

# Reinforcement Learning

Final Exam

11/12/2021

Sanjit K. Kaul

**Instructions:** You have two hours to work on the questions and an additional ten minutes to upload a single PDF containing your work. Please ensure that you give uploading sufficient time. Delayed uploads will not be graded. **Final answers with no supporting steps will receive no credit.** Exam is open book and class notes. No other resources, human or otherwise are allowed. Violation of the policy will be treated as use of unfair means and plagiarism. All academic penalties will apply.

**Question 1. 20 marks** Our environment has one state 0 and a terminal state  $Z$ . In state 0 one can take the two actions of  $-1$  and  $1$ . On taking either action the environment transitions to 0 with probability 0.2 and the agent gets a reward of 0. With probability 0.8, the agent gets a reward 1 and the episode terminates. You want to estimate the value of state 0 when using the policy  $\pi$  that picks action 1 with probability 1. However, you only have access to the policy  $\mu$  that picks either action with probability 0.5. You are allowed to run one episode. Derive the expected value and variance of your estimate of  $v_\pi(0)$ . Assume a discount factor of 1.

**Question 2. 10 marks** Our environment has two states  $s_0$  and  $s_1$  and two terminal states. Five episodes obtained by using a policy  $\pi$  are listed next.

$s_0$  0.5  $s_1$  0,

$s_0$  0.2,

$s_0$  0.3,

$s_1$  0.6,

$s_1$  0.4.

The episodes above are summarised as  $S_t, R_{t+1}, S_{t+1}, R_{t+2}, \dots$ . You are interested in estimating the values of the states. You can use either TD(0) or first visit MC. Derive the estimates both will converge to. Explain your steps.

**Question 3. 20 marks** Our environment has two states  $A$  and  $B$  and actions  $-1$  and  $1$ .  $Z$  is the terminal state. Consider the following episode, where in an episode is summarised as  $(S_t, A_t, R_{t+1}), (S_{t+1}, A_{t+1}, R_{t+2}), \dots$ . The episode is

$(A, 1, 5), (B, -1, 5), (B, 1, 3), (A, -1, -5), (B, 1, 5), Z$ .

Use the episode to derive the action value estimates one would obtain using SARSA and Q-learning.

**Question 4. 20 marks** Consider the following experiences, wherein an experience is in the form  $S_t, A_t, R_{t+1}, S_{t+1}$ .

$A$  and  $B$  are states and  $Z$  is the terminal state.

$A, -1, 3, A,$   
 $A, -1, 5, B,$   
 $A, 1, 2, A,$   
 $B, -1, 4, Z,$   
 $B, 1, 3, A.$

Assume initial action-values of  $Q(A, 1) = 2$ ,  $Q(A, -1) = 3$ ,  $Q(B, 1) = 2$ ,  $Q(B, -1) = 3$ . Use these to associate with each experience a TD error. Calculate the probability with which an experience would be picked when using (a) greedy TD error experience replay, (b) uniform replay, and (c) rank based prioritization. Show all your steps.

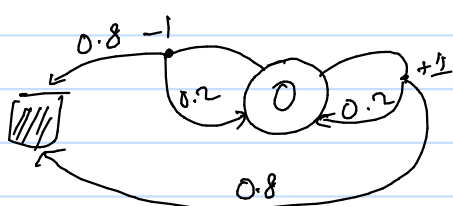
**Question 5. 10 marks** We have states  $1, 2, \dots, n$ . Assume a linear approximation architecture. State  $i$  is mapped to a  $n$ -dimensional feature vector with the  $i^{\text{th}}$  element set to 1 and the other elements set to 0. Suppose we observe an episode with the sequence  $1, 2, \dots, n$  of states. Derive the eligibility trace vector  $z_n$  for the episode. Assume TD(0.3) and a discount factor of 0.8.

**Question 6. 20 marks** Consider  $r(\theta)$ , where  $r$  is the average reward and  $\theta$  parameterizes the policy  $\pi$ . Write down  $r(\theta)$  in terms of the parameterized policy and the true action value function  $q_\pi(s, a)$ . Further write down the gradient of  $r(\theta)$  with respect to  $\theta$ . Assume the steady state distribution is given by  $\mu_\pi(s)$ .

Assume that we approximate the value function (as we would in an actor-critic method) using a linear approximation architecture with the feature vector for  $(s, a)$  given by  $\nabla_\theta \log(\pi(a|s, \theta))$  and weight vector  $w$ . We want to minimize the average squared error in estimating  $r(\theta)$  when using the linear approximation in place of the true action value function. Write down the expression for the error and the necessary first order condition that the error minimizing  $w$  must satisfy.

Derive the gradient of  $r(\theta)$  in terms of the optimal linear approximation and show that it is the same as the gradient that you defined earlier using the true action-value function.

**Question 1. 20 marks** Our environment has one state 0 and a terminal state Z. In state 0 one can take the two actions of -1 and 1. On taking either action the environment transitions to 0 with probability 0.2 and the agent gets a reward of 0. With probability 0.8, the agent gets a reward 1 and the episode terminates. You want to estimate the value of state 0 when using the policy  $\pi$  that picks action 1 with probability 1. However, you only have access to the policy  $\mu$  that picks either action with probability 0.5. You are allowed to run one episode. Derive the expected value and variance of your estimate of  $v_\pi(0)$ . Assume a discount factor of 1.



Your target policy is  $\pi$ , wherein  $\pi(1|0) = 1$ .

The behavior policy is  $\mu$ , wherein  $\mu(1|0) = \mu(-1|0) = 0.5$ .

Note that any episode has a return of 1.

However, when using the behavior policy, an episode will contain a random  $T$  steps.

The mean of our estimate is  $E_\mu [R_{0:T-1} | S_0 = 0]$

$$= E_\mu [R_{0:T-1} | S_0 = 0]$$

$\left\{ \begin{array}{l} S_0 = 0 \text{ is implicit} \end{array} \right.$

$$= \sum_{x=1}^{\infty} E_\mu [R_{0:T-1} | T=x] P[T=x]$$

$$= \sum_{x=1}^{\infty} E_\mu [R_{0:T-1} | T=x] P[T=x]$$

$$P[T=1] = 2(0.5)(0.8) = 0.8$$

$$P[T=2] = (1-0.8)(0.8) = (0.2)(0.8)$$

$$P[T=3] = (0.2)^2(0.8)$$

$\vdots$

$\vdots$

$$P[T=x] = (0.2)^{x-1}(0.8)$$

$$P_{0:x-1} = \frac{\pi(a_0|s_0)}{\mu(a_0|s_0)} \frac{\pi(a_1|s_1)}{\mu(a_1|s_1)} \dots \frac{\pi(a_{x-1}|s_{x-1})}{\mu(a_{x-1}|s_{x-1})}$$

$$= \begin{cases} 0 & \text{otherwise} \\ \frac{1}{(0.5)^x} & a_0 = a_1 = \dots = a_{x-1} = 1. \end{cases}$$

$$E_\mu [P_{0:x-1}] = \frac{1}{(0.5)^x} (0.5)^x = 1.$$

$\left\{ \begin{array}{l} \text{Here we use the fact that} \\ P[a_0 = a_1 = \dots = a_{x-1} = 1] \\ = (0.5)^x \text{ when using } \mu. \end{array} \right.$

$$\therefore E_\mu [R_{0:T-1}] = 1(0.8) + 1(0.2)(0.8) + \dots$$

$$= (0.8)(1 + 0.2 + (0.2)^2 + \dots)$$

$$= (0.8) \left( \frac{1}{1-0.2} \right) = 1.$$

The variance of the estimate is

$$\text{Var}[R_{0:T-1}] = E_\mu [(R_{0:T-1})^2] - (E_\mu [R_{0:T-1}])^2$$

$$E_\mu [(R_{0:T-1})^2] = (1)^2(0.8) + (1)^2(0.2)(0.8) + \dots$$

$$= 1.$$

$$\therefore \text{Var}[\ ] = 0.$$

This makes sense as the behavior policy can only result in a return of 1 at the end of any episode.

**Question 2. 10 marks** Our environment has two states  $s_0$  and  $s_1$  and two terminal states. Five episodes obtained by using a policy  $\pi$  are listed next.

$$s_0 \ 0.5 \ s_1 \ 0,$$

$$s_0 \ 0.2,$$

$$s_0 \ 0.3,$$

$$s_1 \ 0.6,$$

$$s_1 \ 0.4.$$

The episodes above are summarised as  $S_t, R_{t+1}, S_{t+1}, R_{t+2}, \dots$ . You are interested in estimating the values of the states. You can use either TD(0) or first visit MC. Derive the estimates both will converge to. Explain your steps.

Recall that TD(0) converges to certainty equivalent estimates.

Given the episodes above, a transition from  $s_1$  results in the episode being terminated. The returns starting in  $s_1$  are 0, 0.6, 0.4.

$$\therefore V_{\pi}^{\text{TD}(0)}(s_1) \rightarrow \frac{0 + 0.6 + 0.4}{3} = \frac{1}{3}$$

As regards  $s_0$ :

We have  $s_0$  transitioning to  $s_1$  with probability  $\frac{1}{3}$  (as determined from the data).

$$\text{We have } V_{\pi}^{\text{TD}(0)}(s_0) = \frac{1}{3} \left( 0.5 + V_{\pi}^{\text{TD}(0)}(s_1) \right) + \frac{1}{3} (0.2) + \frac{1}{3} (0.3)$$

$$= \frac{1}{3} \left( 0.5 + \frac{1}{3} \right) + \frac{1}{3} \left( \frac{1}{3} \right)$$

$$= \frac{1}{3} + \frac{1}{9} = \frac{4}{9}$$

For MC

$$V_{\pi}^{\text{MC}}(s_0) = \frac{0.5 + 0.2 + 0.3}{3} = \frac{1}{3}$$

$$V_{\pi}^{\text{MC}}(s_1) = \frac{1}{3}$$

**Question 3, 20 marks** Our environment has two states  $A$  and  $B$  and actions  $-1$  and  $1$ .  $Z$  is the terminal state. Consider the following episode, where in an episode is summarised as  $(S_t, A_t, R_{t+1}), (S_{t+1}, A_{t+1}, R_{t+2}), \dots$ . The episode is

$$(A, 1, 5), (B, -1, 5), (B, 1, 3), (A, -1, -5), (B, 1, 5), Z.$$

Use the episode to derive the action value estimates one would obtain using SARSA and Q-learning.

$$\alpha = 0.2 \quad [ \alpha \text{ must be picked from } (0, 1) ]$$

$$r = 1$$

Initial Q-values are all set to zero.

$$Q(A, -1) = 0$$

$$Q(A, 1) = 0$$

$$Q(B, -1) = 0$$

$$Q(B, 1) = 0$$

Q-learning:

$$Q(A, 1) = Q(A, 1) + \alpha (5 + \gamma Q(B, \arg \max_{a'} Q(B, a')))$$

$$= 0 + (0.2)(5 + 0)$$

$$= 1.0.$$

We have  $Q(A, -1) = 0$ ,  $Q(A, 1) = 1$ ,  $Q(B, -1) = 0$ ,  $Q(B, 1) = 0$ .

$$Q(B, -1) = Q(B, -1) + \alpha (5 + \gamma Q(B, \arg \max_{a'} Q(B, a')))$$

$$= 0 + 0.2(5 + 0) = 1.0.$$

We have  $Q(A, -1) = 0$ ,  $Q(A, 1) = 1$ ,  $Q(B, -1) = 1$ ,  $Q(B, 1) = 0$ .

$$Q(B, 1) = Q(B, 1) + \alpha (3 + \gamma \max_{a'} Q(A, a'))$$

$$= 0 + 0.2(3 + 1)$$

$$= 0.2(4) = 0.8.$$

We have  $Q(A, -1) = 0$ ,  $Q(A, 1) = 1$ ,  $Q(B, -1) = 1$ ,  $Q(B, 1) = 0.8$

$$Q(A, -1) = Q(A, -1) + \alpha (-5 + \gamma \max_{a'} Q(B, a'))$$

$$= 0 + 0.2(-5 + 1)$$

$$= 0.2(-4) = -0.8.$$

We have:

$$Q(A, -1) = -0.8, \quad Q(A, 1) = 1, \quad Q(B, -1) = 1, \quad Q(B, 1) = 0.8$$

$$Q(B, 1) = Q(B, 1) + \alpha (5 + \gamma (0))$$

$$= 0.8 + 0.2(5) = 1.8.$$

We have:

$$Q(A, -1) = -0.8$$

$$Q(A, 1) = 1$$

$$Q(B, -1) = 1$$

$$Q(B, 1) = 0.8$$

SARSA

For reference:  $(A, 1, 5), (B, -1, 5), (B, 1, 3), (A, -1, -5), (B, 1, 5), Z.$

$$\text{Again } Q(A, -1) = 0$$

$$Q(A, 1) = 0$$

$$Q(B, -1) = 0$$

$$Q(B, 1) = 0$$

$$Q(A, 1) = Q(A, 1) + \alpha (5 + \gamma Q(B, -1))$$

$$= 0 + 0.2(5) = 1.0.$$

We have:

$$Q(A, -1) = 0, \quad Q(A, 1) = 1.0, \quad Q(B, -1) = 0, \quad Q(B, 1) = 0$$

$$Q(B, -1) = Q(B, -1) + \alpha (5 + \gamma Q(B, 1))$$

$$= 0 + 0.2(5 + 0)$$

$$= 1.$$

We have:

$$Q(A, -1) = 0, \quad Q(A, 1) = 1.0, \quad Q(B, -1) = 1, \quad Q(B, 1) = 0$$

$$Q(B, 1) = Q(B, 1) + \alpha (3 + \gamma Q(A, -1))$$

$$= 0.2(3 + 0) = 0.6.$$

We have:

$$Q(A, -1) = 0, \quad Q(A, 1) = 1.0, \quad Q(B, -1) = 0, \quad Q(B, 1) = 0.6$$

$$Q(A, -1) = Q(A, -1) + \alpha (-5 + \gamma Q(B, 1))$$

$$= 0.2(-5 + 0.6)$$

$$= -1 + 0.12 = -0.88$$

We have:

$$Q(A, -1) = -0.88, \quad Q(A, 1) = 1.0, \quad Q(B, -1) = 0, \quad Q(B, 1) = 0.6$$

$$Q(B, 1) = Q(B, 1) + \alpha (5 + 0)$$

$$= 0.6 + 0.2(5) = 1.6.$$

We have:

$$Q(A, -1) = -0.88, \quad Q(A, 1) = 1.0, \quad Q(B, -1) = 0, \quad Q(B, 1) = 1.6$$

Question 4. 20 marks Consider the following experiences, wherein an experience is in the form  $S_t, A_t, R_{t+1}, S_{t+1}$ .

$A$  and  $B$  are states and  $Z$  is the terminal state.

- $A, -1, 3, A,$
- $A, -1, 5, B,$
- $A, 1, 2, A,$
- $B, -1, 4, Z,$
- $B, 1, 3, A.$

Assume initial action-values of  $Q(A, 1) = 2, Q(A, -1) = 3, Q(B, 1) = 2, Q(B, -1) = 3$ . Use these to associate with each experience a TD error. Calculate the probability with which an experience would be picked when using (a) greedy TD error experience replay, (b) uniform replay, and (c) rank based prioritization. Show all your steps.

Consider the first experience

$A, -1, 3, A$

TD-error  $\delta_1 = R_1 + V Q(A, \arg\max_a Q(A, a)) - Q(A, -1)$   
 $= 3 + Q(A, -1) - Q(A, -1) = 3.$

$A, -1, 5, B$

$\delta_2 = R_2 + V Q(B, \arg\max_a Q(B, a)) - Q(A, -1)$   
 $= 5 + Q(B, -1) - Q(A, -1)$   
 $= 5 + 3 - 3 = 5.$

$A, 1, 2, A$

$\delta_3 = 2 + Q(A, -1) - Q(A, 1)$   
 $= 2 + 3 - 2 = 3.$

$B, -1, 4, Z$

$\delta_4 = 4 + 0 - Q(B, -1)$   
 $= 4 - 3$   
 $= 1.$

$B, 1, 3, A$

$\delta_5 = 3 + Q(A, -1) - Q(B, 1)$   
 $= 3 + 3 - 2 = 4.$

To summarize:

Experience	TD-error
(1) $A, -1, 3, A$	$\delta_1 = 3$
(2) $A, -1, 5, B$	$\delta_2 = 5$
(3) $A, 1, 2, A$	$\delta_3 = 3$
(4) $B, -1, 4, Z$	$\delta_4 = 1$
(5) $B, 1, 3, A$	$\delta_5 = 4$

Rank	$1/\text{Rank}$	$P[\text{Rank-based picking of transition } i] \ (\alpha=0.2)$
3	$1/3$	$(1/3)^\alpha / (\frac{1}{3}^\alpha + 1^\alpha + \frac{1}{3}^\alpha + \frac{1}{4}^\alpha + \frac{1}{2}^\alpha) = 0.18$
1	1	0.23
3	$1/3$	0.18
4	$1/4$	0.18
2	$1/2$	0.20

Greedy:

Pick the experience with the largest error w.p. 1. So we will pick experience (5).

Uniform:  $P[\text{picking any experience}] = 1/5.$

Rank-based:

$P[\text{Sampling transition } i] = \frac{p_i^\alpha}{\sum_{k=1}^5 p_k^\alpha}$

$p_i = \frac{1}{\text{rank}(i)}$

The ranks and the corresponding probabilities are listed in the table above.

**Question 5. 10 marks** We have states  $1, 2, \dots, n$ . Assume a linear approximation architecture. State  $i$  is mapped to a  $n$ -dimensional feature vector with the  $i^{\text{th}}$  element set to 1 and the other elements set to 0. Suppose we observe an episode with the sequence  $1, 2, \dots, n$  of states. Derive the eligibility trace vector  $z_n$  for the episode. Assume TD(0.3) and a discount factor of 0.8.

Initialize  $z_0 = \bar{0}$ , where  $\bar{0}$  is the  $n$ -dimensional zero vector.

We have  $s_1 = 1, s_2 = 2, \dots, s_n = n$

$$z_1 = r_1 z_0 + \nabla \bar{v}(s_1, \bar{w}_1)$$

$$= \nabla \bar{v}(s_1, \bar{w}_1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{n \times 1}$$

(2)

$$z_2 = r_1 z_1 + \nabla \bar{v}(s_2, \bar{w}_2) = r_1 \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$z_3 = r_1 z_2 + \nabla \bar{v}(s_3, \bar{w}_3) = (r_1)^2 \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + (r_1) \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$\vdots$

$$z_n = \begin{bmatrix} (r_1)^{n-1} \\ (r_1)^{n-2} \\ \vdots \\ (r_1) \end{bmatrix} = \begin{bmatrix} (0.24)^{n-1} \\ (0.24)^{n-2} \\ \vdots \\ (0.24) \end{bmatrix}$$

(8)



**Question 6. 20 marks** Consider  $r(\theta)$ , where  $r$  is the average reward and  $\theta$  parameterizes the policy  $\pi$ . Write down  $r(\theta)$  in terms of the parameterized policy and the true action value function  $q_\pi(s, a)$ . Further write down the gradient of  $r(\theta)$  with respect to  $\theta$ . Assume the steady state distribution is given by  $\mu_\pi(s)$ .

Assume that we approximate the value function (as we would in an actor-critic method) using a linear approximation architecture with the feature vector for  $(s, a)$  given by  $\nabla_\theta \log(\pi(a|s, \theta))$  and weight vector  $w$ . We want to minimize the average squared error in estimating  $r(\theta)$  when using the linear approximation in place of the true action value function. Write down the expression for the error and the necessary first order condition that the error minimizing  $w$  must satisfy.

Derive the gradient of  $r(\theta)$  in terms of the optimal linear approximation and show that it is the same as the gradient that you defined earlier using the true action-value function.

$$r(\theta) = \sum_s \mu_\pi(s) \sum_a \pi(a|s, \theta) q_\pi(s, a)$$

The gradient  $\nabla r(\theta)$  is given by the policy gradient theorem for the average reward case. We have:

$$\nabla r(\theta) = \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a|s, \theta)$$

Let the linear approximation be

$$f_w(s, a) = \nabla_\theta \log(\pi(a|s, \theta))^T w$$

The approximation of  $r(\theta)$  is

$$\hat{r}(\theta) = \sum_s \mu(s) \sum_a \pi(a|s, \theta) f_w(s, a)$$

The average squared error in approximation is:

$$\sum_s \mu(s) \sum_a \pi(a|s, \theta) [(q_\pi(s, a) - f_w(s, a))^2]$$

The necessary first order conditions that the  $w$  that minimizes the above must satisfy are

$$\sum_s \mu(s) \sum_a \pi(a|s, \theta) (-2) (q_\pi(s, a) - f_w(s, a)) \nabla_w f_w(s, a) = 0$$

We have

$$\begin{aligned} \sum_s \mu(s) \sum_a \pi(a|s, \theta) q_\pi(s, a) \nabla_w f_w(s, a) \\ = \sum_s \mu(s) \sum_a \pi(a|s, \theta) f_w(s, a) \nabla_w f_w(s, a) \end{aligned}$$

Note that given the linear approximation architecture that  $f_w$  is, we have

$$\begin{aligned} \nabla_w f_w(s, a) &= \nabla_\theta \log \pi(a|s, \theta) \\ &= \frac{\nabla_\theta \pi(a|s, \theta)}{\pi(a|s, \theta)} \end{aligned}$$

Substituting  $\nabla_w f_w(s, a)$  in the equation above corresponding to first order conditions, we get

$$\begin{aligned} \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla_\theta \pi(a|s, \theta) \\ = \sum_s \mu(s) \sum_a f_w(s, a) \nabla_\theta \pi(a|s, \theta) \end{aligned}$$

Which is to say that the  $\nabla r(\theta)$  calculated using  $q_\pi(s, a)$  is the same as that obtained on replacing  $q_\pi(s, a)$  by the approximation architecture  $f_w(s, a)$ . //