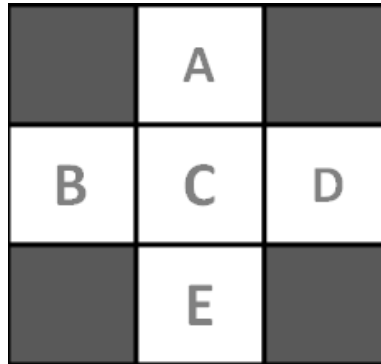


## 1 Learning in Gridworld

Consider the example gridworld that we looked at in lecture. We would like to use TD learning and q-learning to find the values of these states.



Suppose that we have the following observed transitions:

(B, East, C, 2), (C, South, E, 4), (C, East, A, 6), (B, East, C, 2)

The initial value of each state is 0. Assume that  $\gamma = 1$  and  $\alpha = 0.5$ .

1. What are the learned values from TD learning after all four observations?

$$V(B) = 3.5$$

$$V(C) = 4$$

All other states have a value of 0.

2. What are the learned Q-values from Q-learning after all four observations?

$$Q(B, \text{East}) = 3$$

$$Q(C, \text{South}) = 2$$

$$Q(C, \text{East}) = 3$$

All other q-states have a value of 0.

## Q2. 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.

1. 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  $\in \{left, right, up, stay\}$

$$F_p(s, up) = 1 + 2(1) = 3$$

$$F_p(s, left) = 1 + 2(0) = 1$$

$$F_p(s, right) = 1 + 2(0) = 1$$

$$F_p(s, stay) = 1 + 2(0) = 1$$

$$F_g(s, up) = 2 + 0 + 0 = 2$$

$$F_g(s, left) = 2 + 1 + 1 = 4$$

$$F_g(s, right) = 2 + 1 + 1 = 4$$

$$F_g(s, stay) = 2 + 0 + 2 = 4$$

2. 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.

$$Q(s, up) = w_p F_p(s, up) + w_g F_g(s, up) = 100(3) + (-10)(2) = 280$$

$$Q(s, left) = w_p F_p(s, left) + w_g F_g(s, left) = 100(1) + (-10)(4) = 60$$

$$Q(s, right) = w_p F_p(s, right) + w_g F_g(s, right) = 100(1) + (-10)(4) = 60$$

$$Q(s, stay) = w_p F_p(s, stay) + w_g F_g(s, stay) = 100(1) + (-10)(4) = 60$$

3. 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.

$$\begin{aligned}
 Q_{new}(s, a) &= R(s, a, s') + \gamma * \max_{a'} Q(s', a') \\
 &= 250 + 0.5 * \max\{Q(s', down), Q(s', stay)\} \\
 &= 250 + 0.5 * 0 \\
 &= 250
 \end{aligned}$$

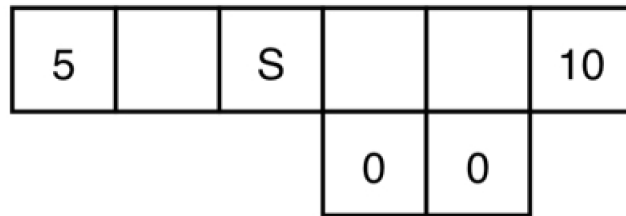
where

$$\begin{aligned}
 Q(s', down) &= w_p F_p(s, down) + w_g F_g(s, down) = 100(0) + (-10)(2) = -20 \\
 Q(s', stay) &= w_p F_p(s, stay) + w_g F_g(s, stay) = 100(0) + (-10)(0) = 0
 \end{aligned}$$

4. With this new estimate and a learning rate ( $\alpha$ ) of 0.5, update the weights for each feature.

$$\begin{aligned}
 w_p &= w_p + \alpha * (Q_{new}(s, a) - Q(s, a)) * F_p(s, a) = 100 + 0.5 * (250 - 280) * 3 = 55 \\
 w_g &= w_g + \alpha * (Q_{new}(s, a) - Q(s, a)) * F_g(s, a) = -10 + 0.5 * (250 - 280) * 2 = -40
 \end{aligned}$$

### Q3. MDPs and RL

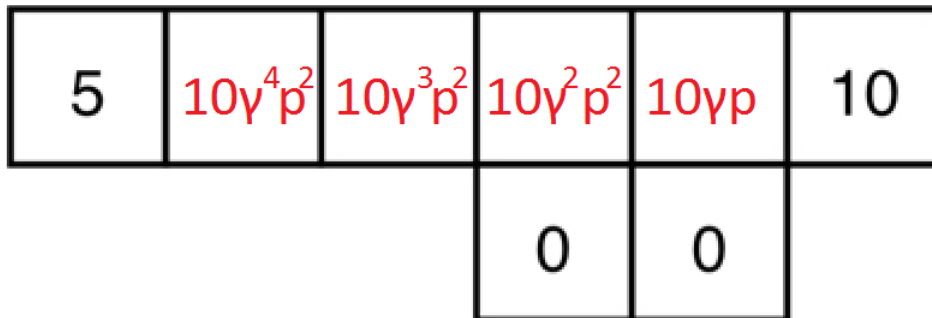


Consider the above gridworld. An agent is currently on grid cell  $S$ , and would like to collect the rewards that lie on both sides of it. If the agent is on a numbered square, its only available action is to Exit, and when it exits it gets reward equal to the number on the square. On any other (non-numbered) square, its available actions are to move East and West. Note that North and South are never available actions.

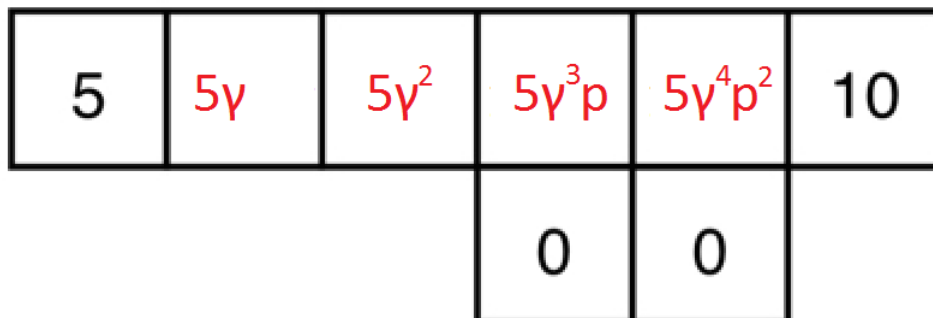
If the agent is in a square with an adjacent square downward, it does not always move successfully: when the agent is in one of these squares and takes a move action, it will only succeed with probability  $p$ . With probability  $1 - p$ , the move action will fail and the agent will instead move downwards. If the agent is not in a square with an adjacent space below, it will always move successfully.

For parts (a) and (b), we are using discount factor  $\gamma \in [0, 1]$ .

- (a) Consider the policy  $\pi_{\text{East}}$ , which is to always move East (right) when possible, and to Exit when that is the only available action. For each non-numbered state  $x$  in the diagram below, fill in  $V^{\pi_{\text{East}}}(x)$  in terms of  $\gamma$  and  $p$ .



- (b) Consider the policy  $\pi_{\text{West}}$ , which is to always move West (left) when possible, and to Exit when that is the only available action. For each non-numbered state  $x$  in the diagram below, fill in  $V^{\pi_{\text{West}}}(x)$  in terms of  $\gamma$  and  $p$ .



(c) For what range of values of  $p$  in terms of  $\gamma$  is it optimal for the agent to go West (left) from the start state ( $S$ )?

We want  $5\gamma^2 \geq 10\gamma^3 p^2$ , which we can solve to get:

$$\text{Range: } p \in [0, \frac{1}{\sqrt{2\gamma}}]$$

(d) For what range of values of  $p$  in terms of  $\gamma$  is  $\pi_{\text{West}}$  the optimal policy?

We need, for each of the four cells, to have the value of that cell under  $\pi_{\text{West}}$  to be at least as large as  $\pi_{\text{East}}$ .

Intuitively, the farther east we are, the higher the value of moving east, and the lower the value of moving west (since the discount factor penalizes far-away rewards).

Thus, if moving west is the optimal policy, we want to focus our attention on the rightmost cell.

At the rightmost cell, in order for moving west to be optimal, then  $V^{\pi_{\text{East}}}(s) \leq V^{\pi_{\text{West}}}(s)$ , which is  $10\gamma p \leq 5\gamma^4 p^2$ , or  $p \geq \frac{2}{\gamma^3}$ .

However, since  $\gamma$  ranges from 0 to 1, the right side of this expression ranges from 2 to  $\infty$ , which means  $p$  (a probability, and thus bounded by 1) has no valid value.

Range:  $\emptyset$

(e) For what range of values of  $p$  in terms of  $\gamma$  is  $\pi_{\text{East}}$  the optimal policy?

We follow the same logic as in the previous part. Specifically, we focus on the leftmost cell, where the condition for  $\pi_{\text{East}}$  to be the optimal policy is:  $10\gamma^4 p^2 \geq 5\gamma$ , which simplifies to  $p \geq \frac{1}{\sqrt{2\gamma^3}}$ . Combined with our bound on any probability being in the range  $[0, 1]$ , we get:

$$\text{Range: } p \in \left[ \frac{1}{\sqrt{2\gamma^3}}, 1 \right], \text{ which could be an empty set depending on } \gamma.$$

Recall that in approximate Q-learning, the Q-value is a weighted sum of features:  $Q(s, a) = \sum_i w_i f_i(s, a)$ . To derive a weight update equation, we first defined the loss function  $L_2 = \frac{1}{2}(y - \sum_k w_k f_k(x))^2$  and found  $dL_2/dw_m = -(y - \sum_k w_k f_k(x))f_m(x)$ . Our label  $y$  in this set up is  $r + \gamma \max_a Q(s', a')$ . Putting this all together, we derived the gradient descent update rule for  $w_m$  as  $w_m \leftarrow w_m + \alpha (r + \gamma \max_a Q(s', a') - Q(s, a)) f_m(s, a)$ .

In the following question, you will derive the gradient descent update rule for  $w_m$  using a different loss function:

$$L_1 = \left| y - \sum_k w_k f_k(x) \right|$$

(f) Find  $dL_1/dw_m$ . Show work to have a chance at receiving partial credit. Ignore the non-differentiable point.

Note that the derivative of  $|x|$  is  $-1$  if  $x < 0$  and  $1$  if  $x > 0$ . So for  $L_1$ , we have:

$$\frac{dL_1}{dw_m} = \begin{cases} -f_m(x) & y - \sum_k w_k f_k(x) > 0 \\ f_m(x) & y - \sum_k w_k f_k(x) < 0 \end{cases}$$

(g) Write the gradient descent update rule for  $w_m$ , using the  $L_1$  loss function.

$$w_m \leftarrow w_m - \alpha dL_1/dw_m \\ \leftarrow \begin{cases} w_m + \alpha f_m(x) & y - \sum_k w_k f_k(x) > 0 \\ w_m - \alpha f_m(x) & y - \sum_k w_k f_k(x) < 0 \end{cases}$$