

# Solutions: Assignment 1

30 July 2025 16:16

MAINAK BISWAS, IISc  
PhD student, PMRF.

## THEORY

1) : - (8 points)

### Markov Chain

(a)  $\pi P = \pi$  : Stationary distribution  $\pi$

where  $\pi$  is a row-vector.

$$\text{let } \pi = [\pi_0, \pi_1, \dots, \pi_N] = [\pi_0, \pi_1, \dots, \pi_N] \begin{bmatrix} p_{00} & p_{01} & \dots & p_{0N} \\ p_{10} & \dots & \dots & p_{1N} \\ \vdots & \vdots & \ddots & \vdots \\ p_{N0} & p_{N1} & \dots & p_{NN} \end{bmatrix}$$

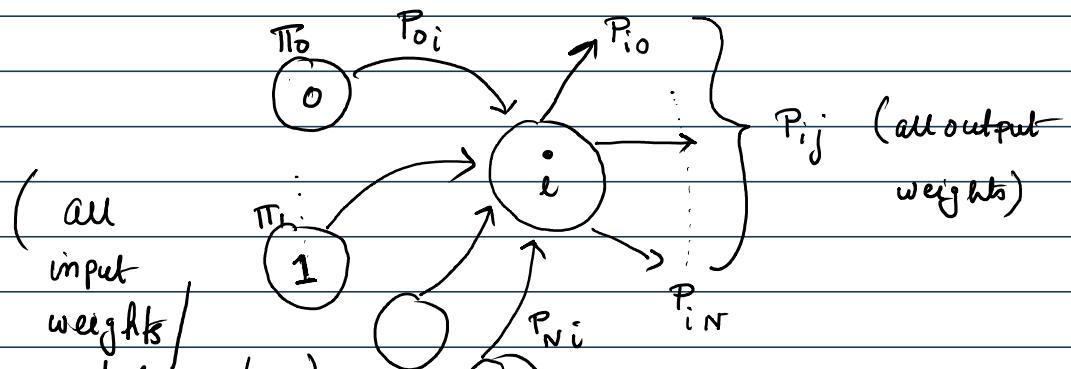
$$\pi_i = \sum_{j=0}^N \pi_j p_{ji}$$

$$\Rightarrow \pi_i (1 - p_{ii}) = \sum_{\substack{j=0 \\ j \neq i}}^N \pi_j p_{ji}$$

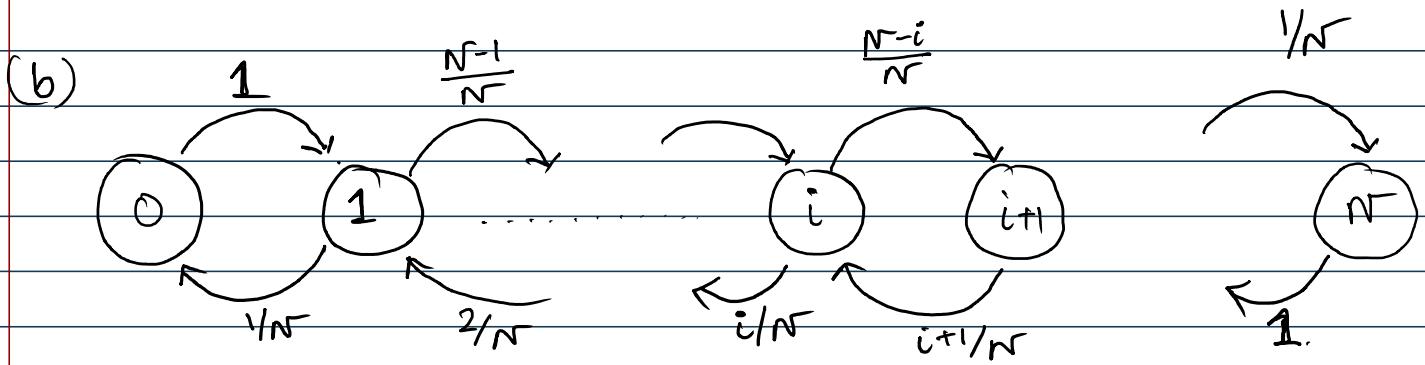
and  $j \neq i$

$$1 - p_{ii} = \sum_{j \neq i} p_{ij}$$

$$\text{Hence: } \pi_i \sum_{j \neq i} p_{ij} = \sum_{\substack{j=0 \\ j \neq i}}^N \pi_j p_{ji}$$



weights /  
(and transitions)



$$P_{xy} = P(X_{t+1} = y | X_t = x)$$

$$= \begin{cases} 1 - x/N & y = x + 1 \\ x/N & y = x - 1 \\ 0 & \text{otherwise.} \end{cases}$$

(c)

Invariant distribution  $\pi$ :  $(\det, \pi = [\pi_0, \pi_1, \dots, \pi_N])$

$$\pi P = \pi$$

Using the balance equation from Part (a)  $[\pi_i \sum_{j \neq i} P_{ij}$

$$= \sum_{j \neq i} \pi_j P_{ji}]$$

State 0:  $\pi_0 \cdot \left(\frac{1}{N}\right) = \frac{1}{N} \pi_1 \Rightarrow (\pi_1 = N \pi_0)$

(prob leaving 0)

State 1:  $\pi_1 \left( \frac{N-1}{N} + \frac{1}{N} \right) = \pi_0 \cdot 1 + \pi_2 \cdot \frac{2}{N}$

$$\Rightarrow \pi_1 = \pi_0 + \frac{2\pi_2}{n}$$

$$\Rightarrow (N-1)\pi_0 = \frac{2\pi_2}{N} \Rightarrow \frac{N(N-1)}{2}\pi_0 = \pi_2$$

State 2:

$$\Rightarrow \frac{N(N-1)}{2} \pi_0 = (N-1) \pi_0 + \pi_3 \cdot \frac{3}{N}$$

$$\pi_3 = \frac{n(n-1)(n-2)}{2 \cdot 3} \pi_0 = \frac{n!}{(n-3)! \cdot 3!} \pi_0 = \binom{n}{3} \pi_0$$

1

•

$$\text{step } i : \quad \pi_i = \frac{n(n-1)(n-2) \dots (n-i+1)\pi_0}{i!} = \frac{n!}{(n-i)!i!} \pi_0$$

n law of total prob:

$$\sum_{i=1}^n \pi_i = 1$$

$$\Rightarrow \pi_0 \sum_{i=0}^N \binom{N}{i} = 1.$$

$$\pi_0 = \frac{1}{2^N} \quad \left. \begin{array}{l} \{\vdots \\ \pi_i = \frac{1}{2^N} \binom{N}{i} \end{array} \right\}$$

$$2) : - (8 \text{ pts})$$

## (Multi-armed Bandits)

(a)

$P(R_t | A_r = a)$  changes with time.

Given.

(b)

At  $(j=0)$   $\rightarrow$  ( $n$  steps in the past)

At ( $j=0$ )  $\rightarrow$  ( $n$  steps in the past)

weight :  $(1-\alpha)^n$

at  $j=1$  ;

weight :  $(1-\alpha)^{n-1}$  ( $n-1$  steps in the past)

↑ exponentially  
decreasing  
weightage

$j=n$  weight : 1 [associated weight with current step]

Theo: the maximum weightage is given to current reward, and it decreases exponentially in the past to discount for rewards from varied reward distribution.

(C) Convergence of Bandits:

$$\theta_{n+1} = \theta_n + \alpha_n \{ R_{n+1} - \theta_n \}$$

$\theta_\infty \rightarrow \theta_*$  (if step sizes satisfy Robin Munro)

$$\text{Consider : } \alpha_n = \frac{\log(n+2)}{n+2} \quad n \geq 0$$

$$(i) \quad \log(n+2) > 0 \quad \forall n \geq 0$$

$\therefore \alpha_n > 0$ . [Condition 1 - positive step sizes]

Condition 2 :

$$(ii) \quad \sum_n \alpha_n = \infty \quad [\text{Asymptotic behaviour of ODE}]$$

To show.  $\sum_n \log(n+2) = \infty$

To show,  $\sum_{n=0}^{\infty} \frac{\log(n+2)}{n+2} = \infty$

(Consider  $\log_2$ ; else scale accordingly)

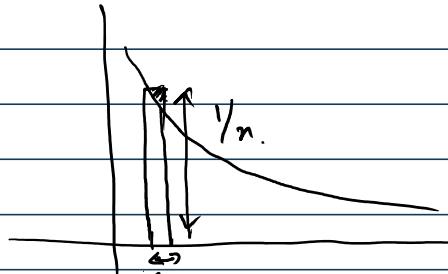
$$\frac{1}{n+2} \leq \frac{\log(n+2)}{n+2}$$

$\forall n \geq 0$ ; as:

$$\log_2(n+2) \geq 1$$

$$\Rightarrow \sum_{n=0}^{\infty} \frac{1}{n+2} \leq \sum_{n=0}^{\infty} \frac{\log_2(n+2)}{n+2}$$

$$\Rightarrow \sum_{n=0}^{\infty} \frac{1}{n+2} \geq \int_0^{\infty} \frac{1}{x+2} dx.$$



$$= \ln(n+2) \Big|_0^{\infty}$$

$$= \infty$$

$$\therefore \infty < \sum_{n=0}^{\infty} \frac{\log(n+2)}{n+2}$$

$$\therefore \left( \sum_{n=0}^{\infty} a_n = \infty \right) \text{ proved}$$

(iii)  $\sum_{n=0}^{\infty} \left[ \frac{\ln(n+2)}{n+2} \right]^2 < \infty$

Consider the sum:  $s_{2k+1} = \sum_{n=2}^{2k+1} \left[ \frac{1}{n^p} \right]$

$$\leq 1 + \sum_{n=1}^K \left[ \frac{1}{(2n)^p} + \frac{1}{(2n+1)^p} \right] \leq \frac{1}{(2n)^p}$$

$$\leq 1 + \sum_{n=1}^{\frac{K}{2}} \frac{2}{(2n)^p}$$

$$= 1 + 2^{1-p} \sum_{n=1}^{\frac{K}{2}} \underbrace{\frac{1}{n^p}}$$

$$\left( \frac{\log x}{x} \right)^2 = \frac{1}{x}.$$

(more terms)

$\downarrow$   
(want work.)

$$\left( \text{But } \frac{1}{x^{1+\epsilon}} \text{ } x > 0 \text{ will} \right) \leq 1 + 2^{1-p} S_{2k+1}$$

$$\Rightarrow S_{2k+1} \leq \frac{1}{1 - 2^{1-p}} < \infty \quad (\forall p > 1) \quad *$$

We know:  $\log x \leq \sqrt{x}$

$$= \lim_{\delta \rightarrow 0} \frac{(\sqrt{x} + \delta)^2 - (\sqrt{x})^2}{\delta}$$

$$\text{Note: } \log x \leq \lim_{\delta \rightarrow 0} x^{\frac{1}{2+\delta}} \quad 2 \left[ \frac{1}{2+\delta} - 1 \right] = \left[ \frac{-(1+\epsilon)}{2+\delta} \right]^2$$

$$\Rightarrow \left( \frac{\log x}{x} \right)^2 \leq \lim_{\delta \rightarrow 0} \left( x^{\frac{-2-\delta}{2+\delta}} - x^{\frac{-\delta}{2+\delta}} \right)$$

$$\leq \frac{1}{x^{1+\epsilon}} \quad (\exists \text{ some } \epsilon)$$

$$\therefore \sum_{x=2}^{\infty} \left( \frac{\log x}{x} \right)^2 \leq \sum_{x=2}^{\infty} \frac{1}{x^{1+\epsilon}} < \infty \quad (\text{from *})$$

QED

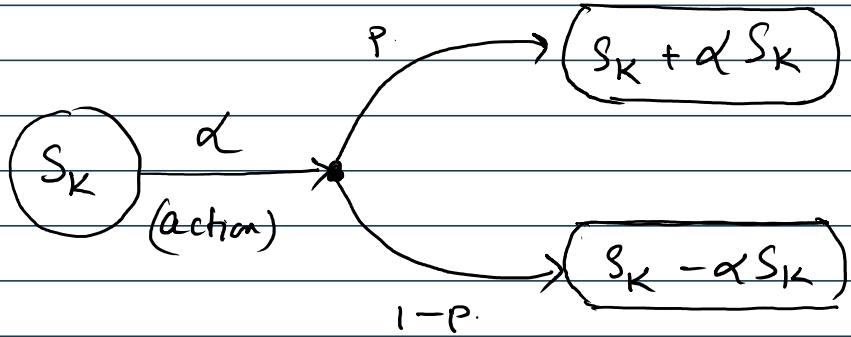
Since:  $\alpha_n = \left(\frac{\log n}{n}\right)^2$  is a valid step schedule.

3) Finite Horizon MDPs ; Model Based Method : (8 pts)

$$\begin{aligned} \text{Prob. of winning: } P \\ \text{Prob. of losing: } q = 1 - P \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{(at every step)}$$

(a) Let  $S_k$  be the amount of money the gambler has at the end of  $k$  round.

Let  $a_k$  be the action at time  $k$ .



Using DP:

Let  $V_k$  be the value function at step  $k$ .

$$V_N(S_N) = \begin{cases} \log S_N & S_N > 0 \\ 0 & \text{otherwise.} \end{cases}$$

(terminal reward)

$$S_{N-1} \rightarrow S_N$$

$$\text{(need } V_0(S_0) = \mathbb{E} \left[ \sum_{i=0}^N R_i \mid S_0 = s_0 \right] \text{)}$$

$$\begin{aligned}
 V_{N-1}(S_{N-1}) &= \mathbb{E} \left[ R_N^0 + v_N(s) \mid S_{N-1}, \alpha \right] \\
 &\quad S \sim P(\cdot \mid S_{N-1}, \alpha_{N-1}) \\
 &= \sup_{\alpha_{N-1}} \left[ p V_N((1+\alpha_{N-1})S_{N-1}) + q V_N((1-\alpha_{N-1})S_{N-1}) \right] \\
 &= \sup_{\alpha_{N-1}} \left[ p \log(1+\alpha_{N-1}) S_{N-1} + (1-p) \log(1-\alpha_{N-1}) S_{N-1} \right]
 \end{aligned}$$

(b)

$$\frac{\partial V_{N-1}}{\partial \alpha_{N-1}} = \left[ \frac{p}{(1+\alpha_{N-1})} - \frac{(1-p)}{(1-\alpha_{N-1})} \right] \frac{1}{S_{N-1}} = 0$$

$$S_{N-1} > 0.$$

$$\Rightarrow \frac{1+\alpha_{N-1}}{1-\alpha_{N-1}} = \frac{p}{1-p}.$$

$$\Rightarrow \frac{1+\alpha_{N-1}}{2} = p.$$

$$\Rightarrow \alpha_{N-1} = 2p - 1 = 2p - p - q$$

$$(\underbrace{\alpha_{N-1} = p - q}_{\text{independent of } N})$$

independent of  $N$ .

$$\frac{\partial^2 V_{N-1}}{\partial \alpha_{N-1}^2} = -\frac{p}{(1+\alpha_{N-1})^2} - \frac{(1-p)}{(1-\alpha_{N-1})^2} < 0$$

$\alpha_{N-1} = p - q$  (maximizes the profit at step  $N-1$ )

$> 0$

$$\begin{aligned}
 v_{n-1}(s_{n-1}) &= P \log((1+p-q)s_{n-1}) + (1-p) \log((1-p+q)s_{n-1}) \\
 &= P \log_2 p s_{n-1} + (1-p) \log_2 q s_{n-1} \\
 &= \log(2q s_{n-1}) + P \log\left(\frac{P}{q}\right) \\
 &= \underbrace{\log 2q + p \log \frac{P}{q}}_{(\text{const independent of } s_{n-1})} + \log s_{n-1} \\
 &= c_{n-1} + \log s_{n-1}
 \end{aligned}$$

Similarly:

$$\therefore v_{n-2}(s_{n-2}) = \sup_{\alpha_{n-2}} \left[ P \left\{ c_{n-1} + \log((1+\alpha)s_{n-2}) \right\} + (1-p) \left\{ c_{n-1} \log((1-\alpha)s_{n-2}) \right\} \right]$$

$$\therefore \frac{\partial v_{n-2}}{\partial \alpha_{n-2}} = 0 \quad (\text{constant drop off})$$

$$\Rightarrow [\alpha_{n-2} = p-q] \quad (\text{Independent of } n)$$

⋮

Similarly: strategy for all  $k$ :  $\boxed{\alpha_k = p-q}$

(C)

$$P < q$$

Now,

$$V_{N-1}(S_{N-1}) = \sup_{\alpha_{N-1}} \left[ P \log \left( \frac{1+\alpha}{n} \right) S_{N-1} + (1-P) \log \left( \frac{1-\alpha}{n} \right) S_{N-1} \right]$$

$$= \sup_{\alpha_{N-1}} \left[ P \log S_{N-1} + P \log (1+\alpha_{N-1}) + q \log (1-\alpha_{N-1}) + q \log S_{N-1} \right]$$

$$= \sup_{\alpha_{N-1}} \left[ \log S_{N-1} + P \log (1+\alpha_{N-1}) + q \log (1-\alpha_{N-1}) \right]$$

$$\frac{\partial V_{N-1}}{\partial \alpha_{N-1}} = \frac{P}{1+\alpha_{N-1}} - \frac{q}{1-\alpha_{N-1}}$$

$$= \frac{(1-\alpha_{N-1})P - (1+\alpha_{N-1})q}{1-\alpha_{N-1}^2}$$

$$= \frac{(P-q) - \alpha_{N-1}(P+q)}{1-\alpha_{N-1}^2}$$

$$= \frac{(P-q) - \alpha_{N-1}^{>0} \xrightarrow{\epsilon[0,1]} 0}{1-\alpha_{N-1}^2} < 0$$

$$\therefore \frac{\partial V_{N-1}}{\partial \alpha_{N-1}} < 0 \quad (\text{decreasing function})$$

[ maximum at  $\alpha_{n+1} = 0$  ] : (Don't bet at any time t)

PROGRAMMING : 16 points)

- 4) Gradient based Bandits }  
5) Value Iteration } Assignment 1.ipynb