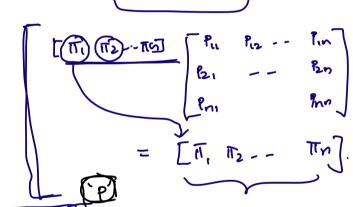Recap :-    Markov chains.
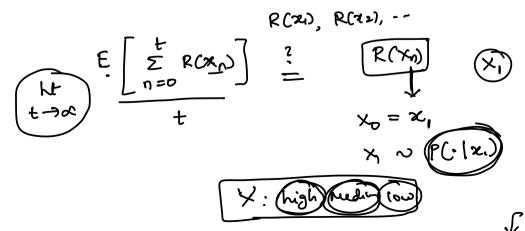
1. Irreducible
2. Periodicity
3. Recurrent (positive).

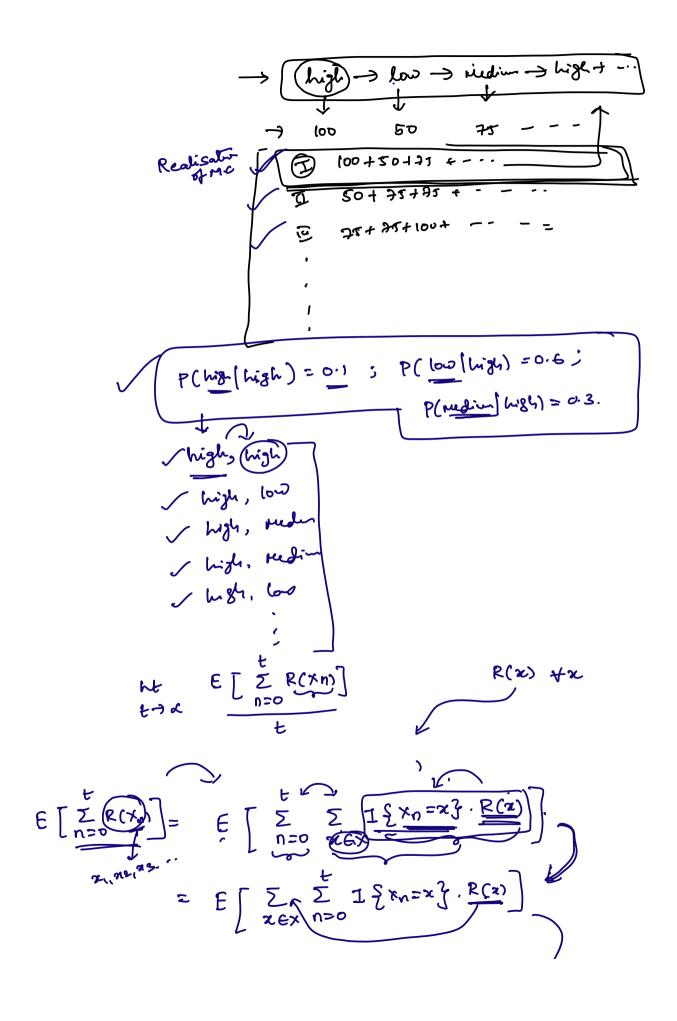• Ergodic M.C    ① Irreducible
                ② Aperiodic
                ③ Positive recurrent.

• Stationary distribution :

$\pi$ S.D w.r.t M.C with $P$ :-
$\pi <$
$\eta \neq 1.$

$P \to n \times n.$

$$P\pi = \pi \longrightarrow$$
$$P^T \pi = \pi$$

$$\begin{bmatrix} \boxed{\begin{matrix} & \\ & \end{matrix}} \end{bmatrix} \left[ \textcircled{$\pi_1$}\, \textcircled{$\pi_2$} \cdots \pi_{(n)} \right] \begin{bmatrix} P_{11} & P_{12} & \cdots & P_{1n} \\ P_{21} & & \cdots & P_{2n} \\ P_{n1} & & & P_{nn} \end{bmatrix}$$

$$= \begin{bmatrix} \pi_1 & \pi_2 & -- & \pi_n \end{bmatrix}.$$

$\textcircled{P}$

(✲). consider an $\boxed{\text{ergodic M.C}}$ with finite states $x_1, x_2, \cdots x_n.$

Define reward function $\underline{R : x \to \mathbb{R}}.$

$R(x_1), R(x_2), \cdots$

$\boxed{\text{lt}\atop t\to\alpha}$  $\dfrac{E\left[ \sum\limits_{n=0}^{t} R(x_n) \right]}{t}$  $\overset{?}{=}$  $\boxed{R(x_n)}$  $\textcircled{$x_1$}$

$x_0 = x_1$

$x_1 \sim \boxed{P(\cdot | x_1)}$

$\boxed{x : \textcircled{high}\ \textcircled{medium}\ \textcircled{low}}$

$\checkmark$

$\rightarrow$ | (high) $\rightarrow$ low $\rightarrow$ medium $\rightarrow$ high + ...

$\rightarrow$ 100     50     75   - - - -

Realisation of MC

(I)   $100 + 50 + 75 + \cdots$

(II)   $50 + 75 + 95 + \cdots$

(iii)   $75 + 75 + 100 + \cdots =$

$\vdots$

$P(\text{high} \mid \text{high}) = 0.1$ ;   $P(\text{low} \mid \text{high}) = 0.6$ ;

$P(\text{medium} \mid \text{high}) = 0.3$.

$\checkmark$ high, (high)

$\checkmark$ high, low

$\checkmark$ high, medium

$\checkmark$ high, medium

$\checkmark$ high, low

$\vdots$

$\lim\limits_{t \to \alpha} \dfrac{E\left[\sum\limits_{n=0}^{t} R(X_n)\right]}{t}$

$R(x) \; \forall x$

$E\left[\sum\limits_{n=0}^{t} R(X_n)\right] = E\left[\sum\limits_{n=0}^{t} \sum\limits_{x \in X} \mathbb{I}\{X_n = x\} \cdot R(x)\right]$

$x_1, x_2, x_3 \ldots$

$= E\left[\sum\limits_{x \in X} \sum\limits_{n=0}^{t} \mathbb{I}\{X_n = x\} \cdot R(x)\right]$

$$= E\left[ \sum_{x \in X} \boxed{R(x)} \cdot \sum_{n=0}^{t} I\{x_n = x\}\right] \swarrow$$

$$= \sum_{x \in X} \underline{R(x)} \sum_{n=0}^{t} \underline{Pr\{x_n = x\}}.$$

---

$$I\{^x A\} = 1, \quad \text{if } x \in A$$
$$= 0, \quad \text{if } x \notin A.$$

$$I\{x_n = x\} = 1, \quad \text{if } x_n = x$$
$$= 0, \quad \text{if not}$$

$$\underline{E\left[ I\{x_n = x\}\right]} = 1 \cdot Pr\{x_n = x\} + 0 \cdot Pr\{x_n \neq x\}$$

$$\quad \rightarrow \quad = \underline{\underline{Pr\{x_n = x\}}}.$$

$$(*). \quad \frac{E\left[\sum_{n=0}^{t} R(x_n)\right]}{t} = \frac{\sum_{x \in X} R(x) \sum_{n=0}^{t} Pr\{x_n = x\}}{t}.$$

$$= \sum_{x \in X} R(x) \cdot \frac{\sum_{n=0}^{t} Pr\{x_n = x\}}{t}$$

$$(*) \quad \lim_{t \to \alpha} \frac{\sum_{n=0}^{t} Pr\{x_n = x\}}{t} = \underline{\underline{\pi(x)}}. \quad \circledast.$$

$$\therefore \quad \lim_{t \to \alpha} \frac{E\left[\sum_{n=0}^{t} R(x_n)\right]}{t} = \sum_{\boxed{x \in X}} \underline{R(x)} \, \underline{\pi(x)};$$

$\square$

stationary distⁿ :-  $\boxed{\pi P = \pi}$   $\left[ \underbrace{\pi(1)} \, \underbrace{\pi(2)} \, \cdots \, \pi(n)\right].$

$\underset{1 \times n \quad n \times n \quad 1 \times n.}{}$

**Multi-arm Bandits:-**

(*). A bandit with "k" arms.

Pay          might    Reward

$50 \rightarrow$ Arm 1 $\rightarrow$ 500 Rs, 0

$\underline{1} \rightarrow$ Arm 2 $\rightarrow$ 2 Rs, 0

0 $\rightarrow$ Arm 3 $\rightarrow$ 1 Rs, 0

$X_3 = 1$
              0

✓ Mode strategy :- Arm 1.

$\rightarrow$ $\dfrac{0+0+0+1+0+\cdots}{n}$ ✓ $\rightarrow Pr\{X_3=1\}$

$\hookrightarrow E[X_3]$

$\rightarrow$ Mean strategy

" Find the arm with highest expected value "

$\hookrightarrow$ optimal arm.

(*) Arm $\hookleftarrow$

$\begin{cases} \boxed{R} \text{ is reward you get when arm 1 is pulled.} \\ \text{Compute } E[R] \overset{?}{=} \end{cases}$

$t=1 ; \gamma_1$

$t=2 ; \gamma_2$

$t=3 ; \gamma_3$  $\underset{n\to\alpha}{Lt} \left(\dfrac{\gamma_n}{n}\right) = E[R]$

$\dfrac{\gamma_n}{n} :$

Naive approach

$\gamma_1, \gamma_2,$

* | $\gamma_1$ | $\gamma_2$ | $\gamma_3$ | $\cdots$ | $\gamma_{100}$ | $\hookleftarrow$

# Running average :-

$Q_n$ : Average of rewards obtained until time `n+1`.

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^{n} r_i$$

$$\vdots$$

$$= Q_n + \frac{1}{n}[r_n - Q_n]$$

$Q_0 = a$ ~~(crossed out)~~
$Q_1$ ~~(crossed out)~~
$Q_2$
$\vdots$

In general:

$$Q_{n+1} = Q_n + \alpha_n [R_n - Q_n].$$

$$\sum_{n=1}^{\alpha} \alpha_n = \infty \quad \& \quad \sum_{n=1}^{\alpha} \alpha_n^2 < \alpha. \qquad \alpha_n > 0$$

$$\sum_{n=1}^{\alpha} \alpha_n^2$$

Seqm

$$\left(\alpha_1^2\right), \quad \left(\alpha_1^2 + \alpha_2^2\right), \quad \left(\alpha_1^2 + \alpha_2^2 + \alpha_3^2\right) \cdots$$

$$\{n\}$$

$$1, 2, 3, 4. --$$

① optimal arm : arm with highest expected reward

②: Efficient scheme to compute expected values.

## Strategies :-

① $\varepsilon$-greedy :-    'n' arms.

t: $\boxed{Q_t(1)}$ ; $\boxed{Q_t(2)} \cdots \boxed{Q_t(n)}.$  $\boxed{Q_t(5) \text{ is highest.}}$

$(t+1)$

$\Bigg\{$ Pull any random arm co P $\underset{}{\epsilon}$ =0.0001 (exploration) $\Bigg\}$

Pull | best arm so far | $1-\epsilon$ (exploitation).

$\rightarrow$ $\underset{a}{\arg\max}\ Q_t(a)$ $\epsilon$

| Intelligent exploration " |

$\Bigg\{$ Arm 1 $\Rightarrow$
Arm 2 $\Rightarrow$

exploration

UCB :- (action selection scheme) :-

$\circledast$ $A_t \doteq \underset{a}{\arg\max}\ \left[ Q_t(a) + C\ \sqrt{\dfrac{\ln t}{N_t(a)}} \right]$ $\rightarrow$ un certainty

$\downarrow$ exploit

$Q_t(a)$
$N_t(a)$ $\Bigg\}$

$\rightarrow$ no. of times until 't' arm 'a' is pulled.

$\textcircled{I}$ : Arm 1 :- $Q_t(1) \approx Q_t(2)$.
Arm 2 : $N_t(1) << N_t(2)$

Arm 1 | why ? |

$\textcircled{II}$ : Arm 1 : $N_t(1) = N_t(2)$
Arm 2 : $Q_t(1) >> Q_t(2)$

| Selection scheme | $\propto Q_t(\cdot)$ $\Bigg\}$
$\propto \dfrac{1}{N_t(\cdot)}$