



Assignment Project Exams Help

<https://powcoder.com>

Add WeChat powcoder

Paolo Turrini



www.dcs.warwick.ac.uk/~pturrini



p.turrini@warwick.ac.uk

Assignment Project Exam Help

Learning in Games

<https://powcoder.com>
(Artificial) Poker Stars

Add WeChat powcoder

Assignment Project Exam Help

We have seen **extensive games of imperfect information**, where players are typically **uncertain** about the current state of the game being played.

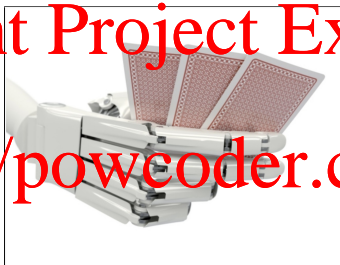
We are going to look at the relevance of this to Artificial Intelligence:

- **Learning through self-play**: the key to many game-playing engines (AlphaGo);
- **Regret**: why it's important to minimise it.

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>



Add WeChat powcoder

“Robots are unlikely to be welcome in casinos any time soon, especially now that a poker-playing computer has learned to play a virtually perfect game — including bluffing.”

(Philip Ball: Game theorists crack poker, *Nature News*, 2015.)

- A difficult game

- Chance, counting odds

- Bluffing, aggressive play

- Still... a game

- An extensive game with imperfect information

- Rational and irrational strategies

Assignment Project Exam Help

<https://powcoder.com>

What is the right solution concept?

Add WeChat powcoder



T.W. Sandholm.

Solving Imperfect-Information Games.

Science, 347(6218):122–123, 2015.

Assignment Project Exam Help

- Everyone who has played poker knows how critical “emotions” are in the game;
- Well... it turns out that there is one emotion that computers use better than anyone else;
- This emotion is regret: how bad I’ve played with respect to how I could have played;
- What is regret really?





<https://powcoder.com>




Add WeChat powcoder

Assignment Project Exam Help

- If you play paper and I play paper, my regret for not playing scissors is 1 (the payoff difference!).
- If you play paper and I play rock, my regret for not playing scissors is 2.



 0	 1	 -1
 -1	 0	 1
 1	 0	 -1

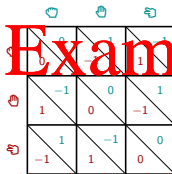
 1	 -1	 0
 -1	 1	 0
 0	 -1	 1






<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

- regret for not playing  in $(\text{blue hand}, \text{blue hand}) = 2$
- regret for not playing  in $(\text{red hand}, \text{blue hand}) = 0$
- regret for not playing  in $(\text{red hand}, \text{red hand}) = -1$



			
	1, -1	0, 0	-1, 1
	-1, 1	1, -1	0, 0

<https://powcoder.com>

Let $\langle N, \mathbf{A}, \mathbf{u} \rangle$ be a normal-form game.

At action profile \mathbf{a} , the regret of player i for not playing a_i is $u_i(a'_i, \mathbf{a}_{-i}) - u_i(\mathbf{a})$.

Avoiding feeling bad (without saying it)

Assignment Project Exam Help

How can we use regret to inform future play?

- The idea is that I want to take actions that I wish I had played in the past;
- Obviously if my opponent knew exactly what I was doing it would not be good;
- My strategy needs to be good and **not exploitable**.

Add WeChat powcoder

Assignment Project Exam Help

<https://powcoder.com>

-1	0	1
1	0	-1
-1	1	0

Add WeChat powcoder

We do this by regret matching: choosing actions at random, with a distribution that is proportional to **positive regrets**.

This means regrets that are proportional to the relative losses one has experienced for not having selected actions in the past.

Assignment Project Exam Help

Games started with (👉, 🤖) followed by (🤖, 🤖):

- regret for not playing 🤖 in (👉, 🤖) = 2
- regret for not playing 🤖 in (👉, 🤖) = 1

	🤖	🤒
👉	-1	1
👊	0	0

<https://powcoder.com>

Add WeChat powcoder

We do this by **regret matching**: choosing actions at random, with a distribution that is proportional to **positive regrets**.

This means regrets that are proportional to the relative losses one has experienced for not having selected actions in the past.

Assignment Project Exam Help

Games started with (♠, ♠) followed by (♠, ♠):

- regret for not playing ♠ in (♠, ♠) = 2
- regret for not playing ♠ in (♠, ♠) = 1

	♠	♥	♣
♠	$\begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$	$\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$
♥	$\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
♣	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

In the next hand, we choose ♠ with probability $\frac{2}{3}$, ♠ with probability $\frac{1}{3}$, ♠ with probability 0.

Notice: positive regrets divided by their sum.

Add WeChat powcoder

We do this by **regret matching**: choosing actions at random, with a distribution that is proportional to **positive regrets**.

This means regrets that are proportional to the relative losses one has experienced for not having selected actions in the past.

Assignment Project Exam Help

- Suppose in the next hand I do play $(2, 3)$... and it turns out to be $(3, 3)$;
- Suppose my opponent plays $(1, 1)$.

	👊	👐	✂️
👊	0	-1	1
👐	1	0	-1
✂️	-1	1	0

<https://powcoder.com>

Add WeChat powcoder

Cumulating regrets

Assignment Project Exam Help

- Suppose in the next hand I do play $(2, 3)$... and it turns out to be $(3, 3)$;

- Suppose my opponent plays $(1, 1)$.

My regret for this hand is 1 for not playing $(1, 1)$,
2 for not playing $(2, 1)$, and obviously 0 for not playing $(3, 1)$.

	1	2	3
1	0/1	-1/0	1/-1
2	1/-1	0/-1	1/0
3	-1/1	1/0	0/1

<https://powcoder.com>

Add WeChat powcoder

Cumulating regrets

Assignment Project Exam Help

- Suppose in the next hand I do play (2, 3) ... and it turns out to be 🖐;
- Suppose my opponent plays 🖐.

My regret for this hand is 1 for not playing 🖐,
2 for not playing 🖐, and obviously 0 for not playing 🖐.

	🖐	🖐	🖐
🖐	0	1	1
🖐	1	0	-1
🖐	-1	1	0

<https://powcoder.com>

Add WeChat powcoder

We add the new regrets to the old ones, and play accordingly.

Cumulating regrets

Assignment Project Exam Help

- Suppose in the next hand I do play $(2, 3)$... and it turns out to be \clubsuit ;

- Suppose my opponent plays \heartsuit .

My regret for this hand is 1 for not playing \heartsuit ,
2 for not playing \clubsuit , and obviously 0 for not playing \spadesuit .

	\clubsuit	\heartsuit	\spadesuit
\clubsuit	0	-1	1
\heartsuit	-1	0	-1
\spadesuit	1	-1	0

We have 2 total regrets for \clubsuit , 2 total regrets for \heartsuit , 2 total regrets for \spadesuit .
Our next strategy is going to be $(\frac{2}{6}, \frac{2}{6}, \frac{2}{6})$

We add the new regrets to the old ones, and play accordingly.

Assignment Project Exam Help

- Cumulative regrets are good, but not very good.
- If our opponent knows that we are using cumulative regrets, then they are always in a position to best respond.

<https://powcoder.com>

Can we do better than this? Yes, we can.

- The key idea, as often the case in modern AI, is to play against ourselves;
- We exploit this 'hypothetical games' to simulate our opponents and strengthen our strategies.

Add WeChat powcoder

Assignment Project Exam Help

Context: you keep playing the same game against the same opponents.

Objective: you want to **learn** their **strategies**.

A good hypothesis might be that the **frequency** with which player i plays action a_i is approximately her probability of playing a_i .

Now suppose you always best-respond to those hypothesised strategies. And suppose everyone else does the same. *What will happen?*

We are going to see that for zero-sum games this process converges to a NE.

This yields a method for **computing** a NE for the (non-repeated) game: just *imagine* players engaging in such “**fictitious play**”.

Assignment Project Exam Help

Given a **history** of actions $H_i^\ell = a_i^0, a_i^1, \dots, a_i^{\ell-1}$ played by player i in ℓ prior plays of game $\langle N, \mathbf{A}, \mathbf{u} \rangle$, fix her **empirical mixed strategy** $s_i^\ell \in S_i$:

$$s_i^\ell(a_i) = \underbrace{\frac{1}{\ell} \cdot \#\{k < \ell : a_i^k = a_i\}}_{\text{relative frequency of } a_i \text{ in } H_i^\ell} \quad \text{for all } a_i \in A_i$$

Add WeChat powcoder

Best Pure Responses

Assignment Project Exam Help

Recall: Strategy $s_i \in S_i$ is a best response for player i to the (partial) strategy profile s_{-i} if $u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i})$ for all $s'_i \in S_i$.

Due to the linearity of expected utilities we get:

Proposition

For any given (partial) strategy profile s_{-i} , the set of best responses for player i must include at least one **pure** strategy.

So we can restrict attention to best pure responses for player i to s_{-i} :

$$a_i^* \in \operatorname{argmax}_{a_i \in A_i} u_i(a_i, s_{-i})$$

Add WeChat powcoder

Fictitious Play

Take any action profile $\mathbf{a}^0 \in A$ for the normal-form game (N, A, u) .

Fictitious play of (N, A, u) starting in \mathbf{a}^0 is the following process:

- In round $\ell = 0$, each player $i \in N$ plays action \mathbf{a}_i^0 .
- In any round $\ell > 0$, each player $i \in N$ plays a **best pure response** to her opponents' **empirical mixed strategies**:

$$\mathbf{a}_i^\ell \in \operatorname{argmax}_{a_i \in A_i} u_i(a_i, \mathbf{s}_{-i}^\ell), \text{ where}$$

$$\mathbf{s}_{i'}^\ell(a_{i'}) = \frac{1}{\ell} \cdot \#\{k < \ell \mid a_{i'}^k = a_{i'}\} \text{ for all } i' \in N \text{ and } a_{i'} \in A_{i'}$$

Assume some deterministic way of breaking ties between maxima.

This yields a sequence $\mathbf{a}^0 \rightarrow \mathbf{a}^1 \rightarrow \mathbf{a}^2 \rightarrow \dots$ with a corresponding sequence of empirical-mixed-strategy profiles $\mathbf{s}^0 \rightarrow \mathbf{s}^1 \rightarrow \mathbf{s}^2 \rightarrow \dots$.

Question: Does $\lim_{\ell \rightarrow \infty} \mathbf{s}^\ell$ exist and is it a meaningful strategy profile?

Example: Matching Pennies

Let's see what happens when we start in the upper lefthand corner HH (and break ties between equally good responses in favour of H):

Assignment Project Exam Help

<https://powcoder.com>

H	-1	1
	1	-1
T	-1	1

Add WeChat powcoder

Any strategy can be represented by a single probability (of playing H).

$$\begin{aligned} HH \left(\frac{1}{1}, \frac{1}{1} \right) &\rightarrow HT \left(\frac{2}{2}, \frac{1}{2} \right) \rightarrow HT \left(\frac{3}{3}, \frac{1}{3} \right) \rightarrow TT \left(\frac{3}{4}, \frac{1}{4} \right) \rightarrow TT \left(\frac{3}{5}, \frac{1}{5} \right) \\ &\rightarrow TT \left(\frac{3}{6}, \frac{1}{6} \right) \rightarrow TH \left(\frac{3}{7}, \frac{2}{7} \right) \rightarrow TH \left(\frac{3}{8}, \frac{3}{8} \right) \rightarrow TH \left(\frac{3}{9}, \frac{4}{9} \right) \\ &\rightarrow TH \left(\frac{3}{10}, \frac{5}{10} \right) \rightarrow HH \left(\frac{4}{11}, \frac{6}{11} \right) \rightarrow HH \left(\frac{5}{12}, \frac{7}{12} \right) \rightarrow \dots \end{aligned}$$

Exercise: Can you guess what this will converge to?

Convergence Profiles are Nash Equilibria

Assignment Project Exam Help

In general, $\lim_{\ell \rightarrow \infty} s^\ell$ does not exist (no guaranteed convergence). But,

Lemma

If fictitious play converges, then it converges to a Nash equilibrium.

Proof: Suppose $s^* = \lim_{\ell \rightarrow \infty} s^\ell$ exists. We need to show that s^* is a NE.

To see that it really is, note that s_i^* is the strategy that player i *seems* to be playing, when in fact she best responds against s_{-i}^* , which she *believes* to be the profile of strategies of her opponents.

Remark: This lemma is true for arbitrary (not just zero-sum) games.

Convergence for Zero-Sum Games

Good news:

Theorem (Robinson, 1951)

For any zero-sum game and initial action profile, fictitious play will converge to a Nash equilibrium.

We know that if FP converges, then to a NE.
Thus, we still have to show that it will converge.

The proof of this fact is difficult and we are not going to discuss it here.



Julia Robinson
(1919–1985)



J. Robinson.

An Iterative Method of Solving a Game.

Annals of Mathematics, 54(2):296–301, 1951.

Playing against ourselves: the procedure

- 1 For each player, initialise cumulative regrets to 0;
- 2 Compute a regret-matching strategy profile;
- 3 Add the strategy profile played to the strategy profile history;
- 4 Select each player action profile according the strategy profile;
- 5 Compute player regrets;
- 6 Add player regrets to the cumulative regrets;
- 7 Repeat, for a fixed number of iterations;
- 8 Return the average strategy profile.

Hart and Mas-Colell (Econometrica, 2000) have shown that this simple procedure converges to a correlated equilibrium, in general, and to the unique NE in two-player zero sum games (like rock-paper-scissors).

Assignment Project Exam Help

- We have all the basics to tackle difficult games;
- The idea extend our basic procedure to extensive games with imperfect information.
- I'm going to present the general procedure

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

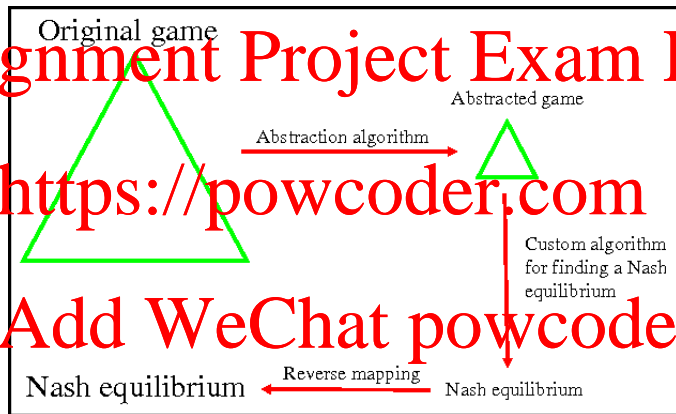
<https://powcoder.com>

Add WeChat powcoder



We are not going to be able to compute the Nash equilibria of Poker.
What can we do instead?

The idea

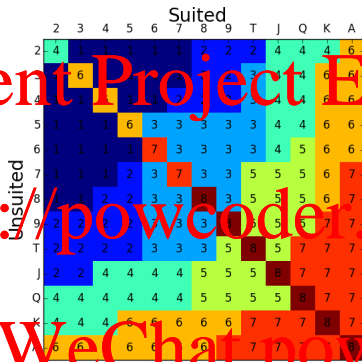


General game-theoretic approach for solving large games.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Preflop hands by strength distance, against various opponents' hands.

 M.Johanson, N.Burch, R.A. Valenzano and M.Bowling

Evaluating state-space abstractions in extensive-form games.

AAMAS, 2013.

Two players, Ann and Bob, are dealt one of the following cards: {A, K, Q}.

Assignment Project Exam Help

Ann	Bob	Ann	outcome
pass	pass		+1 to higher card
pass	bet	pass	+1 to Bob
pass	bet	bet	-2 to higher card
bet	pass		+1 to Ann
bet	bet		+2 to higher card

<https://powcoder.com>

Exercise:

Add WeChat powcoder

Represent Kuhn Poker as an extensive game of perfect information.
How many are the information sets?

Assignment Project Exam Help

The key tool for AI Poker playing engines is Counterfactual Regret Minimization. Basically, CRM ...

- represents the game as a suitably abstracted extensive game of imperfect information;
- uses the regret matching procedure;
- uses behavioural strategies.



M. Zinkevich et al.
Regret Minimization in Games with Incomplete Information.
NIPS, 2012.

Assignment Project Exam Help

We compute our strategy profile iteratively, using the result of regret matching

Let:

- $r_i(h, a)$ be the regret for not taking action a at h
- $r_i^t(h, a)$ be $r_i(h, a)$ at time t
- $r_i^T(h, a)$ be the sum of each $r_i^t(h, a)$ over t , i.e., the cumulative regret.

Let moreover $\mathbf{r}_i(h, a) = \sum_{h' \in h} r_i(h', a)$ be the regret for not taking action a at information set h_i

Finally, let $\mathbf{r}_i^{T+1}(h, a)$ be the corresponding cumulative regret, setting negative regrets to 0.

Assignment Project Exam Help

Now let us define $\sigma_i^{T+1}(h)(a)$, i.e., the probability of taking action a at history h iteratively as follows:

$$\sigma_i^{T+1}(h)(a) := \frac{\mathbf{r}_i^{T+1}(h, a)}{\sum_{a \in \underline{A}(h)} \mathbf{r}_i^{T+1}(h, a)}$$

Requiring that $\sigma_i^{T+1}(h)(a) \rightarrow \frac{1}{|\underline{A}(h)|}$ whenever $\sum_{a \in \underline{A}(h)} \mathbf{r}_i^{T+1}(h, a) = 0$.

Theorem

Add WeChat powcoder

The average strategy profile approaches (Nash) equilibrium as T approaches infinity.

Assignment Project Exam Help

We have seen how players can "learn" their opponents' strategies.

- Self-play as learning in a repeated game
- Convergence to NE if both players do it. Self-play as learning NE!
- Application to Poker

Next: Look at learning in AI and then get back to GT with the new machinery

Add WeChat powcoder