

**Deadline: 23 December 2020**

### Submission Instructions

- We have organized a Github classroom for this assignment. Use this link <https://classroom.github.com/a/3PgCJKfZ> to accept the assignment and to create inside the classroom your own private repository. Commit there your code and report (in PDF form), before the deadline.
- Follow this naming convention for your PDF report: *Lastname\_Firstname\_studentNo\_affiliation.pdf*.
- Late submissions will be penalized (-1 point/day).
- This is an **individual** assignment, sharing code is considered plagiarism and will be punished. Any other online sources used should be declared in the report. Plagiarism will not be tolerated.
- If you have any questions send an email to: [roxana@ai.vub.ac.be](mailto:roxana@ai.vub.ac.be).

## 1 N-Armed Bandit

For this exercise you will have to implement the N-Armed bandit problem using each of the following action selection methods:

- Random
- $\epsilon$ -Greedy with parameter  $\epsilon$
- Softmax with parameter  $\tau$

All  $Q_{a_i}$  are initialized to 0. The rewards are subject to noise according to a normal distribution with mean  $Q_{a_i}^*$  and standard deviation  $\sigma_i$ ,  $\forall i = 1, \dots, N$ .

You are provided with 5 instances of reward distributions (Tables 1-5) for the N-Armed bandit problem. You only have to select **one** of the tables, according to the following criteria:

- Identify the last number of your VUB/ULB enrolment number
- If the number is 0 or 9, select Table 1
- If the number is 1 or 8, select Table 2
- If the number is 2 or 7, select Table 3
- If the number is 3 or 6, select Table 4
- If the number is 4 or 5, select Table 5

Let  $N = 4$ ,  $Q_{a_i}^*$  and  $\sigma_i$  as defined in your selected table.

| Action    | $Q_{a_i}^*$ | $\sigma_i$ |
|-----------|-------------|------------|
| action #1 | 3.4         | 1.0        |
| action #2 | 3.2         | 0.6        |
| action #3 | 2.1         | 0.4        |
| action #4 | 2.3         | 2          |

| Action    | $Q_{a_i}^*$ | $\sigma_i$ |
|-----------|-------------|------------|
| action #1 | 4.0         | 0.8        |
| action #2 | 3.4         | 0.6        |
| action #3 | 4.2         | 1.1        |
| action #4 | 3.3         | 2          |

| Action    | $Q_{a_i}^*$ | $\sigma_i$ |
|-----------|-------------|------------|
| action #1 | 2.2         | 0.4        |
| action #2 | 2.6         | 0.8        |
| action #3 | 1.0         | 2          |
| action #4 | 1.4         | 0.6        |

Table 1: Last number is 0 or 9.

Table 2: Last number is 1 or 8.

Table 3: Last number is 2 or 7.

| Action    | $Q_{a_i}^*$ | $\sigma_i$ |
|-----------|-------------|------------|
| action #1 | 2.2         | 1.1        |
| action #2 | 1.2         | 2          |
| action #3 | 1.5         | 0.5        |
| action #4 | 2.0         | 0.6        |

Table 4: Last number is 3 or 6.

| Action    | $Q_{a_i}^*$ | $\sigma_i$ |
|-----------|-------------|------------|
| action #1 | 3.0         | 0.7        |
| action #2 | 1.7         | 0.8        |
| action #3 | 2.3         | 2          |
| action #4 | 3.2         | 1.2        |

Table 5: Last number is 4 or 5.

## Plotting

For each of the following exercises, run your simulation for 1000 time steps averaged over multiple trials (e.g., 30) and plot the following graphs:

- One combined plot of the cumulative average reward *over time* for each algorithm (One graph).
- One plot per arm showing the  $Q_{a_i}^*$  of that action along with the *actual*  $Q_{a_i}$  estimate *over time* (all action selection methods should be depicted on the same plot!) (One graph per arm).
- For each algorithm, plot a histogram showing the number of times each action is selected (One graph per action selection strategy).

Do **NOT** add any other graphs than the ones specified here.

See Section 2.2 of “Reinforcement Learning, An Introduction” by Sutton and Barto for a more comprehensive explanation of the problem<sup>1</sup>.

### Remarks:

- Since you are running each experiment (i.e., the 1000 time steps for each action selection method) for multiple trials and then averaging over the results, you can also include in your plots an error measurement, such as standard deviation or standard error.
- **Provide comments and explanations for each of the obtained plots/results! Provide answers to all the questions explicitly asked in the following exercises!**

## Exercise 1.1

After selecting the appropriate table, run the following algorithms and provide the requested plots:

- Random
- $\epsilon$ -Greedy with parameter  $\epsilon = 0$
- $\epsilon$ -Greedy with parameter  $\epsilon = 0.1$
- $\epsilon$ -Greedy with parameter  $\epsilon = 0.2$
- Softmax with parameter  $\tau = 1$
- Softmax with parameter  $\tau = 0.1$

**Comment on the results.** Which algorithms arrive close to the best arm after 1000 iterations? Which one arrives there faster? Explain your findings.

## Exercise 1.2

Re-run Exercise 1 doubling the standard deviations of each arm, and comment on the plots and results. **Discuss the results.** Is this problem harder or easier to learn? Does the performance of the algorithms change significantly? Which of the above performs best now?

<sup>1</sup> <http://incompleteideas.net/book/RLbook2018.pdf>

### Exercise 1.3

Re-run Exercise 1 with two additional algorithms, and plot their results. The new algorithms have a time-varying parameter, which depends on the iteration number  $t = 1, \dots, 1000$  as follows:

- $\epsilon$ -Greedy with parameter  $\epsilon(t) = 1/\sqrt{t}$
- Softmax with parameter  $\tau = 4 * \frac{1000-t}{1000}$

**Discuss the results.** Are these algorithms better than the ones with a fixed parameter?

## 2 Stochastic Reward Game

| Player 1 – action $a_1$ |                               |                              |                            | Player 1 – action $a_2$ |                               |                              |                            |
|-------------------------|-------------------------------|------------------------------|----------------------------|-------------------------|-------------------------------|------------------------------|----------------------------|
|                         | $c_1$                         | $c_2$                        | $c_3$                      |                         | $c_1$                         | $c_2$                        | $c_3$                      |
| $b_1$                   | $\mathcal{N}(11, \sigma_0^2)$ | $\mathcal{N}(-30, \sigma^2)$ | $\mathcal{N}(0, \sigma^2)$ | $b_1$                   | $\mathcal{N}(11, \sigma_0^2)$ | $\mathcal{N}(-3, \sigma^2)$  | $\mathcal{N}(0, \sigma^2)$ |
| $b_2$                   | $\mathcal{N}(-30, \sigma^2)$  | $\mathcal{N}(7, \sigma_1^2)$ | $\mathcal{N}(6, \sigma^2)$ | $b_2$                   | $\mathcal{N}(-3, \sigma^2)$   | $\mathcal{N}(7, \sigma_1^2)$ | $\mathcal{N}(6, \sigma^2)$ |
| $b_3$                   | $\mathcal{N}(0, \sigma^2)$    | $\mathcal{N}(0, \sigma^2)$   | $\mathcal{N}(5, \sigma^2)$ | $b_3$                   | $\mathcal{N}(0, \sigma^2)$    | $\mathcal{N}(0, \sigma^2)$   | $\mathcal{N}(5, \sigma^2)$ |

| Player 1 – action $a_3$ |                               |                              |                            |
|-------------------------|-------------------------------|------------------------------|----------------------------|
|                         | $c_1$                         | $c_2$                        | $c_3$                      |
| $b_1$                   | $\mathcal{N}(14, \sigma_0^2)$ | $\mathcal{N}(-50, \sigma^2)$ | $\mathcal{N}(0, \sigma^2)$ |
| $b_2$                   | $\mathcal{N}(-50, \sigma^2)$  | $\mathcal{N}(4, \sigma_1^2)$ | $\mathcal{N}(3, \sigma^2)$ |
| $b_3$                   | $\mathcal{N}(0, \sigma^2)$    | $\mathcal{N}(0, \sigma^2)$   | $\mathcal{N}(3, \sigma^2)$ |

Table 6: Stochastic Climbing Game

Table 6 shows a 3-player stochastic climbing game. It is a cooperative game, in which three agents/players must learn to maximize the score they receive from the game. Each has three available actions, and the rewards given to both agents depend on the actions of all agents, and are stochastic.

After every interaction, for the chosen combination of actions, the reward received by all agents is sampled from a normal distribution with the mean presented in the above table and a standard deviation of  $\sigma_0$  (for the joint actions  $\langle a_i, b_1, c_1 \rangle, i \in 1, 2, 3$ ),  $\sigma_1$  (for the joint actions  $\langle a_i, b_2, c_2 \rangle, i \in 1, 2, 3$ ) or  $\sigma$  (for all the other possible joint actions). For example, for the first agent taking action 3, the second taking action 1 and the third taking action 2 (corresponding to  $\langle a_3, b_1, c_2 \rangle$ ), their reward is sampled from a normal distribution with mean  $-50$  and standard deviation  $\sigma$ .

Implement two types of Joint-Action learners in the above problem. The first type with simple Boltzmann action selection with a fixed temperature of your choice. The second type with a more advanced multi-agent action selection heuristic of your choice (optimistic Boltzmann, FMQ ...). Recall that in Joint-Action learning, you learn about the quality of joint actions (27 in total in this setting), and you often have to maintain beliefs about the other agents' strategies. An explanation of these concepts can be found here: [1, 2]. Further useful references are [3, 4, 5, 6].

### 2.1 Plotting

Consider the following three cases:

- $\sigma_0 = \sigma_1 = \sigma = 2$
- $\sigma_0 = 4$  and  $\sigma_1 = \sigma = 0.5$
- $\sigma_1 = 4$  and  $\sigma_0 = \sigma = 0.5$

For each case, allow agents of the first type to interact with each other for 5000 episodes and do the same for agents of the second type. Average your results over multiple trials (e.g., 30). After learning, plot the following on the same figure, for each of the three cases:

- Average collected reward per episode versus the episode number for each of the two types. Since you are averaging collected reward over multiple trials, you can also include in your plots an error measurement, such as standard deviation or standard error.

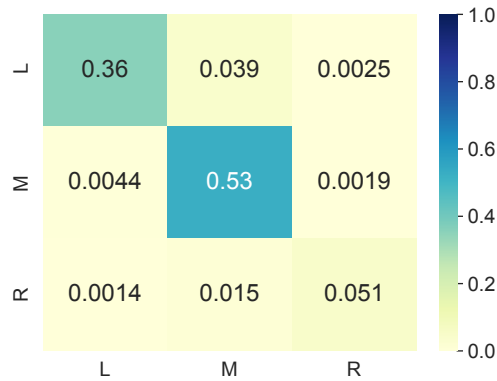


Figure 1: Example of an empirical joint action distribution plot for a game with 2-agents and 3-actions.

- Empirical joint action distribution for the last 10% of the interactions. An example of such a plot is presented in Figure 1. For your case you will need 3 such tables for each of the two types, since you are dealing with a 3-agents game.

Using these plots, you can compare the two algorithms and decide which of the two types learns best in each setting. **Discuss all your results!**

## 2.2 Discussion

Discuss the following:

1. How will the learning process change if we make the agents independent learners?
2. How will the learning process change if we make the agents always select the action that according to them will yield the highest reward (assuming the other agents play the best response)?

NOTE: The choice of the programming language is free. Save yourself work by creating an abstraction for the number of actions and all parameters, so you can re-use your code. **Do not forget to always explain your results!**

## References

- [1] C. Claus and C. Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI/IAAI*, (s 746):752, 1998.
- [2] S. Kapetanakis and D. Kudenko. Reinforcement learning of coordination in heterogeneous cooperative multi-agent systems. In *Adaptive Agents and Multi-Agent Systems II*, pages 119–131. Springer, 2005.
- [3] S. Kapetanakis, D. Kudenko, and M. Strens. Learning of coordination in cooperative multi-agent systems using commitment sequences. *Artificial Intelligence and the Simulation of Behavior*, 1(5), 2004.
- [4] A. Nowé, P. Vrancx, Y. De Hauwere, M. Wiering, and M. van Otterlo. Reinforcement learning: state-of-the-art. *Game Theory and Multi-agent Reinforcement Learning*, pages 441–470, 2012.
- [5] K. Verbeeck, A. Nowé, J. Parent, and K. Tuyls. Exploring selfish reinforcement learning in repeated games with stochastic rewards. *Autonomous Agents and Multi-Agent Systems*, 14(3):239–269, 2007.
- [6] K. Verbeeck, A. Nowé, M. Peeters, and K. Tuyls. Multi-agent reinforcement learning in stochastic single and multi-stage games. In *Adaptive Agents and Multi-Agent Systems II*, pages 275–294. Springer, 2005.