

1. Suppose a game where you choose to flip one of two (possibly unfair) coins. You win \$1 if your chosen coin shows heads and lose \$1 if it shows tails.

- (a) Model this as a K-armed bandit problem: define the action set.

*Ans. Action set : two actions: a1: flip coin 1  
a2: flip coin 2.*

# Worksheet Bandits

CMPUT 397  
January 18, 2021

(b) Is the reward a deterministic or stochastic function of your action?  $\Rightarrow$  The coin is unfair.

Ans: Stochastic because there is probabilities for winning or losing, and you don't know the probabilities.

Ans:  $P_i$  = probability of observing Heads with coin  $i$ .  
 $q_i$  = flip coin  $i$ .



$$\begin{aligned} \text{The true value } q_A(a_A) &= \mathbb{E}[R | A = a_A] \\ &= (1) \cdot P_i + (-1) \cdot (1-P_i) \\ &= 2P_i - 1 \end{aligned}$$

$$\begin{aligned} \text{if we have fair coins, then } q_A(a_A) &= (1) \cdot 0.5 + (-1) \cdot 0.5 \\ &= 0 \\ &= \neq 0. \end{aligned}$$

Mark Holmstrom to Everyone 1:39 PM

MH it's written as being a function of your action, not a function of your belief

I think if  $p = 1$  it would be deterministic

Deepak Ranganatha Sastry Mamillapalli 1:39 PM to Everyone

DR yes

but if you don't know if  $p=1$ , then it's stochastic

Harsh Shah to Everyone 1:41 PM

I think if the action is already chosen and we know one of the trials has coin 2 chosen then it should be deterministic 确定的. regardless of the probability I think

Daniel Asimiakwini to Everyone 1:44 PM

A question: If we do not know  $p_i$  is it inherently stochastic?  $\Rightarrow$  Yes!!

Mohamed Ahmed to Everyone 1:45 PM

MA I missed what was said about deterministic. when would it be deterministic again?  $\Rightarrow$  when we choose coin 2.

Nicholas Rebstock to Everyone 1:47 PM

N We basically never know true probabilities. If you define stochasticity that way, everything is stochastic maybe  $\Rightarrow$  yes.

Cameron Jen to Everyone 1:51 PM

CJ The only situation in which we could say it is deterministic would be 1/0? Yes!!

Cameron Jen to Everyone 1:52 PM

CJ So in all other cases it would be stochastic, regardless of our knowledge of the dist or not  $\Rightarrow$  不管我们是否了解分布 (distribution), 它都是随机的.

# Worksheet Bandits

CMPUT 397  
January 18, 2021

- (c) You do not know the coin flip probabilities. Instead, you are able to view 6 sample flips for each coin respectively: (T,H,H,T,T,T) and (H,T,H,H,H,T). Use the sample average formula (equation 2.1 in the book) to compute the estimates of the value of each action.

Sample Average formula: (formula 2.1)  $\hat{Q}_t(a) = \frac{\text{Sum of rewards when } a \text{ taken prior to } t}{\# \text{ of times } a \text{ taken prior to } t} \approx \hat{q}_{t-1}$   
 $\hat{q}_t$ : Action:  $A_i = \text{flipping the } i^{\text{th}} \text{ coin}$

If observe Heads: win \$1  $P_i \rightarrow$  probability of flipping H by win i.

observe Tails: lose \$1  $(1-P_i) \rightarrow$  probability of flipping T by win i.

$$\hat{Q}_7(\text{flip win 1}) = \frac{(-1) \cdot 4 + (1) \cdot 2}{6} \rightarrow \text{all of the reward @ all of the action we flip win 1.}$$

do the first action  $\approx -0.33$

$$\hat{Q}_7(\text{flip win 2}) = \frac{(1) \cdot 4 + (-1) \cdot 2}{6} \\ \approx 0.33$$

- (d) Decide on which coin to flip next! Assume it's an exploit step.

Exploit step: 行動步驟有執行，並因物次序的 step 的重要-一切狀態。

$\because Q(\text{flip win 1}) = 0.45, \quad Q(\text{flip win 2}) = 0.33$

i. win 1 has the higher estimate value.

$\therefore$  we choose to flip win 2 next.

# Worksheet Bandits

CMPUT 397  
January 18, 2021

2. (Exercise 2.2 from S&B 2nd edition) Consider a  $k$ -armed bandit problem with  $k = 4$  actions, denoted 1, 2, 3, and 4. Consider applying to this problem a bandit algorithm using  $\epsilon$ -greedy action selection, sample-average action-value estimates, and initial estimates of  $Q_1(a) = 0$ , for all  $a$ . Suppose the initial sequence of actions and rewards is  $A_1 = 1, R_1 = -1, A_2 = 2, R_2 = 1, A_3 = 2, R_3 = -2, A_4 = 2, R_4 = 2, A_5 = 3, R_5 = 0$ . On some of these time steps the  $\epsilon$  case may have occurred, causing an action to be selected at random. On which time steps did this definitely occur? On which time steps could this possibly have occurred?

For  $\epsilon$ -greedy action selection, initialize all action values to 0 first:  $Q(a) = 0$   
After  $t$  times we represent action  $a$  to 0:  $N(a) = t$

A simple bandit algorithm					
Initialize, for $a = 1$ to $k$ :					
$Q(a) \leftarrow 0$					
$N(a) \leftarrow 0$					
Loop forever:					
$A \leftarrow \begin{cases} \text{argmax}_a Q(a) & \text{with probability } 1 - \epsilon \text{ (breaking ties randomly)} \\ \text{a random action} & \text{with probability } \epsilon \end{cases}$					
$R \leftarrow \text{bandit}(A)$					
$N(A) \leftarrow N(A) + 1$					
$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$					

$\therefore \text{Actions} = A = \{1, 2, 3, 4\}$

$Q_1(a) = \emptyset$  for all  $a \in A$ .

initialization: could be a greedy step or on a step  $t$  it can choose random actions. i.e.  $\epsilon$  case (3) greedy action

$A_1 = 1, R_1 = -1$ time step 1 we took action 1.	$A_2 = 2, R_2 = 1$	$A_3 = 2, R_3 = -2$ at action 3! we took action 2.	$A_4 = 2, R_4 = 2$	$A_5 = 3, R_5 = 0$
$Q(1) = \emptyset$ $N(1) = 0$	$Q(1) = -1$ $N(1) = 1$	$Q(1) = -1$ $N(1) = 1$	$Q(1) = -1$ $N(1) = 1$	$Q(1) = -1$ $N(1) = 1$
$Q(2) = \emptyset$ $N(2) = 0$	$Q(2) = \frac{1}{1} = 1$ $N(2) = 1$ (action 1) $Q(2) = -\frac{1}{1} = -1$ $N(2) = 1$ (action 2)	$Q(2) = -\frac{1}{2} = -\frac{1}{2}$ $N(2) = 1 + 1 = 2$	$Q(2) = \frac{1 + (-1) + 2}{3} = \frac{1}{3}$ $N(2) = 3$	$Q(2) = \frac{1}{3}$ $N(2) = 3$
$Q(3) = \emptyset$ $N(3) = 0$	$Q(3) = \emptyset$ $N(3) = 0$	$Q(3) = \emptyset$ $N(3) = 0$	$Q(3) = \emptyset$ $N(3) = 0$	$Q(3) = \frac{0}{1} = 0$ $N(3) = 1$
$Q(4) = \emptyset$ $N(4) = 0$	$Q(4) = \emptyset$ $N(4) = 0$	$Q(4) = \emptyset$ $N(4) = 0$	$Q(4) = \emptyset$ $N(4) = 0$	$Q(4) = \emptyset$ $N(4) = 0$

Greedy actions?

1.  $Q(1) = -1$  definitely 1. since we took action 1.

2.  $Q(2) = \frac{1}{3}$  definitely 2. since we took action 2.

3.  $Q(3) = 0$  definitely 3. since we took action 3.

4.  $Q(4) = \emptyset$  possibly 4. since we took action 2.

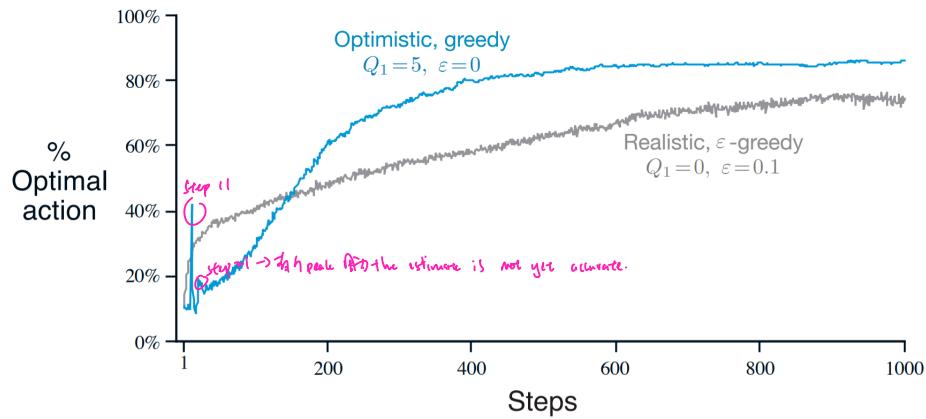
then we choose to take action 2. (the old estimate  $Q(2) = 1$ ),  $\therefore$  the new estimate  $Q(2) = 1 + \frac{1}{2} [-2 - 1] = -\frac{1}{2}$

也可能选择:  $\frac{1 + (-2)}{2} \rightarrow$  从 2 到 -1 的奖励平均值

i. Time step 4 and 5: definitely occur.

Time step 1, 2, 3, 4, 5: possibly occur. ⇒ 只有 4, 5 肯定发生, 1, 2, 3 也可能是可能的。

estimate of action value  $\Rightarrow$  for all action



**Figure 2.3:** The effect of optimistic initial action-value estimates on the 10-armed testbed. Both methods used a constant step-size parameter,  $\alpha = 0.1$ .

4. **Challenge Problem:** (Exercise 2.6 from S&B 2nd edition) The results shown in Figure 2.3 should be quite reliable because they are averages over 2000 individual, randomly chosen 10-armed bandit tasks. Why, then, are there oscillations and spikes in the early part of the curve for the optimistic method? In other words, what might make this method perform particularly better or worse, on average, on particular early steps?

Steps 1~10: Sample all 10 actions at least once. Disappointing reward decrease  $Q(a)$ .  
 optimal action in first 10 steps seen 10% of the time.

the peak:  
 Step 11:  
 first time that we select among all 10 actions.  
 more probable that we select optimal action.

Step 12: Step 11 reward disappointing us, so select second best action...  
 ...

Step 21: once again select among all 10 actions, but our values are less optimistic (as the first time)  
 less probable to select the optimal action.

Why, then, are there oscillations and spikes in the early part of the curve for the optimistic method? In other words, what might make this method perform particularly better or worse, on average, on particular early steps?

那么,为什么乐观的方法在曲线的早期部分会出现振荡和尖峰?换句说话,是什么原因使这种方法在特定的早期步骤上平均表现得特别好或特别差?

1. Suppose  $\gamma = 0.9$  and the reward sequence is  $R_1 = 2, R_2 = -2, R_3 = 0$  followed by an infinite sequence of 7s. What are  $G_1$  and  $G_0$ ?

Reward sequences:  $\frac{R_1}{2}, \frac{R_2}{-2}, \frac{R_3}{0}, \frac{R_4}{7}, \frac{R_5}{7}, \frac{R_6}{7}, \dots \dots$  亂序是 7.

$$\begin{aligned}
 G_t &= R_{t+1} + \gamma G_{t+1} \quad \forall t \geq 0 \\
 &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} \dots \\
 \therefore G_3 &= R_4 + \gamma R_5 + \gamma^2 R_6 + \dots \\
 &= 7 + \gamma 7 + \gamma^2 7 + \dots \\
 &= \frac{7}{1-\gamma} \\
 &= \frac{7}{1-0.9} \\
 &= 70. \\
 \therefore G_0 &= R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 \underbrace{R_4 + \dots}_{G_3} \quad \Rightarrow \text{To compute } G_0, \text{ we have to compute } G_1 \text{ first.} \\
 &= R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 (G_3) \quad G_1 \quad G_2 \\
 &= 2 + \gamma(-2) + 0 + \gamma^2 \times 70 \\
 &= 51.2
 \end{aligned}$$

$G_0 = R_1 + \gamma G_1$   
 $G_1 = R_2 + \gamma G_2$   
 $G_2 = R_3 + \gamma G_3$   
 $G_3 = R_4 + \gamma G_4$   
 $\vdots$

3. Prove that the discounted sum of rewards is always finite, if the rewards are bounded:  $|R_{t+1}| \leq R_{\max}$  for all  $t$  for some finite  $R_{\max} > 0$ .

$$\left| \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| < \infty \quad \text{for } \gamma \in [0, 1)$$

Hint: Recall that  $|a + b| < |a| + |b|$ .

$$\begin{aligned}
 & \because |R_{t+1}| \leq R_{\max} \\
 & \left| \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| = \sum_{i=0}^{\infty} \gamma^i |R_{t+1+i}| \\
 & = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 & \leq R_{\max} + \gamma R_{\max} + \gamma^2 R_{\max} + \dots \\
 & \stackrel{\text{geometric series}}{=} R_{\max} \frac{1 + \gamma + \gamma^2 + \gamma^3 + \dots}{1 - \gamma} \\
 & = \frac{R_{\max}}{1 - \gamma} \\
 & \stackrel{< \infty}{\therefore} \left| \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| < \infty
 \end{aligned}$$

# Worksheet 3:

## Markov Decision Processes

CMPUT 397  
January 21, 2022

6. Write  $r(s, a)$  in terms of  $r(s, a, s')$ .

$$\text{deg: } r(s, a) = \sum_{s'} \sum_r r P(s', r | s, a)$$

这裡的  $r$  是 probability of reward!!

$$= E[r | S_{t-1} = s, A_{t-1} = a]$$

$R_t$  depends on the previous  $a$ ,  $t$  与  $t-1$  重要, watch RT.

$$r(s, a, s') = \sum_r r P(r | s, a, s')$$

由  $r$  与  $s'$ . How do you get the extra condition  $s'$ ?

$$= E[r | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$\Rightarrow r(s, a) = \sum_{s'} \sum_r r P(s', r | s, a)$$

$$= \sum_{s'} \sum_r r P(s', a, s') \cdot P(s' | s, a)$$

already given in the sample space.

$$= \sum_{s'} P(s' | s, a) r(s', a, s')$$

到这一步就行了.

$E[r | s, a, S_t]$  R.V. which makes the whole function a R.V.  
constant R.V. expectation need some random variables! 因为期望是确定的 (outcome), 所以不能有随机变量 (random variable)。

$$r(s, a) = E[r | s, a, S_t]$$

$$\text{R.H.S.} = \sum_{s'} P(s' | s, a) r(s', a, s')$$

$$= \sum_{s'} r(s', a, s') P(r(s', a, s') | s, a) = r(s, a, S_t)$$

$$\text{答: } E[r | s, a] = \sum_{s'} P(s' | s, a) r(s', a, s')$$

R.V. because the input  $X$  is a R.V.

$$= \sum_{s'} f(s') P(X = s')$$

實際上明確指出  $X$  是什麼。

$$E[X] = \sum_{s'} P(X = s')$$

replace  $X$  with  $r(s, a, S_t)$ , 就可得出上面的答案.

Probability of B happen given A

$$\text{答: } P(A, B) = P(B | A) P(A)$$

probability of B happen

$$\therefore r(s, a) = \sum_{s'} \sum_r r P(s', r | s, a)$$

$$= \sum_{s'} \sum_r r P(s', a, s') P(s' | s, a)$$

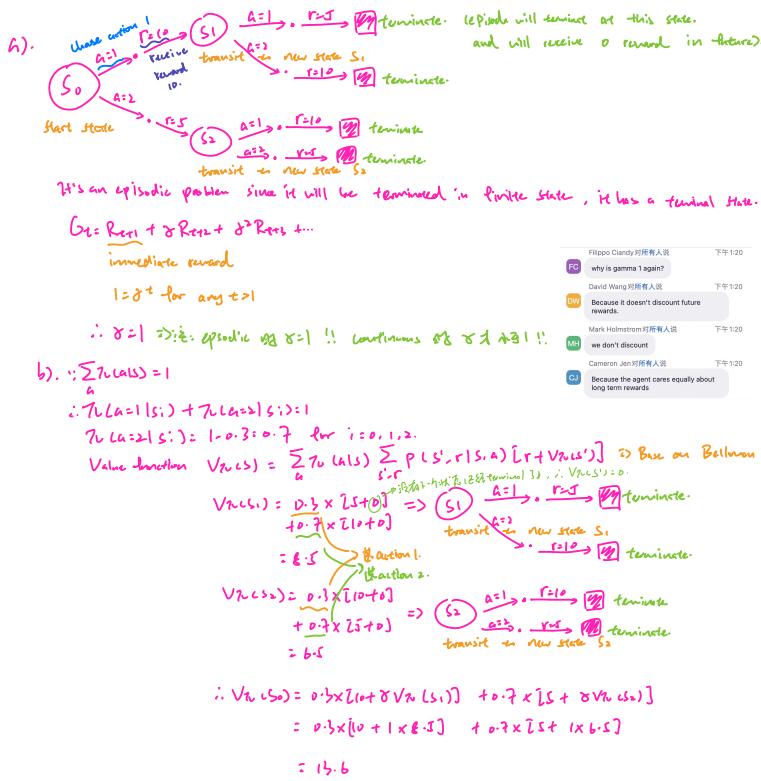


## Worksheet 4

CMPUT 397  
January 25, 2021

2. In this question, you will take a word specification of an MDP, and write the formal terms and determine the optimal policy. Suppose you have a problem with two actions. The agent always starts in the same state,  $s_0$ . From this state, if it takes action 1 it transitions to a new state  $s_1$  and receives reward 10; if it takes action 2 it transitions to a new state  $s_2$  and receives reward 5. From  $s_1$  if it takes action 1 it receives a reward of 5 and terminates; if it takes action 2 it receives a reward of 10 and terminates. From  $s_2$  if it takes action 1 it receives a reward of 10 and terminates; if it takes action 2 it receives a reward of 5 and terminates. Assume the agent cares equally about long term reward as about immediate reward.

- (a) Draw the MDP for this problem. Is it an episodic or continuing problem? What is  $\gamma$ ?
- (b) Assume the policy is  $\pi(a = 1|s_i) = 0.3$  for all  $s_i \in \{s_0, s_1, s_2\}$ . What is  $\pi(a = 2|s_i)$ ? And what is the value function for this policy? In other words, find  $v_\pi(s)$  for all three states.
- (c) What is the optimal policy in this environment?

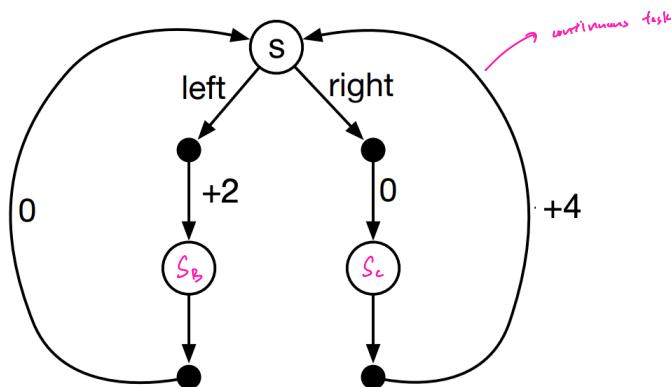


## Worksheet 4

CMPUT 397  
January 25, 2021

4. (Exercise 3.22 in 2<sup>nd</sup> ed.) Consider the continuing MDP shown on the bottom. The only decision to be made is that in the top state, where two actions are available, left and right. The numbers show the rewards that are received deterministically after each action.

- (a) List and describe all the possible policies in this MDP.
- (b) Is the following policy valid for this MDP (i.e. does it fit our definition of a policy): Choose *left* for five steps, then *right* for five steps, then *left* for five steps, and so on? Explain your answer.
- (c) What policy is optimal if  $\gamma = 0$ ? If  $\gamma = 0.9$ ? If  $\gamma = 0.5$ ?
- (d) For each possible policy, what is the value of state  $s$ ? Write down the numeric value to two decimal places. *Hint:* write down the return under each policy starting in state  $s$  (don't forget  $\gamma$ ). Simplify the infinite sum, using the fact that many rewards are zero. Then plug in the rewards and  $\gamma$  and compute the number.



a).  $\pi_a(a=right|s) \in [0, 1]$  a real number between 0 and 1.  
 $\pi_a(a=left|s) = 1 - \pi_a(a=right|s)$

b).  $\pi_a(a|s_B) = 1$   $\pi_a(a|s_C) = 0$

b). Not. The policy is not Markov. 它需要它的上一步来做出决定。2步后才决定。 (需要看前面状态来预测未来, 决策是一步)。

Markov property: The current state includes all information to make a decision.

c). For  $\gamma = 0$ : 只考虑下一步利益, 不考虑长远。

For  $\gamma = 0.9$ : 考虑长远利益。

For  $\gamma = 0.5$ : 考虑长远利益。

The state  $s$  has

action = left	$\text{left: } q_{\pi}(s, \text{left}) = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ $= 2 + 0.9 \times 0 + 0.9^2 \times 2 + \dots = 2 + 0.9^2 \times 2 + 0.9^4 \times 2 + \dots = 2 \cdot \frac{1}{1-0.9^2} = 18.18$
action = right	$\text{right: } q_{\pi}(s, \text{right}) = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ $= 0 + 0.9 \times 4 + 0 + 0.9^2 \times 4 + \dots = 0 + 0.9^2 \times 4 + 0.9^4 \times 4 + \dots = 4 \cdot \frac{0.9^2}{1-0.9^2} = 18.18$

David Wang 对所有人说  
DW  
This policy is not Markov. Its actions depends on its trajectory.

# Worksheet 5: Dynamic Programming

CMPUT 397  
February 2, 2022

1. In iterative policy evaluation, we seek to find the value function  $v_\pi$  for an arbitrary policy  $\pi$ . Using the **Bellman equation** as an update rule, we can generate a sequence of approximate value functions  $\{v_k\}$  that will eventually converge to the value function  $v_\pi$ . The expected update is written as:

$$v_{k+1}(s) \doteq \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')].$$

Write the expected update to generate a sequence of approximate action-value functions  $\{q_k\}$ ?

$$q_{k+1}(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \sum_{a'} \pi(a'|s') q_k(s',a')]$$



∴ approximate action-value functions  $\{q_k\}$  is:  $q_{k+1}(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \sum_{a'} \pi(a'|s') q_k(s',a')]$

# Worksheet 5: Dynamic Programming

CMPUT 397  
February 4, 2022

2. A deterministic policy  $\pi(s)$  outputs an action  $a \in \mathcal{A}(s) = \{a_0, a_1, \dots, a_k\}$ . More generally, a policy  $\pi(\cdot|s)$  outputs the probabilities for all actions  $a \in \mathcal{A}(s)$ :  $\pi(\cdot|s) = [\pi(a_0|s), \pi(a_1|s), \dots, \pi(a_k|s)]$ . How can you write a deterministic policy in this form? (Let  $\pi(s) = a_i$  and define  $\pi(\cdot|s)$ .)

$$\pi(s) = a_i$$

meaning: With probability 1, we take action  $a_i$  in state  $s$ .  $\Rightarrow \pi(a_i|s) = 1$ .

What about the other actions?

$$a \in \mathcal{A}(s) \setminus \{a_i\} \quad \pi(a|s) = 0 \quad \pi(a_i|s) = \begin{cases} 1 & \text{if } i = i \\ 0 & \text{if } i \neq i \end{cases}$$

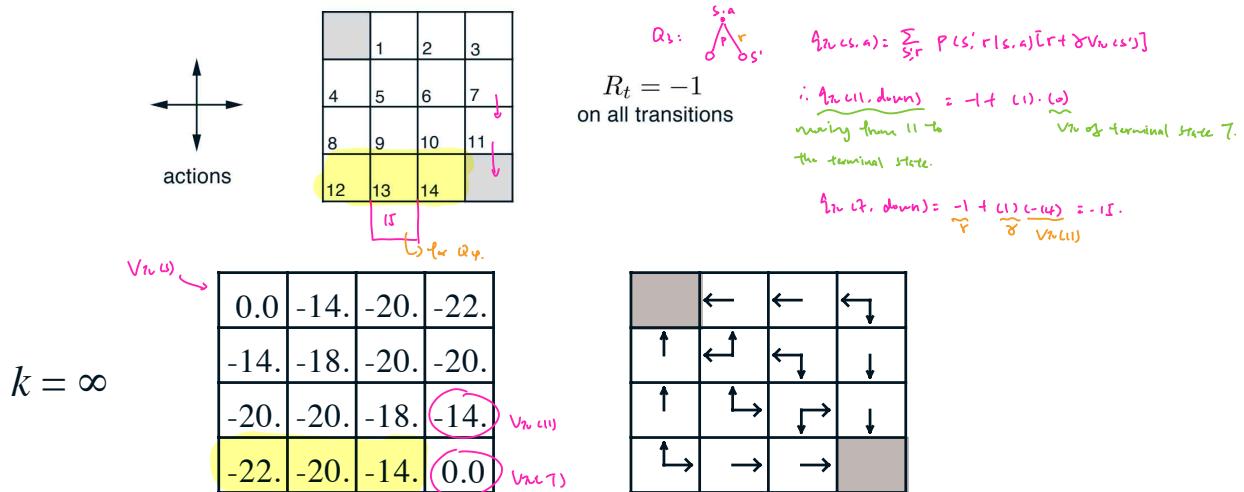
$\{a \mid a \in \mathcal{A}(s) \text{ and } a \neq a_i\}$

$$\therefore \pi(\cdot|s) = [\underbrace{\pi(a_0|s)}, \pi(a_1|s), \dots, \pi(a_k|s)]$$

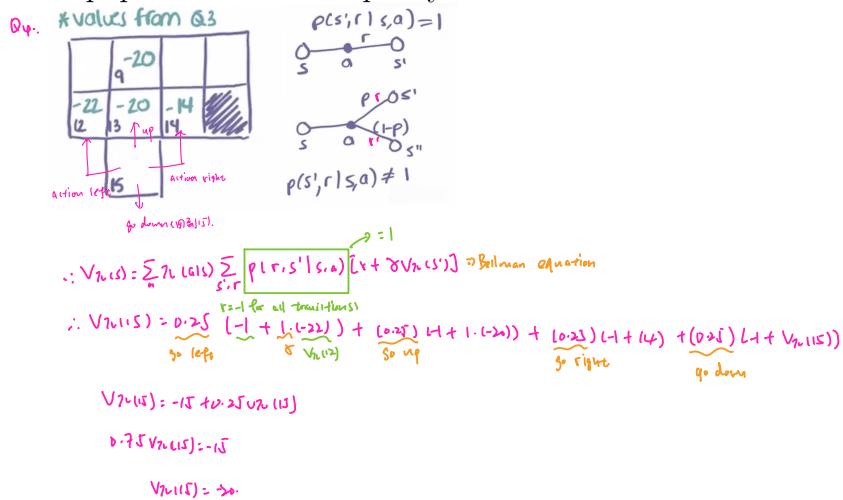
0<sup>th</sup> element (Subscript is 0)

$$\pi(\cdot|s)[i] = \begin{cases} 1 & \text{if } i = i \\ 0 & \text{if } i \neq i \end{cases}$$

3. (Exercise 4.1 S&B) Consider the 4x4 gridworld below, where actions that would take the agent off the grid leave the state unchanged. The task is episodic with  $\gamma = 1$  and the terminal states are the shaded blocks. Using the precomputed values for the equiprobable policy below, what is  $q_\pi(11, \text{down})$ ? What is  $q_\pi(7, \text{down})$ ?  $T = \text{terminal state}$ .

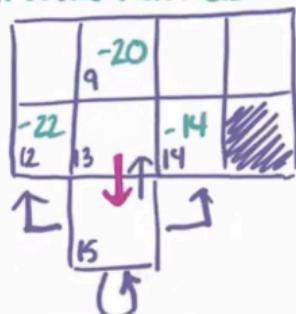


4. (Exercise 4.1 from S&B) Suppose in the above gridworld where a new state 15 is added to the gridworld just below state 13, and its actions, left, up, right, and down, take the agent to the states 12, 13, 14, and 15, respectively. Assume that the transitions from the original states are unchanged. What, then is,  $v_\pi(15)$  for the equiprobable random policy? Now suppose the dynamics of state 13 are also changed, such that action down from state 13 takes the agent to the new state 15. What is  $v_\pi(15)$  for the equiprobable random policy in this case?



4. (Exercise 4.2 S&B) In Question 3, suppose a new state 15 is added to the gridworld just below state 13, and its actions, left, up, right, and down, take the agent to states 12, 13, 14, and 15, respectively. Assume that the transitions from the original states are unchanged. What, then, is  $v_\pi(15)$  for the equiprobable random policy? Now suppose the dynamics of state 13 are also changed, such that action down from state 13 takes the agent to the new state 15. What is  $v_\pi(15)$  for the equiprobable random policy in this case?

\*values from Q3



$$V_{\pi}(13) = 0.25 \left[ -1 + 1 \cdot \underbrace{(-22)}_{V_{\pi}(12)} + (-1 + V_{\pi}(14)) + (-1 + \underbrace{(-14)}_{V_{\pi}(15)}) + (1 + V_{\pi}(15)) \right]$$

$$V_{\pi}(13) = -15 + 0.25 V_{\pi}(15)$$

$$\begin{aligned} V_{\pi}(15) &= (0.25) \left[ (-1 + \underbrace{(-22)}_{V_{\pi}(12)}) + (-1 + V_{\pi}(13)) + (-1 + \underbrace{(-14)}_{V_{\pi}(14)}) + (-1 + \underbrace{(-14)}_{V_{\pi}(15)}) \right] \\ &= -10 + 0.25 V_{\pi}(13) + 0.25 V_{\pi}(14) \\ &= -(10 + 0.25 \left[ (-1 + \underbrace{0.25 V_{\pi}(15)}_{V_{\pi}(13)}) \right] + 0.25 V_{\pi}(14)) \\ &= -13.75 + 0.3125 V_{\pi}(15) \end{aligned}$$

$$V_{\pi}(15) = -20$$

$$\therefore V_{\pi}(13) = -15 + 0.25 (-20) = -20$$

6. **(Challenge Question)** (*Exercise 4.4 S&B*) The policy iteration algorithm on page 80 has a subtle bug in that it may continually switches between two or more policies that are equally good. This is ok for pedagogy, but not for actual use. Modify the pseudocode so that convergence is guaranteed. Note that there is more than one approach to solve this problem.

**Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi_*$**

1. Initialization  
 $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$
2. Policy Evaluation  
 Loop:  $\leftarrow$  value-table & true  
 $\Delta \leftarrow 0$   
 Loop for each  $s \in \mathcal{S}$ :  
 $v \leftarrow V(s)$   
 $V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')] \rightarrow$  update the value by transition probability  
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$  if  $|v - V(s)| > \phi$ , then value-stable & false  
 until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)  
if value-stable = true, stop, return  $\pi$  and  $V$ ; otherwise go to policy improvement.
3. Policy Improvement  
 $policy\text{-stable} \leftarrow true$   
 For each  $s \in \mathcal{S}$ :  
 $old\text{-action} \leftarrow \pi(s)$   
 $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')] \rightarrow$  replace the policy  
 If  $old\text{-action} \neq \pi(s)$ , then  $policy\text{-stable} \leftarrow false$   
 If  $policy\text{-stable}$ , then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

Method: check for a change in the state values.

# Worksheet 6: Monte Carlo

CMPUT 397  
February 8, 2022

5. Let  $\rho_t = \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$ .   
  $\pi(A_t|S_t)$   $\rightarrow$  target policy.  
  $b(A_t|S_t)$   $\rightarrow$  behavior policy.

(a) Verify that  $\mathbb{E}_b[\rho_t|S_t = s] = 1$ .

(b) Verify that  $\mathbb{E}_b[\rho_t R_{t+1}|S_t = s] = \mathbb{E}_\pi[R_{t+1}|S_t = s]$ .

(c) What is the variance of the importance corrected one-step reward,  $\mathbb{V}(\rho_t R_{t+1}|S_t = s)$ ?  
 When would this variance be large?

$$\rho_t = \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$$

a.  $\mathbb{E}_b[\rho_t|S_t = s]$   
 base on behavior policy  
  $= \mathbb{E}_b \left[ \frac{\pi(A_t|S_t)}{b(A_t|S_t)} \mid S_t = s \right]$

All actions from the behavior policy  $b$ .  
  $= \sum_{a \in A} \frac{\pi(a|s)}{b(a|s)} \frac{b(a|s)}{b(a|s)}$  by LOTUS.  
 Action down from the whole action set:  $\mathbb{E}[g(X)] = \sum_{x \in S} g(x)p(x)$ , the law of the unconscious statistician (LOTUS).

$$= \sum_a \pi(a|s)$$

$\therefore 1 \rightarrow$  by definition.

b.  $\mathbb{E}_b[\mathbb{E}_\pi[R_{t+1}|S_{t+1} = s]] = \mathbb{E}_b[\mathbb{E}_\pi[R_{t+1}|S_t = s]]$   
 behavior policy  $\rightarrow$  true target policy.

WHS:  $\mathbb{E}_b[\mathbb{E}_\pi[R_{t+1}|S_t = s]]$   
  $= \mathbb{E}_b \left[ \mathbb{E}_\pi[\mathbb{E}_\pi[R_{t+1}|S_t = s, A_t] \mid S_t = s] \right]$   
 LOTUS

$$= \sum_a b(a|s) \mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] \text{ by LOTUS.}$$

probability of choose action  $a$ . Up to  $\mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a]$ .

$$= \sum_a b(a|s) \sum_{s',r} P(s',r|s,a) \frac{\pi(a|s)}{b(a|s)} \cdot r$$

$$= \sum_a b(a|s) \frac{\pi(a|s)}{b(a|s)} \sum_{s',r} P(s',r|s,a) \cdot r$$

$$= \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) \cdot r$$

$$= \sum_a \pi(a|s) \mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[\mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] \mid S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1}|S_t = s]$$

Stuart Cummings 对所有人说  
is At the same as just At  
下午 1:25  
David Wang 对所有人说  
No, One is an event At  
下午 1:25

DW  
The other is a random variable At

Pr: constant because  $S_t$  and  $A_t$  are given.  
随机数用  $P_{s,a}$   $\frac{\pi(a|s)}{b(a|s)}$

$\therefore \frac{\pi(a|s)}{b(a|s)}$  constant, since we know  $S_t, a$ .  
 $\therefore$  At is not random.

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \Rightarrow \text{when the same variable as an array was taken as an input}$$

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

$$V_{t+1}(S) \doteq V_t(S) \quad \forall S \in S - \{S_t\}, \alpha \in \mathbb{R}$$

## Worksheet 7: TD Learning Methods for Prediction

CMPUT 397  
May 2, 2021

- (Exercise 6.1 S&B) If  $V$  changes during the episode, then

$$G_t - V(S_t) = \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$

only holds approximately; what would the difference be between the two sides? Let  $V_t$  denote the array of state values used at time  $t$  in the TD error and in the TD update. Redo the derivation to determine the additional amount that must be added to the sum of TD errors in order to equal the Monte Carlo error.

$$\therefore \delta_t = R_{t+1} + \gamma V(S_{t+1}) - V_t(S_t)$$

$$\therefore \underbrace{G_t - V_t(S_t)}_{\text{the MC error.}} = R_{t+1} + \gamma G_{t+1} - V_t(S_t) + \gamma V_t(S_{t+1}) - \gamma V_t(S_{t+1})$$

$$= \delta_t + \gamma (G_{t+1} - V_t(S_{t+1}))$$

$$X_t \doteq G_t - V_t(S_t), \quad X_{t+1} \doteq G_{t+1} - V_{t+1}(S_{t+1})$$

$$\rightarrow = \delta_t + \gamma (G_{t+1} - V_t(S_{t+1}) + V_{t+1}(S_{t+1}) - V_{t+1}(S_{t+1}))$$

$$\star X_t = \delta_t + \gamma (V_{t+1}(S_{t+1}) - V_t(S_{t+1})) + \underline{\gamma X_{t+1}}$$

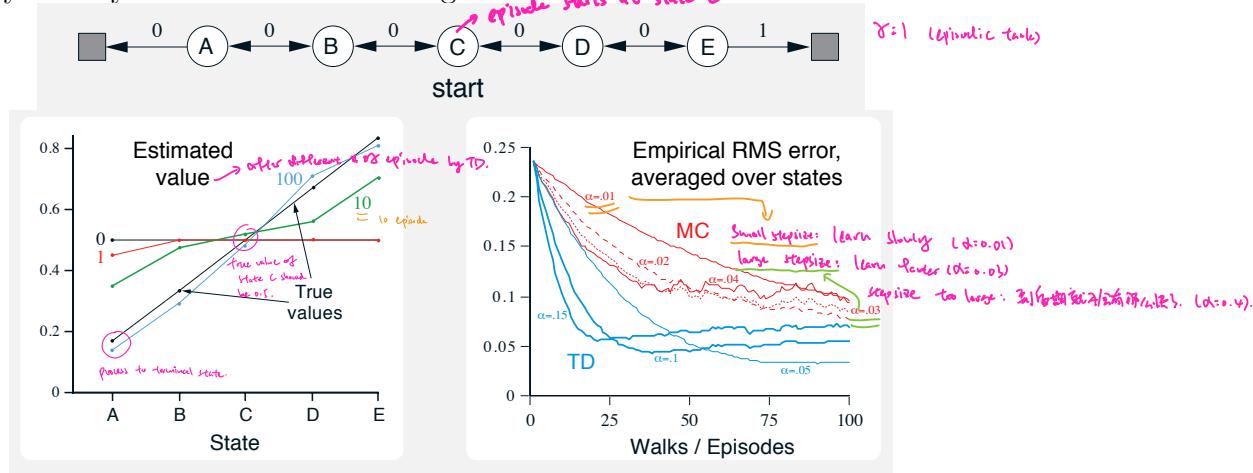
$$= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k + \sum_{k=t}^{T-1} \gamma^{k-t+1} \underbrace{(V_{k+1}(S_{t+1}) - V_k(S_k))}_{\text{if step size is small, } V_{k+1} \approx V_k, \text{ then value is 0.}} \quad \therefore \text{if step size is small, } V_{k+1} \approx V_k, \text{ then value is 0.}$$

$$\rightarrow X_t = \delta_t + \gamma (V_{t+1}(S_{t+1}) - V_t(S_{t+1})) + \underline{\gamma (\delta_{t+1} + \gamma (V_{t+2}(S_{t+2}) - V_{t+1}(S_{t+2}) + \gamma X_{t+2}))}$$

# Worksheet 7: TD Learning Methods for Prediction

CMPUT 397  
May 2, 2021

2. (Exercise 6.3 S&B) From the results shown in the left graph of the random walk example it appears that the first episode results in a change in only  $V(A)$ . What does this tell you about what happened on the first episode? Why was only the estimate for this one state changed? By exactly how much was it changed? *... starts at state C*



MC reward process: transition Schotastic.

by randomness of picking; once the action is chosen, the transition is deterministic.

∴ If ends at state  $A$ , only the state value of state  $A$  will change.

$$\begin{aligned}
 V(LB) &\leftarrow V(LB) + \alpha [d + V(L) - V(LB)] \\
 &\leftarrow 0.5 + 0.1 [d + 0 - 0.5] \\
 &\leftarrow 0.45
 \end{aligned}$$

# Worksheet 7: TD Learning Methods for Prediction

CMPUT 397  
May 2, 2021

5. (Exercise 6.7 S&B) Design an off-policy version of the TD(0) update that can be used with arbitrary target policy  $\pi$  and covering behavior policy  $b$ , using at each step  $t$  the importance sampling ratio  $\rho_{t:t}$  (5.3).

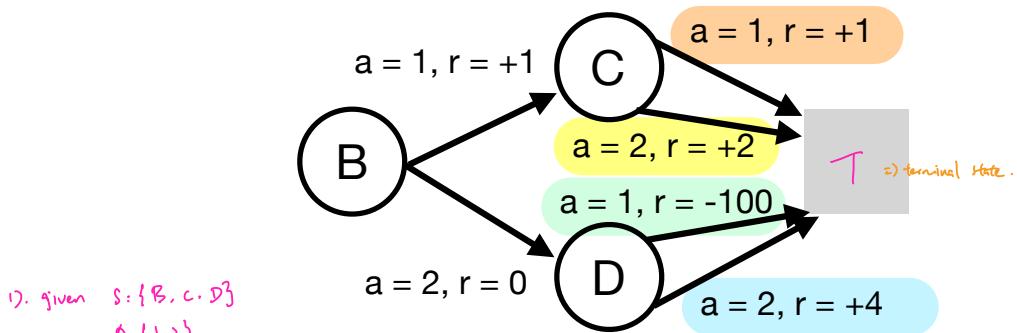
$$\begin{aligned}
 V_{\pi} \approx V_{\pi_b} & \quad \text{A: variable (A是待估常数)} \\
 & \quad \text{Action is drawn from behavior policy } b. \\
 V_{\pi}(s) &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_{\pi}(s')] \quad \text{Bellman equation.} \\
 &= \mathbb{E}_{\pi(a|s)} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \\
 &= \sum_b b(\text{label}) \frac{\pi(a|s)}{b(\text{label})} \sum_{s',r} p(s',r|s,a) [r + \gamma V_{\pi}(s')] \\
 &= \mathbb{E}_{\pi \text{ and } b} [P_{t+1} (R_{t+1} + \gamma V_{\pi}(s_{t+1})) | s_t = s] \\
 \therefore \text{update: } V(s_t) &\leftarrow V(s_t) + \alpha [P_{t+1} (R_{t+1} + \gamma V(s_{t+1})) - V(s_t)] \\
 &\quad \downarrow \\
 &\quad \text{on-policy off-policy update 不需要使用 important sample ratio.} \\
 &\quad \text{off-policy update 使用 target policy } \pi \text{ 和 behavior policy } b \text{ 的重要性采样 ratio!!}
 \end{aligned}$$

# Worksheet 8: TD Learning Methods for Control

CMPUT 397  
May 2, 2021

1. Consider the following MDP, with three states  $B, C$  and  $D$  ( $\mathcal{S} = \{B, C, D\}$ ), and 2 actions ( $\mathcal{A} = \{1, 2\}$ ), with  $\gamma = 1.0$ . Assume the action values are initialized  $Q(s, a) = 0 \forall s \in \mathcal{S}$  and  $a \in \mathcal{A}$ . The agent takes actions according to an  $\epsilon$ -greedy with  $\epsilon = 0.1$ .
  - What is the optimal policy for this MDP and what are the action-values corresponding to the optimal policy:  $q^*(s, a)$ ?
  - Imagine the agent experienced a single episode, and the following experience:  $S_0 = B, A_0 = 2, R_1 = 0, S_1 = D, A_1 = 2, R_2 = 4$ . What are the Sarsa updates during this episode, assuming  $\alpha = 0.1$ ? Start with state  $B$ , and perform the Sarsa update, then update the value of state  $D$ .
  - Using the sample episode above, compute the updates Q-learning would make, with  $\alpha = 0.1$ . Again start with state  $B$ , and then state  $D$ .
  - Let's consider one more episode:  $S_0 = B, A_0 = 2, R_1 = 0, S_1 = D, A_1 = 1, R_2 = -100$ . What would the Sarsa updates be? And what would the Q-learning updates be?
  - Assume you see one more episode, and it's the same on as in 1d. Once more update the action values, for Sarsa and Q-learning. What do you notice?
  - What policy does Q-learning converge to? What policy does Sarsa converge to?

## Deterministic transitions



a).  $\therefore q^*(s, a) = \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} p(s'|s, a) [r + \gamma \max_{a'} q^*(s', a')]$   $\Rightarrow$  By Bellman equation. Since C is a state  $\Rightarrow$  if the transition is deterministic. (即 C 一定会到 terminal state).

the optimal action value: (和  $q^*$  一样是 true value, not estimate)

$q^*(C, 1) = 1$	$q^*(B, 1) = [1 + \gamma q^*(C, 2)]$
$q^*(C, 2) = 2$	$= 1 + 2$ $\Rightarrow$ 有些选择 action 1 是因为 optimal option 是 to move C to action 2.
$q^*(D, 1) = -100$	$= 3$
$q^*(D, 2) = 4$	$q^*(B, 2) = [0 + \gamma q^*(D, 2)]$
	$= 4$

$\pi_{opt}(C|C) = 1, \pi_{opt}(C|D) = 0$  Optimal policy to given state C 的选 action 1 的概率为 0. 而到 state C 的选 action 2 的概率为 1.  $\pi_{opt}(D|C) = 0, \pi_{opt}(D|D) = 1$ .

$\pi_{opt}(D|B) = 1$

b).  $S_0 = B, A_0 = 2, R_1 = 0$   
 $S_1 = D, A_1 = 2, R_2 = 4$   
 $\alpha = 0.1$

$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma Q(s', a') - Q(s, a)]$   $\Rightarrow$  By Sarsa update

$Q(B, 2) = Q(B, 2) + \alpha [0 + \gamma Q(D, 2) - Q(B, 2)]$   $\Rightarrow$  !!!!!!! initialize  $Q(s, a) = 0$ , with the agent hasn't been acting!!

$= 0 + 0.1 [0 + 0 - 0]$   $\Rightarrow$  !!!!!!! agent never choose action 2!! (即 initialize 了  $Q(s, a) = 0$ )

$Q(D, 2) = Q(D, 2) + [4 + \gamma Q(T, 2) - Q(D, 2)]$

$= 0 + 0.1 [4 + 0 - 0]$

$\Rightarrow 0.4$

即 state D 的值是 0.4.

Oh, so for model free methods, we'll always be using R then? Yes!!

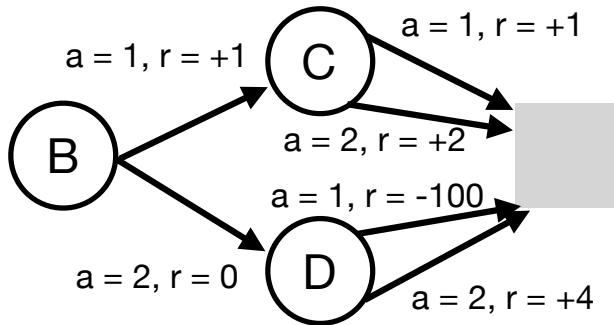
即  $R$  和  $R + \gamma \text{random variable}$  一样

# Worksheet 8: TD Learning Methods for Control

CMPUT 397  
May 2, 2021

1. Consider the following MDP, with three states  $B, C$  and  $D$  ( $\mathcal{S} = \{B, C, D\}$ ), and 2 actions ( $\mathcal{A} = \{1, 2\}$ ), with  $\gamma = 1.0$ . Assume the action values are initialized  $Q(s, a) = 0 \forall s \in \mathcal{S}$  and  $a \in \mathcal{A}$ . The agent takes actions according to an  $\epsilon$ -greedy with  $\epsilon = 0.1$ .
  - What is the optimal policy for this MDP and what are the action-values corresponding to the optimal policy:  $q^*(s, a)$ ?
  - Imagine the agent experienced a single episode, and the following experience:  $S_0 = B, A_0 = 2, R_1 = 0, S_1 = D, A_1 = 2, R_2 = 4$ . What are the Sarsa updates during this episode, assuming  $\alpha = 0.1$ ? Start with state  $B$ , and perform the Sarsa update, then update the value of state  $D$ .
  - Using the sample episode above, compute the updates Q-learning would make, with  $\alpha = 0.1$ . Again start with state  $B$ , and then state  $D$ .
  - Let's consider one more episode:  $S_0 = B, A_0 = 2, R_1 = 0, S_1 = D, A_1 = 1, R_2 = -100$ . What would the Sarsa updates be? And what would the Q-learning updates be?
  - Assume you see one more episode, and it's the same on as in 1d. Once more update the action values, for Sarsa and Q-learning. What do you notice?
  - What policy does Q-learning converge to? What policy does Sarsa converge to?

## Deterministic transitions



c).  $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a') - Q(S, A)] \Rightarrow B_0 Q\text{-learning update}$   
 $Q(B, 2) = Q(B, 2) + 0.1 [0 + \gamma \max_a Q(D, a') - Q(B, 2)]$   
 $= 0 + 0.1 [0 + 0 - 0] = 0.$

DR because the episode generated for sarsa is the greedy policy anyway

$$Q(D, 2) = Q(D, 2) + 0.1 [4 + \gamma \max_a Q(T, a') - Q(D, 2)]$$
 $= 0 + 0.1 [4 + 0 - 0] = 0.4$

d):  $S_0 = B, A_0 = 2, R_1 = 0$   
 $S_1 = D, A_1 = 1, R_2 = -100$

$$Q(B, 2) = Q(B, 2) + 0.1 [0 + \gamma \max_a Q(D, a') - Q(B, 2)]$$

$$= 0 + 0.1 [0 + 0 - 0] = 0$$

$$Q(D, 1) = Q(D, 1) + 0.1 [-100 + \gamma Q(T, \cdot) - Q(D, 1)]$$

$$= 0 + 0.1 [-100 + 0 - 0] = -10.$$

By SARSAs update.

Q-learning

$$Q(B, 2) = Q(B, 2) + 0.1 [0 + \gamma \max_a Q(D, a') - Q(B, 2)]$$

$$= 0 + 0.1 [0 + 0.4 - 0] = 0.04$$

$$Q(D, 1) = -10$$

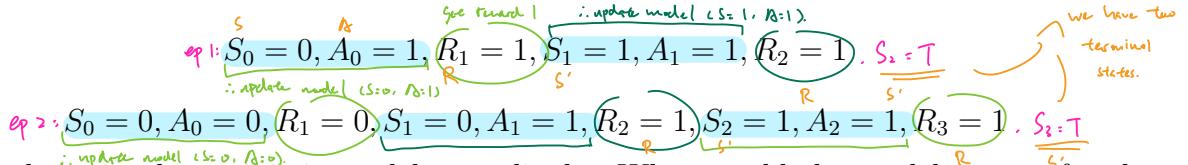
Deepak Ranganatha Sastry Mamillapalli 11:51  
DR q learning converges to optimal policy and sarsa converges to epsilon greedy policy  
 so for  $q(t)$  what is the answer?

# Worksheet 9: Planning, Learning & Acting

CMPUT 397  
March 19, 2021

Dyna-Q

1. An agent observes the following two episodes from an MDP,



and updates its deterministic model accordingly. What would the model output for the following queries:

- (a)  $\text{Model}(S = 0, A = 0)$ :  $\text{None} \rightarrow$   
 (b)  $\text{Model}(S = 0, A = 1)$ :  $\text{None} \rightarrow$   
 (c)  $\text{Model}(S = 1, A = 0)$ :  $\text{None} \rightarrow$   
 (d)  $\text{Model}(S = 1, A = 1)$ :  $\text{None} \rightarrow$
- initializes them all as  $\text{None}$  (RTF, to avoid the empty decimal error!).  
 ep 1:  $t=0 \rightarrow t=1 \rightarrow t=2 \rightarrow \text{None}$   
 ep 2:  $t=0 \rightarrow t=1 \rightarrow t=2 \rightarrow t=3 \rightarrow \text{None}$
- $S \rightarrow R, S' \rightarrow \text{None}$  (step 1)  
 store  $t, 1$  instead of  $(S=0, A=1)$ .  
 None for model  $(S=0, A=1)$ .
- $t=1 \rightarrow t=2 \rightarrow \text{None}$   
 $t=2 \rightarrow t=3 \rightarrow \text{None}$
- $t=0 \rightarrow t=1 \rightarrow \text{None}$   
 $t=1 \rightarrow t=2 \rightarrow \text{None}$   
 $t=2 \rightarrow t=3 \rightarrow \text{None}$
- $t=0 \rightarrow t=1 \rightarrow \text{None}$   
 $t=1 \rightarrow t=2 \rightarrow \text{None}$   
 $t=2 \rightarrow t=3 \rightarrow \text{None}$
- $t=0 \rightarrow t=1 \rightarrow \text{None}$   
 $t=1 \rightarrow t=2 \rightarrow \text{None}$   
 $t=2 \rightarrow t=3 \rightarrow \text{None}$
1. T  $\rightarrow$  output of the model.
- Model( $S, A$ )  $\leftarrow R, S'$   
 corrector from episode 1.

## Tabular Dyna-Q

Initialize  $Q(s, a)$  and  $\text{Model}(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$

Loop forever:

- $S \leftarrow$  current (nonterminal) state
- $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$
- $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- $\text{Model}(S, A) \leftarrow R, S'$  (assuming deterministic environment)
- Loop repeat  $n$  times:
  - $S \leftarrow$  random previously observed state
  - $A \leftarrow$  random action previously taken in  $S$
  - $R, S' \leftarrow \text{Model}(S, A)$
  - $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

# Worksheet 9: Planning, Learning & Acting

CMPUT 397  
March 19, 2021

2. An agent is in a 4-state MDP,  $\mathcal{S} = \{1, 2, 3, 4\}$ , where each state has two actions  $\mathcal{A} = \{1, 2\}$ . Assume the agent saw the following trajectory,

$S_0 = 1, A_0 = 2, R_1 = -1,$   
 Time step 1:  $S_1 = 1, A_1 = 1, R_2 = 1,$   
 $S_2 = 2, A_2 = 2, R_3 = -1,$   
 $S_3 = 2, A_3 = 1, R_4 = 1,$   
 $S_4 = 3, A_4 = 1, R_5 = 100,$   
 $S_5 = 4$   
 Final state.

and uses Tabular Dyna-Q with 5 planning steps for each interaction with the environment.

- (a) Once the agent sees  $S_5$ , how many Q-learning updates has it done with **real experience**? How many updates has it done with **simulated experience**?
- (b) Which of the following are possible (or not possible) simulated transitions  $\{S, A, R, S'\}$  given the above observed trajectory with a deterministic model and random search control?
- $\{S = 1, A = 1, R = 1, S' = 2\} \checkmark \Rightarrow S_1=1, A_1=1, R_2=1, S_2=2.$
  - $\{S = 2, A = 1, R = 1, S' = 3\}$
  - $\{S = 2, A = 2, R = -1, S' = 2\} \checkmark \Rightarrow S_2=2, A_2=2, R_3=-1, S_3=2.$
  - $\{S = 1, A = 2, R = -1, S' = 1\} \checkmark \Rightarrow S_0=1, A_0=2, R_1=-1, S_1=1.$
  - $\{S = 3, A = 1, R = 100, S' = 5\} \times$

a). Tabular Dyna-Q

Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$

Loop forever:

- $S \leftarrow$  current (nonterminal) state
- $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$  ← **real experience** (updated by real experience).
- $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment)
- Loop repeat  $n$  times:
  - $S \leftarrow$  random previously observed state
  - $A \leftarrow$  random action previously taken in  $S$
  - $R, S' \leftarrow Model(S, A)$  ← **simulate step**
  - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$  ← **simulated update** (update 2 by simulated experience).

For each update we do based on the real experience, we do  $n \times$  many updates using the simulated experience.

For each update we do based on the real experience, we do  $n \times$  many updates using the simulated experience.

The question shows that there are 5 time steps of state-action pairs and their subsequent next state and reward that we observed from the real experience, and the question also says that we will do 5 planning steps, so  $n=5$

So 5 updates based on real experience, and for each of these updates we do 5 planning steps

$\rightarrow 5 \times 5 = 25$

For each update we do based on the real experience, we do  $n \times$  many updates using the simulated experience.

对于我们基于真实经验所做的每一次更新，我们使用模拟经验做5次更新。

问题显示，有5个时间步的 state-action 对及其随后的下一个状态和奖励是我们从真实经验中观察到的。问题还说我们将做5个计划步骤，所以 $n=5$

所以有5个基于真实经验的更新，对于每个更新，我们做5个规划步骤

$\rightarrow 5 \times 5 = 25$

# Worksheet 9: Planning, Learning & Acting

CMPUT 397  
March 19, 2021

3. Modify the Tabular Dyna-Q algorithm so that it uses Expected Sarsa instead of Q-learning. Assume that the target policy is  $\epsilon$ -greedy. What should we call this algorithm?

*We can store information regarding  $S$  and  $A$  in the table.*

Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s) \Rightarrow$  initialize model.

Loop forever:

- $S \leftarrow$  current (nonterminal) state
- $A \leftarrow \epsilon\text{-greedy}(S, Q)$
- Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment)
- Loop repeat  $n$  times:

*$\sim$  # of learning steps.*

$S \leftarrow$  random previously observed state

$A \leftarrow$  random action previously taken in  $S$

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

*Call it Soft Actor-Critic (SAC).*

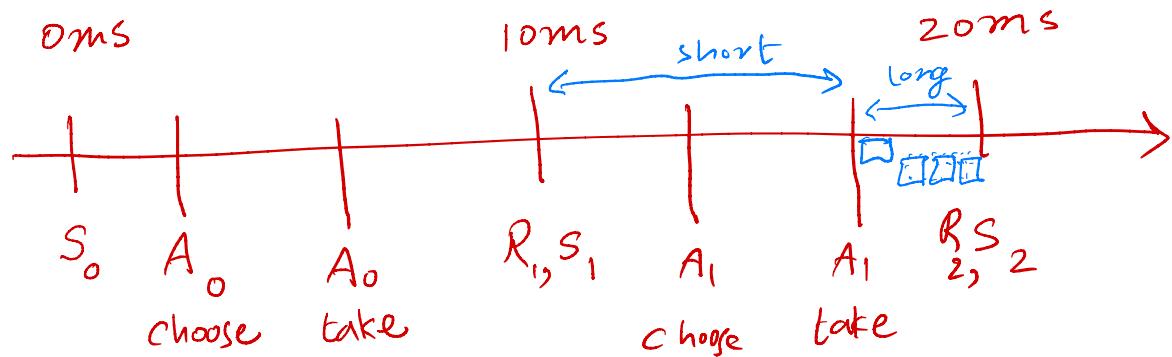
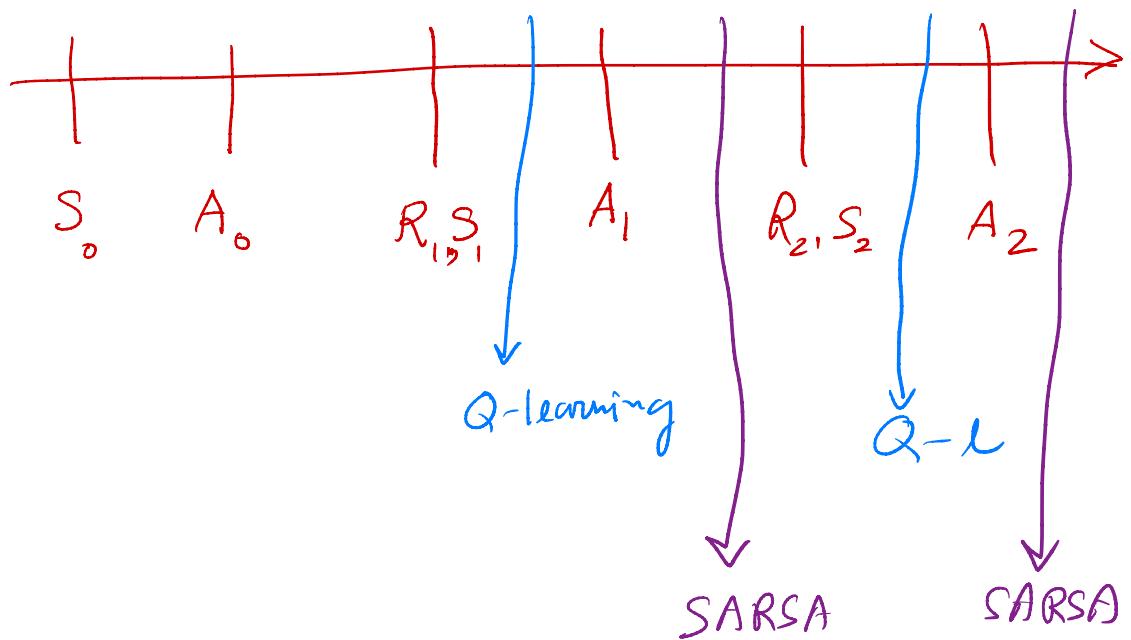
*Model is a table. we access the table by indexing  $S$  and  $A$ . No State 1, Action 1.*

Vova Selin 对所有人说 下午1:13

Does that mean we need to store 2 separate lists of states we've visited and actions we've taken? *yes.*

Siyuan Yu 对所有人说 下午1:15

You would never query for a  $Model(S, A)$  for  $S, A$  that's never seen before in this code, correct? So it doesn't matter how you initialize them



# Worksheet 10:

## On-policy Prediction with Approximation

CMPUT 397  
March 19, 2021

2. Find the gradient of  $f$  <sup>function</sup>

(a) if  $f(x, y, z) = \frac{y^z}{x}$  <sup>scalars.</sup>

(b) if  $f(x) = e^{x^2+5}$  <sup>直接对x求导即可.</sup>

(c) if  $f(\mathbf{x}) = \mathbf{x}^T \beta$  where  $\mathbf{x}$  is a vector in  $\mathbb{R}^N$  and  $\beta$  is a vector of constants in  $\mathbb{R}^N$

(d) if  $f(\mathbf{x}) = (\mathbf{x}^T \beta - y)^2$  where  $\mathbf{x}$  is a vector in  $\mathbb{R}^N$ ,  $\beta$  is a vector of constants in  $\mathbb{R}^N$ , and  $y$  is a scalar in  $\mathbb{R}$

Gradient of  $f$ :  $\nabla f(x, y, z) = \begin{bmatrix} \frac{\partial f(x, y, z)}{\partial x} \\ \frac{\partial f(x, y, z)}{\partial y} \\ \frac{\partial f(x, y, z)}{\partial z} \end{bmatrix} \Rightarrow$  求的是对  $f$  中  $x, y, z$  的偏导数.

(a)  $\nabla f(x, y, z) = \begin{bmatrix} \frac{\partial f(x, y, z)}{\partial x} \\ \frac{\partial f(x, y, z)}{\partial y} \\ \frac{\partial f(x, y, z)}{\partial z} \end{bmatrix}$

$$\begin{aligned} \therefore \frac{\partial f(x, y, z)}{\partial x} &= y^z \frac{\partial (\frac{1}{x})}{\partial x} \Rightarrow x^{-1} \Rightarrow -x^{-2} \\ &= -\frac{y^z}{x^2} \end{aligned}$$

$$\begin{aligned} \frac{\partial f(x, y, z)}{\partial y} &= \frac{1}{x} \frac{\partial y^z}{\partial y} \Rightarrow \text{直接对 } y \text{ 求导.} \\ &= z y^{z-1} \end{aligned}$$

$$\begin{aligned} \frac{\partial f(x, y, z)}{\partial z} &= \frac{1}{x} \cdot \frac{\partial e^{\ln y^z}}{\partial z} \Rightarrow e^{\ln y^z} = y^z \\ &= \frac{1}{x} \cdot \frac{\partial e^{\ln y^z}}{\partial (z \ln y)} \cdot \frac{\partial (z \ln y)}{\partial z} \\ &= \frac{y^z \ln y}{x} \end{aligned}$$

$$\therefore \nabla f(x, y, z) = \begin{bmatrix} -\frac{y^z}{x^2} \\ z y^{z-1} \\ \frac{y^z \ln y}{x} \end{bmatrix}$$

(b)  $f(\mathbf{x}) = \mathbf{x}^T \beta$   $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$

$$\begin{aligned} &= \sum_{i=1}^n x_i \beta_i \quad \text{因为是向量的乘积, 为了做点积, 对 } \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = x_1 \cdot y_1 + x_2 \cdot y_2. \\ \frac{\partial f(\mathbf{x})}{\partial x_j} &= \frac{\partial \sum_{i=1}^n x_i \beta_i}{\partial x_j} \\ &= \frac{\partial x_j \beta_j}{\partial x_j} + \underbrace{\frac{\partial \sum_{i \neq j} x_i \beta_i}{\partial x_j}}_{\text{元素不包含 } x_j, \therefore \text{直接对 } \beta_j \text{ 求导.}} \end{aligned}$$

$$= \beta_j$$

$$\therefore \nabla f(\mathbf{x}) = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} = \beta$$

(c)  $f(\mathbf{x}) = (\mathbf{x}^T \beta - y)^2$   $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$

$$\begin{aligned} \frac{\partial f(\mathbf{x})}{\partial x_j} &= \frac{\partial (\mathbf{x}^T \beta - y)^2}{\partial (\mathbf{x}^T \beta - y)} \cdot \frac{\partial (\mathbf{x}^T \beta - y)}{\partial x_j} \\ &= 2(\mathbf{x}^T \beta - y) \cdot \beta_j \quad \text{因为 } \frac{\partial (\mathbf{x}^T \beta - y)}{\partial x_j} = \beta_j. \end{aligned}$$

$$\therefore \nabla f(\mathbf{x}) = 2 \cdot (\mathbf{x}^T \beta - y) \cdot \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = 2 \cdot (\mathbf{x}^T \beta - y) \cdot \beta$$

# Worksheet 10: On-policy Prediction with Approximation

CMPUT 397  
March 19, 2021

3. (Exercise 9.1 S&B) Show that tabular methods such as presented in Course 2 of the MOOC (and Part I of the book) are a special case of linear function approximation. What would the feature vectors be?

$$S = \{s_1, s_2, \dots, s_n\}$$

State Space

Tabular case: have  $n$  elements in the table;  $V(s_1), V(s_2), \dots, V(s_n)$  to be updated:  $V(s_i) \leftarrow V(s_i) + \alpha [G - V(s_i)]$

linear function approximation:  $\vec{w} = [w_1, w_2, \dots, w_n]^T$

$$\vec{x}(s_k) = [x_1(s_k), x_2(s_k), \dots, x_n(s_k)]^T$$

$$\hat{V}(s_k, \vec{w}) = \sum_{i=1}^n w_i x_i(s_k)$$

$$\frac{\partial \hat{V}(s_k, \vec{w})}{\partial w_i} \Rightarrow \text{the feature of } x_i(s_k).$$

for each  $i \in \{1, \dots, n\}$ ,  $w_i \leftarrow w_i + \alpha [G - \hat{V}(s_k, \vec{w})]$ .  $\frac{\partial \hat{V}(s_k, \vec{w})}{\partial w_i}$   $\Rightarrow$  derive MC method from sample average.

for each  $s_k \in S$ , assign  $w_k$  to approximate its value.

$$\because \text{one-hot state feature } x_i(s_k) = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases}$$

one-hot vector:

vector  $\vec{s}$  State Space  $\rightarrow$   $\vec{s} = (s_1, s_2, \dots, s_n)$   $\vec{s} = (s_1, s_2, \dots, s_n)$   $\vec{s} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{bmatrix}$

$$\hat{V}(s_k, \vec{w}) = w_k$$

$$\therefore \begin{cases} w_i \leftarrow w_i + \alpha [G - w_k] & \text{if } i=k \\ w_i \leftarrow w_i + \alpha [G - w_k] \cdot 0 & \text{if } i \neq k \end{cases} \Rightarrow w_k \leftarrow w_k + \alpha [G - w_k] \quad \text{for the tabular case.}$$

# Worksheet 10: On-policy Prediction with Approximation

CMPUT 397  
March 19, 2021

4. Consider a random walk with three states  $\mathcal{S} = \{1, 2, 3\}$  and a discount rate  $\gamma = 0.9$ . Suppose  $\hat{v}(s; w) = w \cdot s$ , where  $w$  and  $s$  are scalars with  $w$  initialized to 0 and  $s \in \mathcal{S}$ . You observed the following episode:

$$S_0 = 1, R_1 = -7, S_1 = 3, R_2 = 5, S_2 = 1, R_3 = 10.$$

Then if  $S_3$  is the terminal state and  $\alpha = 0.1$ , what is

- (a) the gradient Monte-Carlo estimate for  $\hat{v}(1; w)$ ? Since 1.
- (b) the semi-gradient Linear TD estimate for  $\hat{v}(1; w)$ ?
- (c) the estimate for both gradient Monte-Carlo and semi-gradient Linear TD of  $\hat{v}(1; w)$  after seeing the same episode again?

**Gradient Monte Carlo Algorithm for Estimating  $\hat{v} \approx v_\pi$**

Input: the policy  $\pi$  to be evaluated  
Input: a differentiable function  $\hat{v}: \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$   
Algorithm parameter: step size  $\alpha > 0$   
Initial value-function weights  $w \in \mathbb{R}^d$  arbitrarily (e.g.,  $w = \mathbf{0}$ )  
Loop forever (for each episode):  
  Generate an episode  $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$  using  $\pi$   
  Loop for each step of episode,  $t=0, 1, \dots, T-1$ :  
     $w \leftarrow w + \alpha (G_t - \hat{v}(S_t, w)) \nabla \hat{v}(S_t, w)$

**Semi-gradient TD(0) for estimating  $\hat{v} \approx v_\pi$**

Input: the policy  $\pi$  to be evaluated  
Input: a differentiable function  $\hat{v}: \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$   
Algorithm parameter: step size  $\alpha > 0$   
Initial value-function weights  $w \in \mathbb{R}^d$  arbitrarily (e.g.,  $w = \mathbf{0}$ )  
Loop for each episode:  
  Initialize  $S$   
  Loop for each step of episode:  
    Choose  $A \sim \pi(\cdot | S)$   
    Take action  $A$ , observe  $R, S'$   
     $w \leftarrow w + \alpha [R + \gamma \hat{v}(S', w) - \hat{v}(S, w)] \nabla \hat{v}(S, w)$   
     $S \leftarrow S'$   
  until  $S$  is terminal

a).  $w \leftarrow w + \alpha [G_t - \hat{v}(S_t, w)] \nabla \hat{v}(S_t, w)$  → the gradient of  $\hat{v}(S_t, w)$

$$\hat{v}(S_t, w) = w \cdot s$$

$$G_2 = R_3 = 10$$

$$G_1 = R_2 + \gamma G_2 = 14$$

$$G_0 = R_1 + \gamma G_1 = 3.6$$

when  $t=0$ ,  $w \leftarrow w + \alpha [G_0 - \hat{v}(S_0, w)] \cdot S_0$

$$= 0 + 0.1 [3.6 - (0 \cdot w)] \cdot 1$$

$$= 1.12$$

$t=1$   $w \leftarrow w + \alpha [G_1 - \hat{v}(S_1, w)] \cdot S_1$

$$= 1.12 + 0.1 [14 - 3 \cdot 1.12] \cdot 3$$

$$= 4.31$$

$t=2$   $w \leftarrow w + \alpha [G_2 - \hat{v}(S_2, w)] \cdot S_2$

$$= 4.31 + 0.1 [10 - 1 \cdot 4.31] \cdot 10$$

$$= 4.88$$

b).  $w \leftarrow w + \alpha [R_t + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)] \cdot S_t$

$$\therefore \text{when } t=0, \quad w \leftarrow w + \alpha [R_1 + \gamma \hat{v}(S_2, w) - \hat{v}(S_0, w)] \cdot S_0$$

$$= 0 + 0.1 \times [-7 + 0.9 \times 3 \cdot 10 - (0 \cdot w)] \times 1$$

$$= -0.7$$

$t=1$ ,  $w \leftarrow w + \alpha [R_2 + \gamma \hat{v}(S_3, w) - \hat{v}(S_1, w)] \cdot S_1$

$$= -0.7 + 0.1 [5 + 0.9 \times 10 - 0.7] \times 3$$

$$= 1.24$$

$t=2$ ,  $w \leftarrow w + \alpha [R_3 - \hat{v}(S_2, w)] \cdot S_2$  S<sub>2</sub> is terminal state. ∴ don't have  $\gamma \hat{v}(S_3, w)$

$$= 1.24 + 0.1 [10 - 1 \cdot 1.24] \times 10$$

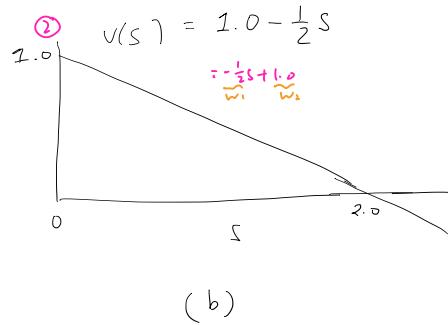
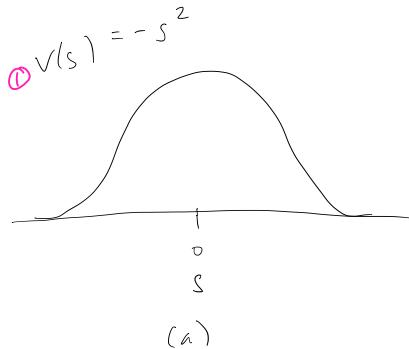
$$= 2.116$$

# Worksheet 11:

## Constructing Features for Prediction

CMPUT 397  
May 1, 2021

1. Consider the following two functions.



- (a) Design features for each function, to approximate them as a linear function of these features. Can you design features to make the approximation exact?
- (b) Can you design **one set of features**, that allows you to represent both functions?

a). ①  $f(s) = s^2$   
 $\hat{v}(s) = w_1 s^2 = w_1 s^2$   
 $\therefore w_1 = 1$

② we have  $s, 1$ ;  
 $\therefore \hat{v}(s) = w_1 s + w_2 1 = 1 - \frac{1}{2}s$   
 $\therefore w_1 = -\frac{1}{2}, w_2 = 1$

b).  $(1, s, s^2) \Rightarrow$  ① ② 指前.

$\therefore$  ①:  $\hat{v}(s) = w_1 s^2 + w_2 s + w_3 1$   
 $w_1 = 1, w_2 = 0, w_3 = 0$ .

②:  $\hat{v}(s) = w_1 s^2 + w_2 s + w_3 1$   
 $w_1 = 0, w_2 = -\frac{1}{2}, w_3 = 1$ .

# Worksheet 11: Constructing Features for Prediction

CMPUT 397  
May 1, 2021

2. Consider the following neural network  $\hat{v}$  with one-hidden layer, relu activation function  $g$ , with weights  $\mathbf{W}^{[0]}, \mathbf{W}^{[1]}, \mathbf{b}^{[0]}, \mathbf{b}^{[1]}$ ,

$$\hat{v}(\mathbf{x}; \mathbf{W}^{[0]}, \mathbf{W}^{[1]}, \mathbf{b}^{[0]}, \mathbf{b}^{[1]}) = \mathbf{W}^{[1]} g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]}) + \mathbf{b}^{[1]}.$$

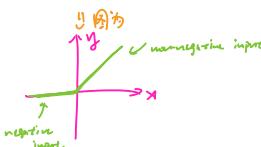
Recall that the following gradients  $\Rightarrow$  input is a vector.  $i$  and  $j$  should also be a vector.

$$\begin{aligned} \frac{\partial \hat{v}}{\partial \mathbf{W}^{[1]}_{ij}} &= g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j \Rightarrow \text{partial derivative of } w_{ij}^{[1]} \\ \frac{\partial \hat{v}}{\partial \mathbf{b}^{[1]}_j} &= 1 \\ \frac{\partial \hat{v}}{\partial \mathbf{b}^{[0]}_i} &= \sum_k \mathbf{W}^{[1]}_{ki} \frac{\partial g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_i}{\partial \mathbf{b}^{[0]}_i} \\ \frac{\partial \hat{v}}{\partial \mathbf{W}^{[0]}_{ij}} &= \sum_k \mathbf{W}^{[1]}_{ki} \frac{\partial g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_i}{\partial \mathbf{W}^{[0]}_{ij}} \end{aligned}$$

(input  $\mathbf{W}^{[0]}$ )

- (a) What are the derivatives specifically for the relu activation  $g$ ?
- (b) We talked about carefully initializing the weights for the NN. For example, each weight can be sampled from a Gaussian distribution. Imagine instead you decided to initialize all the weights to zero. Why would this be a problem? Hint: Consider the derivatives in (a).

A).  $f(a) = \begin{cases} 0 & a < 0 \Rightarrow \text{relu function output is 0 if input is negative.} \\ a & a \geq 0 \Rightarrow \text{relu function output is the input itself if the input is non-negative.} \end{cases}$

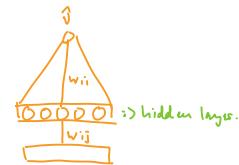


$\therefore$  derivative of relu function  $g$  would be:  $g'(a) = \begin{cases} 0 & a < 0 \\ 1 & a \geq 0 \end{cases}$

$$\frac{\partial \hat{v}}{\partial \mathbf{W}^{[0]}_{ij}} = g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j = \begin{cases} 0 & \text{if } (\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j < 0 \\ (\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j & \text{if } (\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j \geq 0 \end{cases}$$

$$\frac{\partial \hat{v}}{\partial \mathbf{W}^{[0]}_{ij}} = \mathbf{W}^{[1]}_{ij} \frac{\partial g(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j}{\partial \mathbf{W}^{[0]}_{ij}} = \mathbf{W}^{[1]}_{ij} g'(\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j \cdot x_j$$

$$= \begin{cases} 0 & \text{if } (\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j < 0 \\ \mathbf{W}^{[1]}_{ij} \cdot x_j & \text{if } (\mathbf{W}^{[0]}\mathbf{x} + \mathbf{b}^{[0]})_j \geq 0 \end{cases}$$



- b). If all weights are 0, then no matter what the input is (20 or 30), the output will be 0. (因为  $W \cdot x + b = 0$ )

# Worksheet 11:

## Constructing Features for Prediction

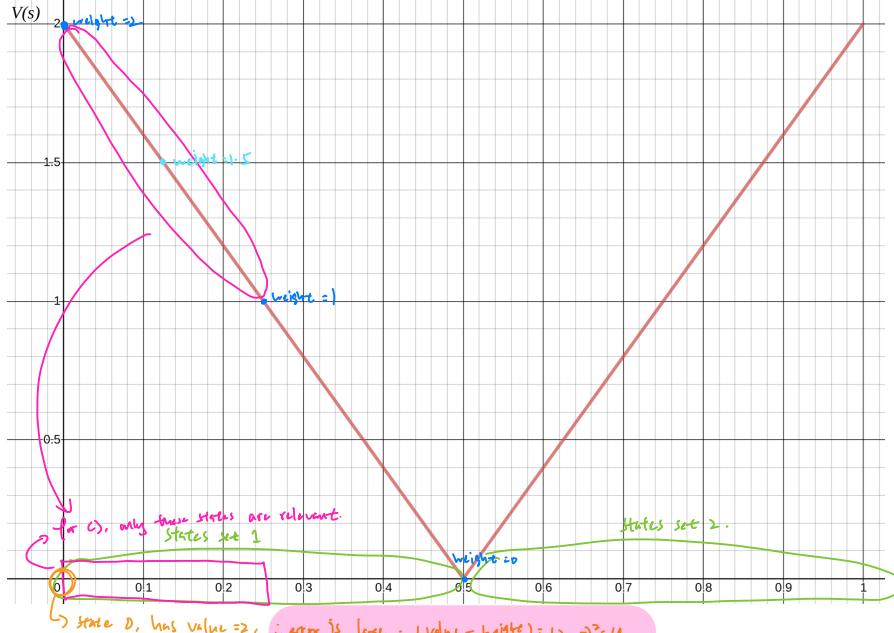
# CMPUT 397

## May 1, 2021

3. Consider a problem with the state space,  $\mathcal{S} = \{0, 0.01, 0.02, \dots, 1\}$ . Assume the true value function is

$$v_\pi(s) = 4|s - 0.5|$$

which is visualized below. We decide to create features with state aggregation, and choose to aggregate into two bins:  $[0, 0.5]$  and  $(0.5, 1]$ .



- (a) What are the possible feature vectors for this state aggregation?

(b) Imagine you minimize the  $\overline{VE}(\mathbf{w}) = \sum_{s \in \mathcal{S}} d(s)(v_\pi(s) - \hat{v}(s, \mathbf{w}))^2$  with a uniform weighting  $d(s) = \frac{1}{101}$  for all  $s \in \mathcal{S}$ . What vector  $\mathbf{w}$  is found?

(c) Now, if the agent puts all of the weighting on the range  $[0, 0.25]$ , (i.e.  $d(s) = 0$  for all  $s \in (0.25, 1]$ ), then what vector  $\mathbf{w}$  is found by minimizing  $\overline{VE}$ ?

6).  $\{1, 0\}^T$

$[0, 1]^7$   $\mathbb{R}^7$  Space  $\mathbb{R}$   $\mathbb{R}^7$   $\mathbb{R}^7$

b). State 1,  $x=0$  and  $y=2$ ; answer is 12.

练习3：阅读下面的短文，完成练习。

AT state 2.5  $x=0.3$  and  $y=0$  the value is 0)  $\therefore$  error :

if we choose weight  $\Rightarrow$ : ~~all states 0~~ state 1: error =  $(2 - 2)^2 = 0$

Section 3: Answers (1-12)

第二步

if we close the weight, in the state 0.5

High  $T_B$  in the mid

∴ 司徒(御率)

∴ the width error we can get =  $(2-1) = 1$

∴ We should choose weight  $m = 7.1$  g.  $\therefore$  weight 5 g is the smallest weight of the beam.  $\therefore$   $1/2$  error is small.

- weight  $\mathbf{w} = [1.5, 1.5]^T$  to give the

# Worksheet 12: Control with Approximation

CMPUT 397  
April 7, 2022

3. (Exercise 10.5 S&B) What equations are needed, in addition to the equation below,

$$\delta_t = R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t) \quad \text{⇒ the TD error for differential TD(0)}$$

to specify the differential version of  $TD(0)$ ?

∴ You have to write the TD update that is missing:

∴  $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha(S_t, \nabla \hat{v}(S_t, \mathbf{w}_t)) \Rightarrow TD \text{ update.}$

Specify that:  $\begin{cases} \text{Sample generated by following policy } \pi_t \text{ and } \hat{v}(S_{t+1}) \approx \hat{v}_{\pi_t}(S_{t+1}) \\ \bar{R}_{t+1} = R_{t+1} + \gamma(\bar{R}_t - \bar{R}_{t+1}) \end{cases}$

The update of  $\bar{R}$

on policy learning  
With off-policy,  $\bar{R}_{t+1}$  importance sampling ratio!

# Worksheet 12: Control with Approximation

CMPUT 397  
April 7, 2022

4. (Exercise 10.6 S&B) Suppose there is an MDP that under any policy produces the deterministic sequence of rewards  $+1, 0, +1, 0, +1, 0, \dots$  going on forever. Technically, this is not allowed because it violates ergodicity; there is no stationary limiting distribution  $\mu_\pi$  and the limit

$$\lim_{t \rightarrow \infty} \mathbb{E}[R_t | S_0, A_{0:(t-1)} \sim \pi]$$

does not exist. Nevertheless, the average reward,

$$\lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t | S_0, A_{0:(t-1)} \sim \pi]$$

is well defined; What is it? Now consider two states in this MDP. From A, the reward sequence is exactly as described above, starting with a  $+1$ , whereas, from B, the reward sequence starts with a  $0$  and then continues with  $+1, 0, +1, 0, \dots$  The differential return

$$G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

A: 1 0 1 0 ...  
B: 0 1 0 1 0 ...

is not well defined for this case as the limit does not exist. To repair this, one could alternately define the value of a state as

$$v_\pi(s) = \lim_{\gamma \rightarrow 1} \lim_{h \rightarrow \infty} \sum_{t=0}^h \gamma^t (\mathbb{E}_\pi[R_{t+1} | S_0 = s] - r(\pi))$$

Under this definition, what are the values of states A and B?

Markov reward process:  $C_1 = \frac{1}{1} = 1$   
 $C_2 = \frac{1+0}{2} = \frac{1}{2}$   
 $C_3 = \frac{1+0+1}{3} = \frac{2}{3}$   
 $C_4 = \frac{1+0+1+0}{4} = \frac{1}{2}$   $\Rightarrow$  even number: all half  
 $C_5 = \frac{1+0+1+0+1}{5} = \frac{3}{5}$  odd number: reducing. (converging to  $\frac{1}{2}$ , i.e. never going below  $0.5$ !!).  
 $C_6 = \frac{1+0+1+0+1+0}{6} = \frac{1}{2}$   
 $C_7 = \frac{4}{7}$   
 $\therefore \lim_{h \rightarrow \infty} C_h = \frac{1}{2}$  is the average reward.

Here we have a period of 2:  $\therefore \underline{V_{\pi}(A)} = \frac{1}{2} (2(R_1 - r(\pi)) + (R_2 - r(\pi))$   
 $\qquad\qquad\qquad = \frac{1}{2} (2(1 - \frac{1}{2}) + (0 - \frac{1}{2}))$   
 $\qquad\qquad\qquad = \frac{1}{2} (1 - \frac{1}{2})$   
 $\qquad\qquad\qquad = \frac{1}{4}$   
 $\underline{V_{\pi}(B)} = \frac{1}{2} (2(R_1 - r(\pi)) + (R_2 - r(\pi)))$   
 $\qquad\qquad\qquad = \frac{1}{2} (2(0 - \frac{1}{2}) + (1 - \frac{1}{2}))$   
 $\qquad\qquad\qquad = -\frac{1}{4}$

Starting from A:  $\overset{R_1 \text{ in } \pi}{\text{B}}$   $\overset{\text{average reward}}{\text{B}}$   $\overset{R_2 \text{ in } \pi}{\text{B}}$   $\overset{\text{average reward}}{\text{B}}$   $\dots$   
 $\therefore V_{\pi}(A) = 1 - \frac{1}{2} + \frac{1}{4} - \frac{1}{8} + \dots$

Starting from B:  $\overset{R_1 \text{ in } \pi}{\text{B}}$   $\overset{\text{average reward}}{\text{B}}$   $\overset{R_2 \text{ in } \pi}{\text{B}}$   $\overset{\text{average reward}}{\text{B}}$   $\dots$   
 $\therefore V_{\pi}(B) = 0 - \frac{1}{2} + \frac{1}{4} - \frac{1}{8} + \dots$

Confirming accuracy through Bellman equation:

$$V_{\pi}(A) = 1 - \frac{1}{2} + \frac{1}{4}$$

$$V_{\pi}(B) = 0 - \frac{1}{2} + \frac{1}{4}$$

$$\therefore V_{\pi}(A) = 1 - \frac{1}{2} + \frac{1}{4}$$

$$= \frac{1}{2} - \frac{1}{4}$$

$$= \frac{1}{4} \Rightarrow \text{correct!!}$$

$V_{\pi}(B) = 0 - \frac{1}{2} + \frac{1}{4}$

$$= -\frac{1}{2} + \frac{1}{4}$$

$$= -\frac{1}{4} \Rightarrow \text{incorrect!!}$$

using Bellman:  $V_{\pi}(B) = \lim_{\gamma \rightarrow 1} (\gamma(1 - \pi(\pi)) + \gamma(0 - \pi(\pi)) + \dots)$   
 $\qquad\qquad\qquad = \lim_{\gamma \rightarrow 1} (1 - \frac{1}{2} + \gamma(0 - \frac{1}{2}) + \gamma^2(1 - \frac{1}{2}) + \dots)$   
 $\qquad\qquad\qquad = \lim_{\gamma \rightarrow 1} (1 - \frac{1}{2} - \frac{1}{2} + \frac{\gamma^2}{2} \dots)$   
 $\qquad\qquad\qquad = \frac{1}{2} \lim_{\gamma \rightarrow 1} (1 - \gamma^2 + \gamma^2 - \gamma^2 \dots)$   
 $\qquad\qquad\qquad = \frac{1}{2} \lim_{\gamma \rightarrow 1} \frac{1}{1-\gamma}$   
 $\qquad\qquad\qquad = \frac{1}{2} \cdot \frac{1}{2} \cdot 4$   
 $\qquad\qquad\qquad = \frac{1}{4}$

# Worksheet 12: Control with Approximation

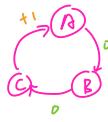
CMPUT 397  
April 7, 2022

5. (Exercise 10.7 SEC) Consider a Markov reward process consisting of a ring of three states A, B, and C, with state transitions going deterministically around the ring. A reward of +1 is received upon arrival in A and otherwise the reward is 0. What are the differential values of the three states using

$$v_\pi = \lim_{\gamma \rightarrow 1} \lim_{h \rightarrow \infty} \sum_{t=0}^h \gamma^t (\mathbb{E}_\pi[R_{t+1}|S_0 = s] - r(\pi))$$

⇒ 這樣你就要 to find the answer equal this way!  
不然就用動力方程來找 answer!!

Markov reward process:



$$\therefore \text{Average reward } r(\pi) = \frac{1+0+0}{3} = \frac{1}{3}$$

⇒ 3/3 得到的 reward.

$$\left. \begin{aligned} v_\pi(C) &= (1 - \frac{1}{3}) + v_\pi(A) \\ v_\pi(B) &= (0 - \frac{1}{3}) + v_\pi(C) \\ v_\pi(A) &= (0 - \frac{1}{3}) + v_\pi(B) \end{aligned} \right\} \text{Solve by Bellman equation}$$

$$v_\pi(B) = \frac{1}{3} (3 (R_1 - r(\pi)) + 2 (R_2 - r(\pi)) + R_3 - r(\pi))$$

$$= \frac{1}{3} (3 (0 - \frac{1}{3}) + 2 (0 - \frac{1}{3}) + 1 (1 - \frac{1}{3}))$$

LL P 的概率, reward 當 0 : 0 → 0 → 1

$$= \frac{1}{3} (-1 - \frac{2}{3} + \frac{2}{3})$$

$$= -\frac{1}{3}$$

∴ can find  $v_\pi(B)$  and  $v_\pi(C)$  by Bellman equation. 但因為用了 Bellman equation 找 answer, 之後就用另外種方法推導.

$$v_\pi(C) = 1 - \frac{1}{3} + v_\pi(A) = \frac{2}{3} - \frac{1}{3} = \frac{1}{3}$$

$$v_\pi(B) = 0 - \frac{1}{3} + v_\pi(C) = -\frac{1}{3} + \frac{1}{3} = 0$$

⇒  $v_\pi(B)$  LL P 當 0, 得到 reward 0, ∵ 0 是 vulnerable 狀態最後, ∴  $v_\pi(B)$  为 0.

$$B \quad \text{2 reward } (2 \times 0) \times 0.1, \therefore B$$

$$\leftarrow \text{2 reward, } \therefore v_\pi(B) = 0.$$

$$C \quad 1 - \text{reward } 0 \times 0.1, \therefore C$$

$$\leftarrow \text{最高, } \therefore v_\pi(C) = 1. \quad \text{最高.}$$

∴ Q.