

## Noncooperative games and strategic behavior : rationality, dominance, iterated dominance, and Nash equilibrium

Crusoe versus Crusoe is not really a game, just two independent decision problems; and we don't need any theory to predict that rational players will choose (T, L).

	L	R
T	2 2	1 1
B	2 1	1 1

**Crusoe vs Crusoe**

**Prisoner's Dilemma** is a game, because players' decisions affect each other's payoffs; but we still don't need a new theory to predict that rational players will choose (Confess, Confess). ("Confess" = "Defect"; "Don't" = "Cooperate".)

(The "s" in Prisoner's Dilemma (not "s")) signals methodological individualism.)

	Don't	Confess
Don't	3 3	5 0
Confess	0 5	1 1

**Prisoner's Dilemma**

A *strictly dominant* decision is a decision that yields a player strictly higher payoff, no matter which decision(s) the other player(s) choose.

E.g. T for Row or L for Column in Crusoe vs Crusoe, or Confess for either player in Prisoner's Dilemma.

	Don't	Confess
Don't	3 3	5 0
Confess	0 5	1 1

**Prisoner's Dilemma**

A rational player must choose a strictly dominant decision if he has one.

A *strictly dominated* decision is a decision that yields a player strictly lower payoff than another feasible decision, no matter which decisions the others choose.

E.g. B for Row or R for Column in Crusoe vs Crusoe, or Don't in Prisoner's Dilemma.

	Don't	Confess
Don't	3 3	0 5
Confess	5 0	1 1

**Prisoner's Dilemma**

A rational player will never play a strictly dominated decision, because there are no beliefs about other players' decisions that make it a best response.

Although it doesn't happen in Crusoe vs Crusoe or Prisoner's Dilemma, there can be dominated decisions without a dominant decision, which makes the notion of dominated decision more useful than the notion of dominant decision.

Because of the way the prisoners' payoffs interact, individually rational decisions yield a collectively suboptimal (Pareto-inefficient, in the prisoners' view) outcome.

Note that what's Pareto-efficient for the prisoners need not be good for society.

	Don't	Confess
Don't	3 3	0 5
Confess	5 0	1 1

**Prisoner's Dilemma**

Prisoner's Dilemma highlights a flaw in libertarianism: an enforceable law against confessing would make both prisoners better off, while limiting their freedom.

(The grain of truth in libertarianism is that it would be beneficial to be the *only one* allowed to break the law.

But that can't yield a universal rule by which to organize society, "universal" as in: "Act only according to that maxim whereby you can, at the same time, will that it should become a universal law."—Kant)

Prisoner's Dilemma's tension between individual rationality and Pareto-efficiency makes it the simplest possible model of incentive problems, which makes it a popular platform for the analyses of institutions that overcome such problems.



The positive flip side of my caveat about modeling situations as games, e.g.:

“If you object to my game analysis on the grounds that players don’t really have to play ‘my’ game, my (only!) remedy is to add to my game’s rules a player’s decision whether to participate, and then to insist that that decision be explained by the same behavioural assumptions as players’ other decisions.”

This insistence is an important constraint on analysis: otherwise there is nothing to pin down the assumptions implicit in speculative “solutions” to problems.

Later in these lectures we will see examples of how repeated interaction can sometimes support cooperation despite incentive problems in the short run.

Yet a Prisoner's Dilemma model of incentive problems is too simple: it ignores the difficulty of coordination and conflicts between different ways to cooperate.

Pigs in a Box, Row (R) is a big (“dominant”) pig and Column (C) a little (“subordinate”) pig. The box is a “Skinner box”, named for B.F. Skinner.

- Pushing a lever at one end of the box yields 10 units of grain at the other. Pushing “costs” either pig 2 units of grain.
- If R pushes while C waits, C can eat 5 units before R comes and shoves C aside.
- If C pushes while R waits, C cannot shove R aside and R gets all but one unit of grain.
- If both push and then arrive at the grain together, C gets 3 units and R gets 7.
- If both wait, both get 0.

	Push	Wait
Push	5      1	3      5
Wait	9      -1	0      0

**Pigs in a Box**

Here rational strategic behavior is more subtle, in that for the first time, it requires at least one player to predict the other’s response to the game.

Its consequences are also a bit surprising:

- R can do anything C can do, which in an individual decision problem would ensure that R does better.
- But in the lab pigs tend to settle down at (R Push, C Wait): C does better!
- In games, evidently, (the right kind of) weakness might be an advantage

Recall that the structure is assumed for simplicity, to be *common knowledge*.



• A player is *rational* (in the decision-theoretic sense) if he maximizes his payoff given *beliefs* (subjective probability distributions) about other players' decisions that are not inconsistent with anything he knows.

• A player who faces uncertainty about the consequences of his decisions is *rational* if he maximizes his expected payoff (*vN-M utility*).

first guess at how to formalize the idea of rational decisions in games is that assuming that players are rational will suffice for a useful theory of behavior.

That guess is correct for games like Crusoe v. Crusoe and Prisoner's Dilemma.

But that guess fails badly in slightly more complex games, such as Pigs in a Box.

A second guess is that assuming that players are rational *and* that that fact is common or at least mutual knowledge is enough to yield a useful theory.

That guess works in some games, in which common knowledge of rationality yields a unique prediction.

But even that guess fails badly in many economically interesting games.

A third guess is that assuming that players are rational and that players' decisions are best responses to *correct* beliefs about others' decisions (which must then be the same for all players) is enough to yield a useful theory.

That guess leads to the idea of Nash equilibrium, which is the leading behavioral assumption in noncooperative game theory; but even it has some drawbacks.

see how common or mutual knowledge of rationality works, imagine that the pigs are as good at reasoning about others' responses to incentives as (some) humans seem to be.

They can then use rationality and knowledge of others' rationality—in this case mutual knowledge is enough—to figure out they should play (R Push, C Wait).

If they have mutual knowledge of rationality, the reasoning goes as follows:

	Push	Wait
Push	5 1	3 5
Wait	9 -1	0 0

**Pigs in a Box**

• No rational C will choose Push.

• Therefore no rational R who knows that C is rational will play Wait.

• (R Push, C Wait) is the only possible outcome.

This incentive effect is what turns R's greater strength into a weakness.



R *might* do better if he can change the game in a way that gives C an incentive to Push, at least some of the time; e.g. by committing himself to giving C more grain if C Pushed. (There's still a coordination problem; more below.)

Understanding which kinds of games commitments help in, and what kinds of commitments help, should help us to understand the usefulness of contracts and other ways to change how relationships are governed.

(As legal “persons”, corporations have the “right” to be sued. This is a “right”, not simply a liability, because it may allow mutually beneficial contracts that would not be in the other party's interests if it could not sue for breach.)

Pigs are probably not really as good as humans at reasoning about others' likely decisions. So why do they still tend to settle down at (R Push, C Wait)?

	Push	Wait
Push	5, 1	3, 5
Wait	9, -1	0, 0

**Pigs in a Box**

- In repeated play, because Push is strictly dominated for player C, it must do worse on average for C than Wait.
- Thus even a C that reacts unthinkingly to rewards will “learn” to choose Wait with higher and higher probability over time.
- Once the probability that C chooses Wait is high enough ( $> 4/7$ ), player R will learn to choose Push with higher and higher probability over time.
- They will eventually settle down at (R Push, C Wait): Learning yields the same outcome in the limit as rationality-based reasoning does: a general result

we can characterize the implications of common knowledge of rationality via

- *Iterated deletion* of strictly dominated decisions (often called “iterated strict dominance”): eliminating strictly dominated decisions for one or both players, then eliminating decisions that become strictly dominated once players' strictly dominated decisions are eliminated, and so on ad infinitum.

If iterated strict dominance reduces the game to a single decision for each player—as in Pigs in a Box eliminating Push for C and then Wait for R reduces the game to (Push, Wait)—the game is said to be *dominance-solvable*.

	<b>Push</b>	<b>Wait</b>
<b>Push</b>	5      1	3      5
<b>Wait</b>	9      -1	0      0

**Pigs in a Box**

Moreover, a rational player with sharply focused beliefs need not choose a weakly dominant decision, and might choose a weakly dominated decision.

Iterated strict dominance is linked to common knowledge of rationality via the notion of *rationalizable* decisions.

- A *rationalizable* decision is one that survives iterated elimination of *never*

(*weak*) *best responses*, those decisions that are not even tied for being a best response to any beliefs.

The set of rationalizable decisions can't be larger than the set that survive iterated strict dominance, because strictly dominated decisions can never be weak best responses (that is why the notion builds on *strict* dominance).

In a two-person game, a rationalizable decision is exactly one that survives iterated strict dominance; well defined because the latter is order-independent.

In games with more than two players, the two notions are not quite the same because players can have correlated beliefs about others' strategies.

eg., M and C are the only rationalizable decisions in (how support them?):

	<b>L</b>	<b>C</b>	<b>R</b>
<b>T</b>	7      0	0      5	0      3
<b>M</b>	5      0	2      2	5      0
<b>B</b>	0      7	0      5	7      3

**Dominance-solvable game**

But any decisions are rationalizable in (how support them?):

	L	C	R
T	7, 0	0, 5	0, 7
M	5, 0	2, 2	5, 0
B	0, 7	0, 5	7, 0

**Unique equilibrium but no dominance**

### Nash Equilibrium

Most economically interesting games have multiple rationalizable outcomes, so players' decisions are not *dictated* by common knowledge of rationality, and the guess that it will yield a useful theory of strategic behavior fails badly.

- To make sharper predictions, noncooperative game theory assumes that

players' decisions are in *Nash equilibrium*, that is, that each player's decision maximizes his payoff or expected payoff, given the others' decisions.

Any equilibrium decision is rationalizable (why?).

It can be shown that an equilibrium always exists in non-pathological games.

Therefore in a dominance-solvable game, players' unique rationalizable decisions are in equilibrium (why?).

non-dominance-solvable games, however, equilibrium also effectively requires that players' decisions are best responses to *correct* beliefs about others' decisions, which must then be the same for all players, e.g.:

	L	C	R
T	7, 0	0, 5	0, 7
M	5, 0	2, 2	5, 0
B	0, 7	0, 5	7, 0

**Unique equilibrium but no dominance**

Nash equilibrium is a kind of "rational expectations" equilibrium, in that if players are rational, and all expect the same decisions and best respond to those beliefs, then their beliefs are self-confirming if and only if they are in Nash equilibrium.





This goes far beyond rationality, or even common knowledge of rationality.

Why might players have correct beliefs about each other's decisions?

There are two possible justifications, which generalize those mentioned in connection with Pig in a Box.

- **Thinking:** If players are rational and have perfect models of each other's

decisions, strategic thinking will lead them to have the same beliefs immediately, and so play an equilibrium even in their initial responses to a game.

- **Learning:** Even without perfect models, if players are rational and repeatedly play analogous games, experience will eventually allow them to predict each others' decisions well enough to play an equilibrium in the game that is repeated.

### mixed strategies

In game theory it is useful to extend the idea of decision, or strategy, from the unrandomized (*pure*) notion to allow randomized (*mixed*) decisions or strategies.

E.g. Matching Pennies has no appealing pure strategies, but there is an appealing way to play using mixed strategies: randomizing 50-50. (Why?)

	Heads	Tails
Heads	1, -1	-1, 1
Tails	-1, 1	1, -1

**Matching Pennies**

Our definitions apply to mixed as well as pure strategies, if the uncertainty mixed strategies cause is handled as for other kinds of uncertainty, by assigning payoffs to outcomes so that rational players maximize their expected payoffs.

Mixed strategies ensure that “well-behaved” games always have rational-expectations strategy combinations: i.e. that *Nash equilibria* always exist

### Nonuniqueness of Equilibrium and Coordination

	Go	Wait
Go	0, 0	1, 1
Wait	1, 1	0, 0

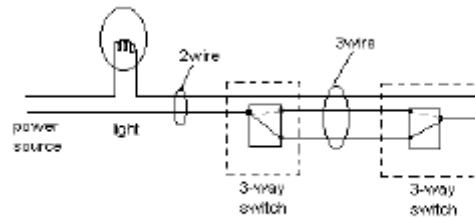
**Alphonse and Gaston**

	Fights	Ballet
Fights	2, 1	0, 0
Ballet	0, 0	1, 2

**Battle of the Sexes**



In the early 1900s Frederick B. Oppen created the comic strip *Alphonse and Gaston*, with two excessively polite fellows saying "after you, my dear Gaston" or "after you, my dear Alphonse" and thus never getting through a doorway. Alphonse and Gaston live on in the dual-control lighting circuits in our homes.



Alphonse and Gaston's problem is that there are *two* good ways to solve their coordination problem...and therefore maybe no good way.

Each way requires them to decide differently, but the setting provides no clue to break the symmetry of their roles.

Battle of the Sexes—the simplest possible bargaining problem—adds to the difficulty of coordination by giving players different preferences about *how* to coordinate, but still no clue about how to break the symmetry.

These games are popular platforms for the analyses of institutions that overcome such problems, e.g. via conventions that use labels to break the symmetry of players' roles, such as "defer to short people" or "defer to women".

In Stag Hunt (Rousseau's story), with two or  $n$  players, there are two symmetric, Pareto-ranked, pure-strategy equilibria, "all-Stag" and "all-Rabbit".

	Other Player	
	Stag	Rabbit
Stag	2, 2	0, 1
Rabbit	1, 0	1, 1

**2-person Stag Hunt**

	All Other Players	
	Stag	Rabbit
Stag	2, 0	0, 0
Rabbit	1, 1	1, 1

**$n$ -Person Stag Hunt**

All-Stag is better for all than all-Rabbit: Kant would have no trouble here.

But Stag is riskier in that unless all others play Stag, a player does better with Rabbit.

Stag Hunt is like a choice between autarky and participating in a highly productive but brittle society, which is more rewarding but riskier because productivity depends on perfect coordination.

Stag Hunt is a special case of Larry Summers's Bank Runs example:

"A crude but simple game, related to Douglas Diamond and Philip Dybvig's [1983 *JPE*] celebrated analysis of bank runs, illustrates some of the issues involved here. Imagine that everyone who has invested \$10 with me can expect to earn \$1, assuming that I stay solvent.



Suppose that if I go bankrupt, investors who remain lose their whole \$10 investment, but that an investor who withdraws today neither gains nor loses. What would you do? Each individual judgment would presumably depend on one's assessment of my prospects, but this in turn depends on the collective judgment of all of the investors.

Suppose, first, that my foreign reserves, ability to mobilize resources, and economic strength are so limited that if any investor withdraws I will go bankrupt. It would be a Nash equilibrium (indeed, a Pareto-dominant one) for everyone to remain, but (I expect) not an attainable one. Someone would reason that someone else would decide to be cautious and withdraw, or at least that someone would reason that someone would reason that someone would withdraw, and so forth. This...would likely lead to large-scale withdrawals, and I would go bankrupt. It would not be a close-run thing. ...Keynes's beauty contest captures a similar idea.

Now suppose that my fundamental situation were such that everyone would be paid off as long as no more than one-third of the investors chose to withdraw. What would you do then? Again, there are multiple equilibria: everyone should stay if everyone else does, and everyone should pull out if everyone else does, but the more favorable equilibria seems much more robust."

—Lawrence Summers, "International Financial Crises: Causes, Prevention, and Cures," (2000 *AER*).

The game Summers describes can be represented by a payoff table as follows:

		Summary statistic	
		In	Out
Representative player	In	1	-10
	Out	0	0
Bank Runs			

In Summers's first example, all investors must stay In to prevent the bank from collapsing, so the summary statistic takes the value In if and only if all but the representative player stay In.

In Summers's second example, two-thirds of the investors need to stay In, so the summary statistic takes the value In if and only if that is the case, adding in the representative player.

In each example there are two pure-strategy equilibria: "all-In" and "all-Out".

In this simplified static model, what happens depends on players' initial responses to the game as shaped by their strategic thinking: specifically, which equilibrium's basin of attraction, "all-In" or "all-Out", the initial responses fall into.

The leading models of initial responses for games like this are Harsanyi and Selten's (1988) notions of *payoff-dominance* and *risk-dominance*.

Payoff-dominance favors equilibria that are Pareto-superior to other equilibria.

Hence here it selects the all-In equilibrium, for any value of the population size  $n$  and deviation cost (here, the -10).

But that seems behaviorally unlikely, even for small  $n$  and "small -10".



Risk-dominance favors the equilibrium with (roughly) the largest “basin of attraction” in beliefs space.

		Summary statistic	
		In	Out
Representative player	In	1	-10
	Out	0	0
		Bank Runs	

In games like this one, that turns out to be the same as selecting the equilibrium that results if each player best responds to a uniform prior over others’ decisions.

Assuming independence of others’ decisions, with these payoffs risk-dominance favors the all-Out equilibrium for any  $n$ , even if only two-thirds need to stay In.

That again seems behaviorally unlikely for small  $n$ .