

Q1. Q-learning

Consider an unknown MDP with three states (A , B and C) and two actions (\leftarrow and \rightarrow). Suppose the agent chooses actions according to some policy π in the unknown MDP, collecting a dataset consisting of samples (s, a, s', r) representing taking action a in state s resulting in a transition to state s' and a reward of r .

s	a	s'	r
A	\rightarrow	B	2
C	\leftarrow	B	2
B	\rightarrow	C	-2
A	\rightarrow	B	4

You may assume a discount factor of $\gamma = 1$.

(a) Recall the update function of Q -learning is:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_{a'} Q(s_{t+1}, a') \right)$$

Assume that all Q -values are initialized to 0, and use a learning rate of $\alpha = \frac{1}{2}$.

(i) Run Q -learning on the above experience table and fill in the following Q -values:

$Q(A, \rightarrow) = \underline{\hspace{10em}}$ $Q(B, \rightarrow) = \underline{\hspace{10em}}$

(ii) After running Q -learning and producing the above Q -values, you construct a policy π_Q that maximizes the Q -value in a given state:

$$\pi_Q(s) = \arg \max_a Q(s, a).$$

What are the actions chosen by the policy in states A and B ?

$\pi_Q(A)$ is equal to:

- $\pi_Q(A) = \leftarrow$.
- $\pi_Q(A) = \rightarrow$.
- $\pi_Q(A) = \text{Undefined}$.

$\pi_Q(B)$ is equal to:

- $\pi_Q(B) = \leftarrow$.
- $\pi_Q(B) = \rightarrow$.
- $\pi_Q(B) = \text{Undefined}$.

(b) Use the empirical frequency count model-based reinforcement learning method described in lectures to estimate the transition function $\hat{T}(s, a, s')$ and reward function $\hat{R}(s, a, s')$. (Do not use pseudocounts; if a transition is not observed, it has a count of 0.)

Write down the following quantities. You may write N/A for undefined quantities.

$\hat{T}(A, \rightarrow, B) = \underline{\hspace{10em}}$ $\hat{R}(A, \rightarrow, B) = \underline{\hspace{10em}}$

$$\hat{T}(B, \rightarrow, A) = \underline{\hspace{10em}} \quad \hat{R}(B, \rightarrow, A) = \underline{\hspace{10em}}$$

$$\hat{T}(B, \leftarrow, A) = \underline{\hspace{10em}} \quad \hat{R}(B, \leftarrow, A) = \underline{\hspace{10em}}$$

(c) This question considers properties of reinforcement learning algorithms for *arbitrary* discrete MDPs; you do not need to refer to the MDP considered in the previous parts.

(i) Which of the following methods, at convergence, provide enough information to obtain an optimal policy? (Assume adequate exploration.)

Model-based learning of $T(s, a, s')$ and $R(s, a, s')$.

Direct Evaluation to estimate $V(s)$.

Temporal Difference learning to estimate $V(s)$.

Q-Learning to estimate $Q(s, a)$.

(ii) In the limit of infinite timesteps, under which of the following exploration policies is Q-learning guaranteed to converge to the optimal Q-values for all state? (You may assume the learning rate α is chosen appropriately, and that the MDP is ergodic: i.e., every state is reachable from every other state with non-zero probability.)

A fixed policy taking actions uniformly at random.

A greedy policy.

An ϵ -greedy policy

A fixed optimal policy.

2 Pacman with Feature-Based Q-Learning

We would like to use a Q-learning agent for Pacman, but the size of the state space for a large grid is too massive to hold in memory. To solve this, we will switch to feature-based representation of Pacman's state.

(a) We will have two features, F_g and F_p , defined as follows:

$$F_g(s, a) = A(s) + B(s, a) + C(s, a)$$
$$F_p(s, a) = D(s) + 2E(s, a)$$

where

- $A(s)$ = number of ghosts within 1 step of state s
- $B(s, a)$ = number of ghosts Pacman touches after taking action a from state s
- $C(s, a)$ = number of ghosts within 1 step of the state Pacman ends up in after taking action a
- $D(s)$ = number of food pellets within 1 step of state s
- $E(s, a)$ = number of food pellets eaten after taking action a from state s

For this pacman board, the ghosts will always be stationary, and the action space is $\{left, right, up, down, stay\}$.

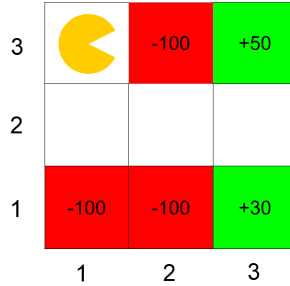


calculate the features for the actions $\in \{left, right, up, stay\}$

- (b) After a few episodes of Q-learning, the weights are $w_g = -10$ and $w_p = 100$. Calculate the Q value for each action $\in \{left, right, up, stay\}$ from the current state shown in the figure.
- (c) We observe a transition that starts from the state above, s , takes action up , ends in state s' (the state with the food pellet above) and receives a reward $R(s, a, s') = 250$. The available actions from state s' are $down$ and $stay$. Assuming a discount of $\gamma = 0.5$, calculate the new estimate of the Q value for s based on this episode.
- (d) With this new estimate and a learning rate (α) of 0.5, update the weights for each feature.

3 Deep inside Q-learning

Consider the grid-world given below and an agent who is trying to learn the optimal policy. Rewards are only awarded for taking the *Exit* action from one of the shaded states. Taking this action moves the agent to the Done state, and the MDP terminates. Assume $\gamma = 1$ and $\alpha = 0.5$ for all calculations. All equations need to explicitly mention γ and α if necessary.



- (a) The agent starts from the top left corner and you are given the following episodes from runs of the agent through this grid-world. Each line in an Episode is a tuple containing (s, a, s', r) .

Episode 1	Episode 2	Episode 3	Episode 4	Episode 5
(1,3), S, (1,2), 0	(1,3), S, (1,2), 0	(1,3), S, (1,2), 0	(1,3), S, (1,2), 0	(1,3), S, (1,2), 0
(1,2), E, (2,2), 0	(1,2), E, (2,2), 0	(1,2), E, (2,2), 0	(1,2), E, (2,2), 0	(1,2), E, (2,2), 0
(2,2), E, (3,2), 0	(2,2), S, (2,1), 0	(2,2), E, (3,2), 0	(2,2), E, (3,2), 0	(2,2), E, (3,2), 0
(3,2), N, (3,3), 0	(2,1), Exit, D, -100	(3,2), S, (3,1), 0	(3,2), N, (3,3), 0	(3,2), S, (3,1), 0
(3,3), Exit, D, +50		(3,1), Exit, D, +30	(3,3), Exit, D, +50	(3,1), Exit, D, +30

Fill in the following Q-values obtained from direct evaluation from the samples:

$$Q((3,2), N) = \underline{\hspace{2cm}} \quad Q((3,2), S) = \underline{\hspace{2cm}} \quad Q((2,2), E) = \underline{\hspace{2cm}}$$

- (b) Q-learning is an online algorithm to learn optimal Q-values in an MDP with unknown rewards and transition function. The update equation is:

$$Q(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q(s_{t+1}, a'))$$

where γ is the discount factor, α is the learning rate and the sequence of observations are $(\dots, s_t, a_t, s_{t+1}, r_t, \dots)$. Given the episodes in (a), fill in the time at which the following Q values first become non-zero. When updating the Q values, You should only go through each transition once, and the order in which you are to go through them is: transitions in ep 1, transitions in ep 2 and so on. Your answer should be of the form **(episode#,iter#)** where **iter#** is the Q-learning update iteration in that episode. If the specified Q value never becomes non-zero, write *never*.

$$Q((1,2), E) = \underline{\hspace{2cm}} \quad Q((2,2), E) = \underline{\hspace{2cm}} \quad Q((3,2), S) = \underline{\hspace{2cm}}$$

- (c) In Q-learning, we look at a window of (s_t, a_t, s_{t+1}, r_t) to update our Q-values. One can think of using an update rule that uses a larger window to update these values. Give an update rule for $Q(s_t, a_t)$ given the window $(s_t, a_t, r_t, s_{t+1}, a_{t+1}, r_{t+1}, s_{t+2})$.

$$Q(s_t, a_t) =$$

$$Q(s_t, a_t) =$$

$$Q(s_t, a_t) =$$