

Mid Term Exam-II Solutions

1 a) $p(\text{head}) = p$

$p(\text{head second time on 5th toss})$

$$p(\text{head on 5th toss}) = p \quad \text{--- (1)}$$

For head to occur second time on 5th toss,
it should have occurred ^{only} once on 1st 4 times.

$p(\text{one head in first 4 tosses})$

$$= 4C_1 p(1-p)^3 \quad \text{--- (2)}$$

↑ ↑ ↑ tail 3 times

Choosing 1 out
of 4

$$\therefore p(\text{head second times on 5th toss}) = 4C_1 p(1-p)^3 * p$$

[from (1) & (2)]

$$= \boxed{4p^2(1-p)^3}$$

16) Dice (fair) rolled repeatedly till atleast 1: 5^{and} One-6

$$P(5 \text{ rolls}) = ?$$

The toll stops if 5^{and} 6 are rolled once.

\therefore 5(or) 6 should be rolled on 5th roll
and 6(or) 5 should be rolled in 1st 4 rolls
respectively

$$\begin{array}{c} \text{---} \\ | \quad \quad \quad | \uparrow \end{array} P(5 \text{ or } 6) = \frac{2}{6}$$

$$\text{One place can be } 5 \text{ or } 6 = \frac{1}{6}$$

$$\text{Choosing 1 place out of 4} = 4C_1$$

Other 3 places can have, 5 values except one in 5th place.

$$= \frac{5}{6}$$

$$\therefore p(5 \text{ rolls}) = 4C_1 \left(\frac{1}{6}\right) \left(\frac{5}{6}\right)^3 \times \frac{2}{6}$$

$$= 4 \times (0.167) \times (0.83)^3 \times (0.33)$$

$$= 0.126.$$

2) 5 Arms

$$\text{Arm 1 } p_1 = \frac{1}{4} \quad R_{t+1} = \begin{cases} 1 - \frac{1}{4} = \frac{3}{4} & \frac{1}{4} : \text{Prob} \\ 0 & \frac{3}{4} \end{cases}$$

$$\text{Arm 2 } p_2 = \frac{1}{3} \quad R_{t+1} = \begin{cases} 1 - \frac{1}{3} = \frac{2}{3} & \frac{1}{3} : \text{Prob} \\ 0 & \frac{2}{3} \end{cases}$$

$$\text{Arm 3 } p_3 = \frac{1}{2} \quad R_{t+1} = \begin{cases} 1 - \frac{1}{2} = \frac{1}{2} & \frac{1}{2} \text{ Prob} \\ 0 & \frac{1}{2} \end{cases}$$

$$\text{Arm 4 } p_4 = \frac{2}{3} \quad R_{t+1} = \begin{cases} 1 - \frac{2}{3} = \frac{1}{3} & \frac{2}{3} \text{ Prob} \\ 0 & \frac{1}{3} \end{cases}$$

$$\text{Arm 5} \quad P_5 = \frac{3}{4}$$

$$R_{t+1} = \begin{cases} 1 - \frac{3}{4} = \frac{1}{4} \\ 0 \end{cases}$$

$\frac{3}{4}$ prob
 $\frac{1}{4}$

a) $E(\text{Arm 1}) = \frac{3}{4} \times \frac{1}{4} + 0 = \frac{3}{16} = 0.19$

$$E(\text{Arm 2}) = \frac{2}{3} \times \frac{1}{3} = \frac{2}{9} = 0.22$$

$$E(\text{Arm 3}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = 0.25$$

$$E(\text{Arm 4}) = \frac{1}{3} \times \frac{2}{3} = \frac{2}{9} = 0.22$$

$$E(\text{Arm 5}) = \frac{1}{4} \times \frac{3}{4} = \frac{3}{16} = 0.19$$

Arm 3 gives maximum expected reward

b) Maximize $E[R_{t+1}^i]$, find value p_i

$$E[R_{t+1}^i] = p_i(1-p_i)$$

$$\therefore R_{t+1}^i = \begin{cases} 1-p_i & \text{with prob } p_i \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial p_i} E[R_{t+1}^i] = \frac{d}{dp_i} (p_i - p_i^2) = 1 - 2p_i = 0$$
$$\Rightarrow p_i = \frac{1}{2}$$

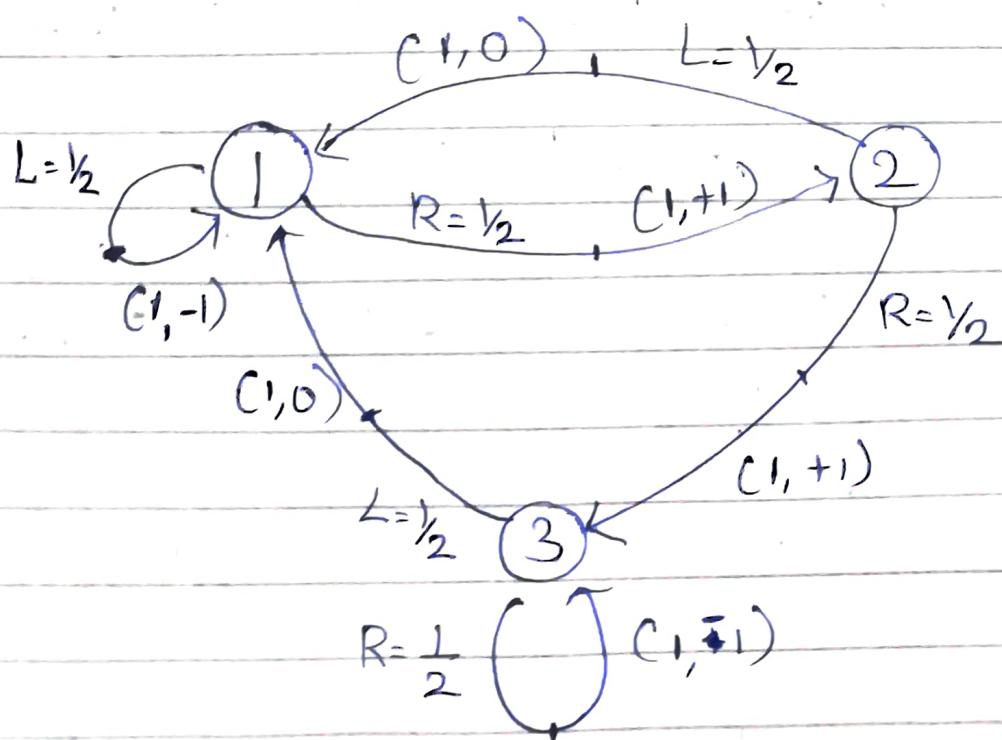
$$\text{Also, } \frac{d^2}{dp_i^2} E[R_{t+1}^i] = -2 < 0 \Rightarrow \text{Maximum}$$

$\therefore p_i = \frac{1}{2}$ that maximises expected reward

3) a)

Current State (s)	Action (a)	Reward (r)	Next state (s')
1	Left	-1	1
1	Right	+1	2
2	Left	0	1
2	Right	+1	3
3	Left	0	1
3	Right	-1	3

b)



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

$$V_{\pi}(1) = \frac{1}{2} [\text{Value from action Left} + \gamma V_{\pi}(1)] + \frac{1}{2} [\text{Value from action Right}] \quad \gamma = 0.9$$

$$V_{\pi}(1) = \frac{1}{2} [(-1 + 0.9 V_{\pi}(1)) + (-1 + 0.9 V_{\pi}(2))]$$

$$2V_{\pi}(1) = 0.9 V_{\pi}(1) + 0.9 V_{\pi}(2)$$

$$2.22 V_{\pi}(1) = V_{\pi}(1) + V_{\pi}(2)$$

$$1.22 V_{\pi}(1) = V_{\pi}(2) \rightarrow ①$$

$$V_{\pi}(2) = \frac{1}{2} [\text{Value from action Left} + \text{Value from action Right}]$$

$$V_{\pi}(2) = \frac{1}{2} [(0 + 0.9 V_{\pi}(1)) + (+1 + 0.9 V_{\pi}(3))]$$

$$2V_{\pi}(2) = 0.9V_{\pi}(1) + 0.9V_{\pi}(3) + 1$$

From(1), $2 \cdot 44 V_{\pi}(1) - 0.9V_{\pi}(1) = 0.9V_{\pi}(3) + 1$

$$(1.54V_{\pi}(1) + 1) / 0.9 = V_{\pi}(3)$$

$$[1.71 V_{\pi}(1) - 1.11 = V_{\pi}(3)] \rightarrow (2)$$

$$V_{\pi}(3) = \frac{1}{2} [\text{Value from action Left} + \text{Value from action Right}]$$

$$V_{\pi}(3) = \frac{1}{2} [(0 + 0.9V_{\pi}(1)) + (-1 + 0.9V_{\pi}(3))]$$

$$2V_{\pi}(3) = 0.9V_{\pi}(1) + 0.9V_{\pi}(3) + 1$$

$$1.1V_{\pi}(3) = 0.9V_{\pi}(1) + 1$$

From(2), $1.1(1.71V_{\pi}(1) - 1.11) = 0.9V_{\pi}(1) + 1$

$$1.881V_{\pi}(1) - 1.221 = 0.9V_{\pi}(1) + 1$$

$$V_{\pi}(1) = 0.221$$

$$\text{In (1)} \quad V_{\pi}(2) = 0.270$$

$$\text{In (2)} \quad V_{\pi}(3) = -0.732$$

$$4) a) Q_{\pi}(s, a) = E_{\pi} [G_t \mid S_t = s, A_t = a]$$

$$\tilde{Q}_{\pi}(s, a) = E_{\pi} [R_{(t+1)} + \gamma R_{(t+2)} + \dots \mid S_t = s, A_t = a]$$

$$\tilde{Q}_{\pi}(s, a) = E_{\pi} [R_{(t+1)} + C + \gamma [R_{(t+2)} + C] + \dots \mid S_t = s, A_t = a]$$

$$\tilde{Q}_{\pi}(s, a) = E_{\pi} [(R_{(t+1)} + \gamma R_{(t+2)} + \gamma^2 R_{(t+3)}) + \dots$$

$$+ C + \gamma [C + \gamma^2 C + \dots] \mid S_t = s, A_t = a]$$

$$\tilde{Q}_{\pi}(s, a) = E_{\pi} \left[\frac{(G_t + C)}{1-\gamma} \right] \mid S_t = s, A_t = a$$

$$\therefore \gamma \in (0, 1)$$

$$\tilde{Q}_\pi(s, a) = Q_\pi(s, a) + \frac{c}{1-\gamma} \quad [\because \frac{c}{1-\gamma} \text{ is constant}]$$

$\therefore Q$ value is changed by $\left(\frac{+c}{1-\gamma} \right)$

4b) If ϵ -greedy used to update policy,

Actions not optimal under π (suggested) will be chosen.
These actions' rewards will not be translated
by ' $+c$ ' as they are not suggested by π .

Thus, actions obtained will differ between two cases, due to partial advantage to selected actions.

4c) If task is discounted and episodic

$$\tilde{Q}(s, a) = E \left(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t} R_T \right)$$

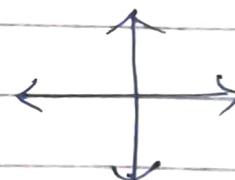
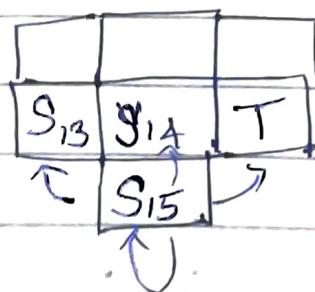
$$+ C + \gamma C + \gamma^2 C + \dots + \gamma^{T-1} C \quad |$$

$$s_t = s, a_t = a$$

$$\tilde{Q}(s, a) = E \left[G_t + \underbrace{C(1-\gamma^T)}_{1-\gamma} \right]$$

$$\tilde{Q}(s, a) = \tilde{Q}(s, a) + \frac{C(1-\gamma^T)}{(1-\gamma)}$$

(5)



Equiprobable

$R_t = -1$ if transitions

$\gamma = 1$

5a) $V_{II}(15) = \frac{1}{4} [$ Value from action Left +
" " " " Top +
" " " " Right +
" " " " Down $]$

$$V_{II}(15) = \frac{1}{4} [-1 + \delta(-20) -1 + \delta(-14) + \delta(0) -1 \\ -1 + \delta V_{II}(15)]$$

$$4V_{II}(15) = -4 - 34 + V_{II}(15)$$

$$3V_{II}(15) = -38 \Rightarrow V_{II}(15) = -12.67$$

5b) If dynamics of S_{14} change, other states values will also be affected.

For simplicity, we consider only value changes of state: S_4 & S_5

	-18	
-20	S_{14}	0
	S_{15}	

$$V(S_{14}) = \frac{1}{4} [-1 - 18 + 0 - 1 - 1 + V(S_{15}) - 1 - 20]$$

$$4V(S_{14}) = V(S_{15}) - 4 - 38$$

$$\boxed{4V(S_{14}) - V(S_{15}) = -42} \rightarrow \textcircled{1}$$

$$V(S_{15}) = \frac{1}{4} [-1 + V(S_{14}) - 1 + 0 + V(S_{15}) - 1 - 1 - 20]$$

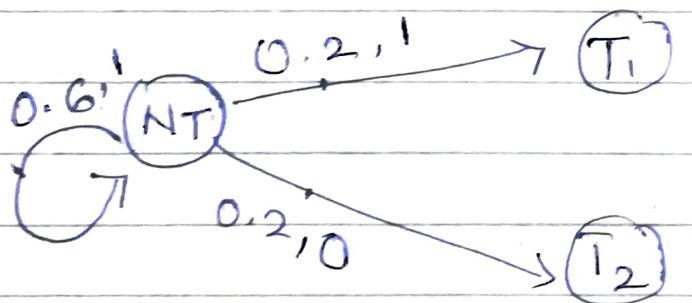
$$4V(S_{15}) = V(S_{14}) + V(S_{15}) - 4 - 20$$

$$[3V(S_{15}) - V(S_{14}) = -24] \rightarrow (2)$$

Solving (1), (2)

$$V(S_{14}) = -13.64 ; V(S_{15}) = -12.55$$

(6)



$$\text{Epi.1} = NT \rightarrow \dots \rightarrow T_1 \quad R=10$$

$$\text{Epi.2} = NT \rightarrow \dots \rightarrow T_2 \quad R=4$$

a) Episode 1.



Episode 2,

NT $\xrightarrow{1}$, NT $\xrightarrow{1}$, NT $\xrightarrow{1}$, NT $\xrightarrow{1}$, NT $\xrightarrow{0}$, T2

b) First Visit

$\frac{G_{\text{Episode 1}} + G_{\text{Episode 2}}}{2}$

Every Visit G_{ij} : i^{th} Episode j^{th} Occurrence of NT

$$(G_{1,1} + G_{1,2} + G_{1,3} + G_{1,4} + G_{1,5} + G_{1,6} + G_{1,7} + G_{1,8} + G_{1,9} + G_{1,10} \\ + G_{2,1} + G_{2,2} + G_{2,3} + G_{2,4} + G_{2,5}) / 15$$

c) First Visit: 2 Estimates ; Every Visit - 15 estimates

d) i) First visit Estimator : $(G_{\text{epi1}} + G_{\text{epi2}}) / 2$

$$= (10 + 4) / 2 = 7$$

ii) Every Visit Estimator

$$(10 + 9 + \dots + 1) + (4 + 3 + 2 + 1) / 15$$

$$= \left(\frac{5}{7} \times 1 \right) + \frac{2}{2} \Bigg) \Bigg|_{15} = 65/15 = 4,33$$

7) G_i has cweight e^{w_i}

$$V_n = \frac{\sum_{k=1}^n e^{w_k} G_k}{\sum_{k=1}^n e^{w_k}}, n \geq 1$$

$$a) V_{n+1} = \frac{\sum_{k=1}^{n+1} e^{w_k} G_k}{\sum_{k=1}^{n+1} e^{w_k}}$$

$$= \frac{\sum_{k=1}^n e^{w_k} G_k + e^{w_{n+1}} G_{n+1}}{\sum_{k=1}^{n+1} e^{w_k}}$$

$$V_{n+1} = \frac{\sum_{k=1}^n e^{lk} G_k}{\sum_{k=1}^{n+1} e^{lk}} + \frac{e^{ln+1} G_{n+1}}{\sum_{k=1}^{n+1} e^{lk}}$$

$$V_{n+1} = \frac{\sum_{k=1}^n e^{lk} G_k}{\sum_{k=1}^{n+1} e^{lk}} * \frac{\sum_{k=1}^n e^{lk}}{\sum_{k=1}^{n+1} e^{lk}} + \frac{e^{ln+1} G_{n+1}}{\sum_{k=1}^{n+1} e^{lk}}$$

$$V_{n+1} = V_n \frac{\sum_{k=1}^{n+1} e^{lk} - e^{ln+1}}{\sum_{k=1}^{n+1} e^{lk}} + \frac{e^{ln+1} G_{n+1}}{\sum_{k=1}^{n+1} e^{lk}}$$

$$V_{n+1} = V_n \left[1 - \frac{e^{ln+1}}{\sum_{k=1}^{n+1} e^{lk}} \right] + \frac{e^{ln+1} G_{n+1}}{\sum_{k=1}^{n+1} e^{lk}}$$

$$V_{n+1} = V_n + \frac{e^{h_{n+1}}}{\sum_{k=1}^{n+1} e^{h_k}} [G_{n+1} - V_n]$$

b) Off-policy MC control Algorithm

Initialize $\pi(s, a)$, $\alpha(s, a)$

$\rightarrow Q(s, a) \in \mathbb{R}$ arbitrarily

$\rightarrow C(s, a) \leftarrow 0$

$\rightarrow \pi(s) \leftarrow \arg\max_a Q(s, a)$

Loop for each episode:

$b \leftarrow$ any soft policy

Generate an episode (length: $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, A_T$)

$G \leftarrow 0$

$H \leftarrow 0$

Loop for each step in episode $t = T-1, T-2, \dots, 0$

$G \leftarrow \gamma G + R_{t+1}$

$C(S_t, A_t) \leftarrow C(S_t, A_t) + H$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{e^W}{C(s_t, a_t)} [G - Q(s_t, a_t)]$$

$$\pi(s_t) \leftarrow \arg \max_a Q(s_t, a)$$

If $A_t \neq \pi(s_t)$ then exit innerloop (Go to next episode)

$$W \leftarrow W + \log \left(\frac{\pi(A_t | s_t)}{b(A_t | s_t)} \right)$$

c)

$$h_i = e_{t_i:T(t)-1}$$

$$h_i = \ln e_{t_i:T(t)-1}$$

$$h_i = \ln \left[\prod_{k=t_i}^{T(t)-1} \frac{\pi(A_k | s_k)}{b(A_k | s_k)} \right]$$

$$h_i = \sum_{k=t_i}^{T(t)-1} \ln \left(\frac{\pi(A_k | s_k)}{b(A_k | s_k)} \right)$$