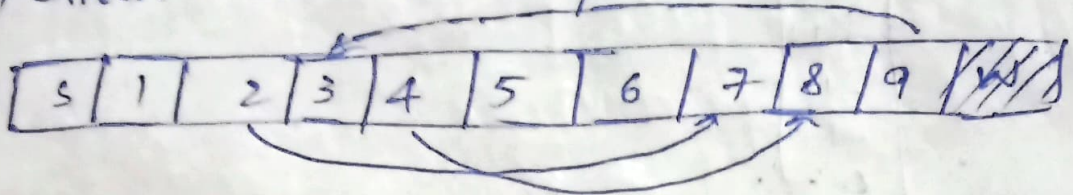- Abburi Venkata Sai Mahesh

- CS18BTECH11001

## Problem 1:

a) Given Markov reward process is



### States:

As 2, 4, 9 are states having either ladder and snake any die roll leading to these places will redirect to othe states. so we can ignore these state (The ans will be same even if we consider these states)

$$\therefore S = \{S, 1, 2, 3, 5, 6, 7, 8, W\}$$

where W is terminal state
others are non terminal states.

### Transition probability matrix

|     | S | 1    | 3    | 5    | 6    | 7    | 8    | W    |
|-----|---|------|------|------|------|------|------|------|
| S   | 0 | 0.25 | 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) Reward function

The reward is $-1$ for non terminal state
$0$ for terminal state

i. $R = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ 0 \end{bmatrix}$

discount factor : 1

Using Bellman equation

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

$$= \begin{bmatrix} -7.08 \ 333 \\ -7 \\ -6.6666 \\ -6.6666 \\ -5.3333 \\ -5.3333 \\ -5.3333 \end{bmatrix}$$

$\left[ \because \text{Considering only} \atop \text{non-terminal states} \right]$

ii. The expected number of die throws $= -7.08$
$\sim 7$

2) a) Given $M = \langle S, A, P, R, \gamma \rangle$

$$R(s,a) = R_1(s,a) + R_2(s,a)$$

Given value action function for policy $\pi$,

$Q_1^\pi(s,a)$ with reward function $R_1(s,a)$

value action function $Q_2^\pi(s,a)$ with

reward function $R_2(s,a)$

So $Q_1^\pi(s,a) = E_\pi \left( \sum_{k>0}^\infty \gamma^k r_{t+k+1} \mid s_t = s, q_t = a \right) \gamma \sim R_1(s,a)$

$Q_2^\pi(s,a) = E_\pi \left( \sum_{k=0}^\infty \gamma^k r_{t+k+1} \mid t = s, q_t = a \right) \gamma \sim R_2(s,a)$

we can
say that $\overline{Q}(s,a) = Q_1^\pi(s,a) + Q_2^\pi(s,a)$

But we cannot exactly say that

$$Q^*(s,a) = Q_1^*(s,a) + Q_2^*(s,a)$$

Because the variation of $Q_1(s,a)$ might not

be propotional to $Q_2(s,a)$ so the optimal

value will not be for same action

So it is not possible to combine

the action value functions in a

simpler manner.

b) Given

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

$f, g : S \times A \to R.$

$(4)\ (s, a) = R(s, a) + \gamma P(s, a) , V_f(s)$

where $\cdot V_f(s) = \max\limits_a f(\cdot, a).$

Consider

$$\| L_f - L_g \|_\infty = \| R(s, a) + \gamma P(s, a) V_f(s) -$$

$$R(s, a) + \gamma P(s, a) V_g(s) \|_\infty$$

$$= \| \gamma P(s, a) [ V_f(s) - V_g(s) ] \|_\infty$$

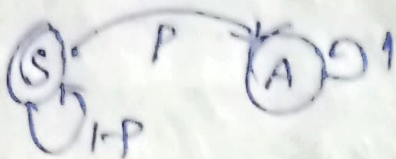$$\leq \gamma \| V_f(s) - V_g(s) \|_\infty \quad [\because P(s, a) \leq 1]$$

Now consider

$$\| V_f - V_g \| = | \max\limits_a f(\cdot, a) - \max\limits_a g(\cdot, a) |$$

$$\leq \max\limits_a \| f(\cdot, a) - g(\cdot, a) |$$

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

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

$$\Rightarrow \| L_f - L_g \|_\infty \leq \gamma \| V_f(s) - V_g(s) \|_\infty$$

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

Hence proved.

3. Given Markov process-transition probabilities

|   | S | A |
|---|---|---|
| S | 1-P | P |
| A | 0 | 1 |



a) This shows that A is a terminal state
   if we start from state s
   then     s → A          sA
            s → s → A       s²A

            s → s → s - s → A   s^k A

   i.e once we hit A we exit
   ∴ we can write general form of trajectory
   as s^k A where k ≥ 1 [∴ s repeated for k
   times and exited with A]

b) We calculate MC for first k trajectories

            S → A          r = 1
            S → S → A       r = 2
            ⋮
            S → S → S..S → A   r = k

   $V(s) = \dfrac{\sum_{k=1}^{k}(k_d)}{k}, \dfrac{K(K+1)}{\frac{2}{k}} ≥ \dfrac{k+1}{2}$

   $$\boxed{\therefore V(s) = \dfrac{k+1}{2}}$$

c) we consider every visit MC for first k trajectories

$$S \rightarrow A \qquad r = 1$$
$$S \rightarrow S \rightarrow A \qquad r = 2 + 1$$
$$S \rightarrow S \rightarrow S \rightarrow A \qquad r = 3 + 2 + 1 \cdots$$

$$\vdots$$

$$S \rightarrow S \rightarrow S \rightarrow \cdots S \rightarrow A \quad r = k + (k-1) + \cdots - 1$$

$$\therefore V(s) = \frac{\sum\limits_{k=1}^{k} \frac{k(k+1)}{2}}{\sum\limits_{k=1}^{k} k} = \frac{\frac{1}{2}\left[\frac{k(k+1)(2k+1)}{6} + \frac{k(k+1)}{2}\right]}{\frac{k(k+1)}{2}}$$

$$= \frac{1}{2}\left[\frac{2k+1}{3} + 1\right]$$

$$\boxed{\therefore V(s) = \frac{k+2}{3}}$$

d) True value of $V(s) = (I-P)^{-1} R$

$$P = \begin{bmatrix} 1-P & P \\ 0 & 1 \end{bmatrix} \qquad R = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} \because g_{Nen} \end{bmatrix}$$

So by removing terminal state

$$P = \begin{bmatrix} 1-P \end{bmatrix} \qquad R = \begin{bmatrix} 1 \end{bmatrix}$$

$$V(s) = \begin{bmatrix} [1] - [1-P] \end{bmatrix}^{-1} \begin{bmatrix} 1 \end{bmatrix}$$

$$= [P]^{-1} [1]$$

$$\boxed{\therefore V(s) \sim \frac{1}{P}}$$

e) Consider the expectation of every visit MC

$$E[V(s)] = E\left[\frac{K+2}{3}\right]$$

$$= \sum_{k=1}^{\infty} \left(\frac{k+2}{3}\right) P_r(s \text{ repeats } k \text{ times})$$

$$= \sum_{k=1}^{\infty} \left(\frac{k+2}{3}\right)(1-p)^{k-1} \cdot p$$

$$= \frac{p}{3} \underbrace{\sum_{k=1}^{\infty} (k+2)(1-p)^{k-1}}_{\text{sum of } \infty \text{ A GP}}$$

$$= \frac{p}{3}\left[\frac{3}{1-(1-p)} + \frac{1(1-p)}{(1-(1-p))^2}\right]$$

$$= \frac{p}{3}\left[\frac{3}{p} + \frac{1-p}{p^2}\right]$$

$$= \frac{p}{3}\left[\frac{1+2p}{p^2}\right]$$

$$\boxed{\therefore E[V(s)] = \frac{1+2p}{3p} \neq V(s)}$$

$$\therefore \text{The MC estimate is biased.}$$

f) The first visit MC has low bias & high variance but the Every visit MC has high bias & low variance.

Both MC converges to the unique $\hat{V}$ as the no. of trajectories goes to $\infty$.

$\therefore$ By law of large numbers both converges.

5) a) for the TD($\lambda$)

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t$$

when $\lambda = 1$

$$G_t^\lambda = 0 + \lambda^{T-t-1} G_t$$

$$= (1)^{T-t-1} G_t$$

$$G_t^1 = G_t$$

Therefore it is Monte carlo method as state, action process goes all the way to the end.

b) T.D(0) is low variance and high bias

~~TD(0) is~~

TD(n) is high variance & low bias

we consider TD($\lambda$) as a trade off between variance and bias.

c) If all rewards are scaled with a positive constant, the expected reward is scaled with that constant. Therefore the best policy is not affected.

d) If the behaviour policy is deterministic the chances for the exploration would decrease and might not perform well for a ~~too~~ stachoistic target policy.

i. It may not be beneficial.

è The convergence take place under following conditions

1. state and action spaces are finite

2. All state-action pairs are visited infinitely often.

3. Robbins-Monroe condition

$$\sum_t \alpha_t = \infty, \quad \sum_t \alpha_t^2 < \infty$$

f) The MC method for policy evaluation is the sample mean for the distribution of rewards. As sample mean is a random variable and is an estimator of population mean, the expected value of sample mean is same as population mean. Therefore policy evaluation MC method is unbiased estimator.

h) The value iteration update for state s is

$$V_{t+1}(s) = \max_{a \in A} R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V_t(s')$$

As the $V_{t+1}(s)$ depends only on values in the, $V_t(s)$ and not on any other entries of $V_{t+1}$, it is possible to parallelize the calculations.