# AI 3000

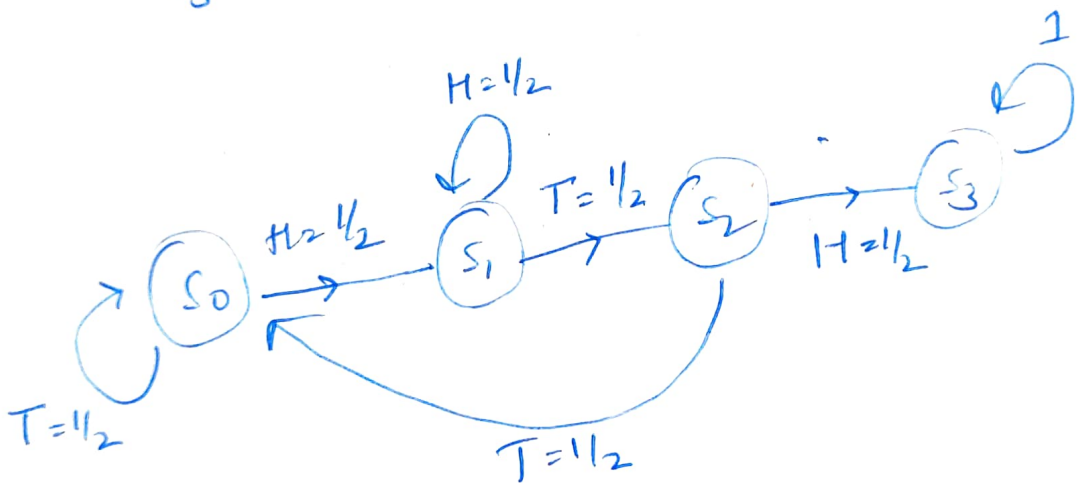## MID -TERM EXAM

EEI8BTECH11026

K. Surya Prakash

**Q1)** State - Transition diagram

$S_0$: Start state $\Big\}$  $S = \{ S_0, S_1, S_2, S_3 \}$
$S_3$: End state



$$P = \begin{bmatrix} 1/2 & 1/2 & 0 & 0 \\ 0 & 1/2 & 1/2 & 0 \\ 0 & 1/2 & 0 & 1/2 \\ 1/2 & 0 & 0 & 1/2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Let **Reward :**

$$R(S) = \begin{cases} -1 & ; \ S \in \{ S_0, S_1, S_2 \} \\ 0 & ; \ S \in \{ S_3 \} \end{cases}$$

Let discount factor : $(\gamma = 1)$

∴ This is a MRP $(S, P, R, \gamma)$

---

* Finding Expected no. of tosses through

Bellman Eqn's

$$V(S) = R(S) + \gamma \sum_{s' \in S} P_{ss'} V(s') \quad \rightarrow (1)$$

Eq(1) for all states:

$$V(S_0) = -1 + \tfrac{1}{2} (V(S_0) + V(S_1))$$
$$V(S_1) = -1 + \tfrac{1}{2} (V(S_1) + V(S_2))$$
$$V(S_2) = -1 + \tfrac{1}{2} (V(S_0) + V(S_3))$$
$$V(S_3) = 0$$

Solving:
$$V(S_2) = -1 + \tfrac{1}{2} (V(S_0))$$
$$V(S_1) = -1 + \tfrac{1}{2} \left( \tfrac{3}{2} V(S_1) - 1 \right)$$

$$\boxed{V(S_1) = -6}$$

$\Rightarrow V(S_0) = -1 + \tfrac{1}{2} (V(S_0) - 6))$

$\Rightarrow \boxed{V(S_0) = -8}$ //

* Since we choose $R(s) = -1 \,\forall\, s \in \{s_0, \cdots s_2\}$

It represents steps counted whenever we move from one state to others

∴ $V(s_0)$ should determine Expected no. of trials to reach the end state ($s_y$).

∴ Averge no. of trails $= \dfrac{8}{2}$

———— x ————

QI — DONE

Q2) Q             **Question 2**

**Q2.a)**      $R(s,a) = R_1(s,a) + R_2(s,a)$

$$\pi_1^*(s) = \arg\max_{a'} Q_1^{\pi_1^*}(s,a')$$

$$\pi_2^*(s) = \arg\max_{a'} Q_2^{\pi_2^*}(s,a')$$

$$Q^{\pi_1^*}(s,a) = R_1(s,a) + \gamma \sum P_{ss'}^{a} V^{\pi_1^*}(s')$$

$$Q^{\pi_2^*}(s,a) = R_2(s,a) + \gamma \sum_{s'} P_{ss'}^{a} V^{\pi_2^*}(s')$$

for $\pi^*$ : Reward: $R(s,a) \supset R_1(s,a), R_2(s,a)$

$$\pi^* = \arg\max_a Q^*(s,a)$$

$$Q^*(s,a) = (R_1(s,a) + R_2(s,a))$$
$$+ \gamma \sum_{s'} P_{ss'}^a \left( V^{\pi^*}(s') \right)$$

∴ It is not possible to formulate $\pi^*$ simply.

## Alternate Explanation:

Since $\pi$ depends on $V$

$$V^{\pi_1^*} = (I - \gamma P^{\pi_1})^{-1} R_1$$

$$V^{\pi_2^*} = (I - \gamma P^{\pi_2})^{-1} R_2$$

$$V^{\pi^*} = (I - \gamma P^{\pi})^{-1} (R_1 + R_2)$$

Cannot simply formulate $V^{\pi^*}$ as $V^{\pi_1} \not= V^{\pi_2^*}$

due to $\underline{\underline{P^\pi}}$ .

$$M = \langle S, A, P, R, \gamma \rangle$$

$$f,g : S \times A \to \mathbb{R} \quad ; \quad \mathcal{L} \to \text{Bellman optimality operator}$$

$$f \pm \quad V_f(s) = \max_a f(s,a)$$

Given;

$$(\mathcal{L}f)(s,a) = R(s,a) + \gamma \, P(s,a) \, V_f(s')$$

Prove that:

$$\| \mathcal{L}f - \mathcal{L}g \|_\infty \leq \gamma \, \| f - g \|_\infty$$

$$\Rightarrow \quad \| \mathcal{L}f - \mathcal{L}g \|_\infty =$$

$$\max_s \max_a \left( | \, \mathcal{L}f(s,a) - \mathcal{L}g(s,a) | \right)$$

$$= \max_s \max_a \left( | \left( R(s,a) + \gamma \, P(s,a) \, V_f(s) \right) \right.$$
$$\left. - \left( R(s,a) + \gamma \, P(s,a) \, V_g(s) \right) | \right)$$

$$= \max_s \max_a \left( | \, \gamma \, P(s,a) \left( V_f(s) - V_g(s) \right) | \right.$$

$$= \| r P(V_f - V_g) \|_\infty$$

$$\leq r \| P \|_\infty \| V_f - V_g \|_\infty \quad \Big\} \quad \begin{array}{l} P \text{ is } \overset{\text{probability}}{\cancel{too}} \\ \text{matrix} \\ \text{all elements} \leq 1 \end{array}$$

$$\leq r \| V_f - V_g \|_\infty \qquad \rightarrow (1)$$

we know that $V_f(s) = \max\limits_a f(\cdot, a)$

or

$$\| V_f - V_g \|_\infty = \max\limits_a \max\limits_s \left| V_f(s, a) - V_g(s, a) \right|$$

$$= \max\limits_a \max\limits_s \left( \max\limits_a (f(s, a)) - \max\limits_a (g(s, a)) \right)$$

$$\leq \max\limits_a \max\limits_s \left( \max\limits_a (f(s, a) - g(s, a)) \right)$$

$$\hookrightarrow \text{IMP!}$$

$$\leq \max\limits_a \max\limits_s \left( f(s, a) - g(s, a) \right)$$

$$\leq \| f - g \|_\infty$$

$$\cdot \gamma \| \nabla_f - \nabla_g \|_\infty \leq \| f - g \|_\infty$$

continuing from (1)

$$\leq \gamma \| \nabla_f - \nabla_g \|_\infty$$

$$\leq \gamma \| f - g \|_\infty$$

$$\boxed{\| L_f - L_g \|_\infty \leq \gamma \| f - g \|_\infty}$$

⌖ Q-2  DONE ⌖

**Q4)** Problem Formulation:

To make it simple:

State of an ingridient depends on:

1) Quanty ( < thresh , > thresholds)

2) Days old ( 0, 1, 2,≥3 )

Each state is defined by quantity of the ingridient ( if less than a certain threshold (1) or not)

6 day olden ($\sigma$, first), $\not\leq 3\not\sigma$, ($\leq 3$, $>3$)
                                                binary.

Action space:

The Owner can have 3 actions

1) Refill: ~~By~~ Buy and add ingridients
           by not discarding alredy non
           expired ingridients

2) Replace: Discard the existing lot &
           replace with a new one

3) Idle: Do nothing.

_____

(b) Reward :-         S: $S_q$, $S_d$

$$\cancel{R(S,a)} = \begin{cases} S_q = 1 & \text{: if } > Thrs \\ S_q = 0 & \text{: else} \end{cases}$$

$$S_d = 0 \text{ ; if } < days$$
$$= 1 \text{ ; if } > days$$

$R(S,a) =) \quad R(S, Refill) =$

# Preferred rewards

| $q$, $d$ | Best action ( High rewards |
|---|---|
| $<T$ ; $<3$ → | Refill > Idle > Replace |
| $<T$ ; $>3$ → | Replace > ~~Idle~~ ~~Refill~~ > Idle |
| $0 >T$; $>3$ → | ~~Idle~~ Replace >> Idle > Refill |
| $>T$; $<3$ → | Idle > Refill > Replace |

(C) Since my state space is binary, will be using non-discounted.

But if no. of days taken ( 0, 1, 2, 3 ... ) discounting would give better results.

→ If d → closer to 3 ; reward for replacing increases.

(d) Since this is an real time problem and accurate formulation of state parameters like (P) are changing. We need to rely on Reinforcement Learning. We can gain more knowledge over time and can make the model better.

e. Since state and action space is finite and less in this case. It will not be harder to traverse through the entire trajectory. We can

$\Rightarrow$ MC methods will be used.

∵ This is because the length of the episode for each ingrndient is finite and less.

(f). Yes. Since the data is based on

Learning is ever changing.

→ An acute function approximator

might lower the variance. We

can model the real world phenomenon

like : likeability of combination

of bread & filling using a

function.

Q-4 _ Done

Q5)                          MISCELLANEOUS

a)   No. MDP formalism requires knowledge
of $\langle S, P, R, \gamma A \rangle$.

→ Some times $|S|$ can be huge and will
be difficult to keep track. Eg. Atari games

→ P is also hard to formulate for larger
environments.

→ This is where we use model-free methods.

_____

b)  For MDP:

Discount $\left.\right\}$ $\gamma \in [0,1]$
factor

→ Modeling MDP requires choosing actions
to maximise total discounted reward

$$\mathbb{E}(G_t | S_t = s):$$

$$G_t = \sum_{k=0}^{\infty} (\gamma^k \, r_{t+k+1})$$

If episodes are long. To avoid $G_t$
from exploding we have to have

$$\gamma \in [0,1].$$

(C) By ordering of policies

if $\pi > \pi' \Rightarrow V^{\pi}(s) \geq V^{\pi'}(s)$    <u>Yes</u>

$\Rightarrow$ In Policy Iteration:

we converge to the best optimal policy $\pi^*$

$\Rightarrow$ ie; $\pi^* > \pi_i$   $\forall$ ~~are~~ possible policies

$$\pi_*(s) = \arg\max_{a \in A} Q_*(s,a)$$

$$V_*(s) = \max_{a \in A} Q_*(s,a) \quad \forall s \in S$$

$$V^{\pi_*}(s) \geq V^{\pi'}(s) \quad ; \quad \pi' \text{ all other policies,}$$

$$\underbrace{\breve{V}(s) \geq V^{\pi'}(s)}$$

The 'Optimal policy achieved through policy iteration attains the ~~to~~ Optimal value function

$\rightarrow$ <u>In Value Iteration</u>:

$$V_{k+1}(s) \leftarrow \max_{a} \left( \sum_{ss'} P^a_{ss'} \left( R^a_{ss'} + \gamma V_k(s') \right) \right)$$

$\Rightarrow$ $V_k$ will converge to $V_*$

$$\boxed{V^*(s) \geq V^{\pi}(s)}$$

(d) In DP-setup, we already modelled the randomnen of the system using 'P'. which we lack in model-free setting.

→ The ε-greedy in the latter tries to inculcate the randomness. ~~by~~

→ This ensures continual exploration which is already present in DP.

(e)

Advantages of MC over DP

i) ~~the~~ Compatible with Non-Markovian domain problems

(ii) Since it is a model-free method, does not require ~~&~~ environment parameters like 'P'.

(f) Evaluating value func. in MRP

$V(S)$ ~~←+~~

$V \leftarrow R + \gamma P V$

$V \rightarrow$ vector of size $|S|$

$R \rightarrow$ " " $|S|$

$P \rightarrow$ matrix : of shape $|S| \times |S|$

Each update consists of $O(|S|^2)$ operations.

∴.

(g) Yes. It is possible.

There might be no action that can jump from $S_1 \rightarrow S_3$.

only possibilities are $S_1 \rightarrow S_2 : a_2$

$S_1 \rightarrow S_1 : a_1$

∴ $a_2$ pushes $S_1 \rightarrow S_2$ : Reward = 2

$a_1$ : $S_1 \rightarrow S_1$ : Reward = -1

∴ If followed a greedy policy

~~π(S)~~ $\pi(S_1) = a_1$

But since, we follow $\varepsilon$-greedy

* 1st action:

  $S_1, a_1, 1, S_1$ }  Random

* 2nd action

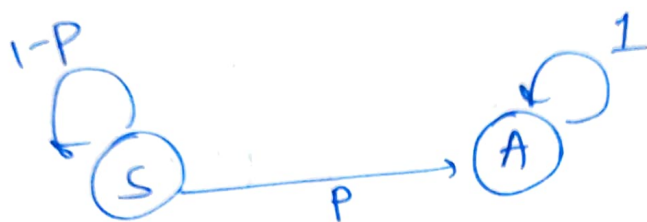  $S_1, a_2, 2, S_2$ }  greedy

  $Q(S_1, a_2) > Q(S_1, a_1)$.

<div align="center">

## Problem 3

</div>

(Q3)  Monte Carlo Methods:

$S = \{S, A\}$  ;  $\tau = 1$ ,  $P \in (0,1)$



a)  Trajectory's  generic  form.

$S \to S \to S \to \cdots \to \underset{\sim}{A}$
first appearance
of 'A'

$S^N A$ ;  where  $A = \{0, 1, 3 \cdots \infty\}$

(b)  $V(S)$  using  first Visit  MC:

The  trajectories  can be

$\{S^N A\}$

Let  prob. of  getting  trajectory :

$$P(S^N A) = P(1-P)^{N-1}$$

## b) First visit MC:

| Trajectory | Reward | Prob. of occuring |
|---|---|---|
| S-A | 1 | $p$ |
| S-S-A | 2 | $p(1-p)$ |
| $S^3 A$ | 3 | $p(1-p)^2$ |
| $S^N A$ | N | $p(1-p)^N$ |

### First visit MC:

average reward:

$V(S)$

$$R = 1(p) + 2p(1-p) + 3p(1-p)^2 + \cdots N \, p(1-p)^{N-1}$$

$$= p(1 + 2(1-p) + 3(1-p)^2 + \cdots N(1-p)^{N-1})$$

- Upon solving

$$R = V(S) = \frac{1 - (N+1)(1-p)^N + N(1-p)^{N+1}}{p}$$

(c) Every Visit MC.

$$S - A \quad = \quad 1 \qquad p$$

$$S - S - A \quad = \quad 1+2 \qquad p(1-p)$$

$$S - S - S - A \quad = \quad 1+2+3 \qquad p^2(1-p)^2$$

$$\therefore \text{Visit} = \sum_{n=1}^{N} \frac{n(n+1)}{2} (1-p)^{n-1} p.$$

(d) True Estimate

$$V(s) = 1 + p(V(A)) + (1-p)(V(s))$$

$$\Rightarrow \quad V(A) = 0 + 1(V(A)) \Rightarrow \boxed{V(A) = 0}$$

$$\Rightarrow \quad V(s) = (1-p)(V(s)) + 1$$

$$\Rightarrow \quad p(V(s)) = 1$$

$$\Rightarrow \quad \boxed{V(s) = 1/p}$$

(e) Yes every visit MC is biased

$$\text{Bias} ( \hat{v}(s)) = E( \hat{v}(s)) - V(s)$$

$$\hat{v}(s) \Rightarrow \text{for } \underline{n \text{ trajectories}}.$$

$$
\left.
\begin{array}{l}
S-A: = 1 \\
S-S-A = 1+2 \\
S-S-S-A = 1+2+3
\end{array}
\right\}
\qquad
V(s) = \sum_{n=1}^{N} \frac{n(n+1)}{2} (1-p)^{n-1}
$$

$$= E \left( \sum_{n=1}^{N} \frac{n(n+1)}{2} (1-p)^{n-1} P \right) - \frac{1}{P} \leq$$

$$\neq 0$$

$\therefore$ Bias exists"

---

(f) $=$ MC convergence is based on the "Law of large numbers"

Each element in sequence should be Identical & independent distributed random variables. This leads to point wise convergence.

— Q3 Done —