Name - T. Vijay

Roll·no - CS17BTECH11040

## Assignment No.1

Q1

a) States S: Tall , Medium , short
$\qquad$ (TT) , (TS or ST) , (SS)

Fig: Transition Probabilities



$$\text{Transition probability Matrix } P = \begin{array}{c} \\ T \\ M \\ S \end{array} \begin{array}{ccc} T & M & S \\ \left[\begin{array}{ccc} 1/2 & 1/2 & 0 \\ 1/4 & 1/2 & 1/4 \\ 0 & 1/2 & 1/2 \end{array}\right] \end{array}$$

b) We know that $n^{th}$ step transition matrix is given by $P^n$

∴ Probabilities of Tall, short, medium offspring belonging to

first generation = second row of P matrix

$$= \quad 1/4 \quad , \quad \frac{1}{2} \quad , \quad 1/4$$
$$\qquad (Tall) \qquad (Medium) \qquad (short)$$

## Second generation

$$P^2 = \begin{bmatrix} 0.375, & 0.5, & 0.125 \\ 0.25, & 0.5, & 0.25 \\ 0.125, & 0.5, & 0.375 \end{bmatrix} \longrightarrow \text{Second row}$$

$\therefore P_{Tall} = 0.25$ , $P_{medium} = 0.5$ , $P_{short} = 0.25$

Similarly for __Third generation__

$P_{Tall} = 0.25$ , $P_{medium} = 0.5$ , $P_{short} = 0.25$

c) Second row of $P^n$ $\forall$ $n \in N$ = $\begin{bmatrix} 0.25, & 0.5, & 0.25 \end{bmatrix}$

$\therefore P_{Tall} = 0.25$ , $P_{medium} = 0.5$ , $P_{short} = 0.25$

**Q2**

@ States : S, 1, 3, 5, 6, 7, 8, W [we can't stay on other states]

Transition Matrix

|   | S | 1 | 3 | 5 | 6 | 7 | 8 | W |
|---|---|---|---|---|---|---|---|---|
| S | 0 | 1/4 | 1/4 | 0 | 0 | 1/4 | 1/4 | 0 |
| 1 | 0 | 0 | 1/4 | 1/4 | 0 | 1/4 | 1/4 | 0 |
| 3 | 0 | 0 | 0 | 1/4 | 1/4 | 1/4 | 1/4 | 0 |
| 5 | 0 | 0 | 1/4 | 0 | 1/4 | 1/4 | 1/4 | 0 |
| 6 | 0 | 0 | 1/4 | 0 | 0 | 1/4 | 1/4 | 1/4 |
| 7 | 0 | 0 | 1/4 | 0 | 0 | 1/4 | 1/4 | 1/4 |
| 8 | 0 | 0 | 1/4 | 0 | 0 | 0 | 1/2 | 1/4 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

b W is the absorbing state as we cannot leave state W after entering it.

© 

We take the reward = -1 at state, so the final value of each state will indicate the no. of ~~steps~~ expected steps but the value will be negative.

From Bellman eqⁿ, we have $V = R + \delta PV$

$$\therefore V = (I - \delta P) R$$

We take $\delta = 1$, $P = $ Probability transision matrix
(Discount factor)  (We will exclude the last row and column from matrix obtained)

$$R = \begin{bmatrix} -1, & -1, & -1, & -1, & -1, & -1, & -1 \end{bmatrix}^T_{1\times7}$$

$P_{7\times7}$ ← Excluding last row and column from P (Q2 a)

$$\therefore V = \begin{bmatrix} 7.083, & 7, & 6.67, & 6.67, & 5.33, & 5.33, & 5.33 \end{bmatrix}^T$$

$\therefore V(i)$ represents the no -of steps from state i

Q3 a)



Actions: D → Driving

ND → not Driving

b) Deterministic policy

$$\pi(s) = \begin{cases} \text{Drive} & , \quad s = \text{top} \\ \text{Drive} & , \quad s = \text{rolling} \\ \text{Drive} & , \quad s = \text{Top Bottom} \end{cases}$$
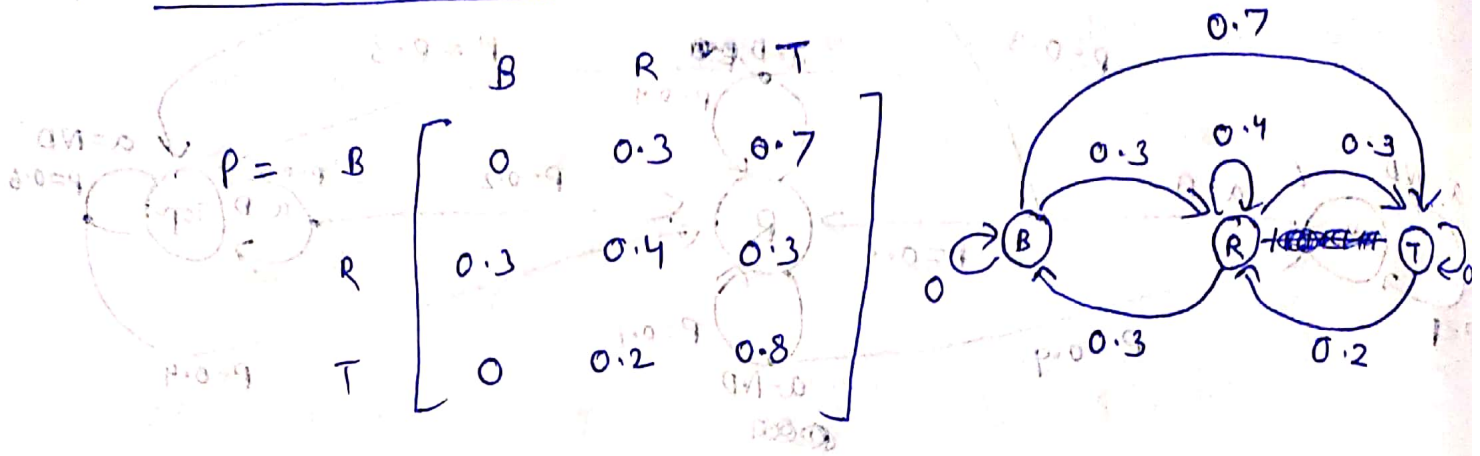
Stochastic policy

$$\pi(a \mid \text{bottom}) = \begin{cases} 0.5 & , \quad a = D \\ 0.5 & , \quad a = ND \end{cases}$$
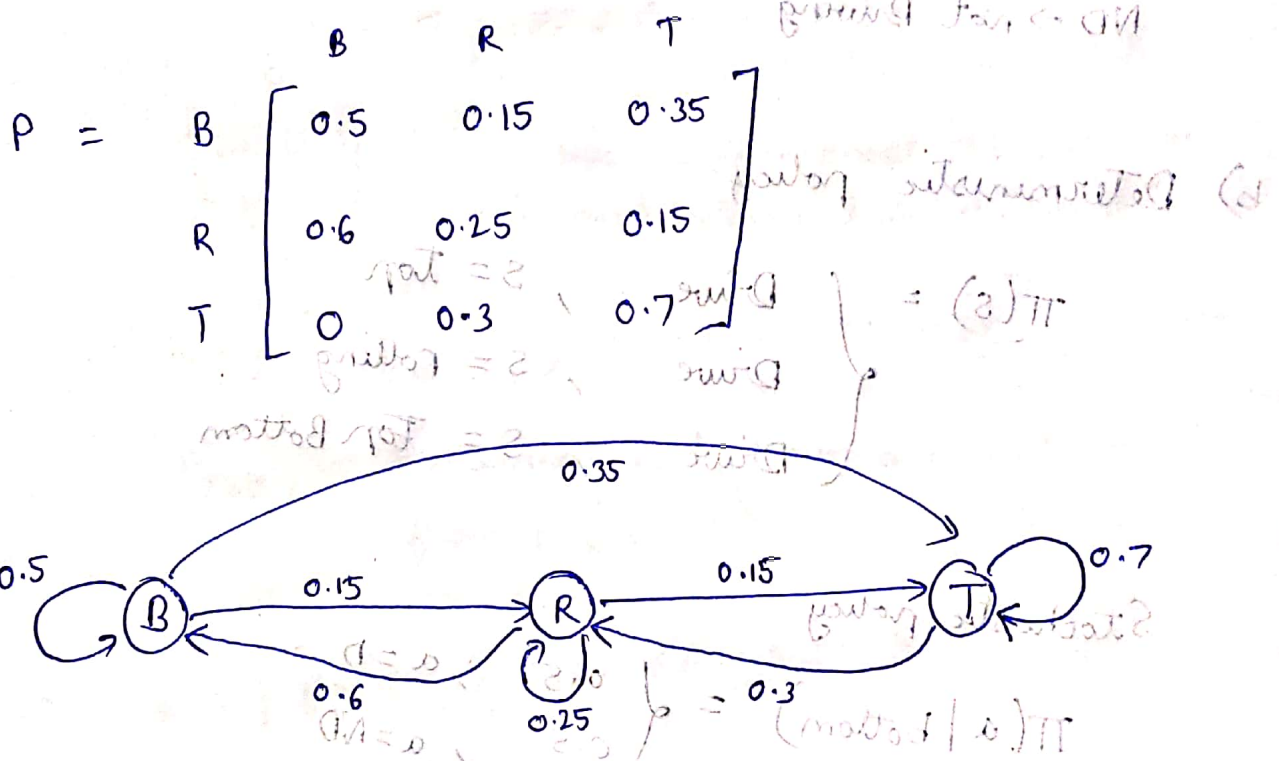
$$\pi(a \mid \text{Top}) = \begin{cases} 0.5 & , \quad a = D \\ 0.5 & , \quad a = ND \end{cases}$$

$$\pi(a \mid \text{Rolling}) = \begin{cases} 0.5 & , \quad a = D \\ 0.5 & , \quad a = ND \end{cases}$$

## c) Deterministic policy

$$P = \begin{array}{c} \\ B \\ R \\ T \end{array} \begin{array}{ccc} B & R & T \\ \left[\begin{array}{ccc} 0 & 0.3 & 0.7 \\ 0.3 & 0.4 & 0.3 \\ 0 & 0.2 & 0.8 \end{array}\right] \end{array}$$



## Stochastic policy

$$P = \begin{array}{c} \\ B \\ R \\ T \end{array} \begin{array}{ccc} B & R & T \\ \left[\begin{array}{ccc} 0.5 & 0.15 & 0.35 \\ 0.6 & 0.25 & 0.15 \\ 0 & 0.3 & 0.7 \end{array}\right] \end{array}$$



d)

$$\pi(a \mid S_t = s, S_{t-1} = s) = \begin{cases} 0.5 & a = \text{Drive} \\ 0.5 & a = \text{Not driving} \end{cases}$$

$$\pi(a \mid S_t = s, S_{t-1} \neq s) = \begin{cases} 0.7 & a = D \\ 0.3 & a = ND \end{cases}$$

Q4 a) $P^{\pi_1} = $

$$\begin{array}{c c} & \begin{array}{c c c c} a & b & c & d \end{array} \\ \begin{array}{c} a \\ b \\ c \\ d \end{array} & \left[ \begin{array}{c c c c} 0 & 0.9 & 0.1 & 0 \\ 0.1 & 0 & 0 & 0.9 \\ 0.9 & 0 & 0 & 0.1 \\ 0 & 0 & 0 & 1 \end{array} \right] \end{array}$$

$P^{\pi_2} = $

$$\begin{array}{c c} & \begin{array}{c c c c} a & b & c & d \end{array} \\ \begin{array}{c} a \\ b \\ c \\ d \end{array} & \left[ \begin{array}{c c c c} 0 & 0.1 & 0.9 & 0 \\ 0.9 & 0 & 0 & 0.1 \\ 0.1 & 0 & 0 & 0.9 \\ 0 & 0 & 0 & 1 \end{array} \right] \end{array}$$

$P^{\pi_3} = $

$$\begin{array}{c c} & \begin{array}{c c c c} a & b & c & d \end{array} \\ \begin{array}{c} a \\ b \\ c \\ d \end{array} & \left[ \begin{array}{c c c c} 0 & 0.42 & 0.58 & 0 \\ 0.1 & 0 & 0 & 0.9 \\ 0.1 & 0 & 0 & 0.9 \\ 0 & 0 & 0 & 1 \end{array} \right] \end{array}$$

$R^{\pi} = \begin{bmatrix} -10 & -10 & -10 & 100 \end{bmatrix}^T$  for all policies
because reward for each
action is same.

Let $\xi = \gamma = 0.9$

$$V^{\pi} = (I - \gamma P^{\pi})^{-1} R^{\pi}$$

a)

$V^{\pi_1} = \begin{bmatrix} 63.2 & 84.69 & 51.08 & 109.89 \end{bmatrix}^T$

$V^{\pi_2} = \begin{bmatrix} 63.2 & 51.08 & 84.69 & 109.89 \end{bmatrix}^T$

$V^{\pi_3} = \begin{bmatrix} 66.49 & 84.99 & 84.99 & 109.89 \end{bmatrix}^T$

b) $\pi_3$ is the best policy because for every state $V^{\pi_3}$ has better values* o than the values of other policies

c) $\pi_1$ and $\pi_2$ are not comparable because, values of $V^{\pi_1}$ is greater than $V^{\pi_2}$ for some states while $V^{\pi_2}$ is greater than $V^{\pi_1}$ for some other states

∴ $\pi_1$ and $\pi_2$ cannot be compared.

We will use value itteration to find the optimal value function

Step I : We initialize $V_i(s)$ for $s \in S_0$ to small value $\varepsilon$

Then for all states, we find

$$V_{k+1}(s) = \max_a \left[ \sum_{s' \in S} P_{ss'}^a \left( R_{ss'}^a + Y V_k(s') \right) \right]$$

If $\left( |V_{k+1}(s) - V_k(s)| < \varepsilon \quad \forall S \in S \right)$

we go to the next step

else we repeat the previous step

Step ② : For all state, $s \in S$

$$\pi_* (s) = \underset{a}{\text{argmax}} \; V_* (s)$$

After using the above process we will get the following optimal function

$$\pi (a|s) = \pi_s (a|s) \quad \forall a \; \forall s$$

<u>Intitution</u>

① The above value function is optimal because for each state S, we consider the corresponding optimal state function and then find an optimal policy that achieve this value function and use the probabilities of action associated at this state

② Another Intitution is that if we apply policy iteration on the above policy, we won't get any better policy

because the current policy achieves the max value at state. Hence it can't be improved.

Q6. a) The optimal policy $\pi$ is to go right at each state.

The optimal value function has value 10 at each state because $\gamma = 1$ (far-sighted).

b) $\gamma = 0.9$     $V = [5.9, 6.56, 7.29, 8.1, 9, 10]$

     $\gamma = 0.5$     $V = [0.31, 0.62, 1.25, 2.5, 5, 10]$

     $\gamma = 0.1$     $V = [10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10]$

The optimal policy still remains the same i.e to go right at each step.

Decrease in Discount factor leads to decrease in optimal value but optimal value of each state increases as we go to the right.

c) For all non-positive values of C, the optimal policy will still remain the same i.e to go right because we don't want to accumulate negative rewards.

When $c \leq 10$, the optimal policy is to move right when ~~more than 10~~ $\gamma$ is small.

when $c > 10$, then the optimal policy is to ~~move~~

~~from $S_1$ to $S_5$ and then~~ track

turn left when the agent reaches state $S_5$.

This will lead to an infinite process.

d) The new policy

$$\hat{V}_B^\pi = (1-\delta p)^{-1} (R+c)$$

$$\text{where } c = [c \quad c \quad c ---]_{1 \times n}^T$$

$$\therefore \hat{V}^\pi = V^\pi + (1-\delta p)^{-1} c$$

where $V^\pi$ is value function for any policy $\pi$.

Q7

a) * When Y is low we need to prefer close exit
   as we are not far sighted

   * when Y is high we will prefer distant exit
   as we are far sighted

   * when $\eta$ is high, we will avoid cliff because
   we don't have the control over the path taken

   * when $\eta$ is low, we can risk the cliff because
   we have control over the our actions

① Prefer close exit but risk cliff
   → $Y = 0.1$ and noise = 0

② Prefer distant exit but risk cliff
   → $Y = 0.9$ and noise = 0.1

③ Prefer close exit by avoiding the cliff
   → $Y = 0.1$ and noise = 0.5

④ Prefer the distant exit by avoiding the cliff
   → $Y = 0.9$ and noise = 0.5

Q 88  Given:  $$L(v) = \max_{a \in A} [R^a + \gamma P^a v] \quad — ①$$

Since, $V_*$ is fixed point of operator $L$, we have

$$L(v^*) = v^*$$

and,

$$|V_{K+1} - V^*|_\infty = |L(v_K) - L(v^*)|_\infty \quad (\text{From } ①)$$

$$= \left| \max_a \{R + \gamma P^a v_K\} - \max_a \{R + \gamma P^a v^*\} \right|_\infty$$

$$\leq \max_a \left| \{R + \gamma P^a v_K\} - \{R + \gamma P^a v^*\} \right|_\infty$$

$$= \gamma |P^a(v_K - v^*)|_\infty$$

$$\leq \gamma |v_K - v^*|_\infty$$

i.e $|V_{K+1} - V^*|_\infty \leq \gamma |v_K - v^*|_\infty$

and recursively $|V_{K+1} - V^*|_\infty$

Hence Proved