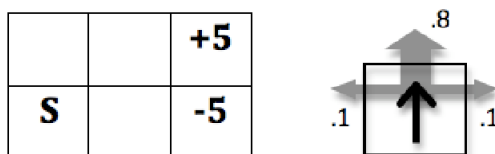


THE HONG KONG UNIVERSITY OF SCIENCE & TECHNOLOGY
Machine Learning
Homework 4

For self-practice. No need to submit your work.
Everyone gets 4 points nonetheless.

Solutions will be provided later.

Question 1: Consider an agent that acts in the gridworld shown below. The agent always starts in state $(1, 1)$, marked with the letter S . There are two terminal goal states, $(3, 2)$ with reward $+5$ and $(3, 1)$ with reward -5 . Rewards are 0 in non-terminal states. (The reward for a state is received as the agent moves into the state.) The transition function is such that the intended agent movement (North, South, West, or East) happens with probability 0.8. With probability 0.1 each, the agent ends up in one of the states perpendicular to the intended direction. If a collision with a wall happens, the agent stays in the same state.



The expected immediate reward function $r(s, a) = \sum_{s'} r(s, a, s')P(s'|s, a)$ is as follows:

$r(s, a)$	N	S	W	E
$(1, 1)$	0	0	0	0
$(1, 2)$	0	0	0	0
$(2, 1)$	-0.5	-0.5	0	-4
$(2, 2)$	0.5	0.5	0	4
$(3, 1)$	0	0	0	0
$(3, 2)$	0	0	0	0

- Assume the initial value function $Q_0(s, a) = 0$ for all states s and actions a . Let $\gamma = 0.9$. The Q-function Q_1 after the first value iteration is the same as $r(s, a)$. What is the Q-function Q_2 after the second value iteration? What is the greedy policy π_2 based on Q_2 . In case of ties, list all tied actions.
- Suppose the agent does not know the transition probabilities and the reward function, and it tries to learn by interacting with the environment. Assume the Q-learning algorithm is used with $Q(s, a) = 0$ initially. Let $\alpha = 0.1$ and $\gamma = 0.9$. Update the Q-function using the following experience tuples. Show the function after each update.

s	a	r	s'
$(2, 2)$	E	5	$(3, 2)$
$(2, 1)$	N	0	$(2, 2)$
$(1, 2)$	E	0	$(2, 2)$
$(1, 1)$	N	0	$(1, 2)$

Give the greedy policy based on the latest Q function. In case of ties, list all tied actions.

Question 2: Here is the parameter update rule for Deep Q-Networks:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} ([r(s, a) + \gamma \max_{a'} Q(s', a'; \theta^-)] - Q(s, a; \theta))^2$$

What do s , a , $r(s, a)$ and s' stand for? What about θ^- ? What is the objective that the update rule is intended to achieve?

Question 3: Here is the update rule for the actor in the Actor-Critic algorithm:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a|s) \hat{A}^{\pi}(s, a),$$

where $\hat{A}^{\pi}(s, a) \leftarrow r + \gamma \hat{V}_{\phi}^{\pi}(s') - \hat{V}_{\phi}^{\pi}(s)$.

What do s , a , $r(s, a)$ and s' stand for? What about $\hat{V}_{\phi}^{\pi}(s)$? Intuitively, what does the update rule try to achieve?

Question 4: (a) In the context of deep image classification, what is adversarial attack?

(b) The CW attack finds an adversarial example \mathbf{x}' for a benign example \mathbf{x} by solving the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{x}'} c \|\mathbf{x} - \mathbf{x}'\|_2^2 + l(\mathbf{x}') \\ & \text{s.t. } \mathbf{x}' \in [0, 1]^n \\ & \text{where } l(\mathbf{x}') = \max\{\max_{i \neq t} Z_i(\mathbf{x}') - Z_t(\mathbf{x}'), -\kappa\}. \end{aligned}$$

What do the terms $Z_t(\mathbf{x}')$, $Z_i(\mathbf{x}')$ stand for? What are the key ideas behind this attack?

Question 5: In some pixel-level explanations, we need to compute $\frac{\partial z_c(\mathbf{x})}{\partial x_i}$. The backpropagation algorithm for the task is given in Lecture 16. Suppose the last layer of the model is a Softmax layer. How would you change the algorithm if we are to compute $\frac{\partial \log P(y=c|\mathbf{x})}{\partial x_i}$?