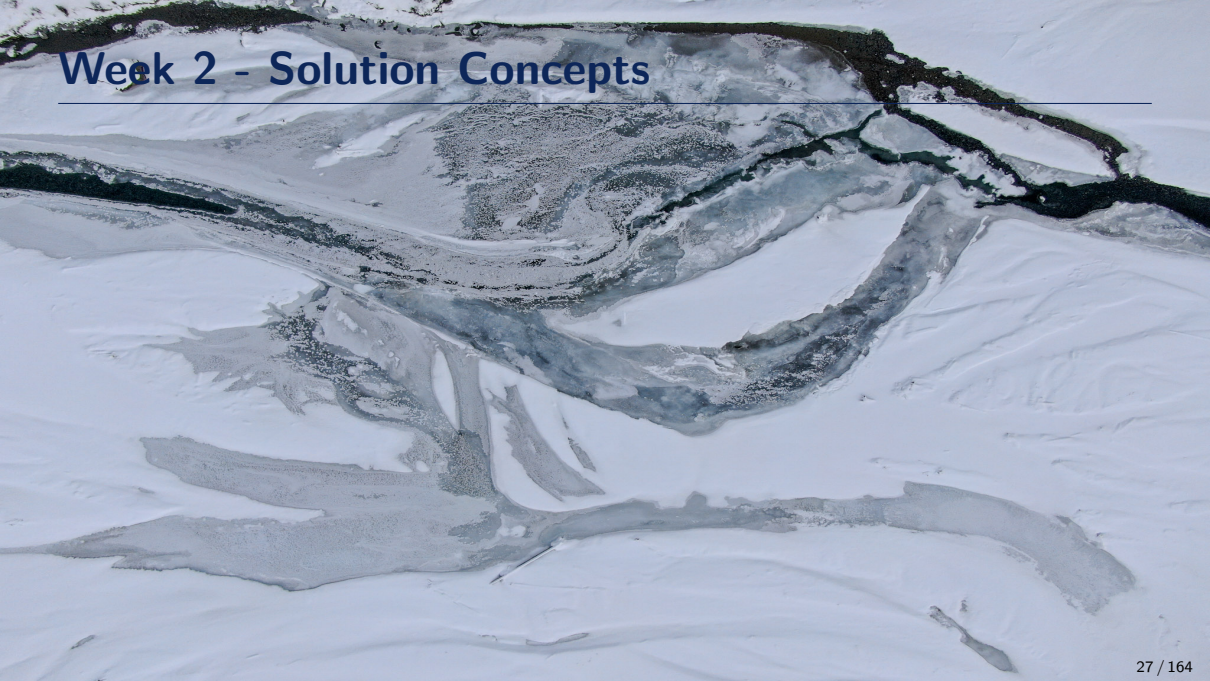


Week 2 - Solution Concepts



Normal Form Games

- Unlike in RL, now the return/reward is a function of all agents in the game
- Fixing the policies of other players results in a single agent environment
- Minimax vs Nash solution concepts

Fixed Strategies

- We assume the agents follow a fixed strategy
 - This does not allow for opponent adaptation or online learning
- The analysis is for tabular, explicit representations
 - However, works for implicit and online algorithms as long as they allow for tabularization¹
- Training and evaluation paradigm - as used in all of the major games AI milestones

¹Sustr M, Schmid M, Moravcik M, Burch N, Lanctot M, Bowling M. Sound search in imperfect information games

Train/Eval Paradigm



(a) Train



(b) Eval

Search

- Online algorithm
- Common idea used in many perfect information games (chess, go, ...)
- Many appealing properties
- We still assume that it is consistent with a fixed, tabular strategy
- Makes only sense in sequential decision making - we will come back to this later

Solution Concepts



(a) Maximin



(b) Nash equilibrium

Maximin

- Maximizing the worst-case scenario
- Assumes everyone else is “out there to get you”

Definition: Maximin Policy

Maximin policy of a player i is:

$$\arg \max_{\pi_i \in \Pi_i} \min_{\pi_{-i} \in \Pi_{-i}} R_i(\pi_i, \pi_{-i}) = \arg \max_{\pi_i \in \Pi_i} BRV_i(\pi_i) \quad (1)$$

Nash equilibrium

- Everyone is happy

Definition: Nash Equilibrium

Strategy profile (π_i, π_{-i}) forms a Nash equilibrium if none of the players benefit by deviating from their policy.

$$\forall i \in N, \forall \pi_i' : R_i(\pi_i, \pi_{-i}) \geq R_i(\pi_i', \pi_{-i})$$

Maximin vs Nash

- One is defined for a strategy, the other for a strategy profile
- We will see some interesting differences
- But we will also see that they are sometimes the same!

Maximin



Maximin in Pure Strategies

	Cooperate	Defect
Cooperate	$(-6, -6)$	$(0, -10)$
Defect	$(-10, 0)$	$(-1, -1)$

Table: Prisoner's dilemma

	Stop	Go
Stop	$(0, 0)$	$(0, 1)$
Go	$(1, 0)$	$(-10, -10)$

Table: The game of Chicken

Maximin in Pure Strategies

	Rock	Paper	Scissors
Rock	0	-1	1
Paper	1	0	-1
Scissors	-1	1	0

Table: Rock Paper Scissors

Maximin

- Let's consider pure strategies
- When we mix, we can do better!
- Opponent does not care about their reward at all!
- How can we find the best mixed strategy?

Optimizing Against Best Response

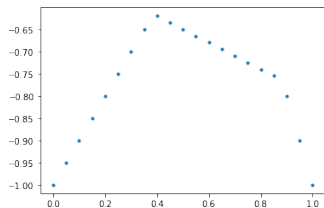
- For two-player zero-sum games, we have

$$\arg \max_{\pi_i \in \Pi_i} \min_{\pi_{-i} \in \Pi_{-i}} R_i(\pi_i, \pi_{-i}) = \arg \max_{\pi_i \in \Pi_i} R_i(\pi_i, br_{-i}(\pi_i))$$

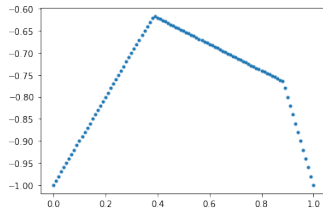
- We are thus optimizing against a best-responding player
- Let's visualise the best-response value function $f(\pi_i) = R_i(\pi_i, br_{-i}(\pi_i))$

Best Response Value Function

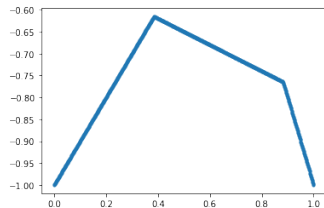
	A	B	C
X	-1	0	-0.8
1-X	1	-1	-0.5



(a) Step size 0.05



(b) Step size 0.01



(c) Step size 0.001

Nash



Nash equilibrium

Nash equilibrium

Strategy profile (π_i, π_{-i}) forms a Nash equilibrium if none of the players benefit by deviating from their policy.

$$\forall i \in N, \forall \pi'_i : R_i(\pi_i, \pi_{-i}) \geq R_i(\pi'_i, \pi_{-i})$$

- Easy to verify - all strategies are best responses
- Brute-force - enumerate all possible pairs and then verify

Nash in Pure Strategies

	Cooperate	Defect
Cooperate	$(-6, -6)$	$(0, -10)$
Defect	$(-10, 0)$	$(-1, -1)$

Table: Prisoner's dilemma

	Stop	Go
Stop	$(0, 0)$	$(0, 1)$
Go	$(1, 0)$	$(-10, -10)$

Table: The game of Chicken

Nash in Pure Strategies

	Rock	Paper	Scissors
Rock	0	-1	1
Paper	1	0	-1
Scissors	-1	1	0

Table: Rock Paper Scissors

Nash equilibrium

- Everyone is happy
- Pure Nash - enumerate, but might not exist!

Mixed Nash Equilibrium

$$\forall i \in N, \forall \pi'_i : R_i(\pi_i, \pi_{-i}) \geq R_i(\pi'_i, \pi_{-i})$$

- Recall that the opponents are best-responding
- We also know that for the best-response strategy, all the actions in the support have the same value

Enumerating Support

- For sparse x, y consider the corresponding elements of $x A_1, A_2 y^T$
- All the elements must correspond to the best-response value from the perspective of the **other** player
- In other words, player needs to mix actions in their support so that the action-values in the opponent's support are best-responding (and thus all the same)

Enumerating Support

	A	B	C
X	(0, 0)	(0, 1)	(-10, -10)
Y	(1, 0)	(-10, -10)	(-10, -10)
Z	(-10, -10)	(-10, -10)	(-10, -10)

Column player needs to mix their actions in the support (A, B) so that the values of the row player for actions in their support (X, Z) are all the same

$$0p(A) + 0p(B) = v_1 \quad (2)$$

$$-10p(A) - 10p(B) = v_1 \quad (3)$$

Same reasoning for the other side

$$0p(X) - 10p(Z) = v_2 \quad (4)$$

$$1p(X) - 10p(Z) = v_2 \quad (5)$$

No solution \Rightarrow no Nash for this support

Enumerating Support

	A	B	C
X	(0, 0)	(0, 1)	(-10, -10)
Y	(1, 0)	(-10, -10)	(-10, -10)
Z	(-10, -10)	(-10, -10)	(-10, -10)

Column player needs to mix their actions in the support (A, B) so that the values of the row player for actions in their support (X, Y) are all the same

$$0p(A) + 0p(B) = v_1 \quad (6)$$

$$1p(A) - 10p(B) = v_1 \quad (7)$$

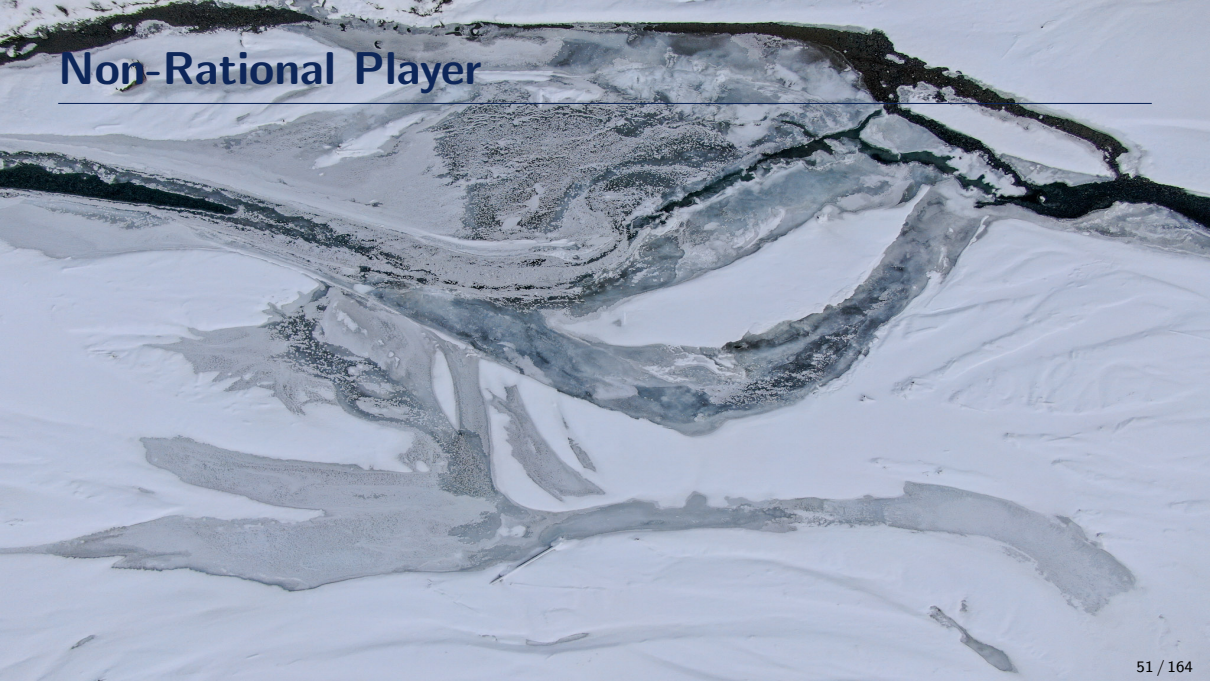
Same reasoning for the other side

$$0p(X) + 0p(Y) = v_2 \quad (8)$$

$$1p(X) - 10p(Y) = v_2 \quad (9)$$

Solution: $p(A) = 0.909091$, $p(B) = 0.090909$, $p(X) = 0.909091$, $p(Y) = 0.090909$

Non-Rational Player



Non-Rational Players - Deviating from Nash

Let's elaborate on the properties of this solution concept!

- What are the implications for the players?
- What are the situations we would/wouldn't play Nash?

Non-Rational Players

Suppose we play versus a stupid opponent

- Non-rational player does not maximize his utility, he can play arbitrarily
- Given Nash equilibrium $\pi = (\pi_1, \pi_2)$, we decide to play π_1 , what do we know?
- Even though π_2 maximizes the utility of the opponent, he can make mistakes and select different (non-equilibristic) strategy π_2'
- Choosing a different strategy than π_2 is no better for the opponent
- But it can be much worse for us! It can be the case that $u_1(\pi_1, \pi_2) \gg u_1(\pi_1, \pi_2')$

Moral of the story: opponent mistakes can hurt us!

Deviating from Nash

	Cooperate	Defect
Cooperate	$(-6, -6)$	$(0, -10)$
Defect	$(-10, 0)$	$(-1, -1)$

Table: Prisoner's dilemma

	Stop	Go
Stop	$(0, 0)$	$(0, 1)$
Go	$(1, 0)$	$(-10, -10)$

Table: The game of Chicken

Rational Players, Multiple Equilibria (I)

Suppose there are two optimal strategy profiles in the game (π_1, π_2^a) and (π_1, π_2^b)

- The opponent is indifferent between his two strategies, he does not care which strategy he chooses (given our strategy π_1), since $u_2(\pi_1, \pi_2^a) = u_2(\pi_1, \pi_2^b)$
- We care!
- $u_1(\pi_1, \pi_2^a) \neq u_1(\pi_1, \pi_2^b)$

Moral of the story: even though both players play optimally, different optimal strategies can lead to different utilities!

Week 2 Homework

1. Draw the best-response value function for a $2 \times N$ matrix game
2. For a two-player matrix game (does not have to be a zero-sum)
 - 2.1 Given a support for both players, construct a set of linear equations and see whether there exists a Nash equilibrium for that support. Use SciPy to solve the constructed linear program²
 - 2.2 Enumerate all possible supports and try to find a Nash equilibrium for each support

²<https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.linprog.html>