

**Deep Learning and Artificial Intelligence**  
WS 2024/25

**Exercise 12: Policy Gradient Methods**

**Exercise 12-1      Softmax Policy**

- (a) Name an example of a situation in which a stochastic policy is better than a deterministic one.

**Detailed Solution**

A stochastic policy is better:

- In cases where it is disadvantageous that an action is predictable (rock-scissors-paper example)
  - If the environment is partially observable and you need to take different actions despite observing the same features (Aliased Gridworld example)
  - In cases with continuous action spaces
- (b) Given a feature vector  $\mathbf{x}(s, a)$  for state  $s$  and action  $a$ , and a weight vector  $\theta$  with  $\mathbf{x}, \theta \in \mathbb{R}^d$ . The Softmax policy  $\pi_\theta$  parameterized by  $\theta$  is defined as:

$$\pi_\theta(a \mid s) = \frac{e^{\mathbf{x}(s,a)^T \theta}}{\sum_{a' \in \mathcal{A}} e^{\mathbf{x}(s,a')^T \theta}}$$

where  $\mathcal{A}$  is the set of all possible actions.

Calculate the corresponding score function  $\nabla_\theta \log \pi_\theta(a \mid s)$ !

**Detailed Solution**

We have

$$\begin{aligned} \log \pi_\theta(a \mid s) &= \log(e^{\mathbf{x}(s,a)^T \theta}) - \log\left(\sum_{a'} e^{\mathbf{x}(s,a')^T \theta}\right) \\ &= \mathbf{x}(s, a)^T \theta - \log\left(\sum_{a'} e^{\mathbf{x}(s,a')^T \theta}\right) \end{aligned}$$

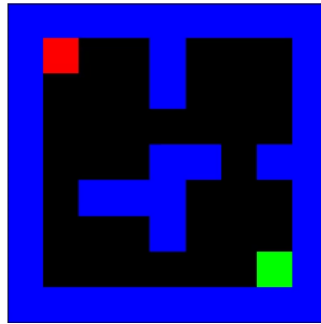
Thus we get:

$$\begin{aligned}
\nabla_{\theta} \log \pi_{\theta}(a \mid s) &= \mathbf{x}(s, a) - \frac{\nabla_{\theta} \sum_{a'} e^{\mathbf{x}(s, a')^T \theta}}{\sum_{a'} e^{\mathbf{x}(s, a')^T \theta}} \\
&= \mathbf{x}(s, a) - \frac{\sum_{a'} \mathbf{x}(s, a') e^{\mathbf{x}(s, a')^T \theta}}{\sum_{a'} e^{\mathbf{x}(s, a')^T \theta}} \\
&= \mathbf{x}(s, a) - \sum_{a'} \pi_{\theta}(a' \mid s) \mathbf{x}(s, a') \\
&= \mathbf{x}(s, a) - \mathbb{E}_{\pi_{\theta}}[\mathbf{x}(s, \cdot)].
\end{aligned}$$

## Exercise 12-2 REINFORCE

In this exercise, you will implement the REINFORCE algorithm, a policy gradient method that uses the return of complete episodes for the updates of the policy parameter.

On the website, you can find a zip file containing the files “rooms.py”, “montecarlo\_main.py” and “REINFORCE.py”. The first file contains a class “RoomsEnv” which simulates the rooms domain depicted in the figure below.



The goal of the agent (red square, upper left) is to find a path to the green goal state (bottom right). The blue squares are walls that cannot be walked through. The reward is 0 at all steps and 1 when reaching the goal state. It is an episodic task and we use a discounting factor  $\gamma$ . The main file already implements the simulation of complete episodes (function “train”) and stores the sampled states, actions and rewards received in separate arrays. After each episode, it calls `agent.update_montecarlo()`.

Task:

In the file “REINFORCE.py”, implement the missing functionality for the class `ReinforceAgent`:

- The method `softmax_policy()` should return the probability  $\pi(a \mid s)$  of choosing an action  $a$  in state  $s$  as given in exercise 13-1.
- The function `score_function()` should calculate the gradient  $\nabla_{\theta} \log \pi_{\theta}(a \mid s)$
- The method `choose_action` should sample an action for state  $s$  according to the probabilities given by the softmax policy.
- Finally, the method `update_montecarlo()` should take the sampled states, actions and rewards arrays and perform the update of the policy parameter `self.theta` for each time step  $t$ :

$$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi_{\theta}(a \mid s),$$

where  $G_t$  is the discounted return at time step  $t$  (for the last time step, it is just the last reward in the rewards array).

Train the agent and look at the results. You can also play around with different feature vectors. By calling `env.save_video()` you can generate a visualization of the agent acting according to the trained policy.

### Exercise 12-3 Policy Gradient with Baseline

In this exercise we will take a look at the effect of the baseline in Policy Gradient algorithms. Assume we have an MDP and a stochastic policy  $\pi_\theta(a|s)$  given. Recall, the goal of policy gradient is to update the parameters  $\theta$  such that the expected return will be maximized. Let us consider the updates for the given state  $s$ . The feature vector  $\mathbf{x}(s, a)$  is defined as follows:

$$\mathbf{x}(s, a) = \begin{cases} (1, 0, 0)^T & \text{if } a = a_1 \\ (0, 1, 0)^T & \text{if } a = a_2 \\ (0, 0, 1)^T & \text{if } a = a_3 \end{cases}$$

The average return  $G$  for action  $a_1, a_2, a_3$  is 101, 110, 104 respectively. Assume that we have values for  $\theta$  such that the softmax-policy  $\pi_\theta(a|s)$  is defined as follows:

$$\pi_\theta(a|s) = \begin{cases} 0.3 & \text{if } a = a_1 \\ 0.6 & \text{if } a = a_2 \\ 0.1 & \text{if } a = a_3 \end{cases}$$

- (a) Compute the mean and the variance for the latter part of the update rule of policy gradient, that is

$$G_t \nabla_\theta \log \pi_\theta(a | s).$$

#### Detailed Solution

Mean:

$$\begin{aligned} & E_{\pi_\theta} [G_t \nabla_\theta \log \pi_\theta(a | s)] \\ &= E_{\pi_\theta} \left[ G_t (\mathbf{x}(s, a) - \sum_{a'} \pi_\theta(a') \mathbf{x}(s, a')) \right] \\ &= 0.3(101 ((1, 0, 0)^T - (0.3 (1, 0, 0)^T + 0.6 (0, 1, 0)^T + 0.1 (0, 0, 1)^T))) \\ &\quad + 0.6(110 ((0, 1, 0)^T - (0.3, 0.6, 0.1)^T)) \\ &\quad + 0.1(104 ((0, 0, 1)^T - (0.3, 0.6, 0.1)^T)) \\ &= (21.21, -18.18, -3.03)^T + (-19.8, 26.4, -6.6)^T + (-3.12, -6.24, 9.36)^T \\ &= (-1.71, 1.98, -0.27)^T \end{aligned}$$

#### Detailed Solution

Variance:

$$\begin{aligned} & E_{\pi_\theta} [(G_t \nabla_\theta \log \pi_\theta(a | s) - \mu)^2] \\ &= E_{\pi_\theta} [(G_t \nabla_\theta \log \pi_\theta(a | s))^2] - E_{\pi_\theta} [G_t \nabla_\theta \log \pi_\theta(a | s)]^2 \\ &\approx (2247.37, 2648.76, 979.23)^T \end{aligned}$$

- (b) What happens if we subtract a baseline  $b(s) = 100$  from the return  $G_t$ ? Compute the mean and variance for the formula

$$(G_t - b(s)) \nabla_{\theta} \log \pi_{\theta}(a | s).$$

### Detailed Solution

Mean:

$$\begin{aligned} & E_{\pi_{\theta}} [(G_t - b(s)) \nabla_{\theta} \log \pi_{\theta}(a | s)] \\ &= E_{\pi_{\theta}} \left[ (G_t - b(s)) (\mathbf{x}(s, a) - \sum_{a'} \pi_{\theta}(a') \mathbf{x}(s, a')) \right] \\ &= 0.3(1 ((1, 0, 0)^T - (0.3 (1, 0, 0)^T + 0.6 (0, 1, 0)^T + 0.1 (0, 0, 1)^T)) \\ &\quad + 0.6(10 ((0, 1, 0)^T - (0.3, 0.6, 0.1)^T)) \\ &\quad + 0.1(4 ((0, 0, 1)^T - (0.3, 0.6, 0.1)^T)) \\ &= (0.21, -0.18, -0.03)^T + (-1.8, 2.4, -0.6)^T + (-0.12, -0.24, 0.36)^T \\ &\approx (-1.71, 1.98, -0.27)^T \end{aligned}$$

### Detailed Solution

Variance:

$$\begin{aligned} & E_{\pi_{\theta}} [((G_t - b(s)) \nabla_{\theta} \log \pi_{\theta}(a | s) - \mu)^2] \\ &= E_{\pi_{\theta}} [((G_t - b(s)) \nabla_{\theta} \log \pi_{\theta}(a | s))^2] - E_{\pi_{\theta}} [(G_t - b(s)) \nabla_{\theta} \log \pi_{\theta}(a | s)]^2 \\ &\approx (2.77, 6.36, 1.83)^T \end{aligned}$$

- (c) Compare the results of (a) and (b). What do you observe? How does the baseline affect the learning of the optimal policy?

### Detailed Solution

A reasonable baseline can reduce the variance of the updates and thus the parameters will converge more stable to the same result.