

Homework 1

Submitted by:

Sayan Chakraborty

Roll No: EE18MTECH11030

Subject: Introduction to
Modern AI

Ans 1:

Derive the Bellman's equation for:

(i) State-value function $V_{\pi}(s)$.

We know,

$$V_{\pi}(s) = E_{\pi} [R_t | S_t = s], \quad R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | S_t = s \right] \quad \text{--- (1)}$$

We use the law of total expectation to derive the Bellman's equation.

The law of total expectation is given as,

$$E[X] = \sum_i E[X | A_i] P(A_i), \quad \{A_i\}_i \text{ is a countable partition of the sample space.}$$

Now,

conditioning (1) over each chosen action at time step t , i.e., a_t

$$V_{\pi}(s) = \sum_a E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | S_t = s, a_t = a \right] P(a_t = a | S_t = s) \quad \text{--- (2)}$$

Now, $P(a_t = a | S_t = s)$ is the policy $\pi(s, a)$

(2) \Rightarrow

$$V_{\pi}(s) = \sum_a \pi(s, a) E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | S_t = s, a_t = a \right] \quad \text{--- (3)}$$

Now, conditioning (3) on the future state s' we get,

$$(3) \Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a, s_{t+1} = s' \right] \cdot \mathbb{P}[s_{t+1} = s' \mid s_t = s, a_t = a]$$

Denote, $\mathbb{P}_{(s,a)}^{s'} = \mathbb{P}[s_{t+1} = s' \mid s_t = s, a_t = a]$

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a, s_{t+1} = s' \right] \mathbb{P}_{(s,a)}^{s'}$$

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} \mathbb{E}_{\pi} [r_{t+1} \mid s, a, s'] + \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s, a, s' \right]$$

Denote, $R_{s,a,s'}^{t+1} = \mathbb{E}_{\pi} [r_{t+1} \mid s, a, s']$

Also by Markov property we have that

$$\mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s, a, s' \right] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s' \right]$$

∴ we have,

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} R_{s,a,s'}^{t+1} + \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s' \right]$$

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} R_{(s,a,s')}^{t+1} + \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} v_{\pi}(s')$$

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} \mathbb{P}_{(s,a)}^{s'} \left[R_{(s,a,s')}^{t+1} + \sum v_{\pi}(s') \right]$$

(iii) Action-value function:

We know,

$$q_{\pi}(s, a) = E_{\pi} [R_t \mid s_t = s, a_t = a], \quad R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$
$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s, a \right] \quad \text{--- (4)}$$

Now, conditioning (4) on the future states s' , we obtain,

$$\Rightarrow q_{\pi}(s, a) = \sum_{s'} E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s, a, s' \right] P(s' \mid s, a)$$

$$\text{denote, } P_{(s,a)}^{s'} = P(s' \mid s, a)$$

Then,

$$\Rightarrow q_{\pi}(s, a) = \sum_{s'} E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s, a, s' \right] P_{(s,a)}^{s'}$$
$$= \sum_{s'} P_{(s,a)}^{s'} E_{\pi} [r_{t+1} \mid s, a, s'] +$$
$$\gamma \sum_{s'} P_{(s,a)}^{s'} E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s, a, s' \right]$$

$$\text{denote, } R_{(s,a,s')}^{t+1} = E_{\pi} [r_{t+1} \mid s, a, s']$$

and observe that, by Markov property,

$$E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s, a, s' \right] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s' \right]$$

$$\Rightarrow q_{\pi}(s, a) = \sum_{s'} P_{(s,a)}^{s'} R_{(s,a,s')}^{t+1} + \gamma \sum_{s'} P_{(s,a)}^{s'} E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s' \right] \quad \text{--- (5)}$$

Now,

conditioning ⑤ on a' (possible actions when on states'),
we have,

$$\begin{aligned} \textcircled{5} \Rightarrow q_{\pi}(s, a) &= \sum_{s'} P_{(s, a)}^{s'} R_{(s, a, s')}^{a+1} + \\ &\quad \gamma \sum_{s'} P_{(s, a)}^{s'} E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s' \right] \\ &= \sum_{s'} P_{(s, a)}^{s'} R_{(s, a, s')}^{a+1} + \\ &\quad \gamma \sum_{s'} P_{(s, a)}^{s'} \sum_{a'} E_{\pi} [R_{t+1} \mid s', a'] P(a' \mid s') \\ &= \sum_{s'} P_{(s, a)}^{s'} R_{(s, a, s')}^{a+1} + \\ &\quad \gamma \sum_{s'} P_{(s, a)}^{s'} \sum_{a'} q_{\pi}(s', a') \pi(s', a') \end{aligned}$$

$$q_{\pi}(s, a) = \sum_{s'} P_{(s, a)}^{s'} \left[R_{(s, a, s')}^{a+1} + \gamma \sum_{a'} q_{\pi}(s', a') \pi(s', a') \right]$$

Ans 2

① (i) state transition table for example 1

Description: A mobile robot has the job of picking cans and placing them on bin. The bot runs on battery. The decision on how to search for a can is made by reinforcement learning agent based on the charge level of the bot's battery. The agent should decide whether the bot should (i) actively search for cans, (ii) wait for someone to bring the cans to it (iii) recharge the battery.

The states are:

$S = \{low, high\}$ (charge levels of bot's battery)

The actions corresponding to states are:

$A(low) = \{search, wait, recharge\}$

$A(high) = \{search, wait\}$

The corresponding rewards are assumed to be:

r_{search} : average reward during search.

r_{wait} : average reward during wait.

-3 : when recharge after running out of battery.

0 : during recharge

State transition table:

s	s'	a	$P(s' s, a)$	$r(a, s, s')$
high	high	search	α	r_{search}
high	low	search	$1-\alpha$	r_{search}
high	high	wait	1	r_{wait}
high	low	wait	0	r_{wait}
low	low	search	β	r_{search}
low	high	search	$1-\beta$	-3
low	low	wait	1	r_{wait}
low	high	wait	0	r_{wait}
low	low	recharge	0	0
low	high	recharge	1	0

(ii) State transition table for example 2

The states are:

$$\mathcal{S} = \{ G, B \}, \quad \begin{array}{l} G = \text{good state} \\ B = \text{bad state} \end{array}$$

The actions corresponding to states are:

$$\mathcal{A}(G) = \{ \text{stay, move} \}$$

$$\mathcal{A}(B) = \{ \text{stay, move} \}$$

The corresponding rewards are assumed to be:

$$r_{GG} = 3 \quad (\text{for staying in state G})$$

$$r_{GB} = -1 \quad (\text{for going to state B but attempting to stay at state G})$$

$$r_{BB} = -1 \quad (\text{for staying in state B})$$

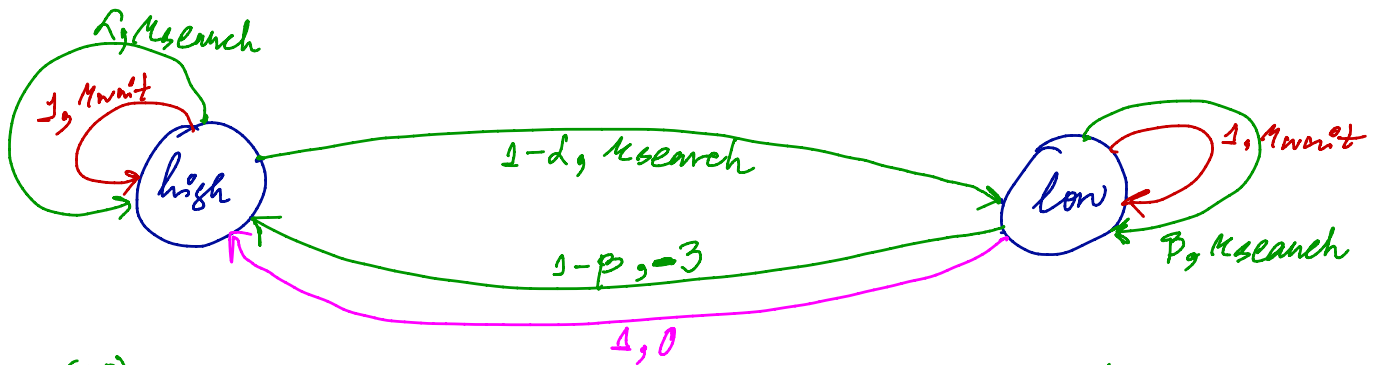
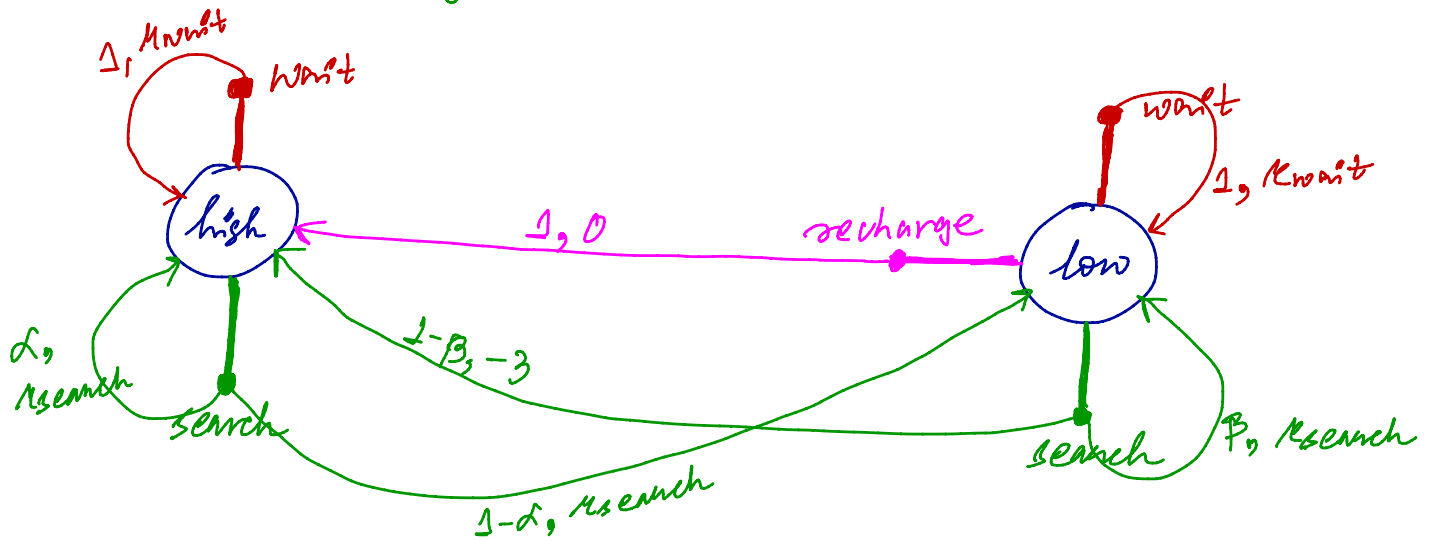
$$r_{BG} = 3 \quad (\text{for moving to state G from state B})$$

$$r_{GB} = -1 \quad (\text{for moving to state B from state G})$$

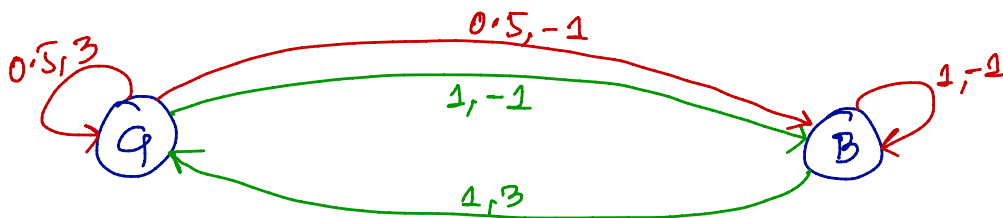
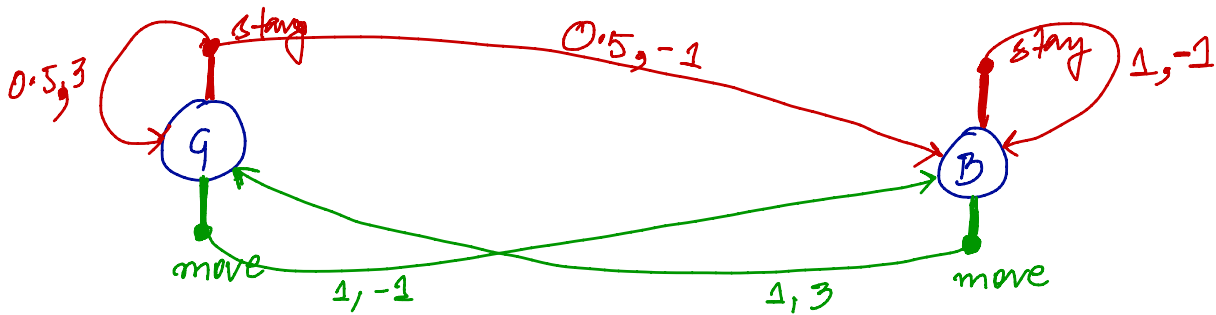
State transition table:

s	s'	a	$p(s' s,a)$	$r(a,s,s')$
G	G	stay	0.5	3
G	B	stay	0.5	-1
G	G	move	0	0
G	B	move	1	-1
B	B	stay	1	-1
B	G	stay	0	0
B	B	move	0	0
B	G	move	1	3

⑥ (i) State-space diagram for example 1



(ii) State-space diagram for example 2



② After implementing the value iteration algorithm for the can collecting robot example, the following results are obtained:
Considering $\Delta = 100$, $E = 1 \times 10^{-9}$,

optimal values

$$v^*(S=High) = 3.176470587$$

$$v^*(S=Low) = 1.99999999$$

Iteration needed = 31

policy

Policy: State = High, Action = Search

Policy: State = Low, Action = Wait

Python file: EE18MTECH11030_Question_2-C_Final.ipynb

③ After implementing the value iteration algorithm for the second MDP example, the following results are obtained:
Considering $\Delta = 100$, $E = 5 \times 10^{-9}$,

optimal values

$$v^*(S=High) = 2.799999999$$

$$v^*(S=Low) = 4.39999999$$

Iteration needed = 30

policy

Policy: State = Good, Action = Stay

Policy: State = Bad, Action = Move

Python file: EE18MTECH11030_Question_2-d_Final.ipynb

Ans: 1)

Deriving the Bellman's equation for state value function as derived in class.

Now,

$$\begin{aligned} v_{\pi}(s) &= E_{\pi}[G_t | s_t] \\ &= E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | s_t\right] \\ &= E_{\pi}\left[R_{t+1} + \sum_{k=0}^{\infty} \gamma^{k+1} R_{t+k+2} | s_t\right] \\ &= E_{\pi}\left[r + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | s_t\right] \\ &= E_{\pi}[r + \gamma G_{t+1} | s_t] \\ &= E_{\pi}[r | s_t] + \gamma E_{\pi}[G_{t+1} | s_t] \quad \text{--- (1)} \end{aligned}$$

Now,

$$E_{\pi}[E_{\pi}[G_{t+1} | s_{t+1}] | s_t] = E_{\pi}[G_{t+1} | s_t]$$

Proof: Let $G_{t+1} = g'$, $s_{t+1} = s'$, $G_t = g$, $s_t = s$

then,

$$\begin{aligned} E_{\pi}[E_{\pi}[g' | s'] | s] &= E_{\pi}[E_{\pi}[g' | s', s]] \\ &= E_{\pi}\left[\sum_{g'} g' p(g' | s', s)\right] \\ &= \sum_{s'} \sum_{g'} g' p(g' | s', s) p(s' | s) \\ &= \sum_{s'} \sum_g g' \times \frac{p(g', s', s)}{p(s', s)} \cdot \frac{p(s' | s)}{p(s)} \\ &= \sum_{s'} \sum_g g' \frac{p(g', s', s)}{p(s)} \end{aligned}$$

$$= \sum_{s'} \sum_{g'} g' \frac{p(g', s', s)}{p(s)}$$

$$= \sum_{g'} g' \frac{p(g', s)}{p(s)}$$

$$= \sum_{g'} g' p(g' | s)$$

$$= E_{\pi}[g' | s]$$

$$= E_{\pi}[G_{t+1} | s_t]$$

$$\begin{aligned} \circ \circ (1) \Rightarrow v_{\pi}(s) &= E_{\pi}[r | s_t] + \gamma E_{\pi}[G_{t+1} | s_t] \\ &= E_{\pi}[r | s_t] + \gamma E_{\pi}[E_{\pi}[G_{t+1} | s_{t+1}] | s_t] \\ &= E_{\pi}[r + \gamma E_{\pi}[G_{t+1} | s_{t+1}] | s_t] \end{aligned}$$

$$v_{\pi}(s) = E_{\pi}[r + \gamma v_{\pi}(s') | s_t]$$