

RL A2

1) arms = 1, 2, ..., 10

$$P(A=a) = \begin{cases} 2p, & a=2, 4, 6, 8, 10 \\ p, & a=1, 3, 5, 7, 9 \end{cases} = \begin{cases} \frac{2}{15}, & a=2, 4, 6, 10 \\ \frac{1}{15}, & a=1, 3, 5, 7, 9 \end{cases}$$

$$5p + (2p) \times 5 = 1 \Rightarrow p = \frac{1}{15}$$

Arm i: $R_i \sim N(\mu=i, \sigma=1)$

No of steps = 10

Expected sum reward

$$E[R^1 + R^2 + \dots + R^{10}]$$

 $R^j \rightarrow$ reward at time
step j $j=1, \dots, 10$

$$E[R_i] = \mu = i$$

$$E[R] = E\left[\sum_{j=1}^{10} R^j\right] = \sum_{j=1}^{10} E[R^j]$$

$$E[R^j] = \sum_{a \in A} r_a P(\text{arm}) \quad r = E[R]$$

$$= \frac{1}{15} (1+3+5+7+9)$$

$$\frac{2}{15} (2+4+6+8+10) = 1.66 + 4$$

$$= 5.66$$

$$E\left[\sum_{j=1}^{10} R^j\right] = 10 \times 5.66 = 56.66$$

(2.) $K = 10$ (no of arms)
 arms - 1, 2, 4, 5, 7, 9, 10 reward = 0
 probability 0.5

~~$R(\text{Reward} = 0 / \text{arm} = 1, 2, 4, 5, 7, 9, 10) = 0.5$~~

arm = 3, 6, 8 reward = 0 $P = 0.3$

$R = 0.2 P = 0.3$

$R = 1 P = 0.4$

$$P(r_i | a_i) = \begin{cases} 0.3 & r=0, a_1 = 3, 6, 8 \\ 0.3 & r=0.2, a_1 = 3, 6, 8 \\ 0.4 & r=1, a_1 = 3, 6, 8 \end{cases}$$

$$P(r_i | a_2) = \begin{cases} 0.5 & r=0 \\ 0.5 & r=1 \end{cases}$$

maximise (E[Reward])

$$\max_{\{a_1, a_2, \dots, a_m\}} E\left[\sum_{i=1}^m R_i\right]$$

$$E[R] = 0 \times 0.3 + 0.2 \times 0.3 + 0.4 \times 1$$

for arm $\{3, 6, 8\}$ = $0.06 + 0.4 = 0.46$

$$E[R] = 0.5 \times 0 + 0.5 \times 1$$

for arm $\{1, 2, 4, 5, 7, 9, 10\}$ = 0.5

Optimal stochastic policies

$$(1) P(A=a) = \begin{cases} \frac{1}{7} & \text{arm } 1, 2, 4, 5, 7, 9, 10 \\ 0 & \text{arm } 3, 6, 8 \end{cases}$$

$$(2) P(A=a) = \begin{cases} \frac{1}{11} & a = 1, 2, 4 \\ \frac{2}{11} & a = 5, 7, 9, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

$$3p + 4(2p) = 1$$

$$3p + 8p = 1$$

$$11p = 1$$

$$(3) P(A=a) = \begin{cases} \frac{1}{11} & a = 1, 5, 7 \\ \frac{2}{11} & a = 2, 4, 9, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

$$(4) P(A=a) = \begin{cases} \frac{1}{11} & a = 9, 10, 2 \\ \frac{2}{11} & a = 1, 5, 7, 4 \\ 0 & a = 3, 6, 8 \end{cases}$$

$$(5) P(A=a) = \begin{cases} 3p = \frac{3}{17} & a = 9, 10 \\ 4p = \frac{4}{17} & a = 5, 7 \\ p = \frac{1}{17} & a = 2, 4, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

$$6p + 8p + 3p = 1$$

$$p = \frac{1}{17}$$

$$(6) \quad P(A=a) = \begin{cases} 3/17 & a = 5, 7 \\ 4/17 & a = 9, 1 \\ 1/17 & a = 2, 4, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

(Q3.) ϵ -greedy explore only non greedy options

$$k=3 \quad R=\{0, 1\}$$

$$Q_1(a) \quad a \in \{1, 2, 3\}$$

$$Q_1(1) = 1$$

$$Q_1(2) = 4$$

$$Q_1(3) = 9$$

$$Q_t(a), A_t, R_t \text{ for } t = 1, 2, 3, 4, 5, 6$$

explore - odd time

$$t = 1, 3, 5$$

exploit - even time

$$t = 2, 4, 6$$

Example Seq:

$t=1$ (explore)

$$A_1 = 1, R_1 = 0$$

$$Q_2(1) = 0$$

$$Q_2(2) = 4$$

$$Q_2(3) = 9$$

$t=2$ (exploit)

$$A_2 = 3, R_2 = 0$$

$$Q_3(1) = 0$$

$$Q_3(2) = 4$$

$$Q_3(3) = 0$$

$t=3$ (explore)

$$A_3 = 1, R_3 = 1$$

$$Q_4(1) = 1$$

$$Q_4(2) = 4$$

$$Q_4(3) = 0$$

$t=4$ (exploit)

$$A_4 = 2, R_4 = 0$$

$$Q_5(1) = 1$$

$$Q_5(2) = 0$$

$$Q_5(3) = 0$$

$t=5$ (explore)

$$A_5 = 3, R_5 = 1$$

$$Q_6(1) = 1$$

$$Q_6(2) = 0$$

$$Q_6(3) = 1$$

$t=6$ (exploit)

$$A_6 = 1, R_6 = 1$$

$$Q_7(1) = 1$$

$$Q_7(2) = 0$$

$$Q_7(3) = 1$$

$$\textcircled{6} \quad P(A=a) = \begin{cases} 3/17 & a = 5, 7 \\ 4/17 & a = 9, 1 \\ 1/17 & a = 2, 4, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

Q3.) ϵ -greedy explore only non greedy options

 $k=3$

$$R = \{0, 1\}$$

$$Q_1(a) \quad a \in \{1, 2, 3\}$$

$$Q_1(1) = 1$$

$$Q_1(2) = 4$$

$$Q_1(3) = 9$$

$Q_t(a), A_t, R_t$ for
 $t = 1, 2, 3, 4, 5, 6$

explore - odd time

$$t = 1, 3, 5$$

exploit - even time

$$t = 2, 4, 6$$

Example seq:

$t=1$ (explore)

$$A_1 = 1, R_1 = 0$$

$$Q_2(1) = 0$$

$$Q_2(2) = 4$$

$$Q_2(3) = 9$$

$t=2$ (exploit)

$$A_2 = 3, R_2 = 0$$

$$Q_3(1) = 0$$

$$Q_3(2) = 4$$

$$Q_3(3) = 0$$

$t=3$ (explore)

$$A_3 = 1, R_3 = 1$$

$$Q_4(1) = 1$$

$$Q_4(2) = 4$$

$$Q_4(3) = 0$$

$t=4$ (exploit)

$$A_4 = 2, R_4 = 0$$

$$Q_5(1) = 1$$

$$Q_5(2) = 0$$

$$Q_5(3) = 0$$

$t=5$ (explore)

$$A_5 = 3, R_5 = 1$$

$$Q_6(1) = 1$$

$$Q_6(2) = 0$$

$$Q_6(3) = 1$$

$t=6$ (exploit)

$$A_6 = 1, R_6 = 1$$

$$Q_7(1) = 1$$

$$Q_7(2) = 0$$

$$Q_7(3) = 1$$

$$\textcircled{6} \quad P(A=a) = \begin{cases} 3/17 & a = 5, 7 \\ 4/17 & a = 9, 1 \\ 1/17 & a = 2, 4, 10 \\ 0 & a = 3, 6, 8 \end{cases}$$

(3) ϵ -greedy explore only non greedy options

$$k=3 \quad R=\{0, 1\}$$

$$Q_1(a) \quad a \in \{1, 2, 3\}$$

$$Q_1(1) = 1$$

$$Q_1(2) = 4$$

$$Q_1(3) = 9$$

$$Q_t(a), A_t, R_t \text{ for}$$

$$t = 1, 2, 3, 4, 5, 6$$

explore - odd time

$$t = 1, 3, 5$$

exploit - even time

$$t = 2, 4, 6$$

Example Seq:

$t=1$ (explore)

$$A_1 = 1, R_1 = 0$$

$$Q_2(1) = 0$$

$$Q_2(2) = 4$$

$$Q_2(3) = 9$$

$t=2$ (exploit)

$$A_2 = 3, R_2 = 0$$

$$Q_3(1) = 0$$

$$Q_3(2) = 4$$

$$Q_3(3) = 0$$

$t=3$ (explore)

$$A_3 = 1, R_3 = 1$$

$$Q_4(1) = 1$$

$$Q_4(2) = 4$$

$$Q_4(3) = 0$$

$t=4$ (exploit)

$$A_4 = 2, R_4 = 0$$

$$Q_5(1) = 1$$

$$Q_5(2) = 0$$

$$Q_5(3) = 0$$

$t=5$ (explore)

$$A_5 = 3, R_5 = 1$$

$$Q_6(1) = 1$$

$$Q_6(2) = 0$$

$$Q_6(3) = 1$$

$t=6$ (exploit)

$$A_6 = 1, R_6 = 1$$

$$Q_7(1) = 1$$

$$Q_7(2) = 0$$

$$Q_7(3) = 0$$

Q4) episodic \rightarrow repeated attempts to balance the pole
 Reward - +1 failure did not occur
 + discounting
 -1 upon failure
 continuous
 return at each time

Q4) $p(s', \tau | s, a)$ $s, a, s', \tau \in p(s', \tau | s, a)$
 3.4 and a row for every 4-tuple for which
 $p(s', \tau | s, a) > 0$

Since the rewards do not have a probability distribution, they are constant depending on the state

$$\therefore p(s', \tau | s, a) = p(s' | s, a)$$

hence table will be same.

Q5) 3.15 3.14

Bellman eq :

$$V_{\pi}(s) = E_{\pi} [G_t | s_t = s]$$

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

Entire state of +0.7

$$\begin{array}{c} 2.3 \\ 0.7 \\ 0.4 \\ -0.4 \end{array}$$

$\forall s \in S$

for $s' = 2, 3$ Action \uparrow

~~$v_\pi(s)$~~

Calculating

for all
actions

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

$$\text{Assuming } P(A=a) = \frac{1}{4}$$

1 state's corresponding to 1 action

Reward -1 off the grid
or 0 others

from A \rightarrow reward +10

B \rightarrow +5

Given: $\gamma = 0.9$

To find $v_\pi(s)$ of centre state

$$a = \uparrow \quad \frac{1}{4} (0 + 0.9(2.3))$$

$$a = \downarrow \quad \frac{1}{4} (0 + 0.9 \times (-0.4)) +$$

$$a = \leftarrow \quad \frac{1}{4} (0 + 0.9 \times 0.7)$$

$$a = \rightarrow \quad \frac{1}{4} (0 + 0.9 \times 0.4)$$

$$0.9 \times 0.25 [2.3 - 0.4 + 0.7 + 0.4]$$

$$= 0.9 \times 0.25 \times 3.0$$

$$= 0.675 \approx 0.7$$

3.15

- ve reward → goal
- ve running into edge
- 0 rest

→ adding constant c to reward

V_c to all value of state

does not affect the relative value of any state

→ V_c in terms of c & γ

$$3.8 \quad G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

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

$$G_t' = (R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \dots$$

(adding constant c to all reward)

$$G_t' = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots +$$

$$c + \gamma c + \gamma^2 c + \dots$$

c - constant $\Rightarrow c + \gamma c + \gamma^2 c + \dots$ is a GP

γ - constant

$$\text{Sum of GP} = \frac{a}{1-\gamma} = \frac{c}{1-\gamma} \quad a=c \quad \gamma=\gamma$$

$\underbrace{\qquad\qquad\qquad}_{\text{constant}}$

$$V_c = \frac{c}{1-\gamma} \quad \gamma < 1$$

3.16

For episodic task

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

$$G_t' = (R_{t+1} + c) + (R_{t+2} + c) + \dots + (R_T + c)$$

adding constant c to all rewards

$$G_t' = R_{t+1} + R_{t+2} + \dots + R_T + [T - (t+1)]c$$

There is a constant added for the time intervals of the episodic task

Since now the value depends on the time interval, the task is changed as the return is incentivized to add additional delay in completing the task.

Q8.) equation for V_* in terms of q_*

$$V_*(s) = \max_{\pi} V_{\pi}(s) = \max_{a \in A} q_*(s, a)$$

$$V_*(s) = \max_{a \in A} q_{aa} \sum_{s', \gamma} p(s', \gamma | s, a)$$

$$[\gamma + \gamma V_*(s')]$$

Q(11) $\text{MDP} \rightarrow \text{PMF}$

R_{t+1} depends on state S_t & Action A_t

Is R_{t+2} dependent on S_t & A_t ? Using PM&Fs used to define MDP. \rightarrow No

Conditional prob of R_{t+2}

According to MDP,

$$p(s', r | s, a) = \Pr \{ S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a \}$$

$$P(R_{t+1} | S_t, A_t) = ? \quad P(R_{t+1} | S_t = s, A_t = a)$$

For R_{t+2}

$$p(s', r | s, a) = \Pr \{ S_{t+2} = s', R_{t+2} = r | S_{t+1} = s, A_{t+1} = a \}$$

In case of MDP, the probability of each possible values for S_t & A_t depends only on the preceding state & action. $\therefore R_{t+2}$ depends only on S_{t+1} & A_{t+1} & and not S_t & A_t .

Q(12) $E[R_{t+2} | S_t = s, A_t = a] = \sum_{r \in R_{t+2}} r p(r | s, a)$

$$P(R_{t+1} | S_t, A_t) = P(R_{t+1} | S_t = s, A_t = a, S_{t+1})$$

$$P(R_{t+2} | S_t, A_t) = \sum_{S' \in S_{t+1}} P(r', s' | S_t, A_t)$$

$$P(R_{t+2} | S_{t+1} = s, A_t = a, S_{t+2} = s')$$

$$\sum_{S'' \in S_{t+2}} P(R_{t+2} = r, S_{t+2} = s'' | S_{t+1} = s', A_t = a)$$

$$\sum_{s'' \in S_{t+2}} p(r_t, s'' | s, a) \times p(s_{t+1} = s' | s_t = s, A_t = a)$$

$$P(A_{t+1} = a | s) \rightarrow \pi(a | s)$$

$$\sum_{s'' \in S_{t+2}} P(R_{t+2} = r_t, S_{t+2} = s'' | S_{t+1} = s, A_{t+1} = a)$$

PMF	$\leq P(R_{t+2} = r_t, S_{t+2} = s'' S_{t+1} = s', A_{t+1} = a')$
$r_t \in R_{t+2}$	
$s_t = s_1$	$P(R_{t+1} = r_t, S_{t+1} = s' S_t = s, A_t = a)$
$a_t \in A_{t+1}$	$\pi(a' s')$

$$E[R_{t+2} | A_t = a, S_t = s] = \sum_{r_t'' \in R_{t+2}} r_t'' P[R_{t+2} = r_t'' | S_t = s, A_t = a]$$

$$= \sum_{\substack{r_t'' \in R_{t+2} \\ s'' \in S_{t+2} \\ r_t'' \in R_{t+1} \\ s' \in S_{t+1} \\ a' \in A_{t+1}}} r_t'' p(r_t, s'' | s', a') p(r_t, s' | s, a)$$

$$Q13.) V_\pi(s) = E[G_t | S_t = s]$$

Bellman eq for $V_\pi(s) \forall s \in S$
 \hookrightarrow in terms of MDP

policy π

$$PMF(p(s', r | s, a))$$

$$V_{\pi}(s) = E \left[G_t \mid S_t = s \right]$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

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

$$V_{\pi}(s) = E_{\pi} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s \right]$$

$$= E_{\pi} [R_{t+1} \mid S_t = s] + \gamma E [R_{t+2} + \gamma R_{t+3} + \dots \mid S_t = s]$$

$$[V_{\pi}(S_{t+1}) = E_{\pi} [R_{t+2} + \gamma R_{t+3} + \dots \mid S_{t+1} = s']]$$

$$V_{\pi}(s) = E_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s]$$

$$= \sum_{r_2, s'} (r_2 + \gamma V_{\pi}(s')) P [R_{t+1} = r_2, S_{t+1} = s' \mid S_t = s]$$

$$= \sum_{r_2, s'} (r_2 + \gamma V_{\pi}(s')) \sum_a P [R_{t+1} = r_2, S_{t+1} = s' \mid S_t = s, A_t = a] \pi(a \mid s)$$

$$V_{\pi}(s) = \sum_{a, r_2, s'} (r_2 + \gamma V_{\pi}(s')) p(r_2, s' \mid s, a) \pi(a \mid s)$$

$$(Q14) R_1 = 2, R_2 = -1, R_3 = 10, R_4 = -3$$

$$\gamma = 0.5$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4}$$

At time step $t=1$: $[t=0]$

$$G_t = 2 + 0.5(-1) + (0.5)^2 10 + (0.5)^3 -3$$

$$G_{t=0} = 1.5 + 2.5 - 0.375$$

At $t=1$

$$G_1 = 2$$

At $t=2$

$$G_2 = R_2 + \gamma G_1$$

At time $t=0$

$$G_{t=0} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4}$$

$$= 13.625$$

At time $t=1$

$$G_{t=1} = R_2 + \gamma R_3 + \gamma^2 R_4 = -1 + 0.5(8.5)$$

$$= 3.25$$

 $t=2$

$$G_{t=2} = R_3 + \gamma R_4 = 10 + 0.5(-3)$$

$$= 8.5$$

 $t=3$

$$G_{t=3} = R_4 = -3$$

$$G_t' = (R_{t+1} + c) + \gamma(R_{t+2} + c) + \gamma^2(R_{t+3} + c) + \dots$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$+ \underbrace{c + \gamma c + \gamma^2 c + \dots}_{\text{GP infinite}}$$

$$c + \gamma c + \gamma^2 c + \dots \rightarrow \text{GP Sum} = \frac{c}{1-\gamma}$$

infinite horizon discounted return

$$G_t' = G_t + \left[\frac{c}{1-\gamma} \right]$$

Q15) $V_*(s) \forall s \in S$

Optimal policy can be found using policy iteration

when $V_{k+1}(s) = V_k(s)$ and policy is optimal at this condition

However since we already have $V_*(s)$

$$\pi_*(s) = \underset{a \in A(s)}{\operatorname{argmax}} E[R_{t+1} + \gamma V_*(s_{t+1}) | s_t = s] \quad A \in \mathbb{A}$$

$\forall s$

$$q_*(s, a)$$

Q16) States: fresh or stale

Action: query or remain silent

$$\begin{aligned} P(\text{fresh} | \text{stale}) &= 0.8 && \left. \begin{array}{l} \text{given} \\ \text{query} \end{array} \right. \\ P(\text{stale} | \text{fresh}) &= 0.1 && \end{aligned}$$

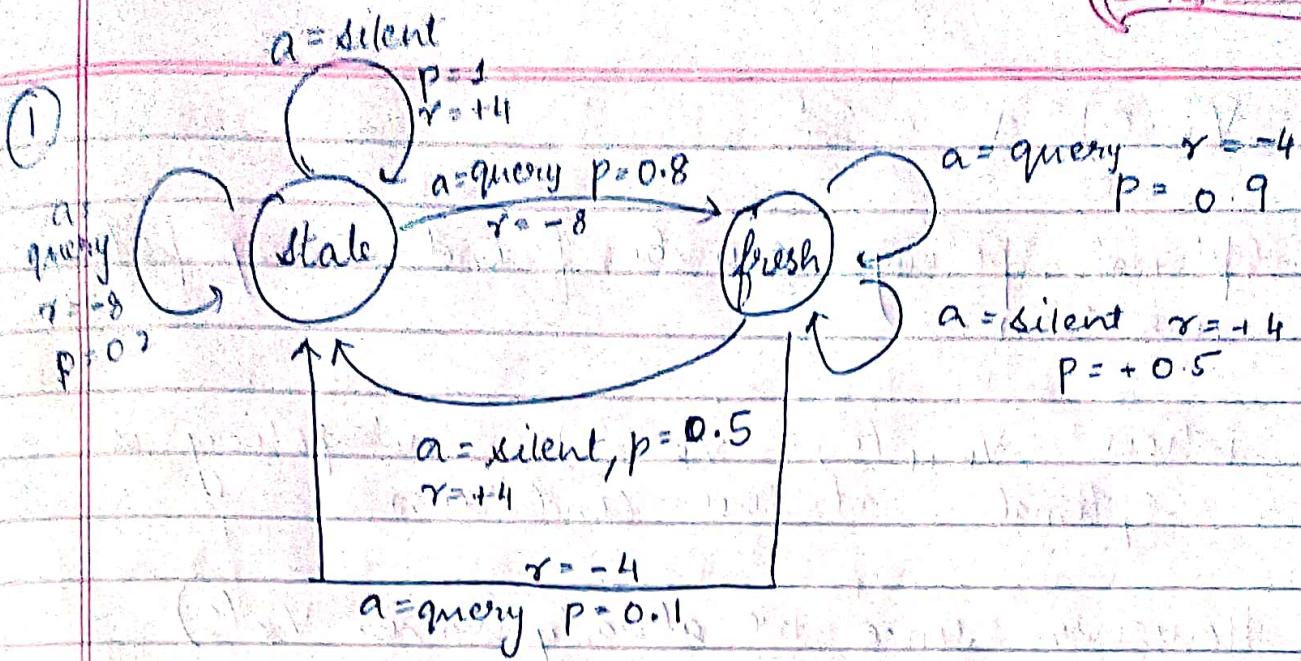
$$P(\text{stale} | \text{stale}) = 1 \quad \& \quad P(\text{stale} | \text{fresh}) = 0.5$$

given silent

Query when current is stale = -8

" " " is fresh = -4

Silent reward = +4



S. s' . a . $p(s'|s, a)$

fresh silent stale +4 0.5

fresh query stale -4 0.1

fresh silent fresh +4 0.5

fresh query fresh -4 0.9

stale silent stale +4 1

stale silent fresh - - [NP]

stale query stale -8 0.2

stale query fresh -8 0.8

For $t = 1$

$$\pi(\text{fresh}) = \max \left\{ \begin{array}{l} q = -4 + 0.5(v_2(\text{fresh}) 0.9 + v_2(\text{stale}) 0.1) \\ s = -4 + 0.5(v_2(\text{fresh}) 0.5 + v_2(\text{stale}) 0.5) \end{array} \right.$$

Action = silent

$$\pi(\text{silent}) = \max \left\{ \begin{array}{l} q = -8 + 0.5(v_2(\text{fresh}) 0.8 + v_2(\text{stale}) 0.2) = -6.5 \\ s = -4 + 0.5(v_2(\text{fresh}) 0.1 + v_2(\text{stale}) 0.1) = 3.5 \end{array} \right.$$

$a = \text{silent}$

$t = 0$

$$\pi(\text{fresh}) = \max \left\{ \begin{array}{l} q = -4 + 0.5(v_1(\text{fresh}) 0.9 + v_2(\text{stale}) 0.1) = -1.68 \\ s = -4 + 0.5(v_1(\text{fresh}) 0.5 + v_1(\text{stab}) 0.5) = 6.0 \end{array} \right.$$

$\Rightarrow \text{Silent}$

$$\pi(\text{stale}) = \max \left\{ \begin{array}{l} q = -8 + 0.5(v_1(\text{stab}) 0.8 + v_1(\text{stale}) 0.2) = -5.75 \\ s = -4 + 0.5(v_1(\text{fresh}) 0.1 + v_1(\text{stale}) 0.1) = -5.7 \end{array} \right.$$

Silent

$$Q17) G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

→ Selecting action

TP: ~~✓~~ policy improvement either improves current policy or current policy is optimal

Let π, π' be any pair of deterministic policies s.t $\forall s \in S$

π - original
 π' - changed

$$\begin{cases} V_{\pi_K}(s) \geq V_{\pi'_K}(s) & \forall s \\ \text{When } V_{\pi_K}(s) = V_{\pi'_K}(s) \Rightarrow \pi_K = \text{Optimal policy} \end{cases}$$

$$\therefore V_{\pi'}(s) \geq V_{\pi}(s) \quad \forall s \quad \text{or} \quad q_{\pi}(s, \underbrace{\pi'(s)}_a) \geq V_{\pi}(s)$$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$= E[R_{t+1} + \gamma V_{\pi}(s_{t+1}) / S_t = s, A_t = \pi'(s)]$$

$$= E_{\pi'}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) / S_t = s]$$

$$\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) / S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma E_{\pi'}[R_{t+2} + \gamma V_{\pi}(s_{t+2}) / S_{t+1}, A_{t+1} = \pi'(s_{t+1})] / S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) / S_t = s]$$

$$\leq V_{\pi'}(s)$$

Policy improvement for all states $s \in S$ have value

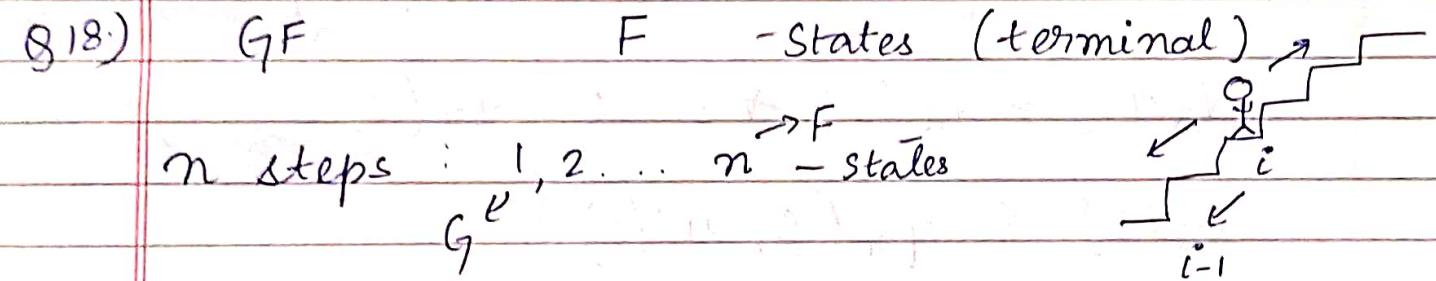
$$V_{\pi'}(s) \geq V_{\pi}(s)$$

In case $V_\pi = V_{\pi'}$ $\forall s \in S$

$$\begin{aligned}V_{\pi'}(s) &= \max_a E[R_{t+1} + \gamma V_{\pi'}(s_{t+1}) / S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi'}(s')]\end{aligned}$$

This same as Bellman optimality
 $\therefore V_{\pi'} \text{ must be } V_\pi$

\therefore both π & π' are optimal



reward 2 $p=0.5$
 0 $p=1-0.5=0.5$

$i \rightarrow i-1$

$1 \leq i \leq n$

reward 1 $\rightarrow G$ reward = 1

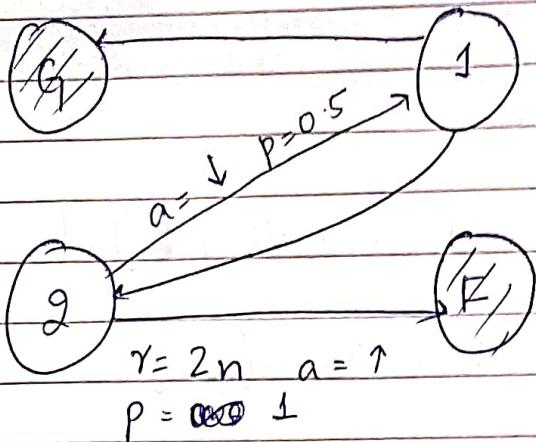
reward -1 \uparrow $i \rightarrow i+1$
 but 2n $n \rightarrow F$

PMF $P(A=a) = \begin{cases} 0.5 & a=\uparrow \\ 0.5 & a=\downarrow \end{cases}$

MDP for $n=2$

$$\gamma = 1$$

States: G 1 2 F



$s \quad a \quad s' \quad r \quad p(s', r | a, s)$

1 \downarrow G 1 1

1 \uparrow 2 -1 1

2 \uparrow F $2n=4$ -1

2 \downarrow 1 2 0.5

2 \downarrow 1 0 0.5

Using current policy

$$\pi(a|s) = \begin{cases} 0.5 & a = \uparrow \\ 0.5 & a = \downarrow \end{cases}$$

Iteration 1 : Policy evaluation

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

$$V_{\pi}(G) = 0$$

$$V_{\pi}(F) = 0$$

$$V_{\pi}(1) = 0.68$$

$$a = \uparrow \rightarrow 0.5 \times 1 [-1 + 1 V_{\pi}(2)]$$

$$+ \\ a = \downarrow \rightarrow 0.5 \times 1 [1 + V_{\pi}(2)]$$

$$V_{\pi}(1) = 0.5 V_{\pi}(2) \quad -(1)$$

$$V_{\pi}(2) =$$

$$a = \uparrow \rightarrow 0.5 \times [1 \times 4] = 2$$

$$+ \\ a = \downarrow \rightarrow 0.5 [0.5 [2 + V_{\pi}(1)] + 0.5 [0 + V_{\pi}(1)]]$$

$$V_{\pi}(2) = 2.5 + 0.5 V_{\pi}(1) \quad -(2)$$

$$\text{Solving eq } 1 \& 2, V_{\pi}(1) = 1.66 \\ V_{\pi}(2) = 3.34$$

Iteration 1: policy iteration improvement

$$\pi(a|s) = \arg\max_{a \in A(s)} q_\pi(s, a)$$

$$= \arg\max_{a \in A(s)} E[R_{t+1} + \gamma v(s_{t+1}) | s_t=s, A_t=a]$$

State $s=1, a=\uparrow \rightarrow -1 + 1 \cdot v_\pi(2)$
 $\rightarrow 2.334$

$s=1, a=\downarrow \rightarrow 1 + 1 \cdot v_\pi(0) = 1$
for $\pi(a|s=1) \rightarrow \uparrow$

State $s=2, a=\uparrow \rightarrow +4 + v_\pi(F) = 4$
 $a=\downarrow \rightarrow 0.5(2 + v_\pi(1)) + 0.5(0 + v_\pi(1))$
 $= 2.667$

 $v_\pi(2) = 1$

Iteration 2: policy evaluation

$$v_\pi(1) \quad a=\uparrow$$

$$= 1 \times 1 [-1 + v_\pi(2)]$$

$$= -1 + v_\pi(2) \quad v_\pi(1) = -1 + 4$$

$$= 3$$

$$v_\pi(2) \quad a=\uparrow = 4$$

Policy improvement

$$\pi(a|s) \quad s=1 \quad a=\uparrow = -1 + v_\pi(2) = 3$$

$$s=1 \quad a=\downarrow = 1$$

$$\pi(a|s=1) = \uparrow$$

$$\text{for } S=2 \quad a=F : 4 + v_{\pi}(F) = 4$$

$$a=J \quad [2 + v_{\pi}(1)] 0.5 + 0.5 [0 + v_{\pi}(1)] \\ = 4$$

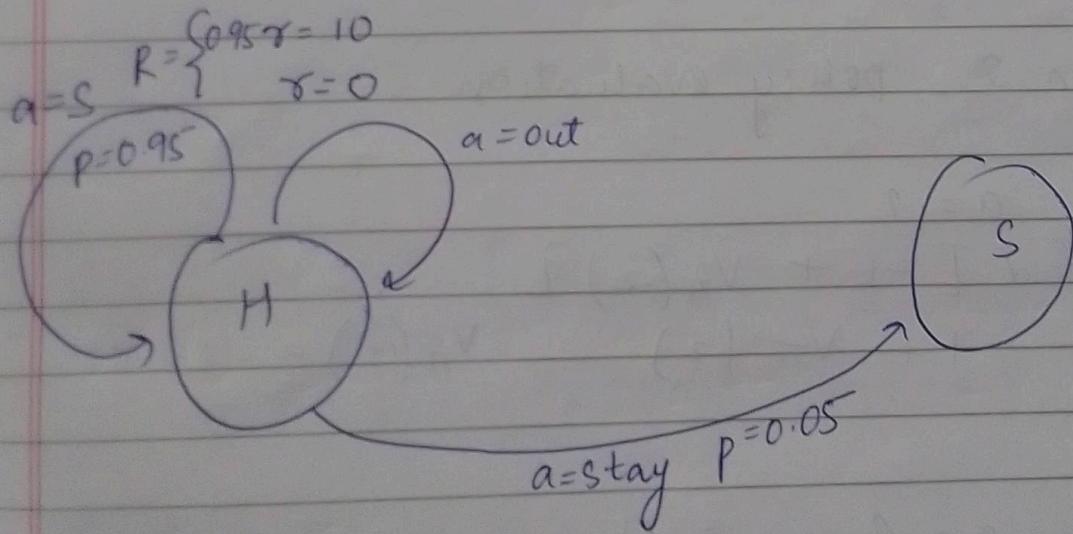
$$\pi(a/S=2) = J$$

Since The policy after iteration 1 & iteration 2 are same, we have found the optimal

The no of iteration will increase with $O(n)$

Q19) State: healthy or sick

Action: Stay home or go out \rightarrow Healthy
med or no med \rightarrow sick



$$\text{for } S=2 \quad a=1 : \quad 4 + v_{\pi}(F) = 4$$

$$a=J \quad [2 + v_{\pi}(I)] 0.5 + 0.5 [0 + v_{\pi}(II)] \\ = 4$$

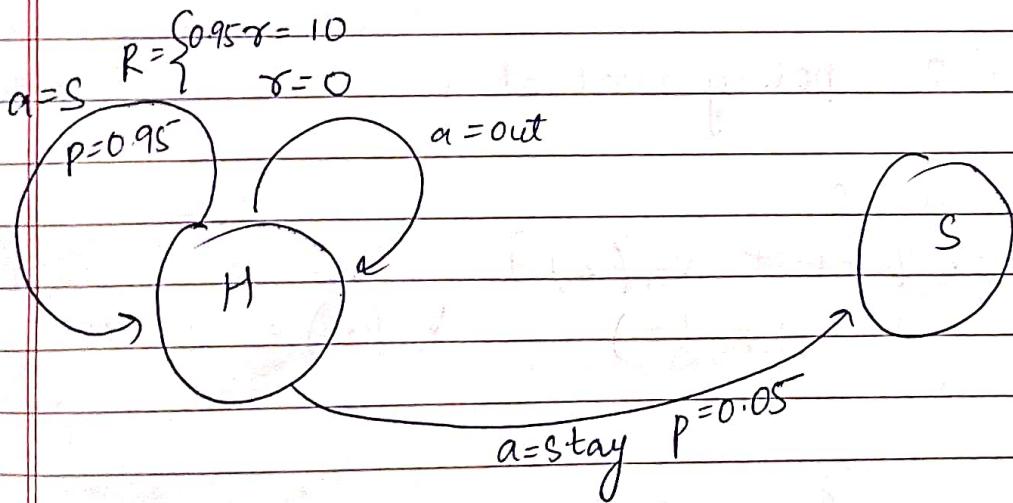
$$\pi(a/S=2) = 1$$

Since The policy after iteration 1 & iteration are same, we have found the optimal policy

The no of iteration will increase with $O(n)$

Q19) State: healthy or sick

Action: Stay home or go out \rightarrow Healthy
med or no med \rightarrow sick



$s \quad a \quad s' \quad r \quad p(s', r | s, a)$

H stay H 10 $0.95 \times 0.95 = 0.9025$

H out

H stay H 0 $0.95 \times 0.05 = 0.0475$

H stay S -10 $0.05 \times 0.9 = 0.045$

H stay S 0 $0.05 \times 0.1 = 0.005$

H out H 20 $0.7 \times 0.9 = 0.63$

H out H 0 $0.7 \times 0.1 = 0.07$

H out S -10 $0.3 \times 0.9 = 0.27$

H out S 0 $0.3 \times 0.1 = 0.03$

S med S -2 0.1

S med H -1 0.9

S nomed S -1 0.4

S nomed H 0 0.6

$$\gamma = 0.9$$

Q20) $\pi(a|s)$ - equal prob.

$V_\pi(s) \neq s$

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

State $\rightarrow H$

$V_\pi(H)$

$$\begin{aligned} a = \text{stay} \Rightarrow & 0.5 \left[0.9025 (10 + 0.9 V_\pi(H)) + \right. \\ & 0.0475 (0 + 0.9 V_\pi(H)) + \\ & \cancel{0.0005} 0.045 (-10 + 0.9 V_\pi(S)) + \\ & \left. 0.005 (0 + 0.9 V_\pi(S)) \right] \\ & + \end{aligned}$$

$$\begin{aligned} a = \text{out} \Rightarrow & 0.5 \left[0.63 (20 + 0.9 V_\pi(H)) + \right. \\ & 0.07 (0 + 0.9 V_\pi(H)) + \\ & 0.27 (-10 + 0.9 V_\pi(S)) + \\ & \left. 0.03 (0 + 0.9 V_\pi(S)) \right] \end{aligned}$$

$$V_\pi(H) = \frac{0.7425 V_\pi(H)}{9.2375} + \frac{0.15775 V_\pi(S)}{-1} \quad -(1)$$

State $\rightarrow S$

$V_\pi(S)$

$$\begin{aligned} a = \text{med} \Rightarrow & 0.5 \left[0.1 (-2 + 0.9 V_\pi(\cancel{H})) + \right. \\ & \left. 0.9 (-1 + 0.9 V_\pi(H)) \right] \end{aligned}$$

$$\begin{aligned} a = \text{no med} \Rightarrow & 0.5 \left[0.6 (0 + 0.9 V_\pi(H)) + \right. \\ & \left. 0.4 (-1 + 0.9 V_\pi(S)) \right] \end{aligned}$$

$$V_{\pi}(s) = 0.675 V_{\pi}(H) + 0.225 V_{\pi}(S) - 0.75 \quad -(2)$$

Solving eq 1 & 2

$$V_{\pi}(\text{healthy}) = 75.641$$

$$V_{\pi}(\text{sick}) = 64.913$$

(21.) $t \quad t+1$ $\pi \rightarrow 0.5 \text{ for all actions}$
 $H \quad S$

$$E[G_t | S_t = \text{healthy}, S_{t+1} = \text{sick}]$$

$$G_t = r_{t+1} + \gamma V_{\pi}(S_{t+1})$$

$$\hookrightarrow E[r_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = \text{healthy}, S_{t+1} = \text{sick}]$$

~~$$r_t + \gamma V_{\pi}(S_{t+1}) - p(r_t | s, s', a) \pi(a | s)$$~~

$$= \sum_{a, s} (r_t + \gamma V_{\pi}(S_{t+1})) \cdot p(r_t | s, s', a) \pi(a | s)$$

act stay ~~act out~~

$$a = \text{stay} \quad 0.5 [0.9 (-10 + 0.9 \times 64.913) + 0.1 (0 + 0.9 \times 64.913)]$$

+

$$a = \text{go out} \quad 0.5 [0.9 (-10 + 0.9 \times 64.913) + 0.1 (0 + 0.9 \times 64.913)]$$

CLASSMATE

Date _____
Page _____

$$2 \times 25 [0.01 \{ 43.57953 + 5.84217 \}]$$

$$= 49.4217$$