

NUID: 002190127

Reinforcement Learning & Sequential Decision Making

HW #02

Sep 28, 2022

Question #01

(a) In this case, state space is coordinate points (x, y)

such that:

$$S = \{(x, y)\}$$

$$x \in [0, 10]$$

$$y \in [0, 10]$$

There are four possible actions, so,

$$A = \{\text{'left', 'right', 'up', 'down'}\}$$

(b) Since, dynamics function is given by:

$$P(s', r | s, a) = \{ \Pr(s_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a) \}$$

so, for given state action pairs:

s	a	r	s'	$P(s', r s, a)$
$(0, 0)$	Down	0	$(0, 0)$ $\xrightarrow{\text{due to noise}}$ $(0, 1)$	α where $\alpha = 0.9$
$(0, 0)$	Down	0	$(1, 0)$	$1 - \alpha$ so, $(1 - \alpha) = 0.1$
$(1, 0)$	Up	0	$(1, 1)$	1
$(1, 0)$	Right	0	$(0, 0)$	$1 - \alpha = 0.1$
$(1, 0)$	Right	1	$(1, 0)$	α ($\because \alpha = 0.9$)

Question #02

(a) Since, unified episodic discounting return is given by:

$$G_t = \sum_{K=t+1}^T \gamma^{K-t-1} R_K$$

so, episodic discounted return will be:

$$G_t = \gamma^0 R_{t+1} + \gamma^1 R_{t+2} + \dots + \gamma^{T-t-1} R_T$$

Here reward = -1 for failure and zero otherwise,

so,

$$G_t = 0 + \gamma^{T-t} (-1)$$

$$G_t = -\gamma^{T-t}$$

where T is timestep when task is failed and hence episode ends.

Similarly,

For continuing case:

$$G_t = -\sum_{K \in K} \gamma^{K-t-1}$$

where K is set of timesteps where task is failed. Note that this reward will increase in long run, irrespective of improved performance. So, designing it as continuous task does not make sense here.

02

Question # 02 (b)

We have designed it such that reward = +1 only when it has exited the maze, and are therefore not incentivising it to learn how to exit the maze faster. To do so, we will need to provide a negative reward proportional to time in maze e.g -1 per timestep.

Question # 03 (a)

↳ The sign of reward is not important. In fact, it is interval between each reward that drives behaviour.

↳ Since eq 3.8 is given by

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

when we add a constant to all rewards, expected return will just receive constant additive term i.e,

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

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

$$G_t = \frac{R_{t+k+1}}{1-\gamma} + \frac{c}{1-\gamma}$$

Since,

$$V_c^{(s)} = E[G_t | S_t = s]$$

$$\begin{aligned} \text{So, } V_c &= E[G_t | S_t = s] \\ &= E\left[\sum_{k=0}^{\infty} \gamma^k R_k\right] \quad \text{so } R_k = c \\ &= E\left[\sum_{k=0}^{\infty} \gamma^k c\right] \\ V_c &= \frac{c}{1-\gamma} \end{aligned}$$

Question # 03 (b)

In episodic task adding a constant to all rewards

does not affect the agent. Here,

$$G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

cumulative reward depends on length of episode.
Time steps that incur positive rewards act to lengthen the episode and vice versa. In maze running example, we may have chosen to give agent

-1 reward at each timestep to ensure it completes task quickly.

So, adding $\epsilon=2$ at every reward such that reward at each timestep is now positive, the agent is now incentivised to not find exit and continue collecting intermediate rewards indefinitely.

Question #04 (a)

Since, Bellman equation is:

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_{\pi}(s')]$$

Here, probability for all 4 actions is equal,

$$\text{so, } \pi(a|s) p(s', r|s, a) = \frac{1}{4} = 0.25$$

where $s \in \{\text{north, south, east, west}\}$

$r = 0$ except for state A & B so,

$$\begin{aligned} \text{so, } V_{\pi}(s) &= (0.25)(0.9 \times 2.3) + (0.25)(0.9 \times 0.7) + \\ &\quad (0.25)(0.9 \times 0.4) + (0.25)(0.9 \times -0.4) \\ &= 0.68 \approx 0.7 \end{aligned}$$

Hence, it satisfies bellman equation.

Question #04 (b)

Since, $V_*(s) = \max_a V_x(s)$

So, for bellman equation

$$V_*(s) = \max_a \left(\sum_{s', r} P(s', r | s, a) [r + \gamma V_*(s')] \right)$$

Here, there are two possible actions for center state, north and west. Both of these action result in same state-value, so, we can choose any of these. Hence, probability of each action is $\frac{1}{2}$.

$$V_*(s) = 0.5 \times (0.9 \times 19.8) + 0.5 (0.9 \times 19.8)$$

$$= 17.8$$

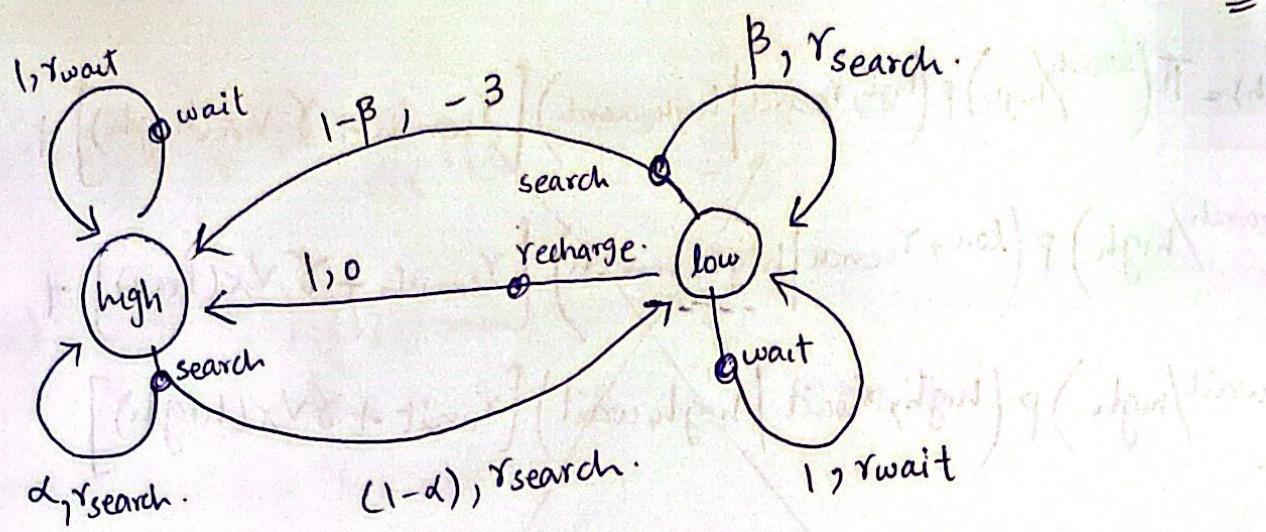
Hence, it satisfies bellman equation.

Question (5) (a)

since, Bellman equation is:

$$V_x = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) [r + \gamma V_x(s')]$$

Now, for 2-state recycling robot,



In tabular form

s	a	s'	r	$P(s', r s, a)$
high	search	high	rsearch	d
high	search	low	rsearch	$1-d$
low	search	high	rsearch	-3
low	search	low	rsearch	B
high	wait	high	rwait	$1-B$
high	wait	low	rwait	d
low	recharge	high	rwait	0

putting values of all these states in bellman equation gives:

so,

\leq

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_{\pi}(s')]$$

$$v_{\pi}(\text{high}) = \sum_a \pi(a|s) \left\{ (1)(r_{\text{wait}} + \gamma v_{\pi}(\text{high})) + (\alpha(r_{\text{search}} + \gamma v_{\pi}(\text{high}))) + ((1-\alpha)(r_{\text{search}} + \gamma v_{\pi}(\text{low}))) \right\}$$

$$= \sum_a \pi(a|s) \left\{ r_{\text{wait}} + \gamma v_{\pi}(\text{high}) + \cancel{\alpha r_{\text{search}}} + \cancel{\alpha \gamma v_{\pi}(\text{high})} + \gamma r_{\text{search}} + \cancel{\gamma v_{\pi}(\text{low})} - \cancel{\alpha r_{\text{search}}} - \cancel{\alpha \gamma v_{\pi}(\text{low})} \right\}$$

$$v_{\pi}(\text{high}) = \sum_a \pi(a|s) \left\{ r_{\text{wait}} + \gamma r_{\text{search}} + (1+\alpha) \gamma v_{\pi}(\text{high}) + (1-\alpha) \gamma v_{\pi}(\text{low}) \right\} \quad (\text{A})$$

and

$$v_{\pi}(\text{low}) = \sum_a \pi(a|s) \left\{ (1)(r_{\text{wait}} + \gamma v_{\pi}(\text{low})) + \beta(r_{\text{search}} + \gamma v_{\pi}(\text{low})) + (1-\beta)(r_{\text{search}} + \gamma v_{\pi}(\text{high})) + (1)(0 + \gamma v_{\pi}(\text{high})) \right\}$$

so,

$$v_{\pi}(\text{low}) = \sum_a \pi(a|s) \left\{ r_{\text{wait}} + \gamma v_{\pi}(\text{low}) + \beta r_{\text{search}} + \beta \gamma v_{\pi}(\text{low}) + r_{\text{search}} + \gamma v_{\pi}(\text{high}) - \beta r_{\text{search}} - \beta \gamma v_{\pi}(\text{high}) + \gamma v_{\pi}(\text{high}) \right\}$$

$$V_{\pi}(\text{low}) = \sum_a \pi(a|s) \left\{ \gamma_{\text{wait}} + \gamma_{\text{search}} + (1+\beta)\gamma V_{\pi}(\text{low}) + (2-\beta)\gamma V_{\pi}(\text{high}) \right\}$$

— (B)

Question # 05 (b)

putting all values in (A)

$$\begin{aligned} V_{\pi}(\text{high}) &= \sum_a \pi(a|s) \left[\gamma_{\text{wait}} + \gamma_{\text{search}} + (1+\alpha)\gamma V_{\pi}(\text{high}) + (1-\alpha)\gamma V_{\pi}(\text{low}) \right] \\ &= (1) \left[3 + 10 + (1+0.8)(0.9) V_{\pi}(\text{high}) + (1-0.8)(0.9) V_{\pi}(\text{low}) \right] \\ &= (1) \left[13 + 1.62 V_{\pi}(\text{high}) + 0.18 V_{\pi}(\text{low}) \right] \\ &\quad + (1) \left[V_{\pi}(\text{high}) - 1.62 V_{\pi}(\text{high}) - 0.18 V_{\pi}(\text{low}) = 13 \right] \\ V_{\pi}(\text{high}) - 1.62 V_{\pi}(\text{high}) - 0.18 V_{\pi}(\text{low}) &= 13 \\ -0.62 V_{\pi}(\text{high}) - 0.18 V_{\pi}(\text{low}) &= 13 \\ 0.62 V_{\pi}(\text{high}) + 0.18 V_{\pi}(\text{low}) &= -13 \quad — (1) \end{aligned}$$

and putting in eq (B):

$$\begin{aligned} V_{\pi}(\text{low}) &= \sum_a \pi(a|s) \left\{ 3 + 10 + (1+0.6)(0.9) V_{\pi}(\text{low}) + (2-0.6)(0.9) V_{\pi}(\text{high}) \right\} \\ &= \sum_a \pi(a|s) \left\{ 13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \right\} \end{aligned}$$

$$\begin{aligned}
 V_{\pi}(\text{low}) &= \pi(\text{wait/low}) \left[13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \right] + \\
 &\quad \pi(\text{recharge/low}) \left[13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \right] \\
 &= 0.5 \left[13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \right] + \\
 &\quad 0.5 \left[13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \right] \\
 V_{\pi}(\text{low}) &= 13 + 1.44 V_{\pi}(\text{low}) + 1.26 V_{\pi}(\text{high}) \\
 -0.44 V_{\pi}(\text{low}) - 1.26 V_{\pi}(\text{high}) &= 13
 \end{aligned}$$

$$\begin{aligned}
 1.26 V_{\pi}(\text{high}) + 0.44 V_{\pi}(\text{low}) &= -13 \quad \text{--- ②} \\
 1.26 V_{\pi}(\text{high}) &= -13 - 0.44 V_{\pi}(\text{low})
 \end{aligned}$$

Solving eq ① and ② gives:

$$V_{\pi}(\text{high}) = -73.47$$

$$V_{\pi}(\text{low}) = 180.86$$

Question #06

(a) The state value function V_x is equal to the expected cumulative return from that state given a distribution of actions. The state-action value function q_{π} is the value of being in a state and taking a deterministic action. Therefore, the state value function is the weighted sum of the state action value function, with weights equal to probabilities of selecting each action:

$$V_{\pi} = \sum_a \pi(a|s) q_{\pi}(s, a)$$

(b) Given action a , the state-action value function is the probability distributions over possible next states & rewards from action times the one-step reward and distributed state value function at next timestep:

$$q_{\pi} = \sum_{s' \in S} \sum_{r \in R} P(s', r | s, a) [r + \gamma V_x(s_{t+1})]$$

$$\begin{aligned}
 (c) \quad q_{\pi}(s, a) &= E_{\pi} \left[G_t | s_t = s, A_t = a \right] \\
 &= E_{\pi} \left[R_{t+1} + \gamma G_{t+1} | s_t = s, A_t = a \right] \\
 &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma E \left[G_{t+1} | s', a' \right] \right] \\
 &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma q_{\pi}(s', a') \right]
 \end{aligned}$$