



COOPERATION: HOW TO MODEL IT, HOW TO FOSTER IT, AND HOW IT MIGHT HAVE EMERGED

GAME THEORY 101 STAGS, PRISONERS AND EQUILIBRIA

Adrian Haret

a.haret@uva.nl

November 9-16, 2023

Game theory is about interactions
between independent, self-interested
agents.

We start with how agents quantify their preferences and take decisions.

DECISIONS, DECISIONS

How should a rational agent make decisions?



SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX





SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX

[...] no matter what Bankman-Fried was doing, he was constantly assessing the odds, costs, and benefits.

• • •

Faux, Z. (2023). *Number Go Up: Inside Crypto's Wild Rise and Staggering Fall*. Crown Currency.



SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX

[...] no matter what Bankman-Fried was doing, he was constantly assessing the odds, costs, and benefits.

Any decision could be boiled down to an “expected value,” [...] whether that was a move in a board-game marathon, a billion-dollar trade, or whether to chat with Bezos at a party.

● ● ●

Faux, Z. (2023). *Number Go Up: Inside Crypto’s Wild Rise and Staggering Fall*. Crown Currency.



SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX

[...] no matter what Bankman-Fried was doing, he was constantly assessing the odds, costs, and benefits.

Any decision could be boiled down to an “expected value,” [...] whether that was a move in a board-game marathon, a billion-dollar trade, or whether to chat with Bezos at a party.

Bankman-Fried’s goal was always to make as much money as possible, so that he could give it to charity.

● ● ●

Faux, Z. (2023). *Number Go Up: Inside Crypto’s Wild Rise and Staggering Fall*. Crown Currency.



SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX

[...] no matter what Bankman-Fried was doing, he was constantly assessing the odds, costs, and benefits.

Any decision could be boiled down to an “expected value,” [...] whether that was a move in a board-game marathon, a billion-dollar trade, or whether to chat with Bezos at a party.

Bankman-Fried’s goal was always to make as much money as possible, so that he could give it to charity.

By this metric, even sleep was an unjustifiable luxury. The expected value of staying awake to trade was too high.

Faux, Z. (2023). *Number Go Up: Inside Crypto’s Wild Rise and Staggering Fall*. Crown Currency.



SAM BANKMAN-FRIED





SAM BANKMAN-FRIED

We should always think about *expected values*.



ZEKE FAUX

[...] no matter what Bankman-Fried was doing, he was constantly assessing the odds, costs, and benefits.

Any decision could be boiled down to an “expected value,” [...] whether that was a move in a board-game marathon, a billion-dollar trade, or whether to chat with Bezos at a party.

Bankman-Fried’s goal was always to make as much money as possible, so that he could give it to charity.

By this metric, even sleep was an unjustifiable luxury. The expected value of staying awake to trade was too high.

Faux, Z. (2023). *Number Go Up: Inside Crypto’s Wild Rise and Staggering Fall*. Crown Currency.



SAM BANKMAN-FRIED

Every minute you spend sleeping is costing you x -thousand dollars, and that directly means you can save this many less lives.

Let's look at a concrete example.

ADRIAN

I'm taking the train from Brussels to Munich.





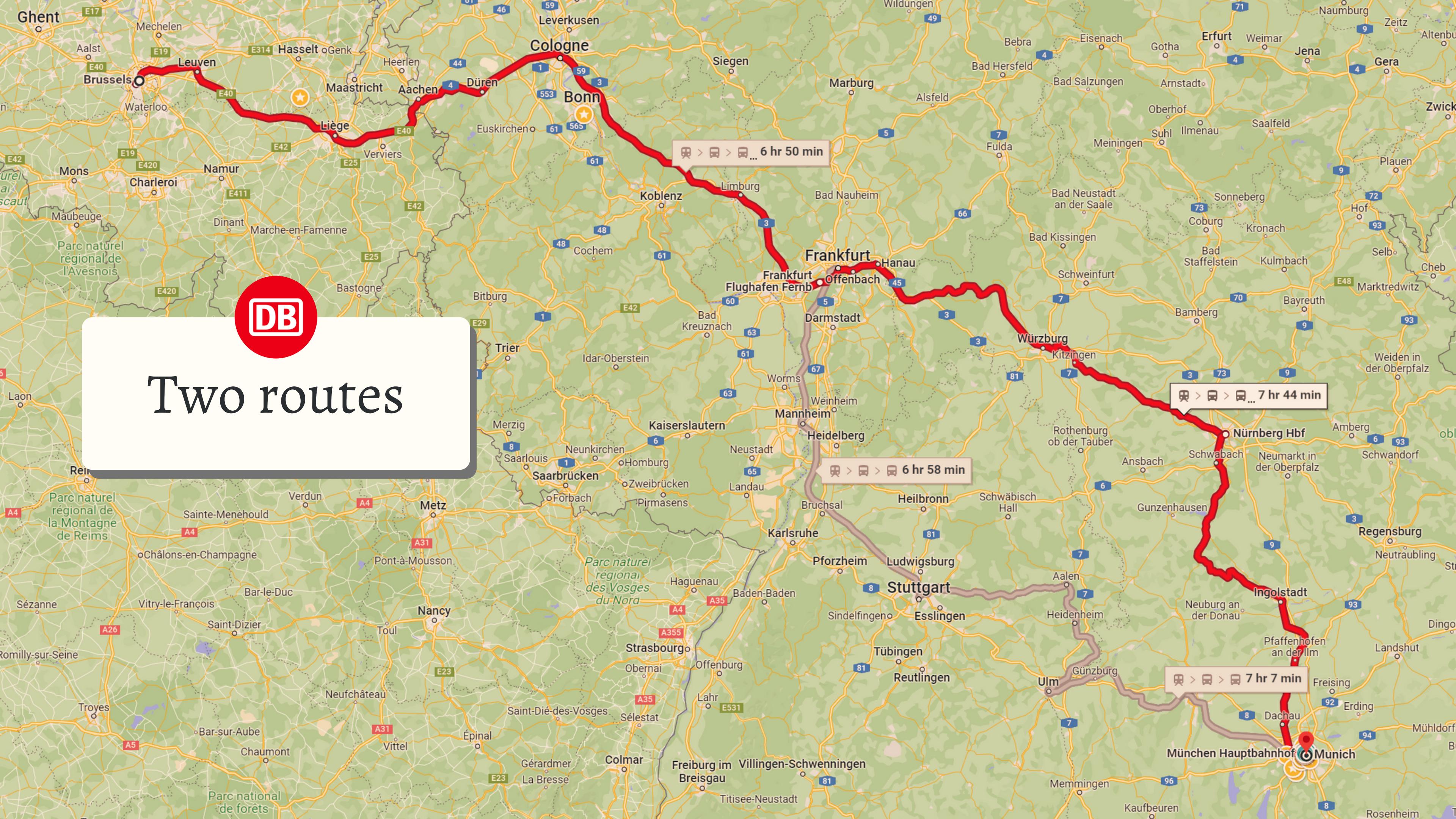
Two routes

6 hr 50 min

7 hr 44 min

6 hr 58 min

7 hr 7 min



ADRIAN

I'm taking the train from Brussels to Munich.



Option 1

Brussels - Frankfurt - München

23:00

Option 2

Brussels - Köln - München

23:20

ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Option 1

Brussels - Frankfurt - München

0

Option 2

Brussels - Köln - München

-20

ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

Option 1

Brussels - Frankfurt - München

0

Option 2

Brussels - Köln - München

-20

ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

Option 1 ✓

Brussels - Frankfurt - München

0

Option 2

Brussels - Köln - München

-20

ADRIAN

I'm taking the train from Brussels to Munich.



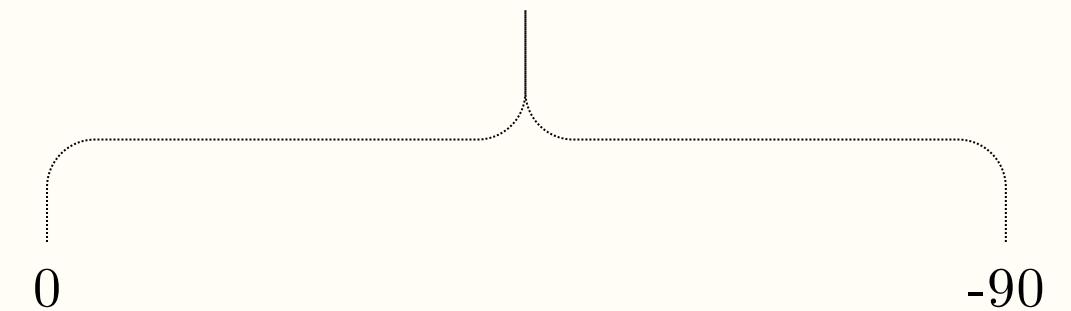
My utility is determined by the arrival time.

Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

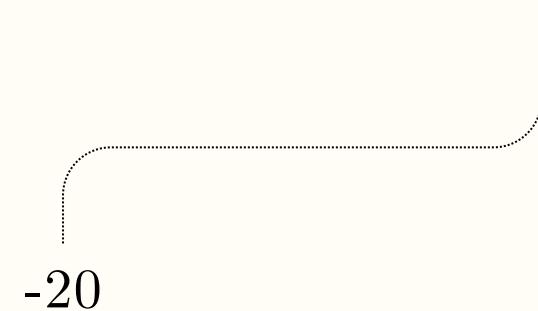
Option 1

Brussels - Frankfurt - (Mannheim?) - München



Option 2

Brussels - Köln - München



ADRIAN

I'm taking the train from Brussels to Munich.

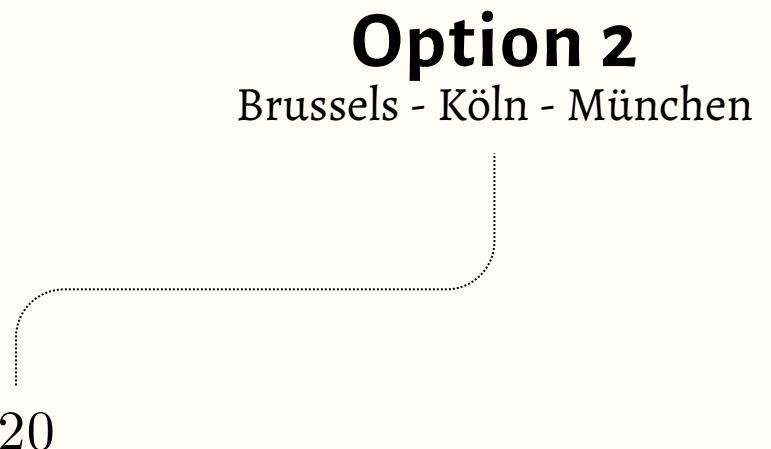
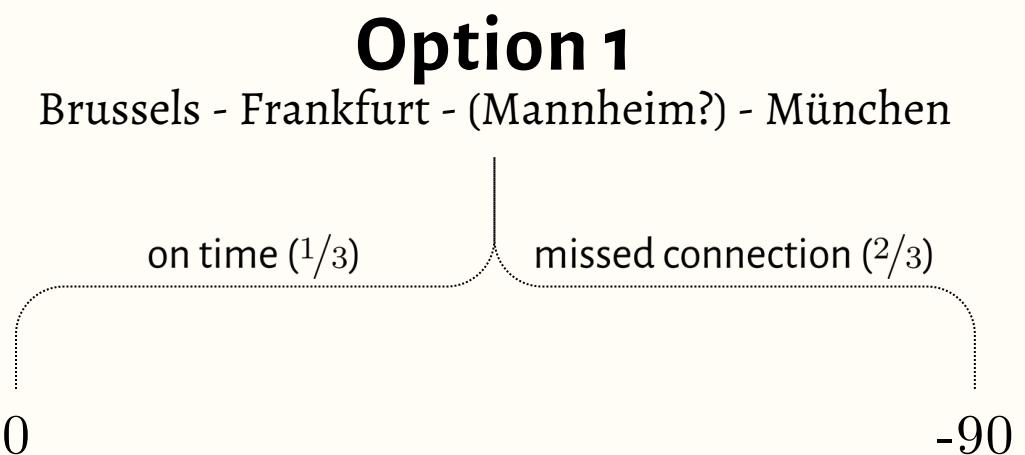


My utility is determined by the arrival time.

Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

This is very likely to happen... So how should we think of this possibility?



ADRIAN

I'm taking the train from Brussels to Munich.



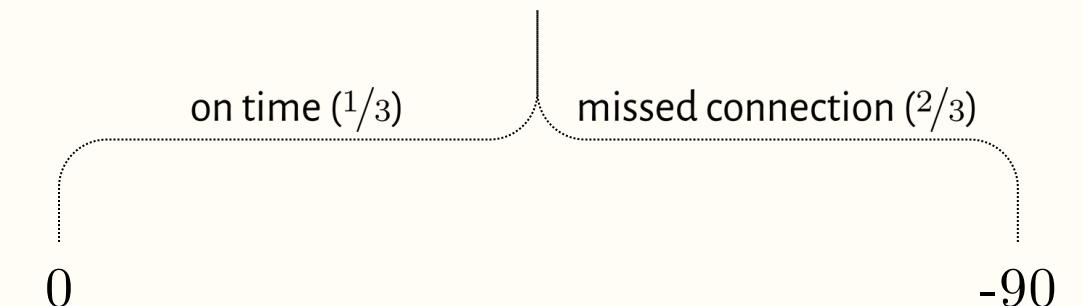
My utility is determined by the arrival time.

Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

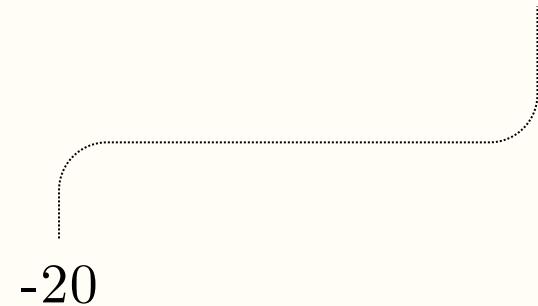
This is very likely to happen... So how should we think of this possibility?

Option 1
Brussels - Frankfurt - (Mannheim?) - München



$$\begin{aligned}\mathbb{E}[\text{Route 1}] &= \Pr[\text{on time}] \cdot 0 + \Pr[\text{missed connection}] \cdot (-90) \\ &= \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-90) \\ &= -60.\end{aligned}$$

Option 2
Brussels - Köln - München



ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

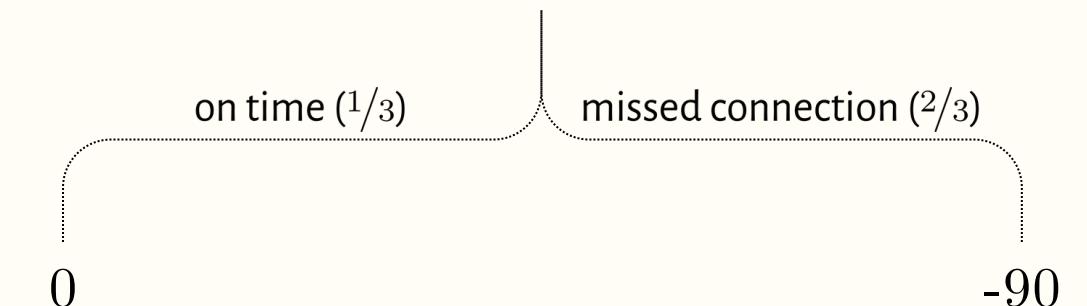
Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

This is very likely to happen... So how should we think of this possibility?

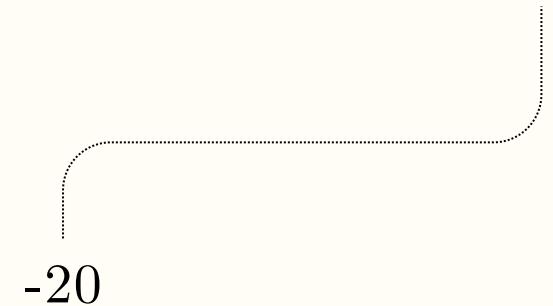
Now the second option seems better.

Option 1
Brussels - Frankfurt - (Mannheim?) - München



$$\begin{aligned}\mathbb{E}[\text{Route 1}] &= \Pr[\text{on time}] \cdot 0 + \Pr[\text{missed connection}] \cdot (-90) \\ &= \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-90) \\ &= -60.\end{aligned}$$

Option 2 ✓
Brussels - Köln - München



ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

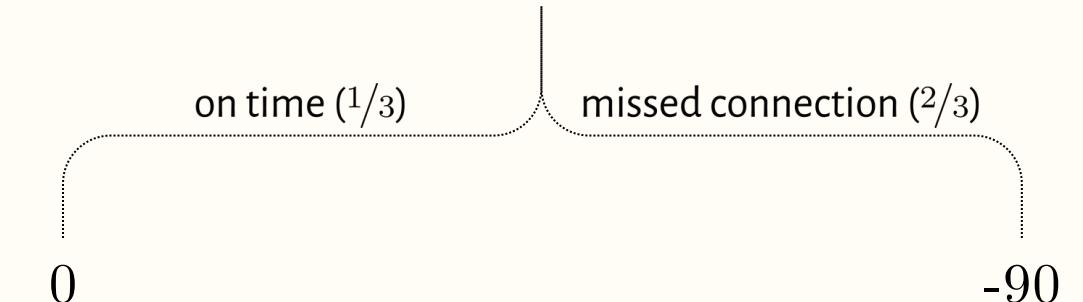
But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

This is very likely to happen... So how should we think of this possibility?

Now the second option seems better.

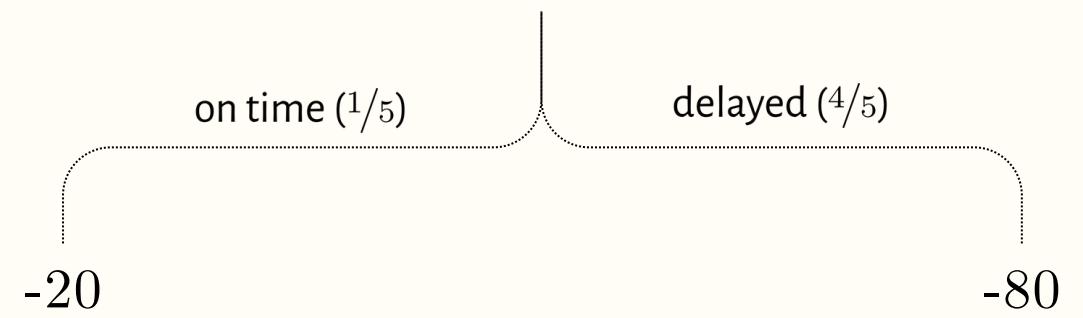
But wait! The second option has some uncertainty too: past experience suggests a likely delay.

Option 1
Brussels - Frankfurt - (Mannheim?) - München



$$\begin{aligned}\mathbb{E}[\text{Route 1}] &= \Pr[\text{on time}] \cdot 0 + \Pr[\text{missed connection}] \cdot (-90) \\ &= \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-90) \\ &= -60.\end{aligned}$$

Option 2
Brussels - Köln - München



ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

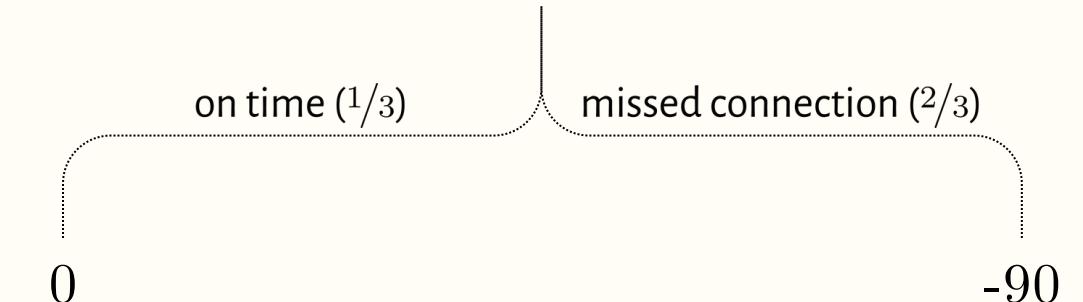
But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

This is very likely to happen... So how should we think of this possibility?

Now the second option seems better.

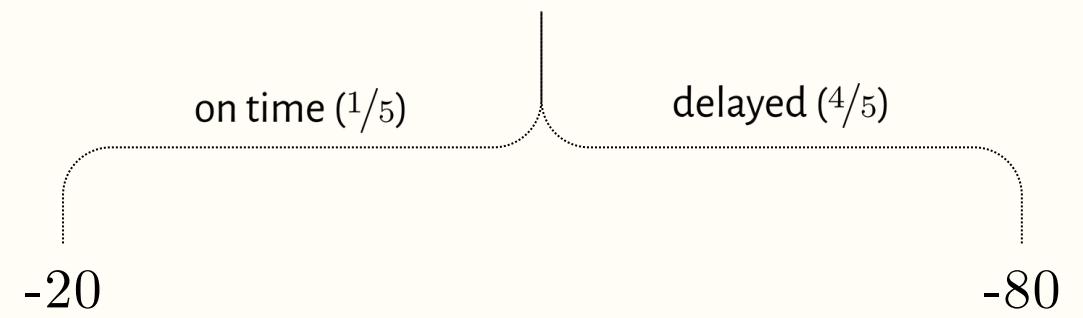
But wait! The second option has some uncertainty too: past experience suggests a likely delay.

Option 1
Brussels - Frankfurt - (Mannheim?) - München



$$\begin{aligned}\mathbb{E}[\text{Route 1}] &= \Pr[\text{on time}] \cdot 0 + \Pr[\text{missed connection}] \cdot (-90) \\ &= \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-90) \\ &= -60.\end{aligned}$$

Option 2
Brussels - Köln - München



$$\begin{aligned}\mathbb{E}[\text{Route 2}] &= \Pr[\text{on time}] \cdot (-20) + \Pr[\text{delayed}] \cdot (-80) \\ &= \frac{1}{5} \cdot (-20) + \frac{4}{5} \cdot (-80) \\ &= -68.\end{aligned}$$

ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

This is very likely to happen... So how should we think of this possibility?

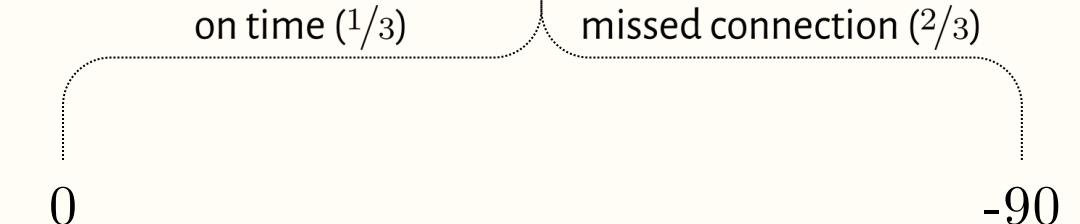
Now the second option seems better.

But wait! The second option has some uncertainty too: past experience suggests a likely delay.

Better to stick with the first option after all...

Option 1

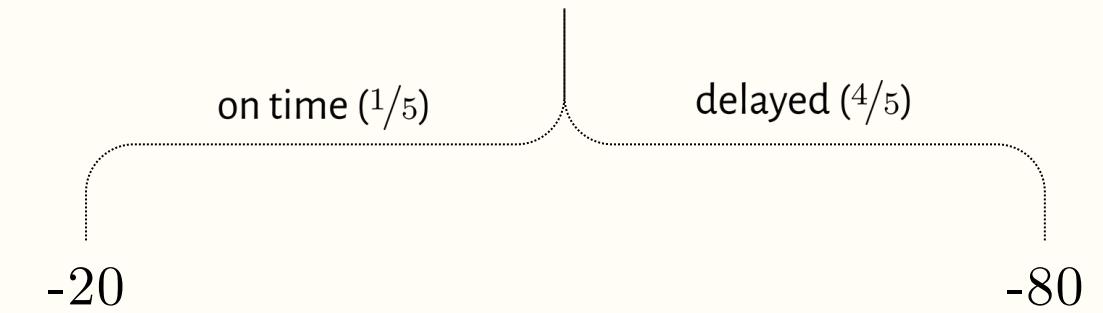
Brussels - Frankfurt - (Mannheim?) - München



$$\begin{aligned}\mathbb{E}[\text{Route 1}] &= \Pr[\text{on time}] \cdot 0 + \Pr[\text{missed connection}] \cdot (-90) \\ &= \frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-90) \\ &= -60.\end{aligned}$$

Option 2

Brussels - Köln - München



$$\begin{aligned}\mathbb{E}[\text{Route 2}] &= \Pr[\text{on time}] \cdot (-20) + \Pr[\text{delayed}] \cdot (-80) \\ &= \frac{1}{5} \cdot (-20) + \frac{4}{5} \cdot (-80) \\ &= -68.\end{aligned}$$

ADRIAN

I'm taking the train from Brussels to Munich.



My utility is determined by the arrival time.

Which option is best?

But with the first option I might miss the Frankfurt connection, meaning and will get home even later.

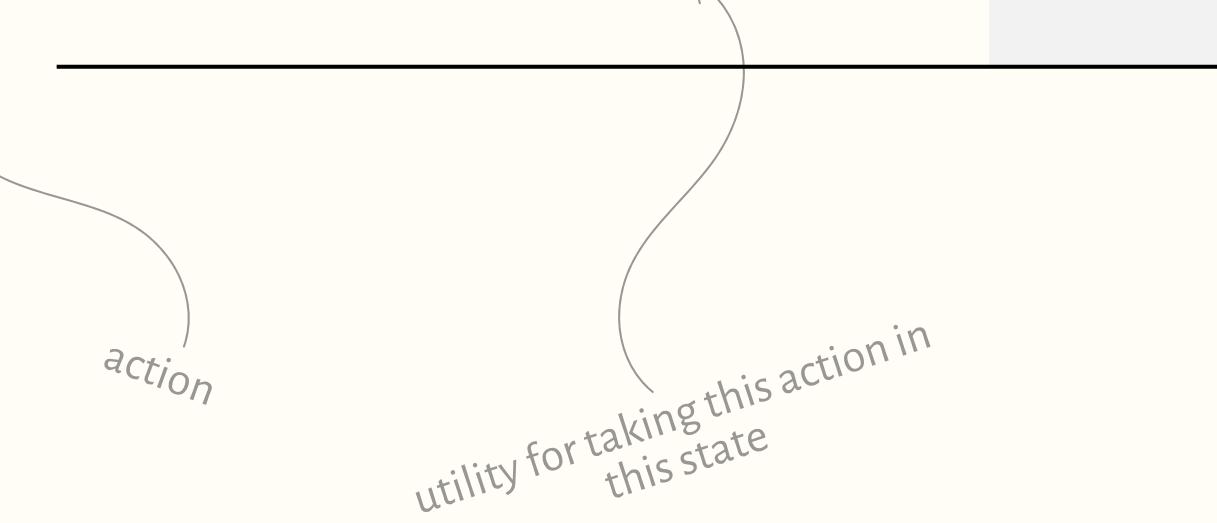
This is very likely to happen... So how should we think of this possibility?

Now the second option seems better.

But wait! The second option has some uncertainty too: past experience suggests a likely delay.

Better to stick with the first option after all...

	state of nature				
	on time	missed connection in Frankfurt	delay on Köln - München line	...	expected utility*
Option 1	0	-90	0	...	-60 
Option 2	-20	0	-80	...	-68



*Table isn't 100% correct: in general, states need to be mutually exclusive

In general, rational agents (aim to) maximize expected utility.

$$\mathbb{E}[u(\text{action})] = \sum_{\text{state}} \left(u(\text{action}, \text{state}) \cdot \Pr[\text{state}] \right)$$

In the previous example utility was defined relative to my immediate self-interest (i.e., getting home as early as possible).

But, in general, it can be whatever we want.

SAM BANKMAN-FRIED

In *effective altruism* we aim to maximize overall expected utility,
i.e., over all (including future) people.



Faux, Z. (2023). *Number Go Up: Inside Crypto's Wild Rise and Staggering Fall*. Crown Currency.

SAM BANKMAN-FRIED

In *effective altruism* we aim to maximize overall expected utility,
i.e., over all (including future) people.



So yeah, I'd take a bet where 51% you double the earth out
somewhere else, 49% it all disappears.

Faux, Z. (2023). *Number Go Up: Inside Crypto's Wild Rise and Staggering Fall*. Crown Currency.

SAM BANKMAN-FRIED

In *effective altruism* we aim to maximize overall expected utility,
i.e., over all (including future) people.



So yeah, I'd take a bet where 51% you double the earth out
somewhere else, 49% it all disappears.

I honestly think it's negative EV for me to cut my hair. I think it's
important for people to think I look crazy.

Faux, Z. (2023). *Number Go Up: Inside Crypto's Wild Rise and Staggering Fall*. Crown Currency.

GAMES IN NORMAL FORM

Now we know how to take optimal decisions, given different states of nature.

But sometimes the ‘states’ are someone else’s decisions.

So what’s best for me to do depends on what you do, and vice-versa.



JOHN VON NEUMANN

We should call that *game theory*.



OSKAR MORGENSTERN





JOHN VON NEUMANN

We should call that *game theory*.



OSKAR MORGENSTERN

And write a classic textbook on it!

JOHN VON NEUMANN





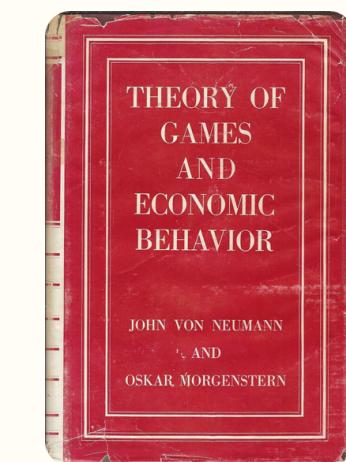
JOHN VON NEUMANN

We should call that *game theory*.



OSKAR MORGENSTERN

And write a classic textbook on it!



JOHN VON NEUMANN

Von Neumann, J., & Morgenstern, O. (1953). *Theory of Games and Economic Behavior*. Princeton University Press.

A game in normal form consists of *players* who can take *actions*, which lead to *payoffs*.

Glossary of Terms

agents, or players	$N = \{1, \dots, n\}$
actions of agent i	A_i
action profile	$\mathbf{a} = (a_1, \dots, a_n)$
set of all action profiles	$\mathbf{A} = A_1 \times \dots \times A_n$
utility (payoff) function of agent i	$u_i: \mathbf{A} \rightarrow \mathbb{R}$
utility profile	$\mathbf{u} = (u_1, \dots, u_n)$
normal-form game	$(N, \mathbf{A}, \mathbf{u})$
pure strategies of agent i	$S_i = A_i$
strategy profile	$\mathbf{s} = (s_1, \dots, s_n)$
set of strategy profiles	$\mathbf{S} = S_1 \times \dots \times S_n$
utility of i with respect to strategy profile \mathbf{s}	$u_i(\mathbf{s}) = u_i(\mathbf{a})$
\mathbf{s} without s_i	$\mathbf{s}_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n)$
\mathbf{s} equivalently	$\mathbf{s} = (s_i, \mathbf{s}_{-i})$

Let's meet our first game!

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

		payoffs	
		Stag	Hare
		Stag	10, 10
		Hare	6, 0
			3, 3

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

action of column hunter

payoffs

action of row hunter

payoff of row hunter
for this combination
of actions

payoff of column
hunter

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

We generally assume that player 1 is the row player and player 2 is the column player.

Players

$$N = \{1, 2\}$$

Actions of player 1

$$\{\text{Stag}, \text{Hare}\}$$

Actions of player 2

$$\{\text{Stag}, \text{Hare}\}$$

Strategies

$$\{(\text{Stag}, \text{Stag}), (\text{Stag}, \text{Hare}), (\text{Hare}, \text{Stag}), (\text{Hare}, \text{Hare})\}$$

Payoffs (aka utilities)

$$u_1(\text{Stag}, \text{Stag}) = 10$$

$$u_2(\text{Stag}, \text{Hare}) = 6$$

...

payoffs

		player 2	Stag	Hare
		Stag	10, 10	0, 6
player 1	Stag	6, 0	3, 3	
	Hare			

OSKAR MORGENSTERN

If we knew what strategies players would play we could go on and compute their utilities, expected utilities and so on.



JOHN VON NEUMANN

But that's not how rational agents behave: strategies change depending on what others do.



OSKAR MORGENSTERN

Indeed! If the column player goes for the hare, the row player will want to do the same.



		Stag	Hare
Stag	Stag	10, 10	0, 6
	Hare	6, 0	3, 3

JOHN VON NEUMANN

We need to reason the other way around: from utilities to strategies.



OSKAR MORGENSTERN

We need to reason about *solution concepts*.



A solution concept describes what strategies we might expect the players will adopt.

And, therefore, the result of the game.

The first solution concept we look at is based on *dominance*.

A player has a dominated strategy if the player could do uniformly better by playing a different strategy.

DEFINITION (STRICT DOMINANCE AMONG STRATEGIES)

Strategy s_i strictly dominates strategy s'_i if $u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$, for any profile s'_i of other agents' strategies.

Strategy s_i is strictly dominant (for player i) if it strictly dominates any other strategy s'_i .

Does T strictly dominate M, for player 1?

	L	C	R
T	3, 0	2, 1	0, 0
M	1, 1	1, 1	5, 0
B	0, 1	4, 2	0, 1

Does T strictly dominate M, for player 1? No!

If player 2 plays L: $3 > 1$ 

If player 2 plays C: $2 > 1$ 

If player 2 plays R: $0 < 5$ 

		L	C	R
		3, 0	2, 1	0, 0
		1, 1	1, 1	5, 0
T	3, 0	2, 1	0, 0	
M	1, 1	1, 1	5, 0	
B	0, 1	4, 2	0, 1	

Does T strictly dominate M, for player 1? No!

If player 2 plays L: $3 > 1$ 

If player 2 plays C: $2 > 1$ 

If player 2 plays R: $0 < 5$ 

Does C strictly dominate L, for player 2?

	L	C	R
T	3, 0	2, 1	0, 0
M	1, 1	1, 1	5, 0
B	0, 1	4, 2	0, 1

Does T strictly dominate M, for player 1? No!

If player 2 plays L: $3 > 1$

If player 2 plays C: $2 > 1$

If player 2 plays R: $0 < 5$

Does C strictly dominate L, for player 2? No!

If player 1 plays T: $1 > 0$

If player 1 plays M: $1 = 1$

If player 1 plays B: $2 > 1$

		L	C	R
		3, 0	2, 1	0, 0
		1, 1	1, 1	5, 0
T				
M				
B				
0, 1		4, 2		0, 1

Does T strictly dominate M, for player 1? No!

If player 2 plays L: $3 > 1$

If player 2 plays C: $2 > 1$

If player 2 plays R: $0 < 5$

Does C strictly dominate L, for player 2? No!

If player 1 plays T: $1 > 0$

If player 1 plays M: $1 = 1$

If player 1 plays B: $2 > 1$

Does C strictly dominate R, for player 2?

	L	C	R
T	3, 0	2, 1	0, 0
M	1, 1	1, 1	5, 0
B	0, 1	4, 2	0, 1

Does T strictly dominate M, for player 1? No!

If player 2 plays L: $3 > 1$

If player 2 plays C: $2 > 1$

If player 2 plays R: $0 < 5$

Does C strictly dominate L, for player 2? No!

If player 1 plays T: $1 > 0$

If player 1 plays M: $1 = 1$

If player 1 plays B: $2 > 1$

Does C strictly dominate R, for player 2? Yes!

If player 1 plays T: $1 > 0$

If player 1 plays M: $1 > 0$

If player 1 plays B: $2 > 1$

		L	C	R
		3, 0	2, 1	0, 0
		1, 1	1, 1	5, 0
T		3, 0	2, 1	0, 0
M		1, 1	1, 1	5, 0
B		0, 1	4, 2	0, 1

There is no point in playing a strictly dominated strategy.

Which means we can successively eliminate any such strategies from a player's arsenal.

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

		payoffs	
		Stag	Hare
Stag	Stag	10, 10	0, 6
	Hare	6, 0	3, 3

strictly dominant strategies

?

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

payoffs

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

strictly dominant strategies

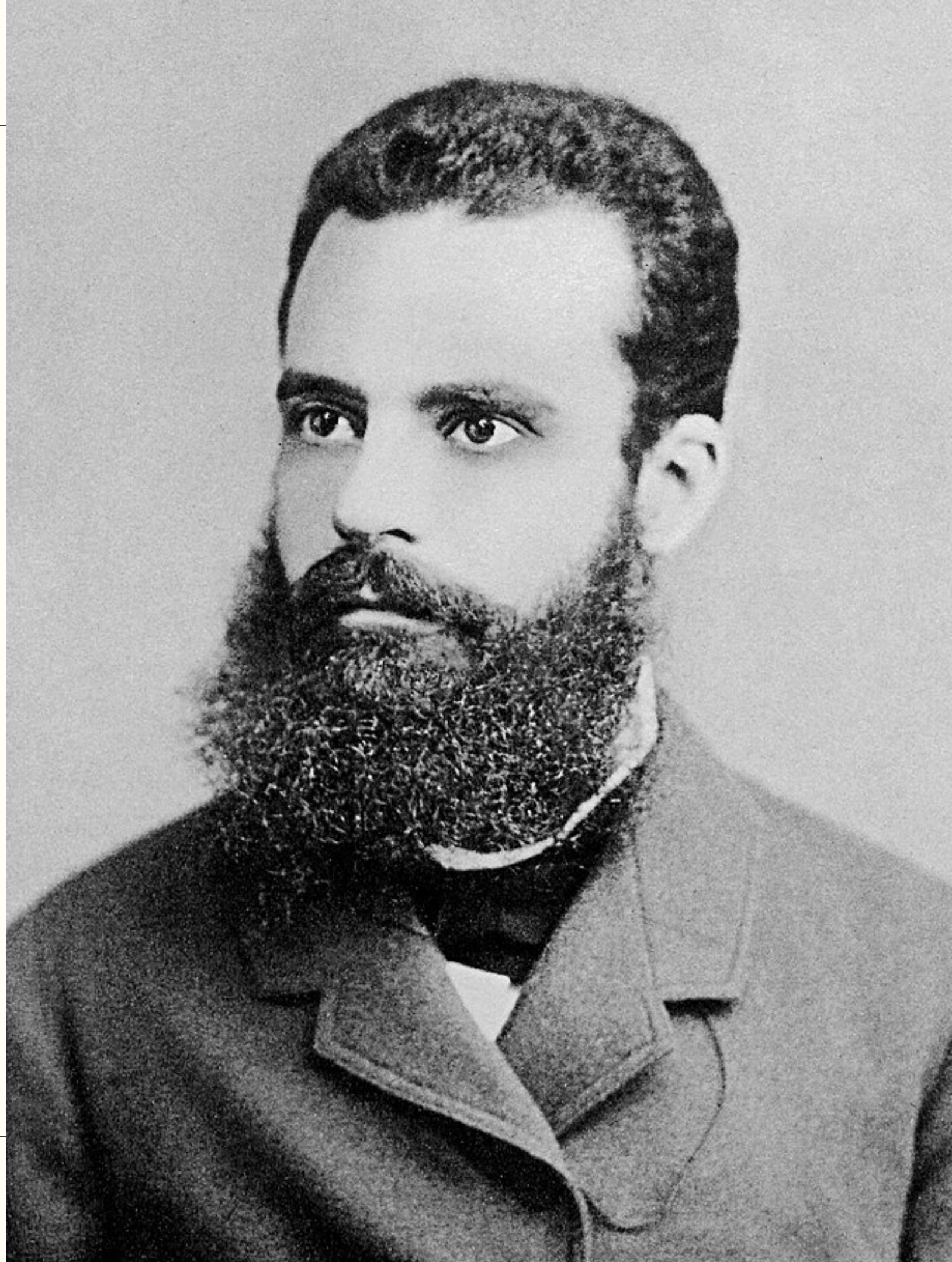
none

Strictly dominant strategies always exist.

Strictly dominant strategies always exist.

Except when they don't: it's a very strong solution concept!

Enter Pareto.



VILFREDO PARETO

Better to look at outcomes where everyone is as well-off as can be.



In a Pareto optimal outcome no one can be made better off without making someone else worse off.

DEFINITION (PARETO DOMINATION)

A strategy profile s *Pareto dominates* strategy profile s' if:

- (i) $u_i(s) \geq u_i(s')$, for every agent i , and
- (ii) there exists an agent j such that $u_j(s) > u_j(s')$.

DEFINITION (PARETO OPTIMALITY)

A strategy profile s is *Pareto optimal* if there is no (other) strategy profile s' that Pareto dominates s .

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

payoffs

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

strictly dominant strategies

none

Pareto optimal strategy profiles

?

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

payoffs

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

strictly dominant strategies

none

Pareto optimal strategy profiles

(Stag, Stag)

Time for a new game!

The Coordination Game



There is a country with no traffic rules.

Two cars are on the road, driving towards each other.

They have to decide what side of the road to take.

If they choose the same side, all is well.

If they choose different sides, they bump into each other.

payoffs

	Left	Right
Left	1, 1	0, 0
Right	0, 0	1, 1

strictly dominant strategies ?

Pareto optimal strategy profiles ?

2/2

The Coordination Game



There is a country with no traffic rules.

Two cars are on the road, driving towards each other.

They have to decide what side of the road to take.

If they choose the same side, all is well.

If they choose different sides, they bump into each other.

		payoffs	
		Left	Right
Left	Left	1, 1	0, 0
	Right	0, 0	1, 1
strictly dominant strategies			
none			
Pareto optimal strategy profiles			
?			

The Coordination Game



There is a country with no traffic rules.

Two cars are on the road, driving towards each other.

They have to decide what side of the road to take.

If they choose the same side, all is well.

If they choose different sides, they bump into each other.

		payoffs	
		Left	Right
Left	Left	1, 1	0, 0
	Right	0, 0	1, 1
strictly dominant strategies			
none			
Pareto optimal strategy profiles			
(Left, Left), (Right, Right)			

dominated by (Left, Left) and (Right, Right)

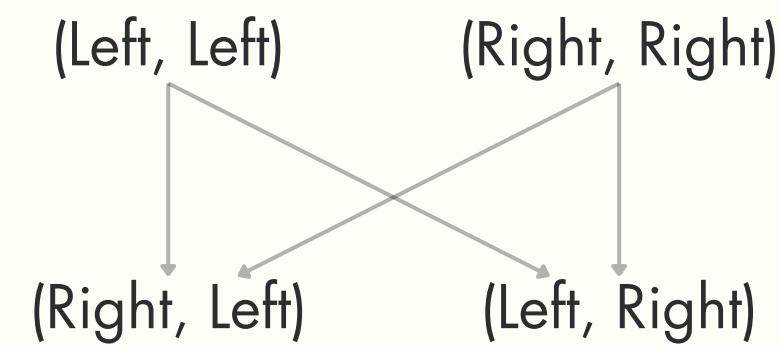
dominated by (Left, Left) and (Right, Right)

2/2

dominated by
(Left, Left) and
(Right, Right)

VILFREDO PARETO

Pareto domination defines a partial order over strategy profiles:



Pareto optimal outcomes always exist.

For real! Check for yourselves if you don't believe it.

May not be unique though.

		Left	Right
Left	Left	1, 1	0, 0
	Right	0, 0	1, 1

DAVID LEWIS

Coordination problems are everywhere in social interactions, and are at the heart of the *conventions* that become social norms.



Rescorla, M. (2019). Convention. Edward N. Zalta (ed.), *The Stanford Encyclopedia of Philosophy* (Summer 2019 Edition).



H. PEYTON YOUNG

Drivers coordinate to avoid collisions on the road.
Economic agents eliminate the need for barter by
coordinating upon a common monetary currency.

Young, H. P. (1996). The Economics of Convention. *The Journal of Economic Perspectives*, 10(2), 105–122.

DAVID LEWIS

Language is a coordination game.



Lewis, D. (2008). *Convention: A Philosophical Study*. Harvard University Press.

VILFREDO PARETO
Pareto optimality doesn't necessarily imply that outcomes are fair.



Just that they're 'efficient', in the sense of not leaving money on the table.

Consider the game on the right, played between the land-owner and the farmers, on how the spoils of the land are divided.

payoffs

		FARMERS		
		Feudalism	Capitalism	Communism
LANDOWNER	Feudalism	90, 10	5, 5	5, 5
	Capitalism	5, 5	70, 30	5, 5
	Communism	5, 5	5, 5	50, 50

strictly dominant strategies
none

Pareto optimal strategy profiles
?

VILFREDO PARETO
Pareto optimality doesn't necessarily imply that outcomes are fair.



Just that they're 'efficient', in the sense of not leaving money on the table.

Consider the game on the right, played between the land-owner and the farmers, on how the spoils of the land are divided.

payoffs

		FARMERS		
		Feudalism	Capitalism	Communism
LANDOWNER	Feudalism	90, 10	5, 5	5, 5
	Capitalism	5, 5	70, 30	5, 5
	Communism	5, 5	5, 5	50, 50

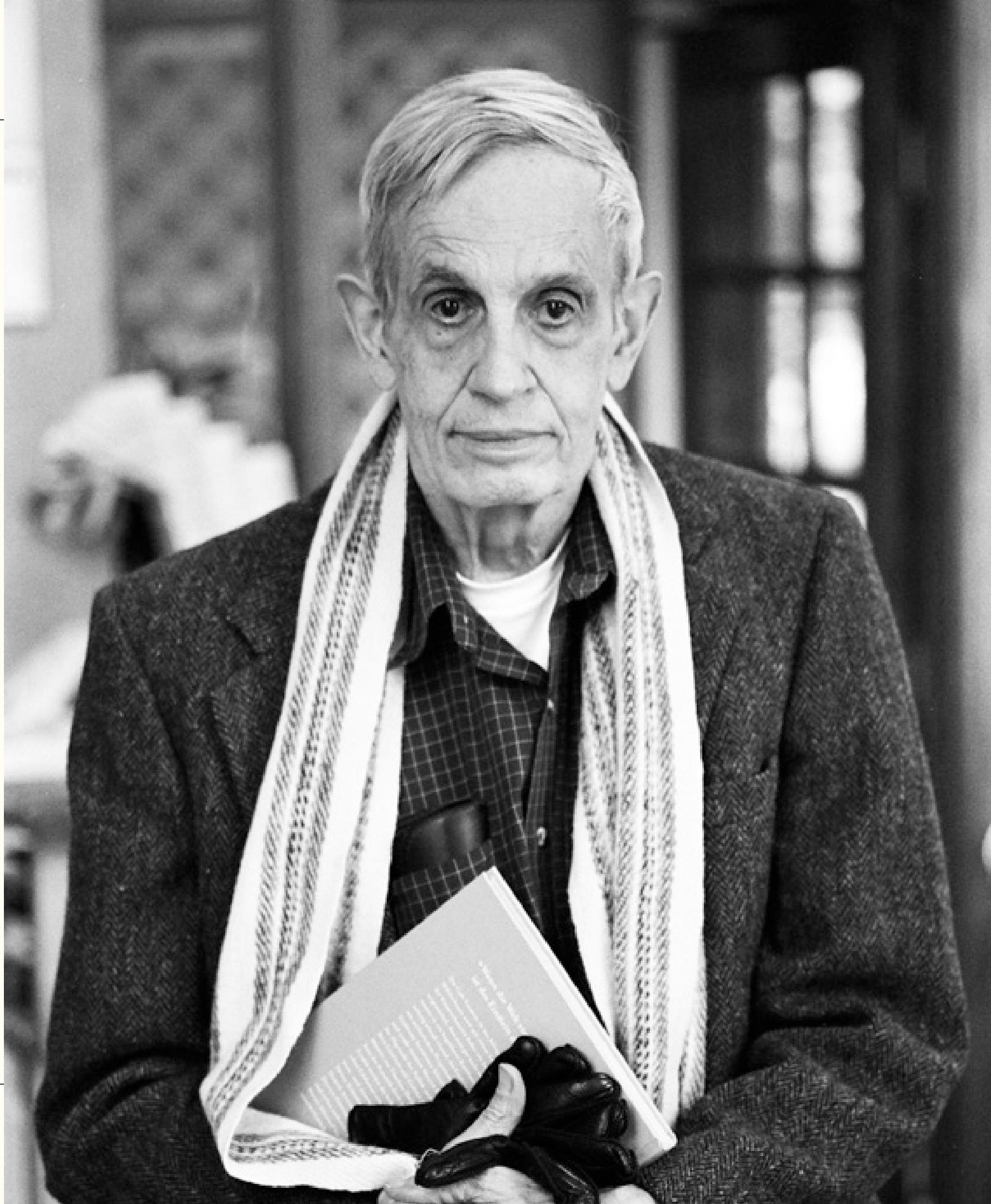
strictly dominant strategies
none

Pareto optimal strategy profiles
(Feudalism, Feudalism), (Capitalism, Capitalism),
(Communism, Communism)

Nonetheless, Pareto optimal is (a minimal requirement on) where we want to be.

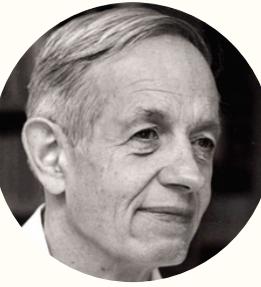
But can we expect that players end up there?

Enter Nash.



JOHN NASH

In a Nash equilibrium no one has an incentive to change their strategy, given the other players' strategies.



DEFINITION (BEST RESPONSE)

Player i 's *best response* to the other players' strategies $s_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n)$ is a strategy s_i^* such that $u_i(s_i^*, s_{-i}) \geq u_i(s_i, s_{-i})$, for any strategy s_i of player i .

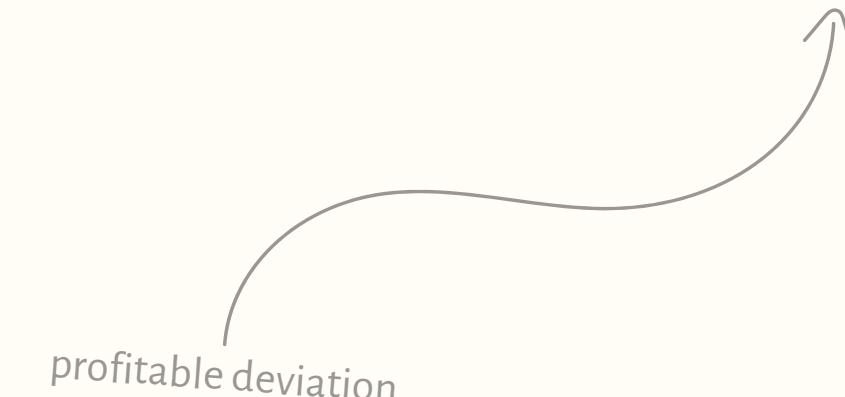
DEFINITION (BEST RESPONSE)

Player i 's *best response* to the other players' strategies $\mathbf{s}_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n)$ is a strategy s_i^* such that $u_i(s_i^*, \mathbf{s}_{-i}) \geq u_i(s_i, \mathbf{s}_{-i})$, for any strategy s_i of player i .

DEFINITION (PURE NASH EQUILIBRIUM)

A strategy profile $\mathbf{s}^* = (s_1^*, \dots, s_n^*)$ is a *pure Nash equilibrium* if s_i^* is a best response to \mathbf{s}_{-i}^* , for every player i .

In other words, given strategy profile \mathbf{s}^* , there is no player i and strategy s'_i such that $u_i(s'_i, \mathbf{s}_{-i}^*) > u_i(s_i^*, \mathbf{s}_{-i}^*)$.



Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

payoffs

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

strictly dominant strategies

none

Pareto optimal strategy profiles

(Stag, Stag)

pure Nash equilibria

?

Stag Hunt



Two hunters have to decide what to hunt: one stag or two hares.

If they hunt together, any catch is divided equally.

A stag is worth a lot (more than both hares combined!), but can only be caught by the two hunters working together. If a hunter goes for the stag alone, they end up with nothing.

A hunter goes for hare while the other for stag, the first gets both hares and does not have to share.

payoffs

	Stag	Hare
Stag	10, 10	0, 6
Hare	6, 0	3, 3

strictly dominant strategies

none

Pareto optimal strategy profiles

(Stag, Stag)

pure Nash equilibria

(Stag, Stag), (Hare, Hare)

And now for the moment
we've all been waiting for.

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

payoffs

	Cooperate	Defect
Cooperate	-20, -20	-100, 0
Defect	0, -100	-50, -50

strictly dominant strategies

Pareto optimal strategy profiles

pure Nash equilibria

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

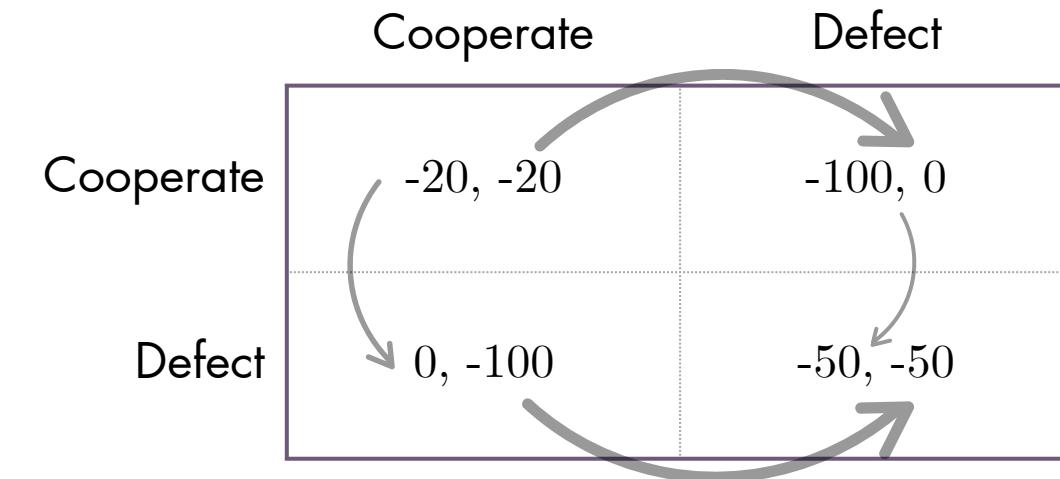
If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

payoffs



strictly dominant strategies

Pareto optimal strategy profiles

pure Nash equilibria

(Defect, Defect)

JOHN NASH
Pure NEs always exist... Except when they don't (think of an example)!



JOHN NASH

Pure NEs always exist... Except when they don't (think of an example)!



In a Nash Equilibrium everyone is as well off as they can be.

JOHN NASH

Pure NEs always exist... Except when they don't (think of an example)!



In a Nash Equilibrium everyone is as well off as they can be.

Except that they're not!

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

payoffs

	Cooperate	Defect
Cooperate	-20, -20	-100, 0
Defect	0, -100	-50, -50

strictly dominant strategies

Pareto optimal strategy profiles

?

pure Nash equilibria

(Defect, Defect)

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

		payoffs	
		Cooperate	Defect
Cooperate	Cooperate	-20, -20	-100, 0
	Defect	0, -100	-50, -50

strictly dominant strategies

Pareto optimal strategy profiles
(Cooperate, Cooperate), (Cooperate, Defect), (Defect, Cooperate)

pure Nash equilibria
(Defect, Defect)

2/2

JOHN NASH

Pure NEs always exist... Except when they don't (think of an example)!



In a Nash Equilibrium everyone is as well off as they can be.

Except that they're not!

In the Prisoner's Dilemma every outcome except the Nash equilibrium is Pareto optimal!

JOHN NASH

Pure NEs always exist... Except when they don't (think of an example)!



In a Nash Equilibrium everyone is as well off as they can be.

Except that they're not!

In the Prisoner's Dilemma every outcome except the Nash equilibrium is Pareto optimal!

In fact defection is an even stronger outcome, in terms of solution concepts we've seen so far.

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

payoffs

		Cooperate	Defect
Cooperate	Cooperate	-20, -20	-100, 0
	Defect	0, -100	-50, -50

strictly dominant strategies ?

Pareto optimal strategy profiles
(Cooperate, Cooperate), (Cooperate, Defect), (Defect, Cooperate)

pure Nash equilibria
(Defect, Defect)

2/2

The Prisoner's Dilemma



You and a friend are at the police station. You are the main suspects in a string of Oktoberfest beer thefts.

You are interrogated at the same time, in separate rooms.

If both of you stick to the common story (Cooperate), you get off with a smallish fine.

But if you tell on your friend (Defect) you get off free, while they get a hefty fine.

Your friend faces the same situation.

If you rat each other out, you split the large fine.

		payoffs	
		Cooperate	Defect
Cooperate	Cooperate	-20, -20	-100, 0
	Defect	0, -100	-50, -50

strictly dominant strategies
Defect, for both players

Pareto optimal strategy profiles
(Cooperate, Cooperate), (Cooperate, Defect), (Defect, Cooperate)

pure Nash equilibria
(Defect, Defect)

JOHN NASH

Pure NEs always exist... Except when they don't (think of an example)!



In a Nash Equilibrium everyone is as well off as they can be.

Except that they're not!

In the Prisoner's Dilemma every outcome except the Nash equilibrium is Pareto optimal!

In fact defection is an even stronger outcome, in terms of solution concepts we've seen so far.

Btw, a strictly dominating strategy profile, if it exists, is a (pure) Nash equilibrium. Though not necessarily the other way around.

How is this relevant to the
problem of cooperation?

JOHN NASH

Note that the numbers in the payoff matrix are not *per se* relevant.



What's important is the *relationship* between them.

That is to say, we should think of the Prisoner's Dilemma as a general scenario in which mutual defection is the equilibrium.

The Prisoner's Dilemma



GENERAL VERSION

There are two players, each with two actions: Cooperate or Defect.

If they both cooperate they both get a payoff of R (the *reward*).

If they both defect, they each get a payoff of P (the *punishment*).

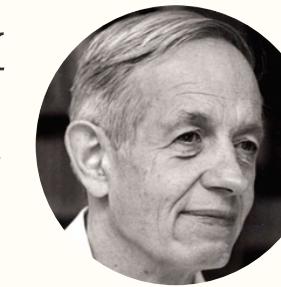
In the case of defection with cooperation, the defector gets T (the *temptation*), while the cooperator gets S (the *sucker's payoff*).

The relationship between the payoffs is T > R > P > S.

		payoffs	
		Cooperate	Defect
Cooperate	Cooperate	R, R	S, T
	Defect	T, S	P, P
strictly dominant strategies			
Defect, for both players			
Pareto optimal strategy profiles			
(Cooperate, Cooperate), (Cooperate, Defect), (Defect, Cooperate)			
pure Nash equilibria			
(Defect, Defect)			

JOHN NASH

Note that the numbers are not *per se* relevant.



What matters are the *relationships* between them.

That is to say, we should think of the Prisoner's Dilemma as a general scenario in which mutual defection is the equilibrium.



MARTIN NOWAK

Things become even clearer when considering a simplified version of the Prisoner's Dilemma: the *Donation Game*.

Nowak, M.A. (2006). *Evolutionary Dynamics*. Belknap Press

The Donation Game



SPECIAL CASE OF PRISONER'S DILEMMA

There are two players, each with two actions: Cooperate or Defect.

A cooperator pays a cost c for the other player to receive a benefit b , with $b > c > 0$.

A defector does not pay any cost, and provides no benefit.

Nowak, M.A. (2006). *Evolutionary Dynamics*. Belknap Press

payoffs

	Cooperate	Defect
Cooperate	$b - c, b - c$	$-c, b$
Defect	$b, -c$	$0, 0$

strictly dominant strategies
Defect, for both players

Pareto optimal strategy profiles
(Cooperate, Cooperate), (Cooperate, Defect), (Defect, Cooperate)

pure Nash equilibria
(Defect, Defect)

Even though cooperation is
overall the better outcome, in a
Prisoner's Dilemma defection
is the rational response!

MOM

These agents are terrible!



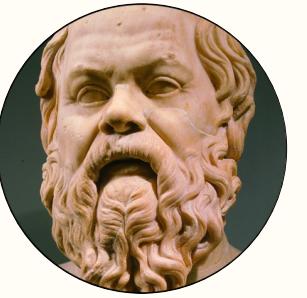
They need better education.



MOM

These agents are terrible!

They need better education.



SOCRATES

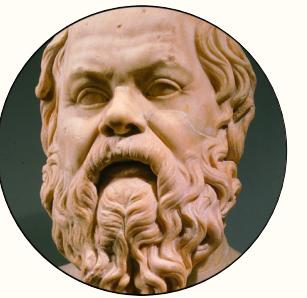
True. All evil is a result of ignorance.



MOM

These agents are terrible!

They need better education.



SOCRATES

True. All evil is a result of ignorance.

ADRIAN

If lack of education means agents are not aware of certain aspects of the games (e.g., payoffs), then ‘educated’ agents should still defect: it’s the dominant action!



If education means acquiring a set of reflexes that keep your selfish impulses in check, then that might work... but we still need to figure out in what situations such reflexes make sense.

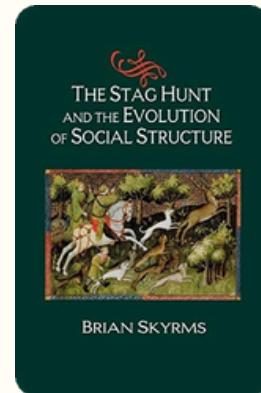
Note that a simple way out of the problem is if the underlying situation is a different game, e.g., Stag Hunt.

BRIAN SKYRMS

The Stag Hunt is a game where the payoffs from cooperating exceed the temptation to defect.



Stag Hunt games are the reason we have nice things, like society.



Skyrms, B. (2003). *The Stag Hunt and the Evolution of Social Structure*. Cambridge University Press.



JOHN NASH

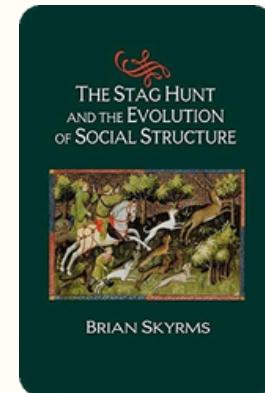


BRIAN SKYRMS

The Stag Hunt is a game where the payoffs from cooperating exceed the temptation to defect.



Stag Hunt games are the reason we have nice things, like society.



Skyrms, B. (2003). *The Stag Hunt and the Evolution of Social Structure*. Cambridge University Press.



JOHN NASH

In the Prisoner's Dilemma the temptation to be selfish and defect is greater than the payoff from pursuing the common good by cooperating.

JOHN W. N. WATKINS

• • •

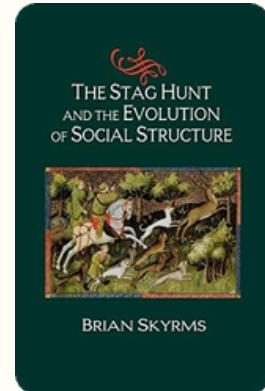


BRIAN SKYRMS

The Stag Hunt is a game where the payoffs from cooperating exceed the temptation to defect.



Stag Hunt games are the reason we have nice things, like society.



Skyrms, B. (2003). *The Stag Hunt and the Evolution of Social Structure*. Cambridge University Press.



JOHN NASH

In the Prisoner's Dilemma the temptation to be selfish and defect is greater than the payoff from pursuing the common good by cooperating.



JOHN W. N. WATKINS

Players in a Prisoner's Dilemma are led, in a way that might have startled Adam Smith, by a malevolent invisible hand to promote an end which was no part of their intention and which none of them wants.

Watkins, J. (1985). Second Thoughts on Self-interest and Morality. In *Paradoxes of Rationality and Cooperation: Prisoner's Dilemma and Newcomb's Problem*. University of British Columbia Press.

VAMPIRE BAT ELDER

Vampire bats face a prisoner's dilemma when having to decide whether to feed their hungry colleagues.



LANCE ARMSTRONG

Sports people too, when deciding whether to take performance enhancing drugs.

Schneier, B. (2006, August 10). [Drugs: Sports' Prisoner's Dilemma](#). *Wired*.

THE UN

Or countries deciding whether to cut down carbon emissions.



VAMPIRE BAT ELDER

Vampire bats face a prisoner's dilemma when having to decide whether to feed their hungry colleagues.



LANCE ARMSTRONG

Sports people too, when deciding whether to take performance enhancing drugs.

Schneier, B. (2006, August 10). [Drugs: Sports' Prisoner's Dilemma](#). *Wired*.

THE UN

Or countries deciding whether to cut down carbon emissions.



MARTIN NOWAK

Indeed, the Prisoner's Dilemma is the paradigmatic game used to study the evolution of cooperation.

Nowak, M.A. (2006). *Evolutionary Dynamics*. Belknap Press

What can we add to our
framework to get cooperation
in prisoner's-dilemma-type
situations?
