y_m \right\}.$$

- (d) *Majority control.* Similarly, we could imagine coalitions can employ the productive of firms they have majority control over. For a coalition S , redefine $F(S)$ to be the set of firms $j \in F$ for which $\sum_{i \in S} s_{ij} > \frac{1}{2}$. Then A_S is as above.
- (e) *Shareholder-scaled production.* We might also imagine that coalitions can use a copy of each firm scaled by their own shareholdings. For a coalition S , Let $s_{Sj} = \sum_{i \in S} s_{ij}$. Then the

$$A_S = \left\{ x \in X_1 \times \cdots \times X_n \mid \sum_{i \in H} x_i \leq \sum_{i \in S} \omega_i + \sum_{j=1}^m s_{Sj} y_j \text{ for some } y \in Y_1 \times \cdots \times Y_m \right\}.$$

It turns out that if $0 \in Y_j$ for each firm $j \in F$ and consumers have locally nonsatiated preferences, then (a), (c) and (e) ensure the competitive equilibrium is in the core for a private ownership economy \mathcal{P} . Moreover, if firms have constant returns to scale technologies, then all of (a)-(e) ensure the competitive equilibrium of \mathcal{P} is in the core. See Kreps (2013), Proposition 15.9. Kreps also has a detailed discussion of how the core converges to the competitive equilibrium in N -replica economies as they grow large enough, a result due to Debreu and Scarf (1963).

9.5 Shapley value

Unlike the previous solution concepts, the Shapley value is a point-valued solution concept. It has several different characterizations. First, Shapley's definition:

Definition 103 (Shapley value). Given a TU game (N, v) , let $\sigma : N \rightarrow N$ denote a permutation of the player set, and let $\Pi(N)$ denote the set of all permutations of the player set. For a given permutation $\sigma \in \Pi(N)$, we define the *set of predecessors of i in σ* by

$$P_\sigma(i) = \{j \in N \mid \sigma^{-1}(j) < \sigma^{-1}(i)\}$$

for all $i \in N$, and we define the *marginal vector m^σ* of σ by

$$m_i^\sigma = v(P_\sigma(i) \cup \{i\}) - v(P_\sigma(i)).$$

The *Shapley value* $\Phi(v)$ of game v is then the average of the marginal vectors of the game, that is,

$$\Phi(v) = \frac{1}{n!} \sum_{\sigma \in \Pi(N)} m^\sigma.$$

Unpacking the definition, imagine that players enter a room in the order according to permutation σ . The i th element of the marginal vector m^σ tells us the marginal contribution of player i entering the room and joining the existing coalition of all players who already entered. The Shapley value is the average across all possible orderings with which players might arrive. Equivalently, suppose players enter the room at random,

and are paid their marginal contribution to the coalition formed from the players already in the room. Then the Shapley value gives the expected payoff for each player.

We can also rewrite the Shapley value as

$$\begin{aligned}\Phi_i(v) &= \frac{1}{n!} \sum_{\sigma \in \Pi(N)} [v(P_\sigma(i) \cup \{i\}) - v(P_\sigma(i))] \\ &= \sum_{S: i \notin S} \frac{|S|!(n-1-|S|)!}{n!} [v(S \cup \{i\}) - v(S)] \\ &= \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [v(S \cup \{i\}) - v(S)]\end{aligned}$$

for each player i . The second and third lines follow because for any given coalition S not containing i , $\frac{1}{n} \binom{n-1}{|S|}^{-1} = \frac{|S|!(n-1-|S|)!}{n!}$ is the fraction of permutations σ for which the set of predecessors of i in σ is S . Probabilistically, then, we can also interpret the Shapley value as follows. Suppose we choose a number $k \in \{0, 1, \dots, n-1\}$ uniformly at random. Given a draw k , we draw a set of size k from the set $N - \{i\}$, with each set S such that $|S| = k$ being drawn with probability $\binom{n-1}{k}^{-1}$. Player i then receives payoff $v(S \cup \{i\}) - v(S)$. The Shapley value gives the expected payoff for i under this procedure.

Example 71 (Glove game). In a glove game, players each have either a left-hand or right-hand glove, and a coalition receives payoff 1 if it can form a left-right pair and 0 otherwise. Suppose $N = \{1, 2, 3\}$, that players 1 and 2 have right-hand gloves and player 3 has a left-hand glove. Then

$$v(S) = \begin{cases} 1 & \text{if } S \in \{\{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}, \\ 0 & \text{otherwise.} \end{cases}$$

It is easy to compute the marginal contributions for each player for each permutation:

| σ | m_1^σ | m_2^σ | m_3^σ |
|----------|--------------|--------------|--------------|
| 1,2,3 | 0 | 0 | 1 |
| 1,3,2 | 0 | 0 | 1 |
| 2,1,3 | 0 | 0 | 1 |
| 2,3,1 | 0 | 0 | 1 |
| 3,1,2 | 1 | 0 | 0 |
| 3,2,1 | 0 | 1 | 0 |

We see that the Shapley value is $\Phi(v) = \left(\frac{1}{6}, \frac{1}{6}, \frac{2}{3}\right)$.

9.5.1 Some axiomatic characterizations

The Shapley value can also be characterized as the unique value satisfying a certain set of attractive axioms.

Definition 104 (Value). We say that a mapping $\psi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is a *value* on \mathcal{G}^N , the space of TU games with player set N .

Definition 105 (Summation of games). Given two TU games (N, v) and (N, w) , we define the sum $(N, v + w)$ by $(v + w)(S) = v(S) + w(S)$ for each $S \in 2^N$.

The Shapley value satisfies the following four axioms:

Axioms.

- (S1) *Efficiency*. $\sum_{i \in N} \psi_i(v) = v(N)$ for all $v \in \mathcal{G}^N$.
- (S2) *Null player property*. In a game (N, v) , call a player i a *null player* if $v(S \cup \{i\}) - v(S) = 0$ for all coalitions S . Under the *null player property*, for all games $v \in \mathcal{G}^N$, if i is a null player in v then $\psi_i(v) = 0$.
- (S3) *Symmetry*. In a game (N, v) , call players i, j *symmetric* if $v(S \cup \{i\}) = v(S \cup \{j\})$ for all coalitions $S \subseteq N - \{i, j\}$. Under symmetry, for all $v \in \mathcal{G}^N$, if i and j are symmetric players then $\psi_i(v) = \psi_j(v)$.
- (S4) *Additivity*. $\psi_i(v + w) = \psi_i(v) + \psi_i(w)$ for all $v, w \in \mathcal{G}^N$ and all $i \in N$.

To prove the following theorem, it is worth noting some bases for \mathcal{G}^N . We denote by $1_T \in \mathcal{G}^N$ the game defined by

$$1_T(S) = \begin{cases} 1 & \text{if } S = T, \\ 0 & \text{if } S \neq T. \end{cases}$$

We denote by $u_T \in \mathcal{G}^N$ the T -unanimity game (see Example 67(c)).

Lemma 24.

- (a) The set $\{1_T \in \mathcal{G}^N \mid T \in 2^N - \{\emptyset\}\}$ is a basis for \mathcal{G}^N .
- (b) The set of unanimity games $\{u_T \in \mathcal{G}^N \mid T \in 2^N - \{\emptyset\}\}$ is a basis for \mathcal{G}^N .

Proof. (a) The set $\mathcal{S} = \{1_T \in \mathcal{G}^N \mid T \in 2^N - \{\emptyset\}\}$ is clearly linearly independent. Consider any game $v \in \mathcal{G}^N$. Define $u(S) = \sum_{T \in 2^N - \{\emptyset\}} v(S) 1_T(S)$. Then we have $u = v$. Thus \mathcal{S} spans \mathcal{G}^N .

- (b) Given (a), we need only show that $\mathcal{U} = \{u_T \in \mathcal{G}^N \mid T \in 2^N - \{\emptyset\}\}$ is linearly independent, for $|\mathcal{U}| = |\mathcal{S}|$ and \mathcal{S} is a basis for \mathcal{G}^N . Towards contradiction, assume there exist scalars c_T , some of which are nonzero, s.t. $v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T = 0$. Fix some U s.t. $c_U \neq 0$ in this sum. Then we must have $\sum_{S: S \supset U} c_S = -c_U$, so at least one $S \supsetneq U$ is s.t. $c_S \neq 0$. Fix a sequence $\{S_k\}$ with $S_0 = U$, $S_{k+1} \supsetneq S_k$, and $c_{S_k} \neq 0$ for all k . Since N is finite, this sequence is of finite length K , and we can choose S_K s.t. there is no $R \supsetneq S_K$ with $c_R \neq 0$. Then $v(S_K) = c_{S_K} \neq 0$, so $v \neq 0$, yielding a contradiction. \square

Theorem 37. *The Shapley value $\Phi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is the unique value on \mathcal{G}^N satisfying the axioms (S1)-(S4).*

Proof. First we show Φ satisfies the axioms:

- (S1) Fix any $v \in \mathcal{G}^N$. For any permutation σ , we have $\sum_{i=1}^n (v(P_\sigma(i) \cup \{i\}) - v(P_\sigma(i))) = v(N)$. Thus,

$$\begin{aligned} \sum_{i=1}^n \Phi_i(v) &= \sum_{i=1}^n \left(\frac{1}{n!} \sum_{\sigma \in \Pi(N)} [v(P_\sigma(i) \cup \{i\}) - v(P_\sigma(i))] \right) \\ &= \frac{1}{n!} \sum_{\sigma \in \Pi(N)} \left[\sum_{i=1}^n (v(P_\sigma(i) \cup \{i\}) - v(P_\sigma(i))) \right] \\ &= \frac{1}{n!} \sum_{\sigma \in \Pi(N)} v(N) = v(N). \end{aligned}$$

- (S2) Suppose i is a null player in v . Then $v(S \cup \{i\}) - v(S) = 0$ for all coalitions S so

$$\Phi_i(v) = \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [v(S \cup \{i\}) - v(S)] = 0.$$

- (S3) If i, j are symmetric players, then $v(S \cup \{i\}) - v(S) = v(S \cup \{j\}) - v(S)$ for all coalitions $S \subseteq N - \{i, j\}$. Furthermore, $v(S \cup \{i, j\}) - v(S \cup \{i\}) = v(S \cup \{i, j\}) - v(S \cup \{j\})$ for any $S \subseteq N - \{i, j\}$, since $v(S \cup \{i\}) = v(S \cup \{j\})$. Thus $\Phi_i(v) = \Phi_j(v)$.

(S4) Fix $v, w \in \mathcal{G}^N$. Now,

$$\begin{aligned}
\Phi_i(v + w) &= \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [(v + w)(S \cup \{i\}) - (v + w)(S)] \\
&= \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [v(S \cup \{i\}) - v(S) + w(S \cup \{i\}) - w(S)] \\
&= \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [v(S \cup \{i\}) - v(S)] + \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [w(S \cup \{i\}) - w(S)] \\
&= \Phi_i(v) + \Phi_i(w).
\end{aligned}$$

Conversely, suppose $\psi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is a value satisfying all four axioms, and fix $v \in \mathcal{G}^N$. By Lemma 24, there are unique numbers c_T s.t. $v = \sum_{T \neq \emptyset} c_T u_T$. By additivity (S4), we must have

$$\psi(v) = \sum_{T \neq \emptyset} \psi(c_T u_T) \quad \text{and} \quad \Phi(v) = \sum_{T \neq \emptyset} \Phi(c_T u_T).$$

Thus it is sufficient to show that $\psi(cu_T) = \Phi(cu_T)$ for all nonempty coalitions T and scalars c . Fix $T \neq \emptyset$ and $c \in \mathbb{R}$. By definition of unanimity games, we have

$$cu_T(S \cup \{i\}) - cu_T(S) = 0$$

for all S and any $i \in N - T$. Thus i is a null player in cu_T . By (S2), we have

$$\psi_i(cu_T) = \Phi_i(cu_T) = 0$$

for all $i \in N - T$. Suppose instead that $i, j \in T$ and $i \neq j$. Then for all $S \subseteq N - \{i, j\}$, we have $cu_T(S \cup \{i\}) = cu_T(S \cup \{j\}) = 0$, and so i and j are symmetric players in cu_T . By symmetry (S3), we have

$$\psi_i(v) = \psi_j(v) \quad \text{and} \quad \Phi_i(v) = \Phi_j(v)$$

for all $i, j \in T$. Together with (S1), it follows that

$$\psi_i(cu_T) = \Phi_i(cu_T) |T|^{-1} c$$

for all $i \in T$. Proof follows. \square

The Shapley value also satisfies stronger axioms than the null player property and symmetry:

Axioms.

(S5) *Dummy player property.* In a game (N, v) , call a player i a *dummy player* if $v(S \cup \{i\}) - v(S) = v(i)$ for all $S \subseteq N - \{i\}$. Under the *dummy player property*, for all games $v \in \mathcal{G}^N$, if i is a dummy player then $\psi_i(v) = v(i)$.

(S6) *Anonymity.* For a permutation $\sigma \in \Pi(N)$, define the game v^σ by $v^\sigma(S) := v(\sigma^{-1}(S))$ for all coalitions $S \in 2^N$. Define $\sigma^* : \mathbb{R}^N \rightarrow \mathbb{R}^N$ by $\sigma^*(x)_{\sigma(k)} = x_k$ for all $x \in \mathbb{R}^N$ and $k \in N$. Under *anonymity*, $\psi(v^\sigma) = \sigma^*(\psi(v))$ for all $v \in \mathcal{G}^N$ and all $\sigma \in \Pi(N)$.

In words, a dummy player contributes only their worth to every coalition of which they are a member. If a value satisfies the dummy player property than a dummy player receives only their worth. A value is anonymous if player labels are irrelevant to the payment players receive under the value – that is, relabelling the players would not change what each is paid.

Proposition 78.

- (i) (S5) implies the null player property axiom (S2), and
- (ii) (S6) implies the symmetry axiom (S3).

Proof. (i) First we show (S5) implies (S5'). Suppose Now to show (S5') implies (S2). A null player i has worth $v(i) = 0$. Hence under (S2), we have $\psi_i(v) = 0$.

- (ii) Let i, j be symmetric players and suppose ψ is anonymous. Consider permutation σ defined by

$$\sigma(k) = \begin{cases} j & \text{if } k = i, \\ i & \text{if } k = j, \\ k & \text{otherwise.} \end{cases}$$

By anonymity, $\psi_i(v^\sigma) = \sigma_i^*(\psi(v)) = \psi_{\sigma(i)}(v) = \psi_j(v)$. Since i, j are symmetric, $v(S \cup \{i\}) = v(S \cup \{j\}) = v^\sigma(S \cup \{i\}) = v^\sigma(S \cup \{j\})$ for all $S \subseteq N - \{i, j\}$. Hence $\psi_i(v) = \psi_i(v^\sigma)$.

□

Proposition 79.

- (i) The Shapley value Φ satisfies the dummy player property, and
- (ii) the Shapley value Φ is anonymous.

Proof. (i) Let i be a dummy player. Then the result follows immediately from $v(S \cup \{i\}) - v(S) = v(i)$ for all $S \subseteq N - \{i\}$ and the definition of the Shapley value written in terms of coalitions S .

- (ii) Note that for all $\rho, \sigma \in \Pi(N)$ and all $v \in \mathcal{G}^N$,

$$\rho^*(m^\sigma(v)) = m^{\rho\sigma}(v^\rho),$$

since

$$\begin{aligned} m^{\rho\sigma}(v^\rho)_{\rho\sigma(i)} &= v^\rho(\{\rho\sigma(1), \dots, \rho\sigma(i)\}) - v^\rho(\{\rho\sigma(1), \dots, \rho\sigma(i-1)\}) \\ &= v(\{\sigma(1), \dots, \sigma(i)\}) - v(\{\sigma(1), \dots, \sigma(i-1)\}) \\ &= m^\sigma(v)_{\sigma(i)} = \rho^*(m^\sigma(v))_{\rho\sigma(i)}. \end{aligned}$$

Take $v \in \mathcal{G}^N$ and $\rho \in \Pi(N)$. Note $\rho \mapsto \rho\sigma$ is a surjection on $\Pi(N)$, and ρ^* is linear. Together, all the above implies that

$$\begin{aligned}\Phi(v^\rho) &= \frac{1}{n!} \sum_{\sigma \in \Pi(N)} m^\sigma(v^\rho) = \frac{1}{n!} \sum_{\sigma \in \Pi(N)} m^{\rho\sigma}(v^\rho) \\ &= \frac{1}{n!} \sum_{\sigma \in \Pi(N)} \rho^*(m^\sigma(v)) = \rho^* \left(\frac{1}{n!} \sum_{\sigma \in \Pi(N)} m^\sigma \right) = \rho^*(\Phi(v)).\end{aligned}$$

□

Another axiomatic characterization of the Shapley value involves only three axioms: efficiency **(S1)**, symmetry **(S3)**, and strong monotonicity.

Axiom.

(S7) Strong monotonicity. For any $i \in N$, $\psi_i(v) \geq \psi_i(w)$ for all $v, w \in \mathcal{G}^N$ for which

$$v(S \cup \{i\}) - v(S) \geq w(S \cup \{i\}) - w(S) \quad \text{for all } S \in 2^N.$$

In words, strong monotonicity is the property that if player i contributes at least as much to any coalition in game v as in game w , then i 's payoff under the Shapley value is weakly greater in v than in w .

To show that the Shapley value uniquely satisfies the three axioms, we make use of the *inclusion-exclusion principle*:

Proposition 80 (Inclusion-exclusion principle). *Let E be a finite set, and let f, g be functions on 2^E such that*

$$f(T) = \sum_{S: S \subseteq T} g(S).$$

Then

$$g(T) = \sum_{S: S \subseteq T} (-1)^{|T-S|} f(S).$$

Proof.

$$\begin{aligned}\sum_{S: S \subseteq T} (-1)^{|T-S|} f(S) &= \sum_{S: S \subseteq T} \sum_{R: R \subseteq S} (-1)^{|T-S|} g(R) \\ &= \sum_{R: R \subseteq T} g(R) \sum_{S: R \subseteq S \subseteq T} (-1)^{|T-S|} = g(T).\end{aligned}$$

□

Theorem 38. *The Shapley value $\Phi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is the unique value on \mathcal{G}^N satisfying **(S1)**, **(S3)** and **(S7)**.*

Proof. From Theorem 37, we have that Φ satisfies (S1) and (S3). Fix player i and suppose v and w are s.t. $v(S \cup \{i\}) - v(S) \geq w(S \cup \{i\}) - w(S)$ for all $S \in 2^N$. Then

$$\begin{aligned}\Phi_i(v) &= \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [v(S \cup \{i\}) - v(S)] \\ &\geq \frac{1}{n} \sum_{S: i \notin S} \binom{n-1}{|S|}^{-1} [w(S \cup \{i\}) - w(S)] = \Phi_i(w),\end{aligned}$$

so Φ satisfies (S7).

Conversely, suppose $\psi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is a value satisfying (S1), (S3) and (S7). Let $0 \in \mathcal{G}^N$ denote the game that is identically zero. All players are symmetric in 0, so (S1) and (S3) imply $\psi(0) = 0$.

Now consider any game v in which i is a null player. By (S7), we have $\psi_i(v) \geq \psi_i(0)$ and $\psi_i(0) \geq \psi_i(v)$. Since $\psi_i(0) = 0$, we have $\psi_i(v) = 0$. Next, let $c \in \mathbb{R}$ and $T \in 2^N - \{\emptyset\}$. Then since any $i \in N - T$ is a null player in the T -unanimity game u_T , it follows that $\psi_i(cu_T) = 0$ for all $i \in N - T$. Together with symmetry (S3) and efficiency (S1), this implies $\psi_i(cu_T) = c|T|^{-1}e^T$ for all $i \in T$, since T consists of symmetric players.

Next, note:

Lemma 25. *For each $v \in \mathcal{G}^N$, v has a unique representation*

$$v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T \quad \text{where } c_T := \sum_{S: S \subseteq T} (-1)^{|T-S|} v(S).$$

Proof. Fix $T \in 2^N$. Now $v(T) = \sum_{S: S \subseteq T} c_S$. Proof follows immediately from the inclusion-exclusion principle (Proposition 80) \square

Thus v has form $v = \sum c_T u_T$. Let $\alpha(v)$ denote the number of terms c_T in this sum for which $c_T \neq 0$.

We proceed by induction. If $\alpha(v) = 0$ then $v = 0$ so $\psi(v) = \Phi(v) = 0$. If $\alpha(v) = 1$, then v is some T -unanimity game u_T and so we must have $\psi(v) = c_T u_T e^T = \Phi(v)$.

Fix $k \geq 2$ and suppose $\psi(w) = \Phi(w)$ for all w s.t. $\alpha(w) < k$. Let v be s.t. $\alpha(v) = k$. Then there exist coalitions T_1, \dots, T_k and numbers $c_{T_1}, \dots, c_{T_k} \neq 0$ s.t. $v = \sum_{s=1}^k c_{T_s} u_{T_s}$. Define $D := \bigcap_{s=1}^k T_s$.

For $i \in N - D$, define $w^i := \sum_{s: i \in T_s} c_s u_{T_s}$. Now, because $\alpha(w^i) < k$, we have, by the induction hypothesis, that $\psi_i(w^i) = \Phi_i(w^i)$. Now, note that for all $S \in 2^N$,

$$\begin{aligned}v(S \cup \{i\}) - v(S) &= \sum_{s=1}^k c_s u_{T_s}(S \cup \{i\}) - \sum_{s=1}^k c_s u_{T_s}(S) \\ &= \sum_{s: i \in T_s} c_s u_{T_s}(S \cup \{i\}) - \sum_{s: i \in T_s} c_s u_{T_s}(S) \\ &= w^i(S \cup \{i\}) - w^i(S).\end{aligned}$$

By strong monotonicity (**S7**), it follows that $\psi_i(v) = \psi_i(w^i) = \Phi_i(w^i) = \Phi_i(v)$. Hence $\psi_i(v) = \Phi_i(v)$ for all $i \in N - D$. By efficiency (**S1**), we have

$$\sum_{i \in D} \psi_i(v) = \sum_{i \in D} \Phi_i(v).$$

Instead consider $i, j \in D$. For every $S \subseteq N - \{i, j\}$, we have

$$0 = v(S \cup \{i\}) = \sum_{s=1}^k c_s u_{T_s}(S \cup \{i\}) = \sum_{s=1}^k c_s u_{T_s}(S \cup \{j\}) = v(S \cup \{j\}).$$

Thus i and j are symmetric players, and so by symmetry (**S3**), $\psi_i(v) = \psi_j(v)$ and $\Phi_i(v) = \Phi_j(v)$. By (**S1**) again, it follows $\psi(v) = \Phi(v)$. \square

One peculiarity about the Shapley value is that in games with three players or more, it does not necessarily assign a point in the core, even when the core is nonempty:

Example 72. Consider the game (N, v) with $N = \{1, 2, 3\}$ and

$$v(S) = \begin{cases} 2 & \text{if } S = \{1\}, \\ 3 & \text{if } S = N, \\ 0 & \text{otherwise.} \end{cases}$$

The core of this game is $C(v) = \text{co}(\{(2, 1, 0), (2, 0, 1), (3, 0, 0)\})$. Now, the marginal contributions are:

| σ | m_1^σ | m_2^σ | m_3^σ |
|----------|--------------|--------------|--------------|
| 1,2,3 | 2 | -2 | 3 |
| 1,3,2 | 2 | 3 | -2 |
| 2,1,3 | 0 | 0 | 3 |
| 2,3,1 | 3 | 0 | 0 |
| 3,1,2 | 0 | 3 | 0 |
| 3,2,1 | 3 | 0 | 0 |

Thus $\Phi(v) = (\frac{5}{3}, \frac{2}{3}, \frac{2}{3})$. Note $\Phi_1(v) < 2$ but $v(1) = 2$. Hence not only is $\Phi(v) \notin C(v)$, but also $\Phi(v) \notin I(v)$, the imputation set!

9.5.2 Harsanyi dividends

Another characterization of the Shapley value involves Harsanyi dividends. The dividend can be considered a measure of synergy – the total extra worth generated by players cooperating.

Definition 106 (Harsanyi dividend). For any TU-game (N, v) , for each coalition $T \subseteq N$, the *dividend* $\Delta_v(T)$ of T is defined recursively by

$$\begin{aligned} \Delta_v(\emptyset) &:= 0, \\ \Delta_v(T) &:= v(T) - \sum_{S: S \subsetneq T} \Delta_v(S) \quad \text{if } |T| \geq 1. \end{aligned}$$

Suppose that for each coalition involving player i , the dividend of the coalition is divided equally between all players. The Shapley value is the sum of these equally divided dividends across all the coalitions involving i :

Theorem 39. *Let Φ be the Shapley value on \mathcal{G}^N . Then*

$$\Phi_i(v) = \sum_{T:i \in T} \frac{\Delta_v(T)}{|T|}.$$

Proof.

Lemma 26. *Let $v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T$ for T -unanimity games u_T and scalars c_T . Then $\Delta_v(T) = c_T$ for all $T \neq \emptyset$.*

Proof. We proceed by induction. Suppose $|T| = 1$, and wlog take $T = \{i\}$. Then $v(i) = \Delta_v(T)$. Now suppose $\Delta_v(T) = c_T$ holds for all $S \subsetneq T$. Then $\Delta_v(T) = v(T) - \sum_{S \subsetneq T} \Delta_v(S) = v(T) - \sum_{S \subsetneq T} c_S = c_T$, for $v(T) = \sum_{S \subseteq T} c_S$. \square

Since $\{u_T \in \mathcal{G}^N \mid T \in 2^N - \{\emptyset\}\}$ is a basis for \mathcal{G}^N (Lemma 24), we can write $v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T$. From the proof of Theorem 37, we have that $\Phi(c_T u_T) = |T|^{-1} c_T e^T$ for all coalitions T . Hence by (ADD), we have $\Phi(v) = \sum_{T \neq \emptyset} c_T |T|^{-1} e^T$, and thus $\Phi_i(v) = \sum_{T:i \in T} c_T |T|^{-1}$. Proof now follows from the lemma. \square

9.5.3 Multilinear extensions

Another alternative characterization is as follows.

Definition 107 (Multilinear function). We call a function $g : \mathbb{R}^N \rightarrow \mathbb{R}$ *multilinear* if g has form

$$g(x) = \sum_{S: S \subseteq N} c_S \left(\prod_{i \in S} x_i \right)$$

for some real numbers c_S .

Proposition 81. *Let (N, v) be a game. Then there is a unique multilinear function $f : [0, 1]^N \rightarrow \mathbb{R}$ such that $f(e^S) = v(S)$ for all $S \in 2^N$. Moreover,*

$$f(x) = \sum_{S \in 2^N} \left(\prod_{i \in S} x_i \prod_{i \in N-S} (1 - x_i) \right) v(S).$$

We call f the multilinear extension of v .

Proof. By definition, f has form

$$f(x) = \sum_{S \subseteq N} c_S \prod_{i \in S} x_i.$$

Now, $f(e^S) = v(S)$ implies that

$$f(e^S) = \sum_{T: T \subseteq S} c_T \quad \text{for all } S \subseteq N.$$

This is a system of linear equations, and so has a unique solution if, when $v(S) = 0$ for all S , we have that $c_T = 0$ for all T is the only solution. Towards contradiction, assume the system

$$\sum_{T: T \subseteq S} c_T = 0 \quad \text{for all } S \subseteq N$$

has a nonzero solution. Let S be s.t. $c_S \neq 0$ but $c_T = 0$ for all $T \subseteq S$. Then $\sum_{T: T \subseteq S} c_T = c_S \neq 0$, yielding a contradiction. Hence the system has no nonzero solutions, and thus the original system has a unique solution.

Now, let

$$f(x) = \sum_{S: S \subseteq N} \left(\sum_{T: T \subseteq S} (-1)^{|S-T|} v(T) \right) \left(\prod_{i \in S} x_i \right).$$

Taking $c_S = \sum_{T: T \subseteq S} (-1)^{|S-T|} v(T)$, we have $f(x) = \sum_{S: S \subseteq N} c_S \left(\prod_{i \in S} x_i \right)$, so f is multilinear. Furthermore, $f(e^S) = v(S)$ for all coalitions S . Now f can equivalently be written as $f(x) = \sum_{S \subseteq 2^N} \left(\prod_{i \in S} x_i \prod_{i \in N-S} (1 - x_i) \right) v(S)$. \square

Note that the set of extreme points of $[0, 1]^N$ is $\text{ext}([0, 1]^N) = \{e^S \mid S \in 2^N\}$, and $f(e^S) = v(S)$. Thus we can also consider f to be the unique multilinear extension of the function $\tilde{v} : \text{ext}([0, 1]^N) \rightarrow \mathbb{R}$ defined by $\tilde{v}(S) := v(S)$ for all coalitions S .

The multilinear extension f of v has several possible interpretations. Probabilistically, we can imagine that each player i independently chooses to cooperate with probability x_i . Then the probability that coalition S forms is $\prod_{i \in S} x_i \prod_{i \in N-S} (1 - x_i)$. Thus $f(x)$ gives the expected value of the worth of the coalition formed through this process. Another interpretation of $x \in [0, 1]^N$ is that $(N, i \mapsto x_i)$ is a fuzzy set (with membership function $i \mapsto x_i$). In this case, x_i is the intensity with which i is available to cooperate.

The payoff to player i under the Shapley value is the integral of $D_i f$ along the main diagonal of $[0, 1]^N$:

Theorem 40. For any $v \in \mathcal{G}^N$ with multilinear extension f ,

$$\Phi_i(v) = \int_0^1 D_i f(t, \dots, t) dt \quad \text{for each } i \in N.$$

Proof. Note

$$\begin{aligned} D_i f(x) &= \sum_{T: i \in T} \left(\prod_{j \in T - \{i\}} x_j \prod_{j \in N - T} (1 - x_j) \right) v(T) - \sum_{S: i \notin S} \left(\prod_{j \in S} x_j \prod_{j \in N - (S \cup \{i\})} (1 - x_j) \right) v(S) \\ &= \sum_{S: i \notin S} \left(\prod_{j \in S} x_j \prod_{j \in N - (S \cup \{i\})} (1 - x_j) \right) [v(S \cup \{i\}) - v(S)]. \end{aligned}$$

We thus have $\int_0^1 D_i f(t, \dots, t) dt = \sum_{S: i \notin S} \left(\int_0^1 t^{|S|} (1-t)^{n-|S|-1} dt \right) [v(S \cup \{i\}) - v(S)]$.

Lemma 27 (Beta-integral formula). *Let m, n be positive integers. Then*

$$\int_0^1 t^{m-1} (1-t)^{n-1} dt = \frac{(m-1)!(n-1)!}{(m+n-1)!}.$$

Proof. Let $B(m, n) = \int_0^1 t^{m-1} (1-t)^{n-1} dt$. Now, for any integer p we can write $(p-1)! = \Gamma(p) := \int_0^\infty t^{p-1} e^{-t} dt$.⁶⁶ Now

$$\begin{aligned} (m-1)!(n-1)! &= \left(\int_0^\infty u^{m-1} e^{-u} du \right) \left(\int_0^\infty v^{n-1} e^{-v} dv \right) \\ &= \int_0^\infty \int_0^\infty u^{m-1} v^{n-1} e^{-u-v} du dv. \end{aligned}$$

Now consider a change of variables with $u = st$ and $v = s(1-t)$. We have

$$\begin{aligned} (m-1)!(n-1)! &= \int_0^\infty \int_0^1 (st)^{m-1} [s(1-t)]^{n-1} e^{-s} s dt ds \\ &= \left(\int_0^\infty e^{-s} s^{m+n-1} ds \right) \left(\int_0^1 t^{m-1} (1-t)^{n-1} dt \right) \\ &= (m+n-1)! \int_0^1 t^{m-1} (1-t)^{n-1} dt. \end{aligned}$$

□

From the Beta-integral formula, it follows that $\int_0^1 t^{|S|} (1-t)^{n-|S|-1} dt = \frac{|S|!(n-|S|-1)!}{n!}$. Thus $\int_0^1 D_i f(t, \dots, t) dt = \sum_{S: i \notin S} \frac{|S|!(n-|S|-1)!}{n!} [v(S \cup \{i\}) - v(S)] = \Phi_i(v)$. □

9.5.4 The potential approach

Another approach to the Shapley value is in terms of potentials, an approach introduced by Hart & Mas-Colell (1989). As we will discuss, this approach can be tied to the notion of (non-cooperative) potential games we discussed in section 2.16.

Let

$$\mathcal{G} = \bigcup_{N: N \subseteq \mathbb{N}, |N| < \infty} \mathcal{G}^N,$$

i.e. \mathcal{G} is the family of all games (N, v) with finite player sets N . Note $(\emptyset, v) \in \mathcal{G}$.

⁶⁶ B is the *beta function* and Γ the *gamma function*. To see that $(p-1)! = \Gamma(p)$, note that in general, for $z \in \mathbb{C}$, $\Gamma(z+1) = \int_0^\infty t^z e^{-t} dt = \left[-t^z e^{-t} \right]_0^\infty + z \int_0^\infty t^{z-1} e^{-t} dt = \lim_{t \rightarrow \infty} [-t^z e^{-t}] + z \int_0^\infty t^{z-1} e^{-t} dt = z\Gamma(z)$. Now, we have $\Gamma(1) = \int_0^\infty e^{-t} dt = \left[-e^{-t} \right]_0^\infty = 1$. By induction, we have $\Gamma(p) = (p-1)!$.

Definition 108 (Potential). A *potential* is a function $P : \mathcal{G} \rightarrow \mathbb{R}$ such that

$$\begin{aligned} P(\emptyset, v) &= 0, \text{ and} \\ \sum_{i \in N} D_i P(N, v) &= v(N) \quad \text{for all } (N, v) \in \mathcal{G}, \\ \text{where } D_i P(N, v) &:= P(N, v) - P(N - \{i\}, v). \end{aligned}$$

We call $D_i P(N, v)$ the *marginal contribution* of player i in the game (N, v) according to P .

That is, P is a potential if the empty game has potential zero and for each game (N, v) , the gradient $\nabla P(N, v) := (D_i P(N, v))_{i \in N}$ of P is an efficient payoff vector for (N, v) .

Theorem 41.

- (i) *There is a unique potential $P : \mathcal{G} \rightarrow \mathbb{R}$;*
- (ii) *For each $v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T$, we have*

$$P(N, v) = \sum_{T \in 2^N - \{\emptyset\}} c_T |T|^{-1};$$

- (iii) $\nabla P(N, v) = \Phi(v)$.

Proof. (i) Note from the gradient condition that

$$P(N, v) = |N|^{-1} \left(v(N) + \sum_{i \in N} P(N - \{i\}, v) \right).$$

If $P(T, v)$ has a unique value for all $T \subseteq N$ s.t. $T = N - \{i\}$ for some $i \in N$, then it follows that $P(N, v)$ has a unique value. Now $P(\emptyset, v) = 0$. Proof follows by induction.

- (ii) Define $Q : \mathcal{G} \rightarrow \mathbb{R}$ by

$$\begin{aligned} Q(\emptyset, v) &:= 0, \\ Q(N, v) &:= \sum_{T \in 2^N - \{\emptyset\}} c_T |T|^{-1} \quad \text{for all } v = \sum c_T u_T \text{ if } N \neq \emptyset. \end{aligned}$$

For each (N, v) , we have

$$D_i Q(N, v) = \sum_{T \in 2^N - \{\emptyset\}} c_T |T|^{-1} - \sum_{T \in 2^{N - \{i\}} - \{\emptyset\}} c'_T |T|^{-1}$$

for all $i \in N$, where $v = \sum_{T \in 2^N - \{\emptyset\}} c_T u_T$ and where $v' = \sum_{T \in 2^{N-\{i\}} - \{\emptyset\}} c'_T u_T$ is the restriction of v to $2^{N-\{i\}}$. For each $S \subseteq N - \{i\}$, we have

$$\sum_{T \in 2^N - \{\emptyset\}} c_T u_T(S) = v(S) = v'(S) = \sum_{T \in 2^{N-\{i\}} - \{\emptyset\}} c'_T u_T(S),$$

and thus by recursion, $c_T = c'_T$ for all $T \in 2^{N-\{i\}} - \{\emptyset\}$. Together with Theorem 39, it follows that

$$D_i Q(N, v) = \sum_{T: i \in T} c_T |T|^{-1} = \Phi_i(v)$$

for all $i \in N$. Since Φ satisfies axiom **(EFF)**, we have

$$\sum_{i \in N} D_i Q(N, v) = \sum_{i \in N} \Phi_i(N, v) = v(N).$$

Thus Q is a potential so by (i), $P = Q$.

(iii) Follows from $P = Q$ and $D_i Q(N, v) = \Phi_i(v)$. □

From (iii), we see that the Shapley value is the discrete gradient of the potential P .

Proposition 82. For each $(N, v) \in \mathcal{G}$,

$$P(N, v) = \sum_{S \subseteq N} \frac{(|S| - 1)! (|N| - |S|)!}{|N|!} v(S).$$

Proof. For each (N, v) , define

$$Q(N, v) := \sum_{S \subseteq N} \frac{(|S| - 1)! (|N| - |S|)!}{|N|!} v(S).$$

Clearly, $Q(\emptyset, v) = 0$. To show that $Q(N, v) = P(N, v)$, we need only show $D_i Q(N, v) = \Phi_i(N, v)$ for all $i \in N$, in view of Theorem 41. We have

$$\begin{aligned} D_i Q(N, v) &= Q(N, v) - Q(N - \{i\}, v) \\ &= \sum_{T \subseteq N} \frac{(|T| - 1)! (|N| - |T|)!}{|N|!} v(T) - \sum_{S \subseteq N - \{i\}} \frac{(|S| - 1)! (|N| - 1 - |S|)!}{(|N| - 1)!} v(S) \\ &= \sum_{S \subseteq N - \{i\}} \frac{|S|! (|N| - |S| - 1)!}{|N|!} v(S \cup \{i\}) + \sum_{S \subseteq N - \{i\}} \frac{(|S| - 1)! (|N| - |S|)!}{|N|!} v(S) \\ &\quad - \sum_{S \subseteq N - \{i\}} \frac{(|S| - 1)! (|N| - 1 - |S|)!}{(|N| - 1)!} v(S) \\ &= \sum_{S \subseteq N - \{i\}} \frac{|S|! (|N| - 1 - |S|)!}{|N|!} (v(S \cup \{i\}) - v(S)) \\ &= \Phi_i(N, v). \end{aligned}$$

□

Proposition 82 has the following interpretation: Suppose the probability that each coalition S forms is $\left(\frac{|N|}{|S|}\right)^{-1} |N|^{-1}$. Then $\frac{P(N,v)}{|N|}$ is the expected normalized worth $\mathbb{E} \left[\frac{v(S)}{|S|} \right]$.

9.5.5 Reduced games

Definition 109 (Reduced game). Let $\psi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ be a value. For any $U \in 2^N - \{\emptyset\}$ and any $(N, v) \in \mathcal{G}$, define the game $(N - U, v_U^\psi)$ by

$$v_U^\psi(S) := \begin{cases} 0 & \text{if } S = \emptyset, \\ v(S \cup U) - \sum_{i \in U} \psi_i(S \cup U, v) & \text{otherwise.} \end{cases}$$

We call v_U^ψ the (U, ψ) -reduced game of v .

Reduced games can be motivated as follows. Suppose a given solution concept is used to allocate payoffs in a game. If any subset of players consider the game arising among themselves and agree to apply the same solution concept, then their payoffs should be identical to those of the original game. Thus the reduced game provides a way to define a certain notion of *consistency* that values ought to satisfy. The standard notion of consistency in TU-games is *Hart-Mas-Colell consistency* (*HM-consistency*).

Definition 110 (HM-consistency). A value $\psi : \mathcal{G}^N \rightarrow \mathbb{R}^N$ is called *HM-consistent* if it satisfies the following property: For all games (N, v) and all $U \in 2^N - \{\emptyset\}$,

$$\psi_i(N - U, v_U^\psi) = \psi_i(N, v) \quad \text{for all } i \in N - U.$$

Lemma 28. Let $(N, v) \in \mathcal{G}$ be a TU game with potential P , and suppose $Q : 2^N \rightarrow \mathbb{R}$ is such that

$$\sum_{i \in S} (Q(S) - Q(S - \{i\}))$$

for all nonempty $S \subseteq N$. Then

$$Q(S) = P(S, v) + Q(\emptyset)$$

for all $S \subseteq N$.

Proof. We proceed by induction. If $|S| = 0$ then $S = \emptyset$, $P(\emptyset, v) = 0$ and so $Q(S) = P(S, v) + Q(\emptyset)$. Now suppose $Q(S) = P(S, v) + Q(\emptyset)$ for all $S \subseteq N$ s.t. $|S| < k$, and

take some $T \subseteq N$ with $|T| = k$. Then

$$\begin{aligned}
Q(T) &= |T|^{-1} \left(v(T) + \sum_{i \in T} Q(T - \{i\}) \right) \\
&= |T|^{-1} \left(v(T) + |T| Q(\emptyset) + \sum_{i \in T} P(T - \{i\}, v) \right) \\
&= Q(\emptyset) + |T|^{-1} \left(v(T) + \sum_{i \in T} P(T - \{i\}, v) \right) \\
&= Q(\emptyset) + P(T, v).
\end{aligned}$$

□

Proposition 83. *The Shapley value Φ is HM-consistent.*

Proof. Fix a TU-game $(N, v) \in \mathcal{G}$. Fix a nonempty coalition U and consider the reduced game v_U^Φ . The reduced game v_U^Φ satisfies

$$\begin{aligned}
v_U^\Phi(S) &= v(S \cup U) - \sum_{i \in U} \Phi_i(S \cup U, v) = \sum_{i \in S} \Phi_i(S \cup U, v) \\
&= \sum_{i \in S} [P(S \cup U, v) - P((S \cup U) - \{i\}, v)],
\end{aligned}$$

where the first equality is by definition of the reduced game, the second is by efficiency, and the third is because $\Phi(N, v) = \nabla P(N, v)$ (Theorem 41), the gradient of the potential. For each $S \in 2^{N-U}$, define $Q(S) := P(S \cup U, v)$. By Lemma 28, this implies that

$$Q(S) = P(S, v_U^\Phi) + Q(\emptyset) = P(S, v_U^\Phi) + P(U, v)$$

for all $S \in 2^{N-U}$. By definition of Q , we have $P(S \cup U, v) = P(S, v_U^\Phi) + P(U, v)$, and so

$$\begin{aligned}
\Phi_i(N - U, v_U^\Phi) &= P(N - U, v_U^\Phi) - P((N - U) - \{i\}, v_U^\Phi) \\
&= P(N, v) - P(N - \{i\}, v) = \Phi_i(N, v).
\end{aligned}$$

Thus for each $U \in 2^N - \{\emptyset\}$, we have $\Phi_i(N - U, v_U^\Phi) = \Phi_i(N, v)$ for all $i \in N - U$. Thus Φ is HM-consistent. □

9.5.6 Myerson value

The Shapley value allocates payoffs to players based on the values of all possible coalitions. Thus it treats all possible coalitions in a symmetric way. In practice, however, there may be factors that prevent certain coalitions from forming. Two people who do not know about each other or who live far apart probably cannot form a coalition, at least not without some intermediary. In many settings, cooperation takes place within

structures such as social structures, networks of business relationships or international agreements, and so on.

Myerson (1977) considers cooperative games where cooperation structures are described by graphs. Fix a set of players N and consider a graph $G = (N, E)$ (c.f. Definition 1). Let \mathcal{N}^N denote the set of all graphs on N .

Definition 111. Consider a TU-game (N, v) and a graph $G = (N, E)$.

- (a) *Connectedness.* Given a coalition S , call a pair of players $i, j \in S$ *connected in S by G* if there is a path connecting i and j that stays within S . Call the coalition S *connected* if every pair of players in S is connected in S by G . That is, S is connected if the subgraph of G induced by S is connected.⁶⁷

For each coalition S , let $S|_G := \{\{i \in S \mid i \text{ and } j \text{ are connected in } S \text{ by } G\} \mid j \in S\}$ be the partition of S into sets of players connected in S by G . Note that S is connected if $S|_G = \{S\}$.

- (b) *Allocation rule.* An *allocation rule* $\psi : \mathcal{N}^N \rightarrow \mathbb{R}^N$ gives, for each graph $G \in \mathcal{N}^N$, an allocation $\psi_i(G)$ for each player $i \in N$, with the property that

$$\sum_{i \in S} \psi_i(G) = v(S)$$

for every coalition $S \in N|_G$, for every graph $G \in \mathcal{N}^N$.

- (c) *Fairness.* We call an allocation rule ψ *fair* if

$$\psi_i(G) - \psi_i(G - ij) = \psi_j(G) - \psi_j(G - ij)$$

for every $ij \in G$ and every graph $G \in \mathcal{N}^N$.

- (d) *Worth of disconnected coalitions.* Given G , the worth of any coalition $S \in 2^N$, possibly disconnected, is given by

$$v|_G(S) = \sum_{T \in S|_G} v(T).$$

Intuitively, fairness requires that in any bilateral relationship between players (represented by the edge $ij \in G$), both players benefit equally.

The *Myerson value* ψ of a game (N, v) under graph G is the Shapley value of the game $(N, v|_G)$.

Theorem 42 (Myerson, 1977). *For any TU-game (N, v) , the unique fair allocation rule ψ is given by $\psi(G) = \Phi(v|_G)$ for every graph $G \in \mathcal{N}^N$, where Φ is the Shapley value.*

⁶⁷A graph $G = (V, E)$ is *connected* if there is a path connecting any pair of vertices $i, j \in V$. The subgraph of G induced by some $U \subset V$ is a graph (U, F) with $ij \in F$ iff $i, j \in U$ and $ij \in E$.

Proof. Suppose ψ and ψ' are two distinct fair allocation rules. Let $G = (N, E)$ be a network with $|E|$ minimal that $\psi(G) \neq \psi'(G)$. Now for all $i, j \in N$, $\psi_i(G - ij) = \psi'_i(G - ij)$ and $\psi_j(G - ij) = \psi'_j(G - ij)$. Since ψ and ψ' are fair, we have $\psi_i(G) - \psi'_i(G) = \psi_j(G) - \psi'_j(G)$. Now, $\sum_{i \in N} (\psi_i(G) - \psi'_i(G)) = 0$, and so $\psi_i(G) - \psi'_i(G) = 0$ for all $i \in N$, and thus $\psi = \psi'$, yielding a contradiction. Hence there is a unique fair allocation rule.

It remains to show that the Myerson value $\psi(G) = \Phi(v|_G)$ is a fair allocation rule. First, we show it is indeed an allocation rule. Fix a coalition T . Now, $T|_G = \bigcup_{S \in N|_G} (T \cap S)|_G$, and thus $v|_G = \sum_{S \in N|_G} u_S$, where u_S is defined by $u_S(T) = \sum_{R \in (T \cap S)|_G} v(R)$ for all coalitions $T \subseteq N$. For u_S , all players not in S are dummy players. Since Φ satisfies the dummy player property (**DPP**), we have $\sum_{i \in S} \Phi_i(u_S) = u_S(N) = v(S)$ and $\sum_{i \in T} \Phi_i(u_S) = 0$ for all $T \in N|_G$ s.t. $T \neq S$. Now by additivity (**ADD**), $\Phi(v|_G) = \sum_{S \in N|_G} \Phi(u_S)$, and thus for any coalition $T \in N|_G$, we have $\sum_{i \in T} \Phi_i(v|_G) = \sum_{S \in N|_G} \sum_{i \in T} \Phi_i(u_S) = u_T(N) = v(T)$. Hence the Myerson value is an allocation rule.

Next, to show the Myerson value is fair. Define a new game $w = v|_G - v|_{G-ij}$. Now, i, j are interchangeable in w and $w(S \cup i) = w(S \cup j) = 0$ for all coalitions $S \not\supseteq \{i, j\}$. By symmetry (**SYM**), $\Phi_i(w) = \Phi_j(w)$. By additivity, we conclude that $\Phi_i(v|_G) - \Phi_i(v|_{G-ij}) = \Phi_j(v|_G) - \Phi_j(v|_{G-ij})$. \square

10 Mathematical appendix

Here we give an overview and prove important results, many of which we relied on throughout. These are not really self-contained but I've stated definitions where they might be less familiar.

10.1 Correspondences

A *correspondence* $F : X \rightrightarrows Y$ is a set-valued function from X into 2^Y .

Call a correspondence $F : X \rightrightarrows Y$ *nonempty-valued* if $F(x) \neq \emptyset$ for all $x \in X$, *convex-valued* if $F(x)$ is convex for all $x \in X$, *closed-valued* if $F(x)$ is closed for all $x \in X$, and *compact-valued* if $F(x)$ is compact for all $x \in X$.

Definition 112 (Hemicontinuity). Let $F : X \rightrightarrows Y$ be a correspondence from a topological space X into a topological space Y .

- (a) *Upper hemicontinuity.* F is said to be *upper hemicontinuous* at a point $x \in X$ if for any open neighbourhood V of $F(x)$, there exists an open neighbourhood U of x such that $F(U) \subseteq V$.

F is called *upper hemicontinuous* if F is upper hemicontinuous at every $x \in X$.

- (b) *Lower hemicontinuity.* F is said to be *lower hemicontinuous* at a point $x \in X$ if for any open set V such that $V \cap F(x) \neq \emptyset$ (i.e. such that V intersects $F(x)$), there exists a neighbourhood U of x such that $F(t) \cap V \neq \emptyset$ for all $t \in U$.

F is called *lower hemicontinuous* if it is lower hemicontinuous at every $x \in X$.

- (c) *Continuity.* F is said to be *continuous* at $x \in X$ if it is both upper and lower hemicontinuous at x . If F is continuous at every $x \in X$, it is said to be *continuous*.

Hemicontinuity is in some sense the correspondence analogue of semicontinuity for functions.⁶⁸

Theorem 43. *Suppose that X and Y are metric spaces, that $F : X \rightrightarrows Y$ is nonempty-valued and that Y is compact.*

- (i) *If for every sequence $\{x_n\}$ in X such that $x_n \rightarrow x$ and for every sequence $\{y_n\}$ in Y such that $y_n \rightarrow y$ and $y_n \in F(x_n)$ for all n , we have $y \in F(x)$, then F is upper hemicontinuous at x .*
- (ii) *If F compact valued, then F is upper hemicontinuous iff for every sequence $\{x_n\}$ in X such that $x_n \rightarrow x$ for some $x \in X$ and for every sequence $\{y_n\}$ in Y such that $y_n \rightarrow y$ and $y_n \in F(x_n)$ for all n , we have $y \in F(x)$.*
- (iii) *F is lower hemicontinuous iff for all $y \in F(x)$ and every sequence $\{x_n\}$ in X such that $x_n \rightarrow x$ for some $x \in X$, there exists a sequence $\{y_n\}$ in Y such that $y_n \in F(x_n)$ for all n and $y_n \rightarrow y$.*

Proof. (i) Suppose otherwise, i.e. the property holds for F at x but F is not upper hemicontinuous. Then there exists an open set $V \supseteq F(x)$ such that for any open set U containing x , there is an $x' \in U$ such that $F(x') \not\subseteq V$. Taking successively smaller U , we can find a sequence $\{x_n\}$ with $x_n \rightarrow x$ and $y_n \in F(x_n)$ for each n but $y_n \notin V$.

Now since V is open, V^c is closed and each $y_n \in V^c$ for all n . Since Y is compact, $\{y_n\}$ has a convergent subsequence. We can thus wlog suppose $y_n \rightarrow y$ (since we can take the subsequence otherwise). Since $y_n \in V^c$, we have that $y := \lim_n y_n \in V^c$, but this implies $y \notin F(x)$, yielding a contradiction.

- (ii) Since (i) proves the implication, we need only prove the converse. Suppose F is compact-valued and upper hemicontinuous, and fix any sequence $\{x_n\}$ s.t. $x_n \rightarrow x$ and any sequence $\{y_n\}$ s.t. $y_n \in F(x_n)$ for each n . Since Y is compact, $\{y_n\}$ contains a convergent subsequence $\{y_{n_k}\}$ i.e. $y_{n_k} \rightarrow y$ for some y . Suppose $y \notin F(x)$. Since $F(x)$ is compact, it is closed, and thus the distance between y and $F(x)$ is strictly positive. Thus there is some closed ϵ -ball B_ϵ containing $F(x)$ s.t. B_ϵ does not contain y . Moreover, B_ϵ contains $F(x_n)$ for sufficiently large n , by upper hemicontinuity, and so $F(x_{n_k})$ lies in B_ϵ for sufficiently large k . But then we must have that $y \in B_\epsilon$, yielding a contradiction.
- (iii) Suppose F is lower hemicontinuous, fix some sequence $\{x_n\}$ s.t. $x_n \rightarrow x$ for some $x \in X$, and fix any $y \in F(x)$. Consider a sequence of $1/k$ -balls $B_{1/k}(y)$ centred on y . Since $y \in F(x)$, $B_{1/k}(y) \cap F(x) \neq \emptyset$. Since F is lower hemicontinuous, for each k there exists a neighbourhood U_k of x s.t. $F(z) \cap B_{1/k}(y) \neq \emptyset$ for each $z \in U_k$.

⁶⁸See e.g. Rudin's RCA, Chapter 2, Definition 8 (p. 37 in the third edition).

Since $x_n \rightarrow x$, $x_n \in U_k$ for each k and n sufficiently large, and thus we can choose a subsequence $\{x_{n_k}\}$ s.t. $x_{n_k} \in U_k$ for each k . Now $y_{n_k} \in B_{1/k}(y) \cap F(x_{n_k})$, and so $y_{n_k} \rightarrow y$.

Conversely, suppose the property holds for F but F is not lower hemicontinuous. Then there exists an open set V s.t. $F(x) \cap V \neq \emptyset$ and every neighbourhood U of x contains some point z with $F(z) \cap V = \emptyset$. Taking $\{U_n\}$ to be a sequence of $1/n$ -balls centred on x , we can choose $x_n \in U_n$ for each n so that $x_n \rightarrow x$ with $F(x_n) \cap V = \emptyset$ for each n . Now any sequence $\{y_n\}$ with $y_n \in F(x_n)$ for each n is contained in the closed set V^c , and so, if convergent, converges in V^c . Hence $\{y_n\}$ cannot converge to any $y \in V \cap F(x)$, yielding a contradiction. \square

Definition 113.

- (a) *Graph.* Given a correspondence $F : X \rightrightarrows Y$ or function $F : X \rightarrow Y$, the *graph* of F is the set $G(F) := \{(x, y) \in X \times Y \mid y \in F(x)\}$.
- (b) *Closed graph property.* Let X and Y be topological vector spaces. We say that a correspondence $F : X \rightrightarrows Y$ has a *closed graph* if $G(F)$ is a closed subset of $X \times Y$.

Theorem 44 (Closed graph theorem). *Suppose X and Y are topological spaces and Y is also a compact Hausdorff space. Then a correspondence $F : X \rightrightarrows Y$ has a closed graph iff F is upper hemicontinuous and closed-valued.*

Proof.

Lemma 29. *If X, Y are topological spaces with Y Hausdorff, and if $F : X \rightrightarrows Y$ is an upper hemicontinuous correspondence, and if F is compact-valued, then F has a closed graph.*

Proof. Suppose $(x, y) \notin G(F)$, i.e. $y \notin F(x)$. Since Y is Hausdorff and F is compact-valued, there exist neighbourhoods V of y and W of $F(x)$ s.t. V and W are disjoint. Define the upper inverse of F by $F^u(A) = \{x \in X \mid F(x) \subseteq A\}$. Now $U = F^u(W)$ is open, and thus $U \times V$ is a neighbourhood of (x, y) and $U \times V$ is disjoint of $G(F)$. Thus $G(F)$ is closed. \square

Now since closed subsets of compact sets are compact, it follows that if F is closed-valued into a compact space then it is compact-valued. Thus if F is closed-valued and upper hemicontinuous, it has a closed graph by the lemma.

For the converse, see the proof of Theorem 17.11 in Aliprentis (2005). \square

Theorem 45 (Berge's theorem of the maximum). *Suppose $f : X \times Y \rightarrow \mathbb{R}$ is a continuous function, where X and Y are metric spaces with Y compact. Then*

- (i) *the function $g : X \rightarrow \mathbb{R}$ defined by*

$$g(x) = \max_{y \in Y} f(x, y)$$

is continuous, and

(ii) the correspondence $F : X \rightrightarrows Y$ defined by

$$F(x) = \arg \max_{y \in Y} f(x, y)$$

is nonempty-valued and has a closed graph.

Proof. See Lemmas 17.29, 17.30, and the proof of Theorem 17.31 in Aliprentis (2005). \square

10.2 Linear programming

Linear programming lies at the heart of much game theory. Indeed, the early history of game theory is effectively a history of linear programming, as von Neumann's minmax theorem illustrates.

10.2.1 Hyperplanes, convex sets and extreme points

Definition 114 (Hyperplane).

- (a) *Hyperplane.* Given a point $p \in \mathbb{R}^n - \{0\}$ and a scalar $c \in \mathbb{R}$, the *hyperplane generated by p and c* is the set

$$H_{p,c} = \{x \in \mathbb{R}^n \mid p \cdot x = c\}.$$

- (b) *Half-space.* Given a hyperplane $H_{p,c}$, the sets

$$\{x \in \mathbb{R}^n \mid p \cdot x \geq c\} \quad \text{and} \quad \{x \in \mathbb{R}^n \mid p \cdot x \leq c\}$$

are called, respectively, the (closed) *half-spaces* above and below $H_{p,c}$.

- (c) *Supporting hyperplane.* Given a set $S \subseteq \mathbb{R}^n$, we say that a hyperplane H is a *supporting hyperplane* of S if
- (i) S is contained in either the half-space above or the half-space below H , and
 - (ii) S has at least one boundary point on H .

There are several variations of the separating hyperplane theorem.

Theorem 46 (Separating hyperplane theorem I). *Let $A \subseteq \mathbb{R}^n$ be a closed convex set and let $x \in \mathbb{R}^n - A$. Then there exists a $y \in \mathbb{R}^n$ with $y \cdot z > y \cdot x$ for all $z \in A$.*

Proof. If A is empty the theorem holds vacuously. Hence suppose A is nonempty.

First note the following lemma:

Lemma 30. *Let C be a (nonempty) closed convex subset of \mathbb{R}^n . Then C contains a unique vector of minimum norm.*

Proof. Define $\delta = \inf\{\|x\| \mid x \in C\}$. Consider any sequence $\{x_i\}$ s.t. $x_i \in C$ for all i and $\|x_i\| \rightarrow \delta$. Since C is convex, $\frac{x_i + x_j}{2} \in C$ for any i, j . Thus $\|x_i + x_j\|^2 \geq 4\delta^2$. Now,

$$\begin{aligned}\|x_i - x_j\|^2 &= \|x_i\|^2 + \|x_j\|^2 - 2x_i \cdot x_j \\ &= 2\|x_i\|^2 + 2\|x_j\|^2 - \|x_i + x_j\|^2 \\ &\leq \|x_i\|^2 + \|x_j\|^2 - 4\delta^2 \rightarrow 0.\end{aligned}$$

Thus $\{x_i\}$ is a Cauchy sequence, so $x := \lim_{i \rightarrow \infty} x_i \in C$. Hence C contains a vector of minimum norm.

For uniqueness, suppose $y \in C$ is also s.t. $\|y\| = \delta$. Then $\|x - y\|^2 \leq 2\|x\|^2 + 2\|y\|^2 - 4\delta^2 = 0$, so $x = y$. \square

By the lemma, there exists some $z' \in A$ s.t. $0 < \|x - z'\| \leq \|x - z\|$ for all $z \in A$.

Let $y = z' - x$ and fix $z \in A$. For any $\alpha \in [0, 1]$, we have that $z' + \alpha(z - z') \in C$ by the convexity of A . Hence

$$\|z' + \alpha(z - z') - x\|^2 \geq \|z' - x\|^2.$$

It follows that

$$2\alpha(z' - x) \cdot (z - z') + \alpha^2\|z - z'\|^2 \geq 0.$$

Dividing through by 2α and taking limits as $\alpha \rightarrow 0$, we see that $(z' - x) \cdot (z - z') \geq 0$. Hence $(z' - x) \cdot z \geq (z' - x) \cdot z' = (z' - x) \cdot x + (z' - x) \cdot (z' - x) > (z' - x) \cdot x$. Since z is arbitrary, the theorem follows. \square

Corollary 13. *In the context of Theorem 46, there exists a real number α such that $y \cdot z > \alpha$ and $y \cdot x < \alpha$ and a real number β such that $y \cdot z > \beta$ and $y \cdot x = \beta$.*

Proof. From the proof of the theorem, $y \cdot z \geq y \cdot z'$ for all $z \in A$, so $y \cdot z'$ is a lower bound on $y \cdot z$. Taking $\alpha = \frac{1}{2}(y \cdot z' + y \cdot x)$, we have $y \cdot z > \alpha$, and since $y \cdot z' > y \cdot x$, we have $y \cdot x < \alpha$. The assertion involving β is trivial. \square

Theorem 47 (Separating hyperplane theorem II). *Let A and B be two disjoint convex subsets of \mathbb{R}^n . Then there exists a vector $y \in \mathbb{R}^n$ and a scalar $c \in \mathbb{R}$ such that*

$$x \cdot y \geq c \quad \text{and} \quad z \cdot y \leq c$$

for all $x \in A$ and all $z \in B$.

Proof. Define

$$C = \{x - z \mid x \in A, z \in B\},$$

i.e. C is the Minkowski sum $C = A + (-B)$. Since $-B$ is convex and Minkowski sums of convex sets are convex, C is convex. The closure \bar{C} of C is also convex. By Lemma 30, \bar{C} contains a (unique) vector y of minimum norm. By convexity of \bar{C} , for any $v \in C$ we have

$$y + t(v - y) \in \bar{C}$$

for all $t \in [0, 1]$. Thus

$$\|y\|^2 \leq \|y + t(v - y)\|^2 = \|y\|^2 + 2ty \cdot (v - y) + t^2\|v - y\|^2.$$

If $t \in (0, 1]$,

$$0 \leq 2y \cdot v - 2\|y\|^2 + t\|v - y\|^2.$$

Taking $t \rightarrow 0$ gives $\|y\|^2 \leq y \cdot v$. Thus $(x - z) \cdot y \geq \|y\|^2$ for any $x \in A$ and $z \in B$. Since $(x - z) \cdot y \geq \inf_{x \in A, z \in B} \{(x - z) \cdot y\} \geq \|y\|^2$ and $\inf_{x \in A, z \in B} \{(x - z) \cdot y\} = \inf_{x \in A} x \cdot y - \sup_{z \in B} z \cdot y$, we have $\inf_{x \in A} x \cdot y \geq \|y\|^2 + \sup_{z \in B} z \cdot y$. Proof follows provided $y \neq 0$.

Extending the argument to the general case, suppose the interior C° of C is nonempty. Construct a sequence of nonempty compact convex sets $\{C_i\}$ as follows:

$$C_i = [-i, i]^n \cap \left\{ x \in C^\circ \mid \|x - x'\| \geq \frac{1}{i} \text{ for all } x' \in (C^\circ)^c \right\}.$$

Then $C_1 \subseteq C_2 \subseteq C_3 \subseteq \dots$, each $C_i \subseteq C^\circ$, and $\bigcup_{i=1}^\infty C_i = C^\circ$. Now $0 \notin C^\circ$, and thus $0 \notin C_i$ for any i . Hence each C_i contains a nonzero vector y_i of minimum norm, by Lemma 30. By the argument of the preceding part, $x \cdot y_i \geq 0$ for all $x \in C_i$. Normalize each y_i s.t. $\|y_i\| = 1$. Letting $S = \{x \in \mathbb{R}^n \mid \|x\| \leq 1\}$, we have $y_i \in S$ for each S . Since S is compact, $\{y_i\}$ contains a convergent subsequence $\{y_{i_k}\}$. Let $y = \lim_{k \rightarrow \infty} y_{i_k}$ and note $y \neq 0$. Now, $x \cdot y \geq 0$ for all $x \in C^\circ$ and, by continuity, for all $x \in C$. Proof follows by the previous arguments.

If the interior of C is empty, then the affine set spanning C has dimension strictly less than n . Thus there is some hyperplane $H_{c,y}$ s.t. $C \subseteq H_{c,y}$. Then $x \cdot y \geq c$ for all $x \in C$. The remainder of the proof precedes as above. \square

A related result is the supporting hyperplane theorem:

Theorem 48 (Supporting hyperplane theorem). *Suppose $A \subseteq \mathbb{R}^n$ is a convex set and $x \notin A^\circ$. Then there exists a $y \in \mathbb{R}^n$ with $y \neq 0$ such that $y \cdot x \geq y \cdot z$ for all $z \in A$.*

Proof. If A is empty the theorem holds vacuously. Hence suppose $A \neq \emptyset$. Suppose $x \notin A$. There exists a sequence $\{x_i\}$ s.t. $x_i \rightarrow x$ and $x_i \notin A$ for all i . By the separating hyperplane theorem (Theorem 46), for each i there exists a $y_i \in \mathbb{R}^n$ $y_i \neq 0$ and a $c_i \in \mathbb{R}$ s.t.

$$y_i \cdot x_i > c_i \geq y_i \cdot z$$

for all $z \in A$. Normalize each y_i s.t. $\|y_i\| = 1$. Let $S = \{x \in \mathbb{R}^n \mid \|x\| \leq 1\}$. Then each $y_i \in S$ and S is compact, so $\{y_i\}$ contains a convergent subsequence $\{y_{i_k}\}$ with $\lim_{k \rightarrow \infty} y_{i_k} =: y \neq 0$ and $\lim_{k \rightarrow \infty} c_{i_k} =: c$. Then

$$y \cdot x \geq c \geq y \cdot z$$

for all $z \in A$. \square

An extreme point of a convex set is a point that is not a convex combination of two distinct points in that set:

Definition 115 (Extreme point). Given a convex set S in a linear space V , a vector $x \in S$ is called an *extreme point* if there are no distinct points $y, z \in S$ and no constant $\lambda \in (0, 1)$ such that $x = \lambda y + (1 - \lambda)z$.

We denote the set of extreme points of S by $\text{ext}(S)$.

There are several equivalent definitions:

Theorem 49. *Let S be a convex set in a linear space V . Then the following are equivalent:*

- (i) $x \in \text{ext}(S)$;
- (ii) for all $y, z \in S$ such that $x = \frac{1}{2}(y + z)$, we have that $x = y = z$;
- (iii) $S - \{x\}$ is convex.

Proof. Clearly (i) immediately implies (ii). Conversely, suppose x satisfies (ii) but there exists a constant $\lambda \in (0, 1)$ s.t. $x = \lambda y + (1 - \lambda)z$ for some distinct $y, z \in S$. If $\lambda > \frac{1}{2}$, let $\alpha = 2\lambda - 1$. Define $z' = \alpha y + (1 - \alpha)z$. Then we have $x = \frac{1}{2}(y + z')$. By convexity of S , $z' \in S$, so we have a contradiction. Similarly, if $\lambda < \frac{1}{2}$, let $\alpha = 2\lambda$ and define $y' = \alpha y + z$. We have $x = \frac{1}{2}(y' + z)$. Likewise, $y' \in S$, yielding a contradiction. Hence (ii) implies (i).

Suppose x satisfies (i) but $S - \{x\}$ is not convex. Then there must be some $y, z \in S - \{x\}$ and some $\lambda \in [0, 1]$ s.t. $\lambda y + (1 - \lambda)z \notin S - \{x\}$. Yet by convexity of S , $\lambda y + (1 - \lambda)z \in S$. Hence we must have $x = \lambda y + (1 - \lambda)z$, but this contradicts (i) [noting $\lambda \in (0, 1)$ since $y, z \neq x$]. Thus (i) implies (iii). Conversely, suppose $S - \{x\}$ is convex. Suppose there is some $\lambda \in (0, 1)$ and distinct points $y, z \in S$ s.t. $x = \lambda y + (1 - \lambda)z$. Then $y, z \neq x$, so $y, z \in S - \{x\}$, but convexity of $S - \{x\}$ would then imply $x \in S - \{x\}$, yielding a contradiction. Hence (iii) implies (i). \square

Theorem 50 (Minkowski-Carathéodory). *Let $K \subseteq \mathbb{R}^n$ be a convex compact set. Then every $x \in K$ can be expressed as a convex combination of at most $n + 1$ extreme points of K .*

Proof. We proceed by induction. Suppose the theorem holds for any convex compact set in \mathbb{R}^{n-1} . Let $K \subseteq \mathbb{R}^n$ be a convex compact set, and fix $x \in K$. Now we can choose some extreme point $y \in K$ and some boundary point $z \in K$ s.t. $x = (1 - \lambda)y + \lambda z$ for some $\lambda \in [0, 1]$. By the supporting hyperplane theorem, K has a supporting hyperplane H at z . Now $H \cap K \subseteq H$ is a convex compact set, and H has dimension $n - 1$. Furthermore, the extreme points of $H \cap K$ are extreme points of K . By hypothesis, we can represent z as a convex combination of at most n extreme points z_1, \dots, z_n , i.e. $z = \sum_{i=1}^n c_i z_i$ for some scalars $c_i \in [0, 1]$ with $\sum_{i=1}^n c_i = 1$. Since x is a convex combination of z and y , it follows that x is a convex combination of at most $n + 1$ extreme points of K .

To complete the proof, note that the induction hypothesis holds for $n = 1$, since any convex set in \mathbb{R} is an interval $[a, b]$, which has 2 extreme points, a and b . For any $x \in [a, b]$, we have $x = (1 - \lambda)a + \lambda b$ for $\lambda = \frac{x - a}{b - a}$, so any $x \in [a, b]$ is a convex combination of at most 2 extreme points. \square

The Krein-Milman theorem establishes that any compact convex subset of an Euclidean space is equal to the (closed) convex hull of its extreme points.⁶⁹ To prove this, we first need a separation lemma:

Lemma 31. *Let C be a nonempty convex subset of \mathbb{R}^n and let $x \in \mathbb{R}^n - C^\circ$. Then there exists a vector $y \in \mathbb{R}^n - \{0\}$ such that $y \cdot x \leq y \cdot z$ for all $z \in C$.*

Proof. First, suppose $x \notin \bar{C}$. Then the lemma follows immediately from the separating hyperplane theorem, Theorem 46, taking $A = \bar{C}$.

Second, suppose $x \in \bar{C}$. Since $x \notin C^\circ$, there is a sequence $\{x_i\}$ with each $x_i \in \mathbb{R}^n - \bar{C}$ and $x = \lim_{i \rightarrow \infty} x_i$. Applying Theorem 46, for each i there exists a $y_i \in \mathbb{R}^n - \{0\}$ s.t. $y_i \cdot x_i \leq y_i \cdot z$ for all $z \in \bar{C}$. Normalize each y_i so that $\|y_i\| = 1$. Since the closed ball $B = \{x \in \mathbb{R}^n \mid \|x\| \leq 1\}$ is compact, $\{y_i\}$ has some convergent subsequence $\{y_{i_k}\}$ with limit $y = \lim_{k \rightarrow \infty} y_{i_k}$. We have $y \cdot x = \lim_{k \rightarrow \infty} y_{i_k} \cdot x \leq \lim_{k \rightarrow \infty} y_{i_k} \cdot z = y \cdot z$ for all $z \in \bar{C}$. \square

Theorem 51 (Krein-Milman). *Let K be a nonempty compact convex set in \mathbb{R}^n . Then $\text{ext}(K)$ is nonempty and $K = \text{co}(\text{ext}(K))$.*

Proof. Since K is compact and $\rho : x \mapsto \|x\|$ is continuous, it follows by the extreme value theorem that ρ attains a maximum on K . Let $x \in \arg \max_{y \in K} \|y\|$. Then $x \in \text{ext}(K)$. For suppose that $x = \frac{1}{2}(x_1 + x_2)$ for some $x_1, x_2 \in K$. Now

$$\|x\| = \left\| \frac{1}{2}(x_1 + x_2) \right\| \leq \frac{1}{2}\|x_1\| + \frac{1}{2}\|x_2\| \leq \frac{1}{2}\|x\| + \frac{1}{2}\|x\|,$$

so $\|x_1\| = \|x_2\| = \left\| \frac{1}{2}(x_1 + x_2) \right\|$. By definition of the Euclidean norm, it follows that $x_1 = x_2 = x$, and thus $x \in \text{ext}(K)$. Thus $\text{ext}(K)$ is nonempty.

Next, we show that $K = \text{co}(\text{ext}(K))$ by induction. First, note that if $\dim K = 0$ then $K = \{x\}$ for some $x \in \mathbb{R}^n$, and thus $\text{ext}(K) = \{x\}$ and $\text{co}(\text{ext}(K)) = \{x\} = K$.

Now fix an integer k and suppose that $\text{co}(\text{ext}(A)) = A$ for all nonempty compact convex $A \subseteq \mathbb{R}^n$ for which $\dim A < k$. Suppose $\dim K = k$. Clearly $\text{co}(\text{ext}(K)) \subseteq K$, so we need only prove $K \subseteq \text{co}(\text{ext}(K))$. Wlog, assume $0 \in K$, and let $W = \bigcap \{E \subseteq \mathbb{R}^n \mid E \text{ is affine and } K \subseteq E\}$. Then $\dim W = k$. Since $\text{ext}(K)$ is nonempty, there is some $x \in \text{ext}(K)$. Fix $y \in K$. If $y = x$ then $y \in \text{co}(\text{ext}(K))$. Hence suppose that $y \neq x$. Let L be the line through x and y . Now $L' = L \cap K$ is a line segment, with one endpoint being x and the other endpoint being some boundary point b of K . By Lemma 31, there is a linear function $f : W \rightarrow \mathbb{R}$ s.t. $f(b) = \min\{f(z) \mid z \in K\}$ and $f \neq 0$. In particular, the lemma shows there is some nonzero vector y so that $f(z) = y \cdot z$ for $z \in W$ has these properties. Now let $A := \{z \in K \mid f(z) = f(b)\}$. Then A is a compact convex subset of K . Note $\text{ext}(A) \subseteq \text{ext}(K)$. Since $f \neq 0$, we have that $\dim A < k$. By the induction hypothesis, $A = \text{co}(\text{ext}(A))$. Thus $b \in A = \text{co}(\text{ext}(A)) \subseteq \text{co}(\text{ext}(K))$. Since $x \in \text{ext}(K)$ and $y \in \text{co}\{x, b\}$, it follows that $y \in \text{co}(\text{ext}(K))$, and thus $K \subseteq \text{co}(\text{ext}(K))$. \square

⁶⁹The theorem actually holds more generally for Hausdorff locally convex topological vector spaces.

Definition 116.

- (a) *Stochastic matrices.* We call an $n \times n$ real matrix D
- (i) *right-stochastic* if $d_{ij} \geq 0$ for all $1 \leq i, j \leq n$ and $\sum_{j=1}^n d_{ij} = 1$ for each row i ;
 - (ii) *left-stochastic* if $d_{ij} \geq 0$ for all $1 \leq i, j \leq n$ and $\sum_{i=1}^n d_{ij} = 1$ for each column j ;
 - (iii) *doubly stochastic* if it is both right-stochastic and left-stochastic.
- (b) *Permutation matrix.* If an $n \times n$ real matrix P is doubly stochastic and each entry $p_{ij} \in \{0, 1\}$, we say that P is a *permutation matrix*.

Theorem 52 (Birkhoff-von Neumann). *Let $\mathcal{D}_{n \times n}$ denote the set of all $n \times n$ doubly stochastic matrices and let $\mathcal{P}_{n \times n}$ denote the set of all $n \times n$ permutation matrices. Then*

- (a) $\text{ext}(\mathcal{D}_{n \times n}) = \mathcal{P}_{n \times n}$ and
- (b) $\mathcal{D}_{n \times n} = \text{co}(\mathcal{P}_{n \times n})$.

Proof.

- (a) First, we show $\text{ext}(\mathcal{D}_{n \times n}) \subseteq \mathcal{P}_{n \times n}$. Let P be an $n \times n$ permutation matrix s.t. $P = \frac{1}{2}(A + B)$ for some matrices $A, B \in \mathcal{D}_{n \times n}$, with ij th entry p_{ij} . We have $p_{ij} = \frac{1}{2}(a_{ij} + b_{ij})$ and $p_{ij} \in \{0, 1\}$. If $p_{ij} = 0$ then since $a_{ij}, b_{ij} \geq 0$ we must have $a_{ij} = b_{ij} = 0$. If $p_{ij} = 1$ then since $a_{ij}, b_{ij} \leq 1$, we must have $a_{ij} = b_{ij} = 1$. Thus $A = B$, and so $P \in \text{ext}(\mathcal{D}_{n \times n})$.

Next, let $D \in \mathcal{D}_{n \times n}$ be s.t. D is not a permutation matrix, with ij th entry d_{ij} . We claim D is not an extreme point. Now, D is not an extreme point if there exists an $n \times n$ matrix $C \neq 0$ s.t. (i) $c_{ij} = 0$ if $d_{ij} \in \{0, 1\}$, (ii) $\sum_{i=1}^n c_{ij} = 0$ for all j s.t. $d_{ij} \neq 1$ for all i , and (iii) $\sum_{j=1}^n c_{ij} = 0$ for all i s.t. $d_{ij} \neq 1$ for all j . If such a matrix C exists, then we have that $D + \epsilon C$ and $D - \epsilon C$ are distinct doubly stochastic matrices for sufficiently small $\epsilon > 0$, and $D = \frac{1}{2}([D + \epsilon C] + [D - \epsilon C])$.

Clearly, for any row or column of D containing a 1, the corresponding row or column of C must contain only 0s. Now, since D is not a permutation matrix, there are $k \geq 2$ rows (and columns) of D that do not contain an entry 1. In each such row, there are at least $2k$ entries d_{ij} s.t. $d_{ij} \notin \{0, 1\}$. The corresponding elements of C must be chosen to satisfy the system of $2k$ homogeneous linear equations described by (ii) and (iii). We can, wlog, assume these equations correspond to the first k rows and columns. If $\sum_{j=1}^k c_{ij} = 0$ for each $i \in \{1, \dots, k-1\}$ and $\sum_{i=1}^k c_{ij} = 0$ for all $j \in \{1, \dots, k\}$, then $\sum_{j=1}^k c_{kj} = 0$, so the last equation is redundant. Hence the system of equations has fewer than $2k$ independent equations and weakly more than $2k$ variables, so has a nonzero solution, C . Thus we have constructed a matrix C satisfying (i)-(iii).

- (b) Fix any $A, B \in \mathcal{D}_{n \times n}$ and any $\lambda \in [0, 1]$. Let $D = (1 - \lambda)A + \lambda B$. Since $a_{ij}, b_{ij} \geq 0$ for all i, j , we have $d_{ij} \geq 0$ for all i, j . For each row i , $\sum_{j=1}^n d_{ij} = \sum_{j=1}^n [(1 - \lambda)a_{ij} + \lambda b_{ij}] = (1 - \lambda) \sum_{j=1}^n a_{ij} + \lambda \sum_{j=1}^n b_{ij} = 1$, and for each column j , $\sum_{i=1}^n d_{ij} = (1 - \lambda) \sum_{i=1}^n a_{ij} + \lambda \sum_{i=1}^n b_{ij} = 1$. Hence $D \in \mathcal{D}_{n \times n}$, so $\mathcal{D}_{n \times n}$ is convex.

Clearly, $\mathcal{D}_{n \times n}$ is bounded since $0 \leq d_{ij} \leq 1$ for each i, j for any $D \in \mathcal{D}_{n \times n}$. Now let A be an accumulation point of $\mathcal{D}_{n \times n}$. Consider any sequence $\{\epsilon^k\}$ s.t. $\epsilon^k > 0$ for all k and $\lim_{k \rightarrow \infty} \epsilon^k = 0$. Since A is an accumulation point, we can construct a sequence $\{D^k\}$ s.t. for each k , $D^k \in \mathcal{D}_{n \times n}$ and $\|D^k - A\| < \epsilon^k$. Then $A = \lim_{k \rightarrow \infty} D^k$. Since $d_{ij}^k \geq 0$ for all i, j and all k , we have that $a_{ij} = \lim_{k \rightarrow \infty} d_{ij}^k \geq 0$. Since $\sum_{j=1}^n d_{ij}^k = 1$ for all i and all k , $\sum_{j=1}^n a_{ij} = \lim_{k \rightarrow \infty} \sum_{j=1}^n d_{ij}^k = 1$ for all i . Likewise, $\sum_{i=1}^n a_{ij} = \lim_{k \rightarrow \infty} \sum_{i=1}^n d_{ij}^k = 1$ for all j . Thus $A \in \mathcal{D}_{n \times n}$. Since $\mathcal{D}_{n \times n}$ contains all its accumulation points, it is closed. Applying the Heine-Borel theorem, we see $\mathcal{D}_{n \times n}$ is compact. Proof of (b) now follows by application of the Krein-Milman theorem. □

10.2.2 Lemmas of the alternative

Lemmas of the alternative are lemmas that describe two systems of linear equations, precisely one of which has a solution.

Lemma 32 (Lemma of the alternative for matrices). *Let A be an $m \times n$ matrix. Then exactly one of the following statements holds:*

- (a) *There exist $y \in \mathbb{R}^n$ and $z \in \mathbb{R}^m$ such that $(y, z) \geq 0$, $(y, z) \neq 0$ and $Ay + z = 0$;*
- (b) *There is an $x \in \mathbb{R}^m$ such that $x > 0$ and $x'A > 0$.*

Proof. First suppose both (a) and (b) hold. Since $Ay + z = 0$, we have that $x'(Ay + z) = x'Ay + x \cdot z = 0$. Since $x'A > 0$ and $y \geq 0$, we have $x'Ay = 0$ iff $y = 0$. Likewise, since $x > 0$ and $z \geq 0$, we have $x \cdot z = 0$ iff $z = 0$. Hence $x'(Ay + z) = 0$ iff $(y, z) = 0$, yielding a contradiction. Hence at most one of (a) and (b) can hold.

Let A_j be the j th column of A and let $e^1, \dots, e^m \in \mathbb{R}^m$ be the standard basis vectors of \mathbb{R}^m . Let $Z = \text{co}(A_1, \dots, A_n, e^1, \dots, e^m)$. In (a), dividing $Ay + z = 0$ by $\sum_{i=1}^n y_i + \sum_{j=1}^m z_j$, we see that 0 is a convex combination of the vectors in Z .

Suppose (a) does not hold. We need only prove that (b) must then have a solution. Since (a) does not hold, $0 \notin Z$. By Theorem 46 and Corollary 13, there exists $x \in \mathbb{R}^m$ and a number $\beta \in \mathbb{R}$ such that $x \cdot p > \beta$ for all $p \in Z$ and $x \cdot 0 = \beta$. Thus $\beta = 0$ and so $x'A > 0$ and $x > 0$ given the columns of A and all e^j all lie in Z . □

The most well-known lemma of the alternative is Farkas' lemma. Proof requires the following:

Lemma 33. *Let A be an $m \times n$ matrix and let*

$$C := \{c \in \mathbb{R}^n \mid \text{there exists an } x \in \mathbb{R}^m, x \geq 0 \text{ s.t. } x'A = c\}.$$

Then C is closed.

Proof. Suppose $\text{rank } A = m$. Then $\text{rank}(AA') = m$.⁷⁰ Hence AA' is invertible. Let $\{c^n\}$ be a sequence in C s.t. $\lim_{n \rightarrow \infty} c^n = c$, and let $c^n = x^{n'}A$ with $x^n \geq 0$ for all n . Since $x^{n'} = x^{n'}(AA')(AA')^{-1}$ for all n , $x^{n'}A \rightarrow c$ implies that $x^{n'} \rightarrow c'A'(AA')^{-1} =: x$. In particular, $x \geq 0$. Since $x^{n'}A \rightarrow x'A$, we have $x'A = c$ and therefore $c \in C$.

Now fix $b \in C - \{0\}$ and choose $x \in \mathbb{R}^m$, $x \geq 0$ with $x'A = b$ s.t. $|S|$ is maximal, where $S := \{i \in \{1, \dots, m\} \mid x_i > 0\}$. We mean to show that the rows of A indexed by $i \in S$ are linearly independent. Suppose otherwise. Then there exists $\mu \in \mathbb{R}^m$ s.t. $\mu_j \neq 0$ for some $j \in S$, $\mu_i = 0$ for all $i \notin S$, and $\mu'A = 0$. Then $(x - t\mu)'A = x'A = b$ for all $t \in \mathbb{R}$. Choose \hat{t} s.t. $x_j - \hat{t}\mu_j \geq 0$ for all $j \in S$, with equality for some $j \in S$. Then $b = (x - \hat{t}\mu)'A \geq 0$ and $|\{i \in \{1, \dots, m\} \mid x_i - \hat{t}\mu_i > 0\}| \leq |S| - 1$, yielding a contradiction.

Now given this, we can write C as

$$C = \bigcup_B \{x'B \mid B \text{ is a } k \times n \text{ submatrix of } A, \text{ rank}(B) = k \leq \text{rank}(A), 0 \leq x \in \mathbb{R}^k\}.$$

By the first part of the proof, each set in the union is closed. Since there are finitely many such B , we therefore have that C is closed. \square

The canonical Farkas' lemma is as follows:

Lemma 34 (Farkas' lemma). *Let A be an $m \times n$ matrix and let $b \in \mathbb{R}^n$. Then exactly one of the following statements holds:*

- (a) *There is an $x \in \mathbb{R}^m$ such that $x'A = b$ and $x \geq 0$;*
- (b) *There is a $y \in \mathbb{R}^n$ such that $Ay \geq 0$ and $b \cdot y < 0$.*

Proof. First, suppose (a) holds. Consider any $y \in \mathbb{R}^n$ s.t. $b \cdot y < 0$. If $Ay \geq 0$ then for $x \geq 0$, we have $x'Ay \geq 0$. But $x'A = b$ and so $x'Ay = b \cdot y < 0$. Hence if (a) holds, (b) cannot hold. Thus at most one of the statements can be true.

Now suppose (a) does not hold. Define

$$C := \{c \in \mathbb{R}^n \mid \text{there exists } x \geq 0 \text{ s.t. } x'A = c\}.$$

By Lemma 33, C is closed.

⁷⁰If an $m \times n$ matrix A has $\text{rank } A = k$ then $\text{rank}(AA') = k$. The kernel of AA' is $\{x \in \mathbb{R}^m \mid AA'x = 0\}$, i.e. the set of $x \in \mathbb{R}^m$ s.t. $\sum_{j=1}^m (\sum_{k=1}^m a_{ik}a_{jk})x_j = 0$ for all i . Follows that $\sum_{i=1}^m x_i \sum_{j=1}^m (\sum_{k=1}^m a_{ik}a_{jk})x_j = 0$ and thus $x'AA'x = x'(AA'x) = x \cdot 0 = 0$. Now, $0 = x'AA'x = (A'x)'(A'x) = \|A'x\|^2$. Hence x must be s.t. $A'x = 0$, so the kernel of AA' and A' are the same. By the rank-nullity theorem, we have that $\text{rank}(AA') = \text{rank}(A') = \text{rank } A = k$.

By Theorem 46 and its corollary, there is a $y \in \mathbb{R}^n$ and a real number α s.t. $y \cdot b < \alpha$ and $y \cdot c > \alpha$ for all $c \in C$. Given $0 \in C$, we have $\alpha < 0$ and so $y \cdot b < \alpha < 0$. We claim $Ay \geq 0$. Suppose otherwise. Then there is some i s.t. $(Ay)_i < 0$. It follows that $e^i Ay < 0$, and so $(Me^i)'Ay \rightarrow -\infty$ as $M \rightarrow \infty$. But for every $M > 0$, $(Me^i)'A \in C$ and therefore $(Me^i)'Ay > \alpha$, yielding a contradiction. Hence if (a) does not hold, (b) must hold. \square

From this, we can derive the following variant of Farkas' lemma:

Lemma 35 (Variant Farkas' lemma). *Let A be an $m \times n$ matrix and let $b \in \mathbb{R}^m$. Then exactly one of the following statements holds:*

- (a) *There is an $x \in \mathbb{R}^m$ such that $x'A \leq b$ and $x \geq 0$;*
- (b) *There is a $y \in \mathbb{R}^n$ such that $Ay \geq 0$, $b \cdot y < 0$ and $y \geq 0$.*

Proof. First, suppose (a) holds. Fix any $y \in \mathbb{R}^n$ s.t. $b \cdot y < 0$ and suppose $Ay \geq 0$. Then $x'Ay \geq 0$ for any $x \geq 0$. But since $x'A \leq b$, we have $x'Ay = (x'A) \cdot y \leq b \cdot y < 0$, yielding a contradiction. Hence if (a) holds, (b) cannot hold, so at most one of the statements holds.

Suppose the system in (a) has no solution. Then the system

$$\begin{pmatrix} x' & \mu' \end{pmatrix} \begin{pmatrix} A \\ I \end{pmatrix} = b$$

has no solution for $x \geq 0$ and $\mu \geq 0$. By Lemma 34, it follows that system

$$\begin{pmatrix} A \\ I \end{pmatrix} y \geq 0, \quad b \cdot y < 0$$

has a solution. Thus (b) holds. \square

10.2.3 Duality theorems

Theorem 53 (Duality theorem of linear programming). *Let A be an $n \times p$ matrix, let $b \in \mathbb{R}^p$ and let $c \in \mathbb{R}^n$. Suppose $V := \{x \in \mathbb{R}^n \mid x'A \geq b, x \geq 0\}$ and $W := \{y \in \mathbb{R}^p \mid Ay \leq c, y \geq 0\}$ are both nonempty sets. Then $\min\{x \cdot c \mid x \in V\} = \max\{b \cdot y \mid y \in W\}$.*

Proof. Note that if $x \in V$ and $y \in W$, then $x \cdot c \geq x'Ay \geq b \cdot y$, since $x'A \geq b$ and $Ay \leq c$.

Next, suppose $\hat{x} \in V$ and $\hat{y} \in W$, and that $\hat{x} \cdot c = \hat{y} \cdot b$. Since $x \cdot c \geq b \cdot y$, for all $x \in V$ and $y \in W$, we have $x \cdot c \geq \hat{x} \cdot c = b \cdot \hat{y}$. Hence $\hat{x} \cdot c = \min\{x \cdot c \mid x \in V\}$. By a similar argument, $b \cdot \hat{y} = \max\{b \cdot y \mid y \in W\}$.

In view of these two results, we need only show that there exists a solution to the system

$$\begin{pmatrix} x' & y' \end{pmatrix} \begin{pmatrix} -A & 0 & c \\ 0 & A' & -b \end{pmatrix} \leq \begin{pmatrix} -b & c & 0 \end{pmatrix}, \quad x \geq 0, \quad y \geq 0.$$

Suppose otherwise. By Lemma 35, there exists a vector $(z, w, t) \in \mathbb{R}^p \times \mathbb{R}^n \times \mathbb{R}$ s.t.

$$\begin{pmatrix} -A & 0 & c \\ 0 & A' & -b \end{pmatrix} \begin{pmatrix} z \\ w \\ t \end{pmatrix} \geq 0, \quad (-b, c, 0) \cdot (z, w, t) < 0, \quad z \geq 0, \quad w \geq 0, \quad t \geq 0.$$

That is, $Az \leq tc$, $w'A \geq tb$, and $c \cdot w < b \cdot z$.

If $t = 0$, then $w'A \geq 0 \leq Az$. Thus we would have, for any $x \in V$ and $y \in W$,

$$b \cdot z \leq x'Az \leq 0 \leq w'Ay \leq w \cdot c,$$

i.e. $b \cdot z \leq w \cdot c$, yielding a contradiction. Hence consider $t > 0$. Since $Az \leq tc$ and $z \geq 0$, $\frac{1}{t}z \in W$. Likewise, since $w'A \geq tb$ and $w \geq 0$, $\frac{1}{t}w \in V$. Given this, we must have $\frac{1}{t}w \cdot c \geq b \cdot \frac{1}{t}z$. Yet this implies $w \cdot c \geq b \cdot z$, which contradicts $c \cdot w < b \cdot z$. Hence our original system must have a solution. \square

A variant duality theorem (which we use to prove the Bondareva-Shapley theorem):

Theorem 54 (Variant duality theorem). *Let A be an $n \times p$ matrix, let $b \in \mathbb{R}^p$ and let $c \in \mathbb{R}^n$. Suppose the sets $V := \{x \in \mathbb{R}^n \mid x'A \geq b\}$ and $W := \{y \in \mathbb{R}^p \mid Ay = c, y \geq 0\}$ are both nonempty. Then $\min\{x \cdot c \mid x \in V\} = \max\{b \cdot y \mid y \in W\}$.*

Proof. Define $B = \begin{pmatrix} A \\ -A \end{pmatrix}$. Note we can write W as

$$W = \{y \in \mathbb{R}^p \mid By \leq (c, -c), y \geq 0\}.$$

Define

$$V' = \{(u, w) \in \mathbb{R}^n \times \mathbb{R}^n \mid (u, w)'B \geq b, (u, w) \geq 0\}.$$

By Theorem 53, we have that

$$\min\{(u, w) \cdot (c, -c) \mid (u, w) \in V'\} = \max\{b \cdot y \mid y \in W\}.$$

Now,

$$\begin{aligned} \min\{(u, w) \cdot (c, -c) \mid (u, w) \in V'\} &= \min\{(u - w) \cdot c \mid (u - w)'A \geq b, (u, w) \geq 0\} \\ &= \min\{x \cdot c \mid x'A \geq b, x \geq 0\} \\ &= \min\{x \cdot c \mid x \in V\}. \end{aligned}$$

\square

We say that a program $\min\{x \cdot c \mid x \in V\}$ or $\max\{x \cdot c \mid x \in V\}$ for some set V and vector c is *infeasible* if the set V is empty. If V is nonempty, we say that the program is *feasible*. We say a program has an *optimal solution* if there is a point in V for which the minimum (or respectively, the maximum) is attained.

Proposition 84. *In Theorems 53 and 54, if one of the two programs is infeasible, then both programs lack an optimal solution.*

10.3 Binary relations, ordered sets and lattices

Definition 117 (Properties of binary relations). Consider a binary relation R on a set X . Note I use “ \neg ” to denote “not”.

- (a) *Basic properties.*
 - (i) *Reflexivity.* We say R is *reflexive* if xRx for all $x \in X$.
 - (ii) *Irreflexivity.* We say R is *irreflexive* if $x \neg Rx$ for all $x \in X$.
 - (iii) *Symmetry.* We say R is *symmetric* if xRy implies yRx for all $x, y \in X$.
 - (iv) *Asymmetry.* We say R is *asymmetric* if xRy implies $y \neg Rx$ for all $x, y \in X$.
 - (v) *Antisymmetry.* We say R is *antisymmetric* if xRy and yRx implies $x = y$ for all $x, y \in X$.
 - (vi) *Transitivity.* We say R is *transitive* if xRy and yRz implies xRz for all $x, y, z \in X$. If R is not transitive then it is called *intransitive*, and it is *antitransitive* if xRy and yRz always implies $x \neg Rz$.
 - (vii) *Negative transitivity.* We say R is *negatively transitive* if $x \neg Ry$ and $y \neg Rz$ implies $x \neg Rz$.
 - (viii) *Acyclicity.* We say R is *acyclic* if for any $x_1, \dots, x_n \in X$, $x_k Rx_{k+1}$ for all $k = 1, \dots, n-1$ implies $x_1 \neq x_n$.
 - (ix) *Connectedness.* We say R is *connected* (or *complete* or *total*) if for all $x, y \in X$, xRy or yRx or $x = y$.
 - (x) *Strong connectedness.* We say R is *strongly connected* if for all $x, y \in X$, xRy or yRx .
 - (xi) *Weak connectedness.* We say R is *weakly connected* if for all $x, y \in X$, $x \neq y$ implies xRy or yRx .
- (b) *Comparability.* Given $x, y \in X$, we say that x and y are *comparable* under R if xRy or yRx . Otherwise, we say x and y are *incomparable*.
- (c) *Preorder.* We say R is a (weak) *preorder* on X (often denoted \leq) if it is reflexive and transitive, and a *strict preorder* on X (often denoted $<$) if it is irreflexive and transitive.
- (d) *Partial order.* We say R is a (weak) *partial order* on X if it is an antisymmetric preorder. Again, we typically denote such a relation by \leq . We say R is a *strict partial order* on x if it is an asymmetric strict preorder, typically denoted as $<$.

If the set X equipped with the partial order \leq , we call (X, \leq) a *partially ordered set* or *poset*. In keeping with the principle that mathematical objects should be represented in the most parsimonious way, we often simply denote this by X .

- (e) *Total order.* We say R is a *total order* or *linear order* on X if it is a strongly connected partial order on X , and a *strict total order* or *strict linear order* if it is a connected strict partial order on X .

If the set X is equipped with total order \leq , we call (X, \leq) a *totally ordered set*.

- (f) *Equivalence relation.* We say R is an *equivalence relation* on X if it is a symmetric preorder on X . This binary relation is typically denoted as \sim .

If \sim is an equivalence relation on X then it partitions X into equivalence classes. Given an element $a \in X$, the *equivalence class* of a is the set $[a]_\sim := \{x \in X \mid x \sim a\}$.

Given an equivalence relation \sim on X , the *quotient set* X/\sim is the set of equivalence classes in X , that is,

$$X/\sim := \{[x]_\sim \mid x \in X\}.$$

- (g) *Subsets of posets.* Consider a poset (X, \leq) .

- (i) *Chain.* We call a subset $C \subseteq X$ a *chain* if C is totally ordered under \leq , i.e. if for all $x, y \in C$, we have $x \leq y$ or $y \leq x$.
- (ii) *Antichain.* We call a subset $C \subseteq X$ an *antichain* if no two distinct elements $x, y \in C$ are comparable under \leq .

Given a partially ordered set (L, \leq) , the *meet* $x \wedge y$ of $x, y \in L$ is the greatest lower bound (infimum) of x and y , i.e. the largest $z \in L$ such that $z \leq x, y$. The *join* $x \vee y$ of $x, y \in L$ is the least upper bound (supremum) of x and y , i.e. the smallest $z \in L$ such that $z \geq x, y$. If $L = \mathbb{R}^k$, then $x \wedge y = (\min\{x_1, y_1\}, \dots, \min\{x_k, y_k\})$ and $x \vee y = (\max\{x_1, y_1\}, \dots, \max\{x_k, y_k\})$. This generalizes to subsets. Given a subset $E \subseteq L$, we denote the meet of E by $\inf E$ and the join of E by $\sup E$.

Definition 118 (Lattices).

- (a) *Semilattice.* A partially ordered set (L, \geq) is a *meet-semilattice* if every pair of points $x, y \in L$ has a meet $x \wedge y$, and a *join-semilattice* if every pair of points has a join $x \vee y$.
- (b) *Lattice.* A partially ordered set (L, \geq) is a *lattice* if every pair of points $x, y \in L$ has a meet $x \wedge y$ and a join $x \vee y$. That is, (L, \geq) is a lattice if it is both a meet-semilattice and a join-semilattice.
- (c) *Sublattice.* A subset $M \subseteq L$ of a lattice (L, \geq) is a *sublattice* if for any $x, y \in P$, we have $x \wedge y \in M$ and $x \vee y \in P$. That is, M is a sublattice of L if (M, \geq) is a lattice.

We say that a sublattice $M \subseteq L$ is *closed* if for every $N \subseteq M$, we have $\inf N \in M$ and $\sup N \in M$.

- (d) *Completeness*. We call a lattice L *complete* if every $E \subseteq L$ has a meet $\inf E \in L$ and a join $\sup E \in L$.

Likewise, we say a sublattice A of a lattice L is *complete* if for all $B \subseteq A$, $\inf B \in A$ and $\sup B \in A$.

- (e) *Chains*. We say that a subset $C \subseteq L$ is a *chain* if C is totally ordered under \leq . That is, if for all $x, y \in C$, we have $x \leq y$ or $y \leq x$.

Lattices are a key object for monotone comparative statics, matching theory, and a number of other areas in game theory.

A useful result that will be needed to prove one of the fixed point theorems:

Lemma 36 (Zorn's lemma). *If every chain in a partially ordered set X has an upper bound, then X has a maximal element.*

We now go through some topological concepts related to lattices.

Definition 119.

- (a) *Intervals*. Suppose L is a lattice. For any $x \in L$, the *principal ideal generated by x* is the set $\downarrow x := \{z \in L \mid z \leq x\}$ and the *principal dual ideal generated by x* is the set $\uparrow x := \{z \in L \mid z \geq x\}$. A subset $I \subseteq L$ is a *closed interval* of L if:
- (i) $I = \emptyset$ or $I = L$, or
 - (ii) $I = \downarrow x$ or $I = \uparrow x$ for some $x \in L$, or
 - (iii) $I = [a, b] := \uparrow a \cap \downarrow b$ for some $a, b \in L$.
- (b) *Order interval topology*. The *order interval topology* of a lattice L is the topology that has the set of all closed intervals of L as its sub-basis. That is, the closed sets in the order interval topology are all intersections of finite unions of closed intervals.
- (c) *Order continuity*. Suppose L is a complete lattice, and consider a function $f : L \rightarrow \mathbb{R}$. Then f is said to be *order upper semicontinuous* if for every chain $C \subseteq L$, we have

$$\limsup_{x \in C, x \downarrow \inf C} f(x) \leq f(\inf C) \quad \text{and} \quad \limsup_{x \in C, x \uparrow \sup C} f(x) \leq f(\sup C).$$

The function f is said to be *order lower semicontinuous* if for every chain $C \subseteq L$, we have

$$\liminf_{x \in C, x \downarrow \inf C} f(x) \geq f(\inf C) \quad \text{and} \quad \liminf_{x \in C, x \uparrow \sup C} f(x) \geq f(\sup C).$$

The function f is said to be *order continuous* if for every chain $C \subseteq L$, we have

$$\lim_{x \in C, x \downarrow \inf C} f(x) = f(\inf C) \quad \text{and} \quad \lim_{x \in C, x \uparrow \sup C} f(x) = f(\sup C).$$

That is, f is order continuous if it is both order upper semicontinuous and order lower semicontinuous.

Proposition 85 (Frink, 1942). *If L is a complete lattice then it is compact in the order interval topology.*

The proof is not particularly illuminating.

Lemma 37. *Let L be a complete lattice. Then any monotone sequence $\{x_n\}$ in L is convergent, with limit $\lim_{n \rightarrow \infty} x_n = \sup_n x_n$ if the sequence is monotonically increasing and limit $\lim_{n \rightarrow \infty} x_n = \inf_n x_n$ if the sequence is monotonically decreasing.*

Proof. Suppose $\{x_n\}$ is monotonically increasing. Let R be the range of $\{x_n\}$ and let $x^* = \sup R$. By monotonicity, x^* is also the supremum of the range of any subsequence of $\{x_n\}$. Fix any neighbourhood V of x^* . Now, V^c is closed and since the closed intervals form a subbasis of the order interval topology, there is some collection $\{E_\alpha\}_{\alpha \in I}$, where I is some index set, s.t. each $E_\alpha := \bigcup_{k=1}^{K_\alpha} [a_k^\alpha, b_k^\alpha]$ is a finite union of closed intervals and $V^c = \bigcap_{\alpha \in I} E_\alpha$.⁷¹ Now $x^* \notin V^c$ so there must be some $\beta \in I$ for which $x^* \notin E_\beta = \bigcup_{k=1}^{K_\beta} [a_k^\beta, b_k^\beta]$. For any $k \in \{1, \dots, K_\beta\}$, there can only be at most finitely many elements of $\{x_n\}$ lying in $[a_k^\beta, b_k^\beta]$, and thus there are at most finitely many elements of $\{x_n\}$ lying in E_β . Thus there is some N s.t. $x_n \geq x_N$ for all $n \geq N$. Since V is an arbitrary neighbourhood, it follows that $x_n \rightarrow x^*$.

The proof if $\{x_n\}$ is monotonically decreasing is analogous. \square

Theorem 55. *Let X be an Euclidean space. Then any nonempty sublattice $L \subseteq X$ is complete iff L is compact.*

Proof. Let $\pi_i : X \rightarrow \mathbb{R}$ be the projection onto the i th coordinate. Suppose $L \subseteq X$ is a compact sublattice. Take $A \subseteq L$. Fix $x^i \in \arg \max\{\pi_i(x) \mid x \in \bar{A}\}$ (which is nonempty by compactness) where \bar{A} is the closure of A . Define x^* so that $x_i^* = x^i$. Since each $x^i \in L$ and L is a sublattice, we have $x^* \in L$ given $x^* = \bigvee_{i=1}^n x^i$. Now, x^* is an upper bound of A since for any $z \in A$, we have $z_i \leq x_i^* = x_i^*$ for all i , and thus $z \leq x^*$. Moreover, if z is an upper bound of A , then it is also an upper bound of \bar{A} , and $x_i^i \in \bar{A}$ so $x_i^i \leq z_i$ for all i . Hence $x^* \leq z$. Thus $x^* = \sup A$. Hence $\sup A \in L$ for all $A \subseteq L$ so L is complete.

Conversely, suppose $L \subseteq X$ is a nonempty complete sublattice. By the Heine-Borel theorem, we need only show L is closed and bounded. Completeness of L immediately implies it is bounded. Take any sequence $\{x_n\}$ in L and define $x^* = \lim_{n \rightarrow \infty} x_n$. Since L is complete, $z_k = \sup\{x_k, x_{k+1}, \dots\} \in L$ for all $k \in \mathbb{N}$ and $x^{**} = \inf\{z_1, z_2, \dots\} \in L$. Note $z_k \leq z_{k+1}$, so $x^{**} = \lim_{k \rightarrow \infty} z_k = \lim_{n \rightarrow \infty} x_n = x^*$. Hence $x^* = x^{**} \in L$. Hence L is closed. \square

10.3.1 Strong set order

When it comes to monotone comparative statics, we will want to order sets. For any space X , the subset relation \subseteq defines an obvious partial order on 2^X . But in monotone comparative statics, the *strong set order* \geq_s is convenient.

⁷¹Since L is complete, all closed intervals of L take the form $[a, b]$ for some $a, b \in L$.

Definition 120 (Strong set order). Fix a lattice (X, \geq) . For $A, B \subseteq X$, we say A is *smaller than B in the strong set order*, denoted $A \leq_s B$, if

$$x \in A \text{ and } y \in B \quad \text{implies} \quad x \wedge y \in A \text{ and } x \vee y \in B.$$

Lemma 38. *The strong set order \geq_s is antisymmetric and transitive on $2^X - \{\emptyset\}$.*

Proof. Let $A \leq_s B$ and $B \leq_s A$. Let $x \in A$ and $y \in B$. Then $x \wedge y \in A$ and $x \vee y \in B$, since $A \leq_s B$. Now $x = (x \vee y) \wedge y \in B$ and $y = y \vee (x \wedge y) \in A$, since $B \leq_s A$. Hence $A = B$.

Let $A \leq_s B$ and $B \leq_s C$. Take $x \in A$ and $y \in C$. Choose arbitrary $z \in B$. Now,

$$x \vee y = x \vee ((y \wedge z) \vee y) = (x \vee (y \wedge z)) \vee y.$$

Since $B \leq_s C$, $y \wedge z \in B$. Hence $x \vee (y \wedge z) \in B$ given $A \leq_s B$. Thus $x \vee y \in C$, since $B \leq_s C$.

Likewise,

$$x \wedge y = (x \wedge (x \vee z)) \wedge y = x \wedge ((x \vee z) \wedge y).$$

Now $x \vee z \in B$, so $(x \vee z) \wedge y \in B$. Then $x \wedge y \in A$. □

Proposition 86. *Suppose (X, \geq) is a lattice and let $L(X)$ denote the set of all nonempty sublattices of X . Then $(L(X), \geq_s)$ is a partially ordered set.*

Proof. □

10.4 Fixed point theorems

There are many fixed point theorems. We only state the most useful here. The proof of Theorem 7 relied on Kakutani's fixed point theorem. Nash's original proof of the existence of Nash equilibrium relied on Brouwer's fixed point theorem.

Recall a function $f : X \rightarrow X$ has a *fixed point* $x^* \in X$ if $f(x^*) = x^*$. The most elementary of the fixed point theorems is as follows:

Theorem 56. *Suppose $f : [0, 1] \rightarrow [0, 1]$ is a continuous function. Then f has a fixed point.*

Proof. Let $g : [0, 1] \rightarrow \mathbb{R}$ be defined by $g(x) = f(x) - x$. Clearly, g is continuous. If $g(0) = 0$ then $f(0) = 0$; if $g(1) = 0$ then $f(1) = 1$. In either of the above cases we are done. If neither of these two cases hold then $g(0) = f(0) > 0$ and $g(1) = f(1) - 1 < 0$. By the intermediate value theorem, we have that there is some $x \in [0, 1]$ s.t. $g(x) = 0$, which implies $f(x) = x$. □

The above theorem is a special case of a more general, and much more harder to prove, result:

Theorem 57 (Brouwer's fixed point theorem). *Let $K \subseteq \mathbb{R}^n$ be a nonempty compact convex subset, and suppose $f : K \rightarrow K$ is a continuous function. Then f has a fixed point.*

Kakutani's fixed point theorem generalizes Brouwer's fixed point theorem to correspondences.

Theorem 6 (Kakutani's fixed point theorem). *Let $K \subset \mathbb{R}^n$ be a nonempty, compact, convex set and suppose $F : K \rightrightarrows K$ is a correspondence satisfying*

- (i) $F(x)$ is nonempty valued;
- (ii) $F(x)$ is convex valued;
- (iii) $F(x)$ is upper hemicontinuous.

Then F has a fixed point.

Kakutani's fixed point theorem is sufficient for proving e.g. the Debreu-Glicksberg-Fan existence theorem (Theorem 7) and thus Nash's existence theorem (Theorem 8). It is also sufficient to prove the existence of most kinds of other equilibria we care about, under the right conditions on the strategy sets and so on.

This said, Kakutani's fixed point theorem has a generalizations that sometimes proves useful in game theory and economics generally:

Theorem 58 (Cellina's fixed point theorem). *Let $K \subset \mathbb{R}^n$ be a nonempty, compact, convex set and consider a correspondence $F : K \rightrightarrows K$. If there exists a compact, convex set $L \subset \mathbb{R}^m$, a correspondence $G : K \rightrightarrows L$ having a closed graph, and a continuous function $f : K \times L \rightarrow K$ such that for each $x \in K$,*

$$F(x) = \{f(x, y) \mid y \in G(x)\},$$

then F has a fixed point.

Another useful fixed point theorem is Banach's fixed point theorem, often known as the contraction mapping theorem. It is particularly important in the theory of dynamic programming.

Definition 121 (Lipschitz continuity). Let (X, ρ_X) and (Y, ρ_Y) be a metric space. We call a function $f : X \rightarrow Y$ *Lipschitz continuous* if there exists a $K > 0$ such that

$$\rho_Y(f(x), f(y)) \leq K \rho_X(x, y) \quad \text{for all } x, y \in X.$$

We call any such value K a *Lipschitz constant* of f .

Definition 122 (Contraction mapping). Let (X, ρ) be a metric space. We call an operator $T : X \rightarrow X$ a *contraction mapping* with *modulus* β if there exists a $\beta \in [0, 1)$ such that

$$\rho(Tx, Ty) \leq \beta \rho(x, y) \quad \text{for all } x, y \in X.$$

That is, T is a contraction mapping if it has a Lipschitz constant $\beta < 1$.

A contraction mapping is simply a Lipschitz continuous function with a Lipschitz constant less than one.

Theorem 59 (Banach's fixed point theorem). *Let X be a nonempty complete metric space and let $T : X \rightarrow X$ be a contraction mapping. Then T has a unique fixed point x^* . Furthermore, for any $x_0 \in X$, $x^* = \lim_{n \rightarrow \infty} T^n(x_0)$.*

Proof. Fix any point $x_0 \in X$, and let $x_n = T^n(x_0)$ for each $n \in \mathbb{N}$. Since T is a contraction mapping with modulus, say, $\beta < 1$, we have $\rho(x_n, x_{n+1}) \leq \beta \rho(x_{n-1}, x_n)$ for all $n \in \mathbb{N}$. Applying this inequality repeatedly, we have $\rho(x_n, x_{n+1}) \leq \beta^n \rho(x_0, x_1)$ for each $n \in \mathbb{N}$. Since $\beta < 1$, this implies that for any $\epsilon > 0$, there is some $N \in \mathbb{N}$ s.t. $\rho(x_n, x_m) \leq \rho(x_n, x_{n+1}) \leq \epsilon$ for all $n, m \geq N$, and thus $\{x_n\}$ is a Cauchy sequence, so $x^* = \lim_{n \rightarrow \infty} x_n$ exists by completeness. By construction, we have $Tx^* = x^* = \lim_{n \rightarrow \infty} T^n(x_0)$.

Suppose x^*, y^* are two fixed points of T . Since T is a contraction mapping $\rho(x^*, y^*) = \rho(Tx^*, Ty^*) \leq \beta \rho(x^*, y^*)$ with $\beta < 1$. This holds iff $\rho(x^*, y^*) = \rho(Tx^*, Ty^*) = 0$. Hence $x^* = y^*$. \square

There is a generalization of Banach's fixed point theorem to correspondences – Nadler's fixed point theorem – though we lose uniqueness.

This requires us to think about distances between sets.

Definition 123 (Hausdorff distance). Given a metric space X with metric ρ , we define the distance between a set $A \subseteq X$ and a point $x \in X$ by

$$\rho(x, A) = \inf_{y \in A} \rho(x, y).$$

We define the *Hausdorff distance* $\rho_H : 2^X \times 2^X \rightarrow [0, \infty)$ as

$$\rho_H(A, B) := \max \left\{ \sup_{x \in A} \rho(x, B), \sup_{y \in B} \rho(y, A) \right\}$$

for any two sets $A, B \subseteq X$.⁷²

Definition 124 (Contractiveness). Let (X, ρ) be a metric space. We call a correspondence $F : X \rightrightarrows X$ *contractive* if there exists a $\beta \in (0, 1)$ such that

$$\rho_H(F(x), F(y)) \leq \beta \rho(x, y) \quad \text{for all } x, y \in X.$$

⁷²This actually has a game-theoretic interpretation. It's just after New Year, and you've put on a few kilos (or so-called "pounds") over the holidays. Your opponent is an overly-paternalistic friend who is worried about your health. The opponent picks a point in any one of the sets A, B . You need to walk to the other set from the point they choose, but you've become pretty lazy over the holidays and want to walk the shortest distance possible to reach the other set (say you walk distance d , then your payoff is $-d$.) Your opponent thinks you need the exercise and so wants to maximize the distance you are forced to walk (their payoff is d). The Hausdorff distance is the most that your opponent can force you to walk. This is a zero-sum game, and so the Hausdorff distance is quite similar to the maximin (the sets are not necessarily closed though.) You could also interpret the opponent as being sadistic and taking pleasure in seeing you suffer. Your call.

Theorem 60 (Nadler’s fixed point theorem). *Let X be a nonempty complete metric space and let $F : X \rightrightarrows X$ be a nonempty-valued, compact-valued contractive correspondence. Then F has a fixed point x^* .*

Proof. For any point $x \in X$ and any compact set $K \subseteq X$, we have that there exists $y \in K$ s.t. $y = \arg \min_{y' \in K} \rho(x, y')$ and thus $\rho(x, y) = \rho(x, K)$. We therefore have that for any two compact sets $A, B \subseteq X$, there exists a pair of points $x \in A, y \in B$ s.t. $\rho(x, y) = \rho_H(A, B)$. Denote such a pair of points $\hat{x}(A, B), \hat{y}(A, B)$.

Fix any pair of points $x_0, x_1 \in X$. For each $n \in \mathbb{N}$, define $x_{2n} = \hat{x}(F(x_{2n-2}), F(x_{2n-1}))$ and $x_{2n+1} = \hat{y}(F(x_{2n-2}), F(x_{2n-1}))$. Since F is compact-valued, such a pair of points always exists. This yields a sequence $\{x_n\}$, with $x_{n+1} \in F(x_n)$ for each $n \in \mathbb{N}$. Since F is contractive, there is some positive $\beta < 1$ s.t. $\rho(x_{2n}, x_{2n+1}) = \rho_H(F(x_{2n-2}), F(x_{2n-1})) \leq \beta \rho(x_{2n-2}, x_{2n-1})$ for all $n \in \mathbb{N}$. Iterating this inequality gives $\rho(x_{2n}, x_{2n+1}) \leq \beta^n \rho(x_0, x_1)$ for each $n \in \mathbb{N}$. Since $\beta < 1$, it follows that for any $\epsilon > 0$, there is some $N \in \mathbb{N}$ s.t. $\rho(x_n, x_m) \leq \rho(x_n, x_{n+1}) \leq \epsilon$ for all $n, m \geq N$ and thus $\{x_n\}$ is a Cauchy sequence. By completeness of X , $\{x_n\}$ thus converges to a limit $x^* = \lim_{n \rightarrow \infty} x_n$. By construction, $x^* \in F(x^*)$. □

The Knaster-Tarski and Tarski fixed point theorems give nice results for monotonic functions. This is particularly useful because unlike the preceding fixed point theorems, we do not need continuity or “pre-continuity” in any sense. The Knaster-Tarski fixed point theorem applies to monotone functions on partially ordered sets in which every chain has a supremum. Tarski’s fixed point theorem applies to monotone functions on complete lattices. Tarski’s fixed point theorem is often important in the matching literature and in the study of supermodular games.

Theorem 61 (Knaster-Tarski fixed point theorem). *Let (X, \geq) be a partially ordered set such that every chain in X has a supremum. Let $f : X \rightarrow X$ be a monotonically increasing function with respect to \geq , and suppose there is some $a \in X$ such that $a \leq f(a)$. Then f has a fixed point and the set of fixed points of f has a maximum.*

Proof. Let $P = \{x \in X \mid x \leq f(x)\}$. Since $a \in P$, P is nonempty. Suppose C is a chain in P and $b = \sup C$. Since $c \leq b$ for all $c \in C$, we have $f(c) \leq f(b)$. Since $C \subseteq P$, $c \leq f(c)$ for all $c \in C$, and thus $f(b)$ is an upper bound for C . Since b is the least upper bound for C , we conclude $b \leq f(b)$, and thus $b \in P$. Hence the supremum of any chain in P lies in P , and so by Zorn’s lemma (Lemma 36), P has a maximal element, say x_0 .

Since $x_0 \in P$, $x_0 \leq f(x_0)$, and since f is monotonically increasing, $f(x_0) \leq f(f(x_0))$. Thus $f(x_0) \in P$. Since x_0 is maximal in P , it follows that $f(x_0) = x_0$. Hence f has a fixed point. Now, any other fixed point of f lies in P , and so x_0 is the maximum of the fixed points of f . □

Unfortunately, we cannot make a similar statement for monotonically decreasing functions on X , even if every chain in X has an infimum.

Theorem 62 (Tarski's fixed point theorem). *Let (L, \geq) be a complete lattice and suppose $f : L \rightarrow L$ is a monotonically increasing function with respect to \geq . Then the set of fixed points of f is a nonempty complete lattice under \geq .*

Proof. Let $P = \{x \in L \mid x \leq f(x)\}$ and let E be the set of fixed points of f . Since (L, \geq) is complete, there is some minimal element x_0 of L and so $x_0 \leq f(x_0)$, so P is nonempty. By the argument of the proof of the Knaster-Tarski theorem (Theorem 61), $\bar{x} := \sup P$ is a fixed point of f . Since $E \subseteq P$, we have that $\bar{x} = \sup E$. Let $Q = \{x \in L \mid x \geq f(x)\}$. Again Q is nonempty since L has a maximal element. A similar argument shows that $\underline{x} := \inf Q$ is a fixed point of f , and since $E \subseteq Q$, $\underline{x} = \inf E$.

Now we claim (E, \geq) is a complete lattice. Fix any nonempty subset $A \subseteq E$, and let \bar{a} be the supremum of A in L . Let I denote the order interval $I = [\bar{a}, 1] := \{x \in L \mid \bar{a} \leq x\}$. Then I is a complete lattice. To see this, consider any subset $B \subseteq I$. Since L is a complete lattice, B has a supremum and infimum in L . Since $B \subseteq I$, $x \geq \bar{a}$ for all $x \in B$, so \bar{a} is a lower bound of B . Hence $\inf B \geq \bar{a}$ so $\inf B \in I$, and so also $\sup B \in I$.

We claim that $f(I) \subseteq I$. If $x \in A$, then $x \leq \bar{a}$, and thus $f(x) \leq f(\bar{a})$. Yet since $A \subseteq E$, $x = f(x)$, and so $x \leq f(\bar{a})$. Thus $f(\bar{a})$ is an upper bound of A , so $\bar{a} \leq f(\bar{a})$. Since $z \geq \bar{a}$ for all $z \in I$, we have that $\bar{a} \leq f(\bar{a}) \leq f(z)$, and thus $f(z) \in I$ for all $z \in I$.

Now let \hat{f} be the restriction of f to I , and let \hat{E} be the set of fixed points of \hat{f} . Since I is a complete lattice, \hat{E} is nonempty. From the above, $\underline{z} := \inf \hat{E}$ is a fixed point of \hat{f} , and therefore of f . Since $\underline{z} \in I$, it is an upper bound of A lying in E . Moreover, it is the least upper bound of A , since \underline{z} is the infimum of \hat{E} , and thus for any other upper bound $b \in I$ of A such that b is a fixed point of f , we have that $\underline{z} \leq b$. Hence A has a least upper bound in E . A similar argument shows it has a greatest lower bound in E . Thus E is a complete lattice. \square

Finally, there is a generalization of Tarski's fixed point theorem to correspondences, due to Zhou (1994).

Definition 125 (Monotonically increasing correspondences). Given partially ordered sets X and Y , we say a correspondence $F : X \rightrightarrows Y$ is *monotonically increasing* if for any $x, x' \in X$ such that $x \geq x'$, for any $y \in F(x)$ and any $y' \in F(x')$, we have that $y \vee y' \in F(x)$ and $y \wedge y' \in F(y')$. That is, F is monotonically increasing in the strong set order, i.e. $x \geq x'$ implies $F(x) \geq_s F(x')$.

Theorem 63 (Zhou's fixed point theorem). *Let (L, \geq) be a complete lattice and suppose $F : L \rightrightarrows L$ is a monotonically increasing correspondence with respect to \geq . If $F(x)$ is a closed sublattice for every $x \in L$, then the set of fixed points of F is a nonempty complete lattice under \geq .*

Proof. Let E denote the set of fixed points of F . Let $P = \{x \in L \mid \text{there exists } y \in F(x) \text{ s.t. } y \geq x\}$. Since (L, \geq) is complete, it has some minimal element x_0 and $x_0 \leq y$ for all $y \in F(x_0)$, given F is monotonically increasing. Let $\bar{x} := \bigvee P$. Now, for any $x \in P$, there exists a $y_x \in F(x)$ with $x \leq y_x$. Since F is monotonically increasing, and $x \leq \bar{x}$, there exists a $z_x \in F(\bar{x})$ such that $x \leq y_x \leq z_x$. Define $z := \bigvee \{z_x \mid x \in P\}$.

Then $\bar{x} \leq z$, because $\bar{x} = \bigvee E \leq \bigvee \{z_x \mid x \in P\} = z$, given $x \leq z_x$ for all $x \in P$. Since $F(\bar{x})$ is a closed sublattice, $z \in F(\bar{x})$. Since F is monotonically increasing, there is some $t \in F(z)$ s.t. $t \geq z$ and so $z \in P$. Since \bar{x} is the supremum of P , $z \leq \bar{x}$. Hence $\bar{x} = z \in F(\bar{x})$, so \bar{x} is a fixed point of F . Moreover, since $E \subseteq P$, $\bar{x} = \bigvee E$.

Let $Q = \{x \in L \mid \text{there exists } y \in F(x) \text{ s.t. } y \leq x\}$. Again, this is nonempty because L has some maximal element. Now a similar argument to the above shows $\underline{x} := \bigwedge Q$ is a fixed point of F and since $E \subseteq Q$, $\underline{x} = \bigwedge E$.

Finally, we claim E is a complete lattice. Fix any nonempty subset $A \subseteq E$, and let \bar{a} be the supremum of A in L . For any $x \in A$, since $x \in F(x)$ and F is monotonically increasing, there exists $y_x \in F(\bar{a})$ s.t. $y_x \geq x$. Let $y = \bigwedge \{y_x \mid x \in A\}$. Then $y \geq \bar{a}$, and $y \in F(\bar{a})$ since $F(\bar{a})$ is a closed sublattice of L . Since F is monotonically increasing, there exists an $x_t \in F(t)$ s.t. $x_t \geq \bar{a}$ for all $t \geq a$. Let I denote the order interval $I = [\bar{a}, 1] := \{x \in L \mid \bar{a} \leq x\}$, and let \hat{F} be defined on I by $\hat{F}(t) := F(t) \cap I$ for all $t \in I$. Then by the above, $\hat{F}(t)$ is nonempty for all $t \in I$. Since for any $t \in I$, $F(t)$ and I are both closed sublattices of L , we have that $\hat{F}(t)$ is a closed sublattice of I . Define $G(t) = I$ for all $t \in I$. Then G is a monotonically increasing correspondence. Since $\hat{F}(t) = F(t) \cap G(t)$ for all $t \in I$ and F and G are monotonically increasing, \hat{F} is monotonically increasing. Let \hat{E} be the set of fixed points of \hat{F} . By the argument of the preceding paragraphs, \hat{E} is nonempty, and $\underline{z} := \bigwedge \hat{E}$ is a fixed point of \hat{F} , and thus of F . Since $\underline{z} \in I$, it is an upper bound of A lying in E , and indeed it is the least upper bound of A , since it is the infimum of \hat{E} , so A has a least upper bound in E . A similar argument shows it has a greatest lower bound in E . Thus E is a complete lattice. \square

Again, there is unfortunately no converse for “descending” correspondences.

10.5 Envelope theorems

In parameterized optimization problems, the envelope theorems are important tools for comparative statics. A relatively general “traditional” envelope theorem is the following.

Theorem 64. *Let $X \subseteq \mathbb{R}^n$ be compact, let $\Theta \subseteq \mathbb{R}^m$ be an open set, and let $f : X \times \Theta \rightarrow \mathbb{R}$ be continuous on $X \times \Theta$. Moreover, suppose $f(x, \theta)$ is continuously differentiable in θ for each fixed x and that the gradient of f with respect to θ is continuous on $X \times \Theta$. Define*

$$v(\theta) = \sup_{x \in X} f(x, \theta)$$

for each $\theta \in \Theta$.⁷³ If, for a given $\hat{\theta}$, there is a unique maximizer $\hat{x} \in \arg \max_{x \in X} f(x, \hat{\theta})$, then v is differentiable at $\hat{\theta}$ and

$$\frac{\partial v(\hat{\theta})}{\partial \theta_i} = \frac{\partial f(\hat{x}, \hat{\theta})}{\partial \theta_i}.$$

⁷³We call v the *upper envelope* of the objective function f . It is also of course the value function of the optimization problem.

Proof. By Berge's theorem (Theorem 45), v is continuous on Θ and

$$x^*(\theta) := \arg \max_{x \in X} f(x, \theta) = \{x \in X \mid f(x, \theta) = v(\theta)\}$$

is non-empty-valued and has a closed graph, so is upper hemicontinuous by the closed graph theorem (Theorem 44). From upper hemicontinuity, it follows that if $\{\hat{\theta}^n\}$ is any sequence s.t. $\hat{\theta}^n \rightarrow \hat{\theta}$ and if $\hat{x}^n \in x^*(\hat{\theta}^n)$ for each n , then $\hat{x}^n \rightarrow \hat{x}$. We mean to show that

$$\lim_{n \rightarrow \infty} \frac{v(\hat{\theta}^n) - v(\hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta})}{\|\hat{\theta}^n - \hat{\theta}\|} = 0,$$

where $\nabla_{\theta} f(\hat{x}, \hat{\theta})$ denotes the gradient of f in θ at $(\hat{x}, \hat{\theta})$.

Now, note that $v(\hat{\theta}^n) = f(\hat{x}^n, \hat{\theta}^n) \geq f(\hat{x}, \hat{\theta}^n)$, and $v(\hat{\theta}) = f(\hat{x}, \hat{\theta}) \geq f(\hat{x}^n, \hat{\theta})$. Hence

$$v(\hat{\theta}^n) - v(\hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta}) = f(\hat{x}^n, \hat{\theta}^n) - f(\hat{x}, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta})$$

for each n and,

$$\begin{aligned} f(\hat{x}^n, \hat{\theta}^n) - f(\hat{x}^n, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta}) &\geq v(\hat{\theta}^n) - v(\hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta}) \\ &\geq f(\hat{x}, \hat{\theta}^n) - f(\hat{x}, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta}). \end{aligned}$$

Now, by the mean value theorem, $f(\hat{x}^n, \hat{\theta}^n) - f(\hat{x}^n, \hat{\theta}) = (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}^n, \tilde{\theta}^n)$, where $\tilde{\theta}^n := \lambda \hat{\theta}^n + (1 - \lambda) \hat{\theta}$ with $\lambda \in (0, 1)$ is some convex combination of $\hat{\theta}^n$ and $\hat{\theta}$. Whence

$$\begin{aligned} \left| f(\hat{x}^n, \hat{\theta}^n) - f(\hat{x}^n, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta}) \right| &= \left| (\hat{\theta}^n - \hat{\theta}) \cdot (\nabla_{\theta} f(\hat{x}^n, \tilde{\theta}^n) - \nabla_{\theta} f(\hat{x}, \hat{\theta})) \right| \\ &\leq \left\| \nabla_{\theta} f(\hat{x}^n, \tilde{\theta}^n) - \nabla_{\theta} f(\hat{x}, \hat{\theta}) \right\| \cdot \left\| \hat{\theta}^n - \hat{\theta} \right\| \end{aligned}$$

This gives

$$\lim_{n \rightarrow \infty} \left| \frac{f(\hat{x}^n, \hat{\theta}^n) - f(\hat{x}^n, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta})}{\|\hat{\theta}^n - \hat{\theta}\|} \right| \leq \lim_{n \rightarrow \infty} \left\| \nabla_{\theta} f(\hat{x}^n, \tilde{\theta}^n) - \nabla_{\theta} f(\hat{x}, \hat{\theta}) \right\| = 0,$$

given $(\hat{x}^n, \hat{\theta}^n) \rightarrow (\hat{x}, \hat{\theta})$ and $\nabla_{\theta} f$ is continuous on $X \times \Theta$. An almost identical argument establishes that

$$\lim_{n \rightarrow \infty} \left| \frac{f(\hat{x}, \hat{\theta}^n) - f(\hat{x}, \hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta})}{\|\hat{\theta}^n - \hat{\theta}\|} \right| \leq \lim_{n \rightarrow \infty} \left\| \nabla_{\theta} f(\hat{x}, \tilde{\theta}^n) - \nabla_{\theta} f(\hat{x}, \hat{\theta}) \right\| = 0.$$

Hence by the sandwich theorem,

$$\lim_{n \rightarrow \infty} \frac{v(\hat{\theta}^n) - v(\hat{\theta}) - (\hat{\theta}^n - \hat{\theta}) \cdot \nabla_{\theta} f(\hat{x}, \hat{\theta})}{\|\hat{\theta}^n - \hat{\theta}\|} = 0.$$

□

This is but one variation of a traditional envelope theorem. The traditional envelope theorems all have a particular form. For ease of exposition, let us restrict to $\Theta = [0, 1]$. Take X to be arbitrary, and consider $f : X \times \Theta \rightarrow \mathbb{R}$. Define

$$v(\theta) = \sup_{x \in X} f(x, \theta),$$

$$X^*(\theta) = \{x \in X \mid f(x, \theta) = v(\theta)\}.$$

The traditional envelope theorems give sufficient conditions for v to be differentiable and show the derivative observes first order condition formula $v'(\theta) = \frac{\partial f(x, \theta)}{\partial \theta}$ for all $\theta \in \Theta$.

Milgrom & Segal (2002) note that the traditional envelope theorems can be extended to the case where X is an arbitrary choice set. This is the “modern envelope theorem”.

Theorem 65 (Milgrom-Segal, 2002). *Let $f : X \times [0, 1] \rightarrow \mathbb{R}$ and define v and X^* as above. Fix $\theta \in [0, 1]$ and $x^* \in X^*(\theta)$ and suppose $\frac{\partial f(x^*, \theta)}{\partial \theta}$ exists. Then*

- (i) *if $\theta > 0$ and v is left-differentiable at x^* , then $v'(\theta-) \leq \frac{\partial f(x^*, \theta)}{\partial \theta}$;*
- (ii) *if $\theta < 1$ and v is right-differentiable at x^* , then $v'(\theta+) \geq \frac{\partial f(x^*, \theta)}{\partial \theta}$;*
- (iii) *if $\theta \in (0, 1)$ and v is differentiable at x^* then*

$$v'(\theta) = \frac{\partial f(x^*, \theta)}{\partial \theta}.$$

Proof. By definition of v and X^* , for any $\theta' \in [0, 1]$ we have that

$$f(x^*, \theta') - f(x^*, \theta) \leq v(\theta') - v(\theta).$$

Taking $\theta' \in (\theta, 1)$, dividing through by $\theta' - \theta > 0$, and taking limits gives $\frac{\partial f(x^*, \theta)}{\partial \theta} \leq \lim_{\theta' \downarrow \theta} \frac{v(\theta') - v(\theta)}{\theta' - \theta} = v'(\theta+)$, assuming the latter exists.

Taking $\theta' \in (0, \theta)$, dividing through by $\theta - \theta' > 0$ and taking limits gives $\frac{\partial f(x^*, \theta)}{\partial \theta} \geq \lim_{\theta' \uparrow \theta} \frac{v(\theta') - v(\theta)}{\theta - \theta'} = v'(\theta-)$, assuming the latter exists.

Finally, if $\theta \in (0, 1)$ and v is differentiable at x^* then $v'(\theta) = v'(\theta+) = v'(\theta-) = \frac{\partial v(x^*, \theta)}{\partial \theta}$. □

This is a powerful result, which generalizes most of the preceding envelope theorems. For example, the Benveniste-Scheinkman theorem (1979), which is widely applied in infinite horizon consumption problems, is a special case of Theorem 65 when the objective function $f(x, \theta)$ is concave in both x and θ and X is convex.

Theorem 65 applies only if the value function is sufficiently well-behaved. The next two theorems give some sufficient conditions for v to be sufficiently well-behaved.

Definition 126.

- (a) *Absolute continuity.* Let I denote an interval in \mathbb{R} . A function $g : I \rightarrow \mathbb{R}$ is *absolutely continuous* on I if for every $\epsilon > 0$, there is a $\delta > 0$ such that for every sequence of pairwise disjoint subintervals $\{(a_k, b_k)\}_{k=0}^{\infty}$ of I such that $\sum_{k=0}^{\infty} (b_k - a_k) < \delta$, we have that

$$\sum_{k=0}^{\infty} |g(b_k) - g(a_k)| < \epsilon.$$

- (b) *Equidifferentiability.* Let $f : X \times [0, 1] \rightarrow \mathbb{R}$. The family of functions $\{f(x, \cdot)\}_{x \in X}$ is *equidifferentiable* at $\theta \in [0, 1]$ if $\frac{f(x, \theta') - f(x, \theta)}{\theta' - \theta}$ converges uniformly as $\theta' \rightarrow \theta$.

Theorem 66 (Milgrom-Segal envelope theorem, 2002). *Suppose $f(x, \cdot)$ is absolutely continuous for all $x \in X$. Suppose there exists a Lebesgue integrable function $g : [0, 1] \rightarrow \mathbb{R}$ such that $\left| \frac{\partial f(x, \theta)}{\partial \theta} \right| \leq g(\theta)$ for all $x \in X$ almost everywhere on $[0, 1]$. Moreover, suppose $f(x, \cdot)$ is differentiable for all $x \in X$ and $X^*(\theta)$ is nonempty almost everywhere. Then for any selection $x^*(\theta) \in X^*(\theta)$,*

$$v(\theta) = v(0) + \int_0^{\theta} \frac{\partial f(x^*(\bar{\theta}), \bar{\theta})}{\partial \theta} d\bar{\theta}.$$

Proof. By definition of v , for any $\theta', \theta'' \in [0, 1]$ s.t. $\theta' < \theta''$, we have

$$\begin{aligned} |v(\theta'') - v(\theta')| &\leq \sup_{x \in X} |f(x, \theta'') - f(x, \theta')| \\ &= \sup_{x \in X} \left| \int_{\theta'}^{\theta''} \frac{\partial f(x, \theta)}{\partial \theta} d\theta \right| \leq \int_{\theta'}^{\theta''} \sup_{x \in X} \left| \frac{\partial f(x, \theta)}{\partial \theta} \right| d\theta \leq \int_{\theta'}^{\theta''} g(\theta) d\theta. \end{aligned}$$

It follows that v is absolutely continuous, so differentiable a.e. and $v(\theta) = v(0) + \int_0^{\theta} v'(\bar{\theta}) d\bar{\theta}$. By hypothesis, $f(x, \theta)$ is differentiable in θ and thus $v'(\bar{\theta}) = \frac{\partial f(x^*(\bar{\theta}), \bar{\theta})}{\partial \theta}$ by Theorem 65. Substituting into the integral gives the result. \square

Theorem 66 is often referred to as the “integral form” of the envelope theorem. It is a very useful result – see Milgrom & Segal (2002) for a wide range of applications.

Milgrom & Segal also provide sufficient conditions for the left- and right-derivatives of v to exist:

Theorem 67. *Suppose the family of functions*

$$\left\{ \frac{\partial f(x, \cdot)}{\partial \theta} \right\}_{x \in X}$$

is equidifferentiable at $\theta_0 \in [0, 1]$, suppose

$$\sup_{x \in X} \left| \frac{\partial f(x, \cdot)}{\partial \theta} \right| < +\infty,$$

and suppose $X^*(\theta)$ is nonempty for all $\theta \in [0, 1]$. Then v is left- and right-differentiable at θ_0 . For any selection $x^* \in X^*(\theta)$, the directional derivatives are

$$v'(\theta+) = \lim_{\theta \downarrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta} \text{ for } \theta_0 < 1, \text{ and}$$

$$v'(\theta-) = \lim_{\theta \uparrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta} \text{ for } \theta_0 > 0.$$

Moreover, v is differentiable at θ_0 iff $\frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta}$ is continuous in θ at $\theta = \theta_0$.

Proof. Given the definition of v and finiteness of $\sup_{x \in X} \left| \frac{\partial f(x, \cdot)}{\partial \theta} \right|$, equidifferentiability implies

$$|v(\theta) - v(\theta_0)| \leq \sup_{x \in X} |f(x, \theta) - f(x, \theta_0)| \leq \sup_{x \in X} \left| \frac{\partial f(x, \theta_0)}{\partial \theta} \right| \cdot |\theta - \theta_0| + o(\theta - \theta_0) \rightarrow 0$$

as $\theta \rightarrow \theta_0$. Hence $\{f(x, \cdot)\}_{x \in X}$ is equicontinuous at θ_0 and v is continuous at θ_0 . Fix $\theta_0 < \theta' < \theta''$. By definition of v , we have

$$\frac{f(x^*(\theta'), \theta'') - f(x^*(\theta'), \theta')}{\theta'' - \theta'} \leq \frac{v(\theta'') - v(\theta')}{\theta'' - \theta'} \leq \frac{f(x^*(\theta''), \theta'') - f(x^*(\theta''), \theta')}{\theta'' - \theta'}.$$

Using equicontinuity and the continuity of v at θ_0 , we get

$$\limsup_{\theta' \downarrow \theta_0} \frac{f(x^*(\theta'), \theta'') - f(x^*(\theta'), \theta_0)}{\theta'' - \theta_0} \leq \frac{v(\theta'') - v(\theta_0)}{\theta'' - \theta_0} \leq \frac{f(x^*(\theta''), \theta'') - f(x^*(\theta''), \theta_0)}{\theta'' - \theta_0}.$$

By equidifferentiability, this implies

$$\limsup_{\theta' \downarrow \theta_0} \frac{\partial f(x^*(\theta'), \theta_0)}{\partial \theta} + \frac{o(\theta'' - \theta_0)}{\theta'' - \theta_0} \leq \frac{v(\theta'') - v(\theta_0)}{\theta'' - \theta_0} \leq \frac{\partial f(x^*(\theta''), \theta_0)}{\partial \theta} + \frac{o(\theta'' - \theta_0)}{\theta'' - \theta_0}.$$

Similarly, taking the limit inferior of the bounds as $\theta'' \downarrow \theta_0$ gives $\limsup_{\theta \downarrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta} \leq \liminf_{\theta \downarrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta}$, so the two must be equal and thus $\lim_{\theta \downarrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta}$ exists.

Applying the sandwich theorem, we get that $v'(\theta+) = \lim_{\theta \downarrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta}$ for $\theta_0 < 1$.

The proof of the claim that $v'(\theta-) = \lim_{\theta \uparrow \theta_0} \frac{\partial f(x^*(\theta), \theta_0)}{\partial \theta}$ for $\theta_0 > 0$ is similar.

Finally, v is differentiable at $\theta_0 \in (0, 1)$ iff $v'(\theta_0-) = v'(\theta_0+) = \frac{\partial f(x^*(\theta_0), \theta_0)}{\partial \theta}$, by Theorem 65. From the limit results we just established, it follows that $\frac{\partial f(x^*(\theta), \theta)}{\partial \theta}$ is continuous at $\theta = \theta_0$. \square

The intuition behind the traditional envelope theorems is that it follows from the first order condition. This is not true of the modern envelope theorems due to Milgrom & Segal (2002), because they extend to settings where the first order condition is not necessarily well-defined. Sinander (2022) proves a “converse envelope theorem” that reestablishes a first order condition intuition behind the envelope theorem.

10.6 Monotone comparative statics

Monotone comparative statics concern the large class of economic problems where endogenous variables in an optimization problem are monotone in exogenous variables. This is important in many applications – games of strategic complements or substitutes, auctions, etc.

10.6.1 Supermodularity and increasing differences

Supermodularity and submodularity need a little lattice theory – we detail some basics in section 10.3. We discuss the fixed point theorems for lattices in section 10.4.

Definition 127 (Increasing differences). Suppose X and Θ are partially ordered sets.

- (a) *Increasing differences.* A function $f : X \times \Theta \rightarrow \mathbb{R}$ is said to have *increasing differences* in θ , if for all $x, x' \in X$ such that $x' \geq x$ and all $\theta, \theta' \in \Theta$ such that $\theta' \geq \theta$, we have that

$$f(x', \theta') - f(x, \theta') \geq f(x', \theta) - f(x, \theta),$$

i.e. $f(x', \theta) - f(x, \theta)$ is nondecreasing in θ .

We say that f has *strictly increasing differences* if this condition holds with strict inequality for all $\theta' > \theta$.

- (b) *Decreasing differences.* A function $f : X \times \Theta \rightarrow \mathbb{R}$ is said to have *decreasing differences* in θ if for all $x, x' \in X$ such that $x' \geq x$ and all $\theta, \theta' \in \Theta$ such that $\theta' \geq \theta$, we have that

$$f(x', \theta') - f(x, \theta') \leq f(x', \theta) - f(x, \theta),$$

i.e. $f(x', \theta) - f(x, \theta)$ is nonincreasing in θ .

We say that f has *strictly decreasing differences* if this condition holds with strict inequality for all $\theta' > \theta$.

Definition 128 (Supermodularity). Suppose (X, \geq) is a lattice and Θ is a partially ordered set.

- (a) *Supermodularity.* A function $f : X \rightarrow \mathbb{R}$ is said to be *supermodular* in x if for all $x, x' \in X$, we have

$$f(x) + f(x') \leq f(x \wedge x') + f(x \vee x').$$

If this inequality is strict whenever x, x' cannot be compared in the order \geq , then we say f is *strictly supermodular*.

We say a function $f : X \times \Theta \rightarrow \mathbb{R}$ is (strictly) *supermodular in x* ($x \in X$) if $x \mapsto f(x, \theta)$ is (strictly) supermodular for all $\theta \in \Theta$.

- (b) *Submodularity*. A function $f : X \rightarrow \mathbb{R}$ is said to be (strictly) *submodular* if $-f$ is (strictly) supermodular. Likewise, $f : X \times \Theta \rightarrow \mathbb{R}$ is (strictly) submodular in x if $-f$ is (strictly) submodular in x .

Supermodularity is a cardinal property: while affine transformations of supermodular functions are supermodular, but in general, monotone transformations are not.

In the context of a game $G = (\mathcal{I}, (S_i, u_i)_{i \in \mathcal{I}})$, for a utility function u_i , we take $X = S_i$ and $\Theta = S_{-i}$. If the utility function u_i has the increasing differences property in s_i , then the incremental gain to choosing a higher value of s_i is weakly greater if s_{-i} is higher (with respect to our given partial order). Supermodularity captures the notion of complementary inputs and submodularity captures the notion of substitutable inputs.

It turns out that supermodularity and increasing differences are equivalent for real functions on \mathbb{R}^k :

Proposition 87. *A function $f : \mathbb{R}^k \rightarrow \mathbb{R}$ is supermodular iff f has increasing differences.*

Proof. The implication follows by definition. For the converse, we have that for any $x, y \in \mathbb{R}^k$,

$$\begin{aligned} f(x \vee y) - f(y) &= \sum_i [f(x_1 \vee y_1, \dots, x_i \vee y_i, y_{i+1}, \dots, y_k) \\ &\quad - f(x_1 \vee y_1, \dots, x_{i-1} \vee y_{i-1}, y_i, \dots, y_k)] \\ &= \sum_i [f(x_1 \vee y_1, \dots, x_{i-1} \vee y_{i-1}, x_i, y_{i+1}, \dots, y_k) \\ &\quad - f(x_1 \vee y_1, \dots, x_{i-1} \vee y_{i-1}, x_i \wedge y_i, y_{i+1}, \dots, y_k)] \\ &\geq \sum_i [f(x_1, \dots, x_{i-1}, x_i, x_{i+1} \wedge y_{i+1}, \dots, x_k \wedge y_k) \\ &\quad - f(x_1, \dots, x_{i-1}, x_i \wedge y_i, x_{i+1} \wedge y_{i+1}, \dots, x_k \wedge y_k)] \\ &= f(x) - f(x \wedge y). \end{aligned}$$

□

An immediate corollary is:

Corollary 14. *A function $f : \mathbb{R}^k \rightarrow \mathbb{R}$ is submodular in $x \in \mathbb{R}^k$ iff f has decreasing differences in x .*

If $f : X \times \Theta \rightarrow \mathbb{R}$ is twice continuously differentiable, then supermodularity is equivalent to f_i having a nonnegative cross partial derivative:

Theorem 68 (Topkis' characterization theorem). *Suppose the function $f : X \times \Theta \rightarrow \mathbb{R}$ is twice continuously differentiable, where $X \subseteq \mathbb{R}^k$ and $\Theta \subseteq \mathbb{R}^\ell$. Then f is supermodular in (x, θ) iff $\frac{\partial f(x_i, \theta_j)}{\partial x_i \partial \theta_j} \geq 0$ for all $(x, \theta) \in X \times \Theta$, $i = 1, \dots, k$ and $j = 1, \dots, \ell$.*

The most important result in monotone comparative statics is the monotonicity theorem (Topkis, 1978):

Theorem 69 (Topkis' monotonicity theorem). *Let X be a lattice and Θ be a partially ordered set. Suppose $f : X \times \Theta \rightarrow \mathbb{R}$ is supermodular in $x \in X$ and has increasing differences in both $x \in X$ and $\theta \in \Theta$. Then the set of maximizers $M(\theta) := \arg \max_{x \in X} f(x, \theta)$ is a monotonically increasing correspondence.*

Proof. Since $x' \in M'$, we have $f(x', \theta') - f(x \wedge x', \theta') \geq 0$. There are two cases:

- (i) $x' \geq x$. Then this becomes $f(x', \theta') - f(x, \theta') \geq 0$. Since $x' = x \vee x'$, we have $f(x \vee x', \theta') - f(x, \theta') \geq 0$.
- (ii) $x' \leq x$. Then $x = x \vee x'$, so $f(x' \vee x, \theta') - f(x, \theta') = f(x, \theta') - f(x, \theta) = 0$.

From the two cases, we establish that $f(x \vee x', \theta) - f(x, \theta) \geq 0$. By supermodularity, we thus have $f(x \vee x', \theta) - f(x', \theta') \geq 0$, implying $x \vee x' \in M$. Since $x \in M$, we must have $f(x, \theta) - f(x \vee x', \theta) \geq 0$, equivalent to $f(x \vee x', \theta) - f(x, \theta) \leq 0$. By supermodularity, $f(x \vee x', \theta') - f(x, \theta') \leq 0$. Again, there are two cases:

- (i) $x' \geq x$. Then this becomes $f(x', \theta') - f(x, \theta') \leq 0$ and $x = x \wedge x'$ so $f(x', \theta') - f(x \wedge x', \theta') \leq 0$.
- (ii) $x' \leq x$. Then $x = x \wedge x'$, so $f(x', \theta') - f(x \wedge x', \theta') = 0$.

From the two cases, we have that $f(x', \theta') - f(x \wedge x', \theta') \leq 0$. Thus $x \wedge x' \in M'$. \square

Milgrom & Roberts (1990) prove the following:

Theorem 70 (Milgrom-Roberts). *Let L be a complete lattice and $f : L \rightarrow \mathbb{R}$ is an order upper semicontinuous, supermodular function, then the set of maximizers of f is a nonempty complete sublattice of L .*

Proof. Let M be the set of maximizers of f . We do not prove M is nonempty – see Theorem 1 in Milgrom & Roberts (1990) for such a proof. Now Topkis' monotonicity theorem (Theorem 69) implies M is a sublattice of L . Fix an arbitrary subset $A \subseteq L$. Let $\hat{M} = M \cap \{x \mid x \leq \sup A\}$. By order upper semicontinuity, every chain $C \subseteq \hat{M}$ is s.t. $\sup C \in \hat{M}$. By Zorn's lemma (Lemma 36), \hat{M} has a maximal element \hat{x} , and \hat{M} is a sublattice so $\hat{x} = \sup \hat{M}$. Now by construction, $\sup \hat{M} = \sup A \in \hat{M} \subseteq M$. A similar argument gives $\inf A \in M$. Thus M is complete. \square

10.6.2 Supermodular capacities and probability measures

Definition 129 (Capacity). Given a finite set Ω , a *capacity* is a monotone increasing function $\nu : 2^\Omega \rightarrow \mathbb{R}$ such that $\nu(\emptyset) = 0$ and $\nu(\Omega) = 1$.

- (i) *Convexity*. We call a capacity *convex* if it is supermodular.
- (ii) *Core*. Given a capacity ν on Ω , the *core* of ν is the set

$$C(\nu) = \{p \in \Delta(\Omega) \mid p \geq \nu\},$$

where $\Delta(\Omega)$ denotes the set of probability measures on Ω .⁷⁴

Like a probability measure, a capacity is a nonnegative set function such that $\nu(\Omega) = 1$ and $\nu(A) \leq \nu(B)$ whenever $A \subseteq B$. However, a capacity does not need to be countably (or finitely) additive.

One can think of a capacity as representing a probability assessment. The relaxation of additivity allows for the flexibility that agents might approach different situations with different priors.⁷⁵

Let $N(\Omega)$ be the set of capacities on Ω . For each $A \subseteq \Omega$, define the difference operator δ_A by $(\delta_A \nu)(B) = \nu(A \cup B) - \nu(B - A)$.

Proposition 88. *A capacity ν is convex iff $\delta_A(\delta_B \nu) \geq 0$ for all $A, B \subseteq \Omega$.*

Proof. Fix $A, B \subseteq \Omega$. For any $C \subseteq \Omega$, we have

$$\begin{aligned} (\delta_A \delta_B \nu)(C) &= (\delta_B \nu)(A \cup C) - (\delta_B \nu)(C - A) \\ &= \nu(A \cup C \cup B) - \nu((A \cup C) - B) - (\nu(C - A) \cup B) + \nu((C - A) - B). \end{aligned}$$

Fix $D, E \subseteq \Omega$. Take $C = D \cap E$, $A = D - C$ and $B = E - C$. Then $A \cup B \cup C = D \cup E$, $(A \cup C) - B = D$, $(C - A) \cup B = C$, and $(C - A) - B = D \cap E$. Hence $(\delta_A \delta_B \nu)(C) \geq 0$ iff $\nu(D \cup E) - \nu(D) - \nu(C) + \nu(D \cap E) \geq 0$. \square

Proposition 89. *Let ν be a capacity on Ω . Then ν is convex iff the core $C(\nu)$ is nonempty.*

Proof. We only prove the only if direction. Suppose ν is convex. Let $\{\omega_1, \dots, \omega_K\}$ be an enumeration of Ω and define a probability measure $p \in \Delta(\Omega)$ as follows: take $p(\{\omega_1\}) = \nu(\{\omega_1\})$ and $p(\{\omega_k\}) = \nu(\{\omega_1, \dots, \omega_k\}) - \nu(\{\omega_1, \dots, \omega_{k-1}\})$ for $k = 2, \dots, K$, and for any $A \subseteq \Omega$ let $p(A) = \sum_{\omega_k \in A} p(\{\omega_k\})$. Then p is clearly additive, nonnegative and has $p(\emptyset) = 0$. Now, $p(\Omega) = \sum_{k=1}^K p(\{\omega_k\}) = \nu(\{\omega_1\}) + \sum_{k=2}^K [\nu(\{\omega_1, \dots, \omega_k\}) - \nu(\{\omega_1, \dots, \omega_{k-1}\})] = \nu(\Omega) = 1$. Hence $p \in \Delta(\Omega)$.

Fix any nonempty $A \subseteq \Omega$, and let ω_k be the first element of the enumeration of Ω for which $\omega_k \notin A$. Then $A \cup \{\omega_1, \dots, \omega_k\} = A \cup \{\omega_k\}$ and $A \cap \{\omega_1, \dots, \omega_k\} = \{\omega_1, \dots, \omega_{k-1}\}$. Hence

$$p(\{\omega_k\}) = \nu(\{\omega_1, \dots, \omega_k\}) - \nu(\{\omega_1, \dots, \omega_{k-1}\}) \leq \nu(A \cup \{\omega_k\}) - \nu(A),$$

giving $p(A) + p(\{\omega_k\}) \leq \nu(A \cup \{\omega_k\}) - \nu(A) + p(A)$, or equivalently, $p(A \cup \{\omega_k\}) - \nu(A \cup \{\omega_k\}) \leq p(A) - \nu(A)$. Repeating the argument for the first subsequent element ω_ℓ not in $A \cup \{\omega_k\}$ and so on gives us that

$$\begin{aligned} 0 &= p(\Omega) - \nu(\Omega) \leq p(A \cup \{\omega_k, \omega_\ell\}) - \nu(A \cup \{\omega_k, \omega_\ell\}) \\ &\leq p(A \cup \{\omega_k\}) - \nu(A \cup \{\omega_k\}) \\ &\leq p(A) - \nu(A). \end{aligned}$$

⁷⁴As usual, $p \geq \nu$ reads as $p(A) \geq \nu(A)$ for all $A \subseteq \Omega$. Note the similarity to the definition of the core in cooperative games.

⁷⁵If we want to generalize von Neumann-Morgenstern preferences to allow for more general uncertainty aversion, then this is one way to go about this. See Schmeidler (1989).

Since A was arbitrary, we have $p(A) \geq \nu(A)$ for all $A \subseteq \Omega$. Thus $p \in C(\nu)$. \square

Proposition 90. *If ν is a convex capacity then $\nu = \inf C(\nu)$.*

Proof. Recall that $C(\nu) = \{p \in \Delta(\Omega) \mid p \geq \nu\}$. Hence ν is a lower bound of $C(\nu)$. If $\nu \neq \inf C(\nu)$ then there is some other capacity μ that is a lower bound for $C(\nu)$ s.t. $\mu(A) \geq \nu(A)$ for all $A \subseteq \Omega$ with strict inequality for some A .

Fix any nonempty $A \subseteq \Omega$. Recall the construction of p in Proposition 89 ensured $p \in C(\nu)$, for any arbitrary enumeration $\{\omega_1, \dots, \omega_K\}$ of Ω . Consider an enumeration $\{\omega_1, \dots, \omega_K\}$ such that $A = \{\omega_1, \dots, \omega_{K_A}\}$, i.e. the elements of A are the first K_A elements of the enumeration, and construct p as in Proposition 89. Then

$$p(A) = \sum_{k=1}^{K_A} p(\{\omega_k\}) = \nu(\{\omega_1\}) + \sum_{k=2}^{K_A} [\nu(\{\omega_1, \dots, \omega_k\}) - \nu(\{\omega_1, \dots, \omega_{k-1}\})] = \nu(A).$$

Hence for any $A \subseteq \Omega$, there is a probability measure $p \in C(\nu)$ s.t. $p(A) = \nu(A)$. Thus if $\nu \leq \mu \leq p$ for all $p \in C(\nu)$ we must have that $\nu = \mu$, a contradiction. \square

We can define Lebesgue-like integrals with respect to capacities. For any function $f : \Omega \rightarrow \mathbb{R}$ and any capacity ν on Ω , we can define

$$\int f \, d\nu = \int_{-\infty}^0 (\nu\{\omega \mid f(\omega) \geq t\} - 1) \, dt + \int_0^{\infty} \nu\{\omega \mid f(\omega) \geq t\} \, dt.$$

If ν is convex, then we equivalently have

$$\int f \, d\nu = \min \left\{ \int f \, dp \mid p \in C(\nu) \right\}.$$

Briefly turning to probability measures:

Proposition 91. *Let $\Delta(A)$ denote the set of Borel probability measures on $A \subseteq \mathbb{R}$. For any $p \in \Delta(A)$ and any bounded measurable function $h : A \rightarrow \mathbb{R}$, the mapping*

$$p \mapsto \int_A h \, dp$$

is supermodular (and submodular).

Proof. For $p \in \Delta(A)$, let F_p denote the cdf of p . Let $A_p := \{x \in A \mid F_p(x) \leq F_q(x)\}$ and let $A_q = A - A_p$. We have

$$\begin{aligned} \int_A h \, dF_{p \vee q} + \int_A h \, dF_{p \wedge q} &= \int_{A_p} h \, dF_p + \int_{A_q} h \, dF_q + \int_{A_p} h \, dF_q + \int_{A_q} h \, dF_p \\ &= \int_A h \, dF_p + \int_A h \, dF_q. \end{aligned}$$

\square

10.6.3 Quasi-supermodularity and single crossing properties

Recall supermodularity is a cardinal property, yet optimization relies on ordinal properties: if an agent's preferences admit a supermodular utility representation, they will also admit non-supermodular utility representations, but the representation does not matter for the optimal choices of the agent. Our lives are made much simpler if we have an ordinal theory of supermodularity. Fortunately, Milgrom and Shannon (1994) have developed such a theory.

Definition 130 (Single crossing property). A function $f : X \times \Theta \rightarrow \mathbb{R}$ satisfies the *single crossing property* in (x, θ) if for all $x, x' \in X$ and all $\theta, \theta' \in \Theta$, we have

$$\begin{aligned} f(x, \theta) \leq f(x', \theta) & \text{ implies } f(x, \theta') \leq f(x', \theta'), \text{ and} \\ f(x, \theta) < f(x', \theta) & \text{ implies } f(x, \theta') < f(x', \theta') \end{aligned}$$

whenever $x \leq x'$ and $\theta \leq \theta'$. Moreover, we say f satisfies the *strict single crossing property* if $f(x, \theta) \leq f(x', \theta)$ implies $f(x, \theta') < f(x', \theta')$ whenever $x \leq x'$ and $\theta \leq \theta'$.

The single crossing property can be thought of as an ordinal generalization of increasing differences:

Proposition 92. *If $f : X \times \Theta \rightarrow \mathbb{R}$ has (strict) increasing differences in (x, θ) then it satisfies the (strict) single crossing property.*

Proof. If f has increasing differences then for any $x \leq x'$ and $\theta \leq \theta'$, we have

$$f(x', \theta') - f(x, \theta') \geq f(x', \theta) - f(x, \theta),$$

If the rhs is nonnegative (positive) then clearly the lhs must be nonnegative (positive), completing the proof. Under strict increasing differences, the above inequality holds strictly and so if the rhs is nonnegative then the lhs must be positive. \square

Similarly, quasisupermodularity is the ordinal generalization of supermodularity.

Definition 131 (Quasisupermodularity). Let (X, \geq) be a lattice. A function $f : X \rightarrow \mathbb{R}$ is called *quasisupermodular* if for all $x, x' \in X$,

$$\begin{aligned} f(x \wedge x') \leq f(x') & \text{ implies } f(x) \leq f(x \vee x'), \text{ and} \\ f(x \wedge x') < f(x') & \text{ implies } f(x) < f(x \vee x'). \end{aligned}$$

Proposition 93. *If f is supermodular then f is quasisupermodular.*

Proof. If $f : X \rightarrow \mathbb{R}$ is supermodular then for all $x, x' \in X$,

$$f(x) + f(x') \leq f(x \wedge x') + f(x \vee x').$$

If $f(x \wedge x') \leq f(x')$ then $\epsilon := f(x') - f(x \wedge x') \geq 0$, and so the inequality becomes $f(x) \leq f(x \vee x') - \epsilon \leq f(x \vee x')$, with the final inequality holding strictly for $\epsilon > 0$. \square

Chambers & Echenique (2009) prove that quasisupermodularity is the appropriate generalization of supermodularity:

Theorem 71 (Chambers & Echenique, 2009). *Let (X, \geq) be a finite lattice. Then a function $f : X \rightarrow \mathbb{R}$ is monotonically increasing and quasisupermodular iff there is some strictly increasing function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $g \circ f : X \rightarrow \mathbb{R}$ is monotonically increasing and supermodular.*

Proof. Suppose there is a strictly increasing function g s.t. $h = g \circ f$ is monotonically increasing and supermodular. Monotonicity of h implies f must be monotonically increasing given g is strictly increasing. Moreover, monotonicity of h implies $h(x \wedge x') \leq h(x')$ and $h(x) \leq h(x \vee x')$, and since g is strictly increasing, it follows that $f(x \wedge x') \leq f(x')$ and $f(x) \leq f(x \vee x')$. If $f(x \wedge x') < f(x')$, then because g is strictly increasing, $h(x \wedge x') < h(x')$, and since $h(x) + h(x') \leq h(x \wedge x') + h(x \vee x')$, we have $h(x) < h(x \vee x')$, and again since g is strictly increasing we conclude $f(x) < f(x \vee x')$. Hence f is quasisupermodular.

Suppose $f : X \rightarrow \mathbb{R}$ is monotonically increasing and quasisupermodular, and, given X is finite, let $f(X) = \{z_1, \dots, z_K\}$ with $z_k < z_{k+1}$. The function $g : f(X) \rightarrow \mathbb{R}$ defined by $g(z_k) = 2^{k-1}$ is strictly increasing, and hence $h = g \circ f$ is monotonically increasing. Moreover, h is quasisupermodular. Take $x, y \in X$. If x, y are ordered in \geq , then $h(x) + h(y) \leq h(x \wedge y) + h(x \vee y)$, so we have nothing to prove. Hence suppose x, y cannot be compared in the order \geq . Then monotonicity implies $h(x \wedge y) \leq h(x)$ and $h(y) \leq h(x \vee y)$. If $h(x \wedge y) = h(x)$ then it follows from $h(y) \leq h(x \vee y)$ that $h(x) + h(y) \leq h(x \wedge y) + h(x \vee y)$. If $h(y) = h(x \vee y)$, quasisupermodularity of h gives us that $h(x \wedge y) = h(x)$, so $h(x) + h(y) = h(x \wedge y) + h(x \vee y)$. Now suppose $h(x \wedge y) < h(x)$ and $h(y) < h(x \vee y)$. Choose k s.t. $f(x \vee y) = z_k$. Then

$$h(x \wedge y) + h(x \vee y) \geq 2^{k-1} = 2^{k-2} + 2k - 2 \geq h(x) + h(y).$$

Hence h is supermodular. □

Corollary 15. *Let X be a finite lattice. If $f : X \rightarrow \mathbb{R}$ is strictly increasing then there exists a function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $g \circ f$ is strictly increasing and supermodular.*

Proof. If f is strictly increasing then it is quasisupermodular, and proof follows from Theorem 71. □

Milgrom & Shannon (1994) generalize Topkis' monotonicity theorem (Theorem 69) to this setting.

Theorem 72 (Milgrom-Shannon monotonicity theorem). *Let X be a lattice and Θ be a partially ordered set, and let $M : \Theta \times 2^X \rightrightarrows \mathbb{R}$ be the correspondence $M(\theta, S) := \arg \max_{x \in S} f(x, \theta)$. Then $M(\theta)$ is a monotonically increasing correspondence iff $f(x, \theta)$ is quasisupermodular in x and f satisfies the single crossing property in (x, θ) .*

When we say that M is a monotonically increasing correspondence here, we mean that if $\theta \leq \theta'$ and $S \leq_s S'$ (in the strong set order \leq_s) and if $x \in M(\theta, S)$ and $x' \in M(\theta', S')$, then $x \vee x' \in M(\theta', S')$ and $x \wedge x' \in M(\theta, S)$.

Proof. First suppose $f(x, \theta)$ is quasisupermodular in x and f satisfies the single crossing property in (x, θ) . Let $S \leq_s S'$, $\theta \leq \theta'$, $x \in M(\theta, S)$ and $x' \in M(\theta', S')$. Since $x \in S$ and $S \leq_s S'$, we have $f(x, \theta) \geq f(x \wedge x', \theta)$. Quasisupermodularity implies $f(x \vee x', \theta) \geq f(x', \theta)$, and by the single crossing property, it follows that $f(x \vee x', \theta') \geq f(x', \theta')$. Hence $x \vee x' \in M(\theta', S')$. Similarly, since $x' \in M(\theta', S')$ and $S \leq_s S'$, we have $f(x', \theta') \geq f(x \vee x', \theta')$ so $f(x \vee x', \theta') - f(x', \theta') \leq 0$, and the single crossing property then implies $f(x \vee x', \theta) - f(x', \theta) \leq 0$. Now quasisupermodularity gives $f(x \wedge x', \theta) \geq f(x, \theta)$, and thus $x \wedge x' \in M(\theta, S)$.

Conversely, suppose M is monotonically increasing. Take θ fixed and $x, x' \in X$. Let $S = \{x, x \wedge x'\}$ and $S' = \{x, x \vee x'\}$. Then $S \leq_s S'$. If $f(x, \theta) \geq f(x \wedge x', \theta)$, then $x \in M(S, \theta)$ and thus $f(x \vee x', \theta) \geq f(x', \theta)$, with strict inequality if the former inequality is strict. Hence f is quasisupermodular in x . Next, take $S = \{x, \tilde{x}\}$ for some $\tilde{x} \geq x$. Now, $f(\tilde{x}, \theta) - f(x, \theta) \geq 0$ implies $\tilde{x} \in M(\theta, S) \leq_s M(\tilde{\theta}, S)$ for $\tilde{\theta} \geq \theta$, and so $f(\tilde{x}, \tilde{\theta}) - f(x, \tilde{\theta}) \geq 0$ for all $\tilde{\theta} \geq \theta$, with strict inequalities implied if $f(\tilde{x}, \theta) - f(x, \theta) > 0$. Hence f has the single crossing property. \square

10.6.4 Infinite supermodularity

Sometimes it is useful to generalize the notion of supermodularity beyond comparison of pairs of points in a lattice X to comparison across all finite sets of points. This motivates *infinite supermodularity*.

Definition 132 (Infinite supermodularity). Let (X, \geq) be a lattice. A function $f : X \rightarrow \mathbb{R}$ is said to be *infinitely supermodular* if for all $n \geq 2$ and all $x_1, \dots, x_n \in X$, we have

$$\sum_{\emptyset \neq I \subseteq \{1, \dots, n\}} (-1)^{|I|+1} f\left(\bigwedge_{i \in I} x_i\right) \leq f\left(\bigvee_{i \in I} x_i\right).$$

For $n = 2$ and any $x_1, x_2 \in X$, this is $f(x_1) + f(x_2) - f(x_1 \wedge x_2) \leq f(x_1 \vee x_2)$, which is exactly supermodularity. Hence infinite supermodularity is a much stronger property. Infinite supermodularity allows us to strengthen Theorem 71:

Theorem 73 (Chateauneuf, Vergopoulos & Zhang, 2017). *Let (X, \geq) be a finite lattice. Then a function $f : X \rightarrow \mathbb{R}$ is monotonically increasing and quasisupermodular iff there is some strictly increasing function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $g \circ f : X \rightarrow \mathbb{R}$ is monotonically increasing and infinitely supermodular.*

Index

- Agent-normal form, 55
- Aumann's agreement theorem, 23
- Bach or Stravinsky, 63
- Bargaining
 - Kalai-Smorodinsky solution, 214
 - Nash solution, 209
 - Rubinstein model, 170
 - Ståhl model, 167
- Battle of the Sexes, 63
- Bayes Nash equilibrium, 95
- Bayesian games, 93
- Bayesian persuasion, 186
- Behavioural strategy, 14
- Berge's theorem of the maximum, 253
- Bertrand competition, 92
- Best response dynamics, 76
- Best response graph, 76
- Birkhoff-von Neumann theorem, 259
- Capacity, 280
 - Convexity of, 281
 - Core of, 280
- Cheap talk, 184
- Closed graph theorem, 253
- Concave games, 51
- Congestion games, 89
- Core, 219
 - And competitive equilibrium, 230
 - Conditions for nonemptiness, 226
 - Noncooperative implementation, 227
- Correlated equilibrium, 63
 - Correlated equilibrium distributions, 66
 - Correlating mechanisms, 64
 - Direct mechanisms, 65
 - Epistemic foundations, 76
 - Subjective correlated equilibrium, 71
- Cournot competition, 86, 96
- Economics of identity, 149
- Envelope theorems
 - Classical envelope theorem, 273
 - Converse envelope theorem, 277
 - Milgrom-Segal, 275
- Equilibrium selection
 - Payoff dominance, 62
 - Risk dominance, 62
- Essentially strict equilibrium, 102
- Ex post equilibrium, 99
- Extensive form games, 5
- Farkas' lemma, 261
- Fictitious play, 78
- Fixed point theorems
 - Banach's, 270
 - Brouwer's, 268
 - Cellina's, 269
 - Kakutani's, 47, 269
 - Knaster-Tarski, 271
 - Nadler's, 271
 - Tarski's, 272
 - Zhou's, 272
- Focal points, 61
- Forward induction, 160
 - Burning money, 163
 - In signalling games, 180
- Game tree, 5
- Global games, 113
- Groves mechanisms, 125
- Harsanyi purification, 100
 - Harsanyi's purification theorem, 102
 - Payoff-perturbed game, 102
- Imputation, 216
- Iterated strict dominance, 35
- Iterated weak dominance, 39
- Keynesian beauty contest, 40

- Knowledge
 - Axioms, 20
 - Common knowledge, 21
 - Mutual knowledge, 21
- Knowledge operator, 19
- Krein-Milman theorem, 258
- Kuhn's theorem, 16
- Level-k rationality, 35
- Market entry, 11, 94, 139, 149, 153, 160
- Matching pennies, 33, 47, 77, 80, 82
- Matrix games, 29
- Minimax theorems, 31
- Mixed strategy, 13
- Monotone comparative statics
 - Increasing differences, 278
 - Infinite supermodularity, 285
 - Milgrom-Shannon monotonicity theorem, 284
 - Quasisupermodularity, 283
 - Single crossing, 283
 - Supermodularity, 278
 - Topkis' characterization theorem, 279
 - Topkis' monotonicity theorem, 280
- Nash equilibrium, 41
 - Nash existence theorem, 49
 - Debreu-Glicksberg-Fan existence theorem, 48
 - Efficiency, 47
 - Epistemic foundations, 74
 - Interpretations, 42
 - Rosen's uniqueness theorem, 53
 - Upper hemicontinuity, 50
 - Weakly dominated Nash, 46
- No trade theorems, 24
- One-shot deviation principle, 150
- Paradox of the absent-minded driver, 17
- Pareto efficiency, 27, 72
- Perfect Bayesian equilibrium, 153
 - Strong perfect Bayesian equilibrium, 155
- Perfect recall, 15
- Potential games
 - Isomorphism with congestion games, 89
 - Ordinal potential, 83
 - Potential, 83
- Price of anarchy, 28
- Price of stability, 28
- Prisoner's dilemma, 10, 77, 86
- Proper equilibrium, 59
- Pure strategy, 9
- Quasi-strict equilibrium, 101
- Rationalizability, 36
 - Coincidence with iterated strict dominance, 38
- Regular equilibrium, 101
- Revelation principle
 - In Bayes Nash equ'm strategies, 120
 - In correlating mechanisms, 66
 - In dominant strategies, 118
 - In ex post equ'm strategies, 119
- Rock-paper-scissors, 80
- Rubinstein's electronic mail game, 111
- Self-confirming equilibrium, 81
- Selten's horse, 82
- Sequential equilibrium, 158
 - Existence, 160
 - Upper hemicontinuity, 159
- Sequential rationality, 139
 - Backward induction, 140
 - Zermelo's theorem, 142
- Shapley value, 234
 - Myerson value, 249
- Signalling games, 173
 - Beer-quiche game, 181
 - Divine equilibrium, 183
 - Hybrid equilibrium, 179
 - Intuitive criterion, 180
 - Pooling equilibrium, 179

- Separating equilibrium, [177](#)
- Spence model, [176](#)
- Social choice, [129](#)
 - Arrow's impossibility theorem, [133](#)
 - Borda rule, [131](#)
 - Condorcet's paradox, [131](#)
 - Majority rule, [131](#)
 - Scoring rules, [130](#)
 - Wilson's theorem, [135](#)
- Stable equilibria, [164](#)
- Stable set, [223](#)
- Stackelberg competition, [144](#)
- Strict dominance, [34](#)
 - In Bayesian games, [98](#)
- Strong Nash equilibrium, [72](#)
- Subgame perfect equilibrium, [146](#)
- Submodular games, [93](#)
- Supermodular games, [90](#)
- Supermodularity of probability measures, [282](#)
- Taxation principle, [123](#)
- Transferable utility games, [215](#)
- Trembling hand perfection, [54](#)
 - Perturbed games, [54](#)
- Type spaces, [105](#)
 - Belief hierarchies, [107](#)
 - Canonical embedding, [110](#)
 - Universal type space, [107](#)
- Upper hemicontinuity, [251](#)
- Weak dominance, [39](#)
 - In Bayesian games, [98](#)
- Zero sum games, [29](#)