

Ex 2

Adhar Banerjee

CS 5180

Wednesday, September 28, 2022 6:31 PM

Q.1

a)

State space = { set of all possible states from  
 $(0,0)$  to  $(10,10)$   
 except the Wall cells }  
 $(x,y)$  + all  $x,y$   $(0,10)$   
 except where  $(x,y)$  is wall

Action Space = The action chosen by the State  
 { UP, DOWN, LEFT, RIGHT }

b)

Dynamic:  $p(s', r | s, a)$

As in the Ex(0) the 0.2 Error to  
 go in 1 direction

$$\{ (0,0), \text{DOWN} \} = \begin{aligned} &P((0,0) | (0,0), \text{DOWN}) = 0.8 \\ &P((0,0) | (0,0), \text{LEFT}) = 0.1 \\ &P((0,1) | (0,0), \text{RIGHT}) = 0.1 \end{aligned}$$

$$\{ (1,5), \text{UP} \} = \begin{aligned} &P((1,5) | (1,5), \text{UP}) = 0.8 \\ &P((1,4) | (1,5), \text{LEFT}) = 0.1 \\ &P((1,6) | (1,5), \text{RIGHT}) = 0.1 \end{aligned}$$

$$\{ (9,10), \text{RIGHT} \} = \begin{aligned} &P((10,10) | (9,10), \text{RIGHT}) = 0.8 \\ &P((9,10) | (9,10), \text{UP}) = 0.1 \\ &P((9,9) | (9,10), \text{DOWN}) = 0.1 \end{aligned}$$

②

(a) Suppose you treated pole-balancing as an episodic task but also used discounting, with all rewards zero except for -1 upon failure. What then would the return be at each time? How does this return differ from that in the discounted, continuing formulation of this task? (Derive expressions for both the episodic and continuing cases.)

Episodic Case

$$R_{\text{terminal state}} = -1$$

Cumulative Reward

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T \quad - (3.7)$$

with Discounting

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T \\ &= \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} = \boxed{-\gamma^{T-t-1}} \end{aligned}$$

for Continuous Case

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = -\gamma^K$$

The difference to the episodic with discounting & Continuing case is that

episodic with discounting will have one failure & return the  $G_t$ . whereas with Continuous Case it will have multiple failure & will be keep on updated.

b)

Imagine that you are designing a robot to run a maze. You decide to give it a reward of +1 for escaping from the maze and a reward of zero at all other times. The task seems to break down naturally into episodes – the successive runs through the maze – so you decide to treat it as an episodic task, where the goal is to maximize expected total reward (Equation 3.7). There is no discounting, i.e.  $\gamma = 1$ . After running the learning agent for a while, you find that it is showing no improvement in escaping from the maze. What is going wrong? Have you effectively communicated to the agent what you want it to achieve?

The problem here is that no matter the time agent spend in the maze, the reward maximum will be +1 for escaping the maze, as there is no discount factor  $\gamma$ .

8.3

(a) In the gridworld example (Figure 3.2 in RL2e, see below), rewards are positive for goals, negative for running into the edge of the world, and zero the rest of the time. Are the signs of these rewards important, or only the intervals between them? Prove, using Equation 3.8, that adding a constant  $c$  to all the rewards adds a constant,  $v_c$ , to the values of all states, and thus does not affect the relative values of any states under any policies. What is  $v_c$  in terms of  $c$  and  $\gamma$ ?

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad \text{--- (3.8)}$$

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right]$$

Adding Constant  $c$

$$V_{\pi}(s)(c) = E_{\pi}\left[\sum_{i=0}^{\infty} \gamma^i (R_{t+i+1} + c) \mid S_t = s\right]$$

$$= E_{\pi}\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i+1} \mid S_t = s\right] + E_{\pi}\left[\sum_{i=0}^{\infty} \gamma^i c \mid S_t = s\right]$$

$$V_{\pi}(s)(c) = V_{\pi}(s) + \frac{c}{1-\gamma}$$

(b) Now consider adding a constant  $c$  to all the rewards in an episodic task, such as maze running. Would this have any effect, or would it leave the task unchanged as in the continuing task above? Why or why not? Give an example.

In the episodic task, such as maze running, the sign of reward is important, as the -ve reward makes the agent to leave the maze fast. Therefore adding constant  $c$  to the reward, if it keep negative reward remains intact the agent will find the exit from the maze but the urgency to find it may vary & if the reward turns +ve after adding constant, the agent might not leave the maze, as be in the maze would be

most beneficial.

Q.4

(a) The Bellman equation (Equation 3.14) must hold for each state for the value function  $v_\pi$  shown in Figure 3.2 (right) (see above) of Example 3.5. Show numerically that this equation holds for the center state, valued at +0.7, with respect to its four neighboring states, valued at +2.3, +0.4, -0.4, +0.7. (These numbers are accurate only to one decimal place.) Note that Figure 3.2 (right) (see above) is the value function for the equiprobable random policy.  $\gamma = 0.9$

$$\begin{aligned}
 V_\pi(s) &\equiv \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [\gamma + \gamma V_\pi(s')] \quad \forall s \in S \quad (3.14) \\
 &= \frac{1}{4} \left( \frac{1}{4} (0 + 0.9 (2.3)) + \frac{1}{4} (0 + 0.9 (0.4)) + \frac{1}{4} (0 + 0.9 (-0.4)) \right) \\
 &\quad + \frac{1}{4} (0 + 0.9 (0.7)) \\
 &= \frac{1}{4} \times 0.9 (2.3 + 0.4 - 0.4 + 0.7) \\
 V_\pi(s) &= \frac{2.7}{4} = \frac{0.675}{1} \\
 &= 0.7 \quad (\text{for one decimal pt})
 \end{aligned}$$

(b) The Bellman equation holds for all policies, including optimal policies. Consider  $v_*$  and  $\pi_*$  shown in Figure 3.5 (middle, right respectively) (See below). Similar to the previous part, show numerically that the Bellman equation holds for the center state, valued at +17.8, with respect to its four neighboring states, for the optimal policy  $\pi_*$  shown in Figure 3.5 (right) (see below).

from Bellman equation

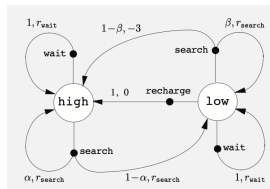
from Centre State +17.8, Can take two state  $\leftarrow \uparrow$  by those action from policy  $\pi_*$

$$\begin{aligned}
 V_\pi(s) &= \frac{1}{2} \left( \frac{1}{2} (0 + 0.9 \times 19.8) + \frac{1}{2} (0 + 0.9 \times 19.8) \right) \\
 &= 0.9 \times 19.8 = 17.82 \\
 \boxed{V_\pi(s) = 17.8} &\quad (\text{one decimal place})
 \end{aligned}$$

Q.5

a)

2 state recycling robot from Example 3.3

Charge level  $S = \{ \text{high, low} \}$  $A(\text{high}) = \{ \text{search, wait} \}$  $A(\text{low}) = \{ \text{search, wait, recharge} \}$ Bellman eq<sup>n</sup>

$$V_\pi(s) \equiv \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [\gamma + \gamma V_\pi(s')] \quad \forall s \in S$$

$$\begin{aligned}
 V(\text{high}) &= \pi(\text{search/high}) \left[ \alpha [\gamma r_{\text{search}} + \gamma V(s'=\text{high})] + (1-\alpha) [\gamma r_{\text{search}} + \gamma V(s'=\text{low})] \right] \\
 &\quad + \pi(\text{wait/high}) \left[ 1 \cdot [\gamma r_{\text{wait}} + \gamma V(s'=\text{high})] + 0 \right]
 \end{aligned}$$

$$\begin{aligned}
 v(\text{low}) = & \pi(\text{search}|\text{low}) [1 - \beta [-3 + \gamma(v(\text{high}))] + \beta (r_{\text{search}} + \gamma(v(\text{low})))] \\
 & + \pi(\text{wait}|\text{low}) [0 + 1(r_{\text{wait}} + \gamma v(\text{low}))] \\
 & + \pi(\text{recharge}|\text{low}) [\gamma v(\text{high})]
 \end{aligned}$$

b)

(b) You should now have two linear equations involving two unknowns,  $v(\text{high})$  and  $v(\text{low})$ , as well as involving the policy  $\pi(\text{als})$ ,  $\gamma$ , and the domain parameters. Let  $\alpha = 0.8$ ,  $\beta = 0.6$ ,  $\gamma = 0.9$ ,  $r_{\text{search}} = 10$ ,  $r_{\text{wait}} = 3$ . Consider the policy  $\pi(\text{search}|\text{high}) = 1$ ,  $\pi(\text{wait}|\text{low}) = 0.5$ , and  $\pi(\text{recharge}|\text{low}) = 0.5$ . Find the value function for this policy, i.e., solve the equations for the values of  $v(\text{high})$  and  $v(\text{low})$ . Check that your solution satisfies the Bellman equation.

$$v(\text{high}) = 1 \cdot [0.8 [10 + 0.9 v(\text{high})] + 0.2 [10 + 0.9 v(\text{low})]]$$

$$\text{Let } v(\text{high}) = x, \quad v(\text{low}) = y$$

$$x = 8 + 0.72x + 2 + 0.18y$$

$$0.28x - 0.18y = 10$$

$$28x - 18y = 1000$$

$$14x - 9y = 500 \quad - \textcircled{1}$$

$$v(\text{low}) = 0.5 [3 + 0.9 v(\text{low})] + 0.5 [0.9 v(\text{high})]$$

$$y = 1.5 + 0.45y + 0.45x$$

$$0.45x - 0.55y + 1.5 = 0$$

$$45x - 55y + 150 = 0$$

$$9x - 11y + 30 = 0 \quad - \textcircled{2}$$

$$\textcircled{1} \times 11 - \textcircled{2} \times 9$$

$$(14 \times 11)x - 99y = 500 \times 11$$

$$- (9 \times 9)x + 99y = +30 \times 9$$

$$(154 - 81)x = 5500 + 270$$

$$73x = 5770$$

$$v(\text{high}) = 79.04$$

$$v(\text{low}) = 67.39$$

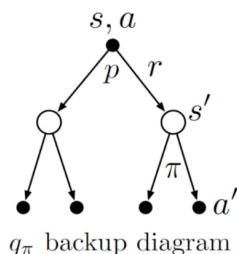
Verify by substituting them back in using Bellman equation

$$\begin{aligned}
 v(\text{high}) &= 0.8 \times (10 + 0.9 \times 79.04) + 0.2 (10 + 0.9 \times 67.39) \\
 &= 64.90 + 14.13
 \end{aligned}$$

$$= 79.03$$

$$\begin{aligned} V(l_{100}) &= 0.5 \times 1 \times (3 \times 0.9 \times 67.39) + 0.5 (0.9 \times 79.04) \\ &= 31.83 + 35.57 \\ V(l_{100}) &= 67.40 \end{aligned}$$

Q.6



$$\begin{aligned} a) \quad V_\pi(s) &= \mathbb{E}_\pi [G_t / S_t = s] \\ &= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} / S_t = s \right] \\ q_\pi(s, a) &= \mathbb{E}_\pi [G_t / S_t = s, A_t = a] \\ &= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} / S_t = s, A_t = a \right] \\ \therefore V_\pi(s) &= \sum_a \pi(a/s) q_\pi(s, a) \end{aligned}$$

$$\begin{aligned} b) \quad q_\pi(s, a) &= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} / S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r / s, a) \left[ r + \gamma V_\pi(s') \right] \quad \text{--- eq (b)} \end{aligned}$$

$$\begin{aligned} c) \quad V_\pi(s') &= \sum_a \pi(a/s') q_\pi(s', a) \quad \text{--- put this in eq (b)} \\ q_\pi(s, a) &= \sum_{s', r} p(s', r / s, a) \left[ r + \gamma \sum_a \pi(a/s') q_\pi(s', a) \right] \end{aligned}$$