

Multi-Agent Systems

VU AI MSc

Final Exam

17 December 2019, 8:45 – 11:30

General Remarks

BEFORE YOU START

- Check if your version of the exam is complete. Your copy should have 13 printed pages, this one included. The last page is blank and can be used as scrap paper.
- Write down your **name and student ID number** on each sheet.
- **Do NOT remove the staple!**
- The blank space provided for each question should be (more than) sufficient for your answer. You can also use the blank last page, if necessary.
- Your mobile phone has to be switched off and in your coat or bag. Your coat and bag must be under your table.
- The use of a calculator is allowed (but isn't really necessary).

PRACTICAL MATTERS

- You are obliged to identify yourself at the request of the examiner (or his representative) with a proof of your enrollment or a valid ID.
- During the examination it is not permitted to visit the toilet, unless the invigilator gives permission to do so.
- 15 minutes before the end, you will be warned that the time to hand in is approaching.

GOOD LUCK!

1 Copying from Wikipedia for homework

Student life is hectic, and there are many essential life skills to be acquired in limited time: throwing and enjoying great parties, conducting (more or less) profound philosophical discussions into the morning hours, exploring the teeming metropolitan bio-sphere, ... to name just a few. It is therefore completely understandable that homework assignments are seen as an unwelcome distraction and need to be dealt with as efficiently as possible. Fortunately, quite often you can simply copy the relevant answers from Wikipedia, saving a lot of valuable time. Unfortunately, the TAs, who in a recent past used to be students themselves, are aware of these time-saving practices and are prone to check the homeworks for plagiarism. It requires more effort on their part, but they get a lot of satisfaction from catching cheating students. In fact, this situation can be interpreted as a simultaneous game with the following pay-off matrix:

		<i>Student</i>	
		<i>Honest</i>	<i>Wikipedia</i>
<i>TA</i>	<i>Check</i>	5, 0	7, -20
	<i>No_check</i>	10, 0	2, 10

Questions

- (4pts) Determine all the Nash equilibria (NE) for this game.
- (2pts) For each of the NE, compute the expected utility for both student and TA.
- (4pts) What can be done (e.g. in terms of pay-offs) to reduce the probability that a student will cheat?

Solution page (continued)

2 Markov Decision Processes (MDP)

Consider an MDP with a finite number of states s_1, s_2, \dots, s_n and actions a_1, a_2, \dots, a_k . For this MDP we define a policy π that specifies the conditional probabilities $\pi(a|s)$. The state value function \mathbf{v}_π satisfies the matrix form of the Bellman equation:

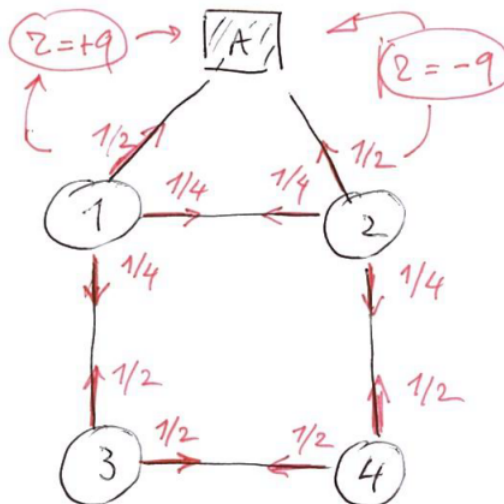
$$\mathbf{v}_\pi = \gamma P \mathbf{v}_\pi + \mathbf{r}$$

where

- $P(s, s') = \sum_a \pi(a|s) p(s'|s, a)$
- $\mathbf{r}(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) r(s, a, s')$,

Questions:

- (2pts) Explain in words the meaning of $P(s, s')$ and $\mathbf{r}(s)$;
 P: ?? **r: reward sum of current step till the end**
- (2pts) Explain in words the meaning of the *product* $P(s, s')P(s', s'')$. How is this different from $P^2(s, s'')$, i.e. **the (s, s'') entry of the matrix $P^2 = P \cdot P$?**
 1. -> P(s, s'') **2.**
- (2pts) Consider the MDP depicted in the figure below. State A is **absorbing**. Transition to A from state 1 yields an immediate reward of 9. Transition to A from state 2 yields an immediate reward of -9. **All other transitions incur a reward of -1.** On this MDP we consider a policy π that assigns transition probabilities as indicated in the figure below. E.g.: $\pi(\text{move to A} | \text{currently in state 1}) = 1/2$ and $\pi(\text{move to 4} | \text{currently in state 2}) = 1/4$, etc. Transitions are deterministic (i.e. each action maps a state s to a unique successor state s').
What are P and \mathbf{r} in this concrete case? Make sure to include the absorbing state A in both P and \mathbf{r} .
- (2pts) Assuming γ is sufficiently small (e.g. $\gamma = 0.1$). How would you calculate an *approximate* solution for \mathbf{v}_π ?
- (2pts) Determine the optimal state value function \mathbf{v}^* assuming $\gamma = 1/3$.

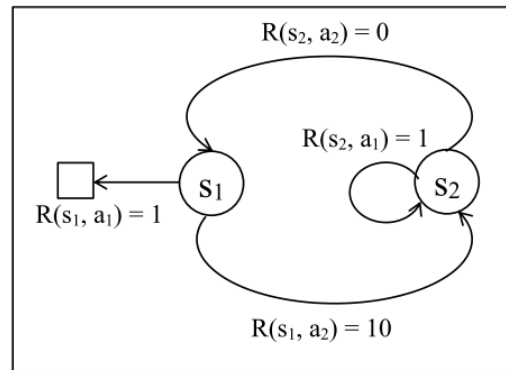


Solution page (continued)

Solution page (continued)

3 MDP 2

Consider the following 2 state MDP.



In both states (s_1 and s_2), there are two possible actions (a_1 and a_2). The actions result in **deterministic transitions**. Taking action a_1 in state s_1 results in a reward of 1, and ends the episode. Taking action a_2 in state s_1 results in a reward of 10, and brings the agent to state s_2 . In state s_2 action a_1 results in a reward of 1 and the agent stays in state s_2 . Action a_2 results in a reward of 0 and brings the agent to state s_1 . **The agent will act (and continues to receive rewards) until the episode ends.**

Questions

- Is this an **infinite horizon**, or a finite horizon problem?
- For a discount factor $\gamma = 0.9$, what is the optimal policy π^* ? Provide the corresponding optimal value $v^*(s_2) = v_{\pi^*}(s_2)$ in state s_2 . Please explain your reasoning and provide your derivation.
- Is it possible to adjust the discount factor γ in such a way that the optimal policy changes? Explain how you would decide whether this is possible or not. If the answer is affirmative, provide an example γ , the corresponding optimal policy, and its corresponding optimal value-function. If you think the answer is negative, provide argument(s).

PS: provide the correct formulae for your answers, even if you can't compute the corresponding numerical result.

Solution page (continued)

4 Vickrey auction

1. (6pts) A Vickrey auction is a sealed bid, second price auction. Explain why truth-telling is a dominant strategy for this auction.
2. (4pts) Would these properties also hold for a sealed bid, third price auction? If affirmative, why would an auctioneer prefer a Vickrey auction to a third price auction? If your answer is negative, explain the difference?

Solution page (continued)

5 Q-learning and SARSA

(10pts) Consider the MDP with a **linear state space**, i.e. **all the states are positioned along a horizontal line**. In each state there are two possible actions: move left ($a = L$) or right ($a = R$). After a number of iteration steps, some of the action values, immediate rewards and current q -values are given by the tabel below. Consider a policy π that picks actions L and R according to the probabilities $\pi(a | s)$ listed in the table below. Furthermore, assume throughout a learning rate $\alpha = 0.9$ and discount factor $\gamma = 2/3$.

$state(s)$	$action(a)$	$next\ state(s')$	$reward(r)$	$q(s, a)$	$\pi(a s)$
2	R	3	-1	5	1/4
2	L	1	0	4	3/4
3	R	4	1	6	2/3
3	L	2	-2	8	1/3

- (2pts) Using this table, what is the current estimate for $v_\pi(2)$ and $v_\pi(3)$ of the state value function v_π .
- (4pts) Compute the next value for $q_\pi(2, R)$ under one **Q-learning** iteration (i.e. only update this state-action pair). Specify both the update formula you're using, and the numerical value that you obtain.
- (4pts) **Expected SARSA** is a variation on SARSA which computes the update using the following formula:

$$q_\pi(S_t, A_t) \leftarrow q_\pi(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \sum_a \pi(a | S_{t+1}) q_\pi(S_{t+1}, a) - q_\pi(S_t, A_t) \right]$$

Compute the next value for $q_\pi(2, R)$ under one iteration step of expected SARSA (using the policy π specified above).

Solution question 5 (continued)

Scrap paper