

# Mathematics for Machine Learning (AI 512)

Amit Chattopadhyay

IIIT-Bangalore

Module 3



# Syllabus

## Topics

- M3: ✓ Markov Chains, Stationary Property and Random Walk
- M3: Page Rank Algorithm
- M3. Markov Chain Monte Carlo (MCMC) Sampling
- M4. Multivariate Distributions, K-Means
- M4. EM Algorithm
- M4. EM Algorithm for GMM.

## Reference Books

- 1. Foundations of Data Science. By John Hopcroft, Ravindran Kannan.
- ✓ 2. Finite Markov Chains and Algorithmic Applications. By OLLE HÄGGSTRÖM.
- ✓ 3. Introduction to Probability Models. By S.M. Ross.
- 4. Pattern Recognition and Machine Learning. By Christopher M. Bishop.
- ✓ 5. Mathematics for Machine Learning. By Marc Peter Deisenroth et al.

# **Markov Chain**

# Introduction: Stochastic Process

## Definition

A stochastic process indexed by an index set  $\mathbb{I}$ , is a collection of random variables  $\{X_t : t \in \mathbb{I}\}$  on a probability space  $(\Omega, \Delta, P)$ .

✓  $\Omega$  : event space/sample space

$\Delta$  :  $\sigma$ -algebra

$P$  : probability function

$$P: \Delta \rightarrow \mathbb{R}$$

## Examples:

1.  $X_n$  : The stock price of a commodity at the  $n$ -th week
2.  $X_t$  : Number of accidents in a city till time  $t$

# Introduction: Stochastic Process

- **State space:** Spectrum of  $X_n$  -  $\{s_1, s_2, \dots, s_k\} \subseteq \mathbb{R}$
- **Stage:** Index set  $\mathbb{I}$
- **Chain:** Both the index set and state space are discrete

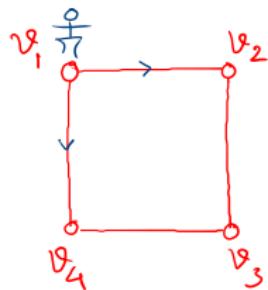
# Example 1

## A “random walker” in a small town

A random walker in a small town with four streets and four street corners  $v_1, v_2, v_3$  and  $v_4$ .

He flips a fair coin and moves forward (clockwise) if ‘H’ occurs and backward (or counter-clockwise) if ‘T’ occurs.

✓  $X_n$ : r.v. denoting the index of the street corner at which the walker stands at time  $n$ .



$$P(X_0 = 1) = 1, \quad P(X_0 = 2) = \dots = P(X_0 = 4) = 0$$

$$P(X_1 = 1) = 0, \quad P(X_1 = 2) = \frac{1}{2}, \quad P(X_1 = 3) = 0, \quad P(X_1 = 4) = \frac{1}{2}$$

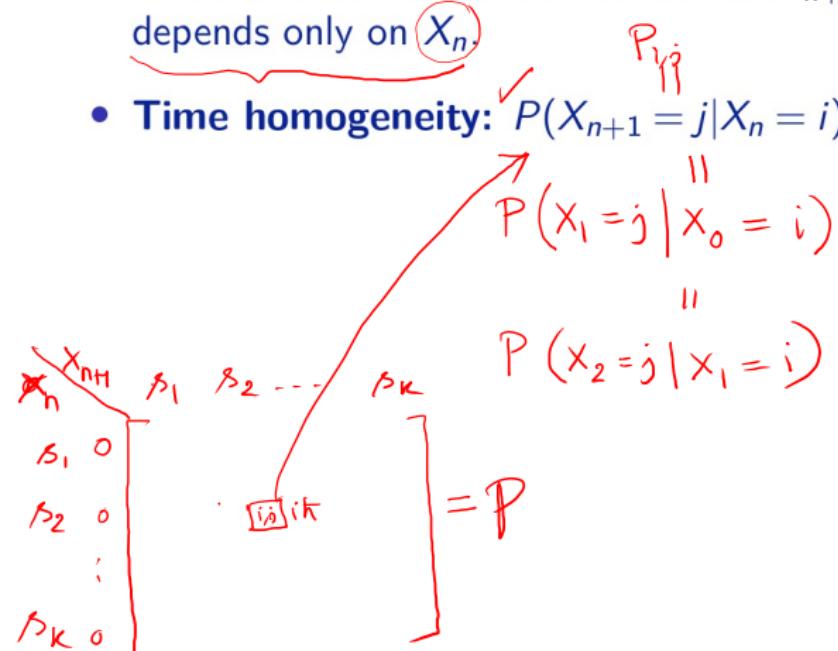
$$P(X_2 = 1 | X_0 = i_0, X_1 = 2) = \frac{1}{2}$$

$$P(X_2 = 3 | X_0 = i_0, X_1 = 2) = \frac{1}{2}$$

# Observations

Mankov Property

- **Memoryless property:** The coin flip at time  $n+1$  is independent of all previous coin flips and hence independent of  $X_0, X_1, \dots, X_n$ . Therefore, conditional distribution of  $X_{n+1}$  given  $(X_0, X_1, \dots, X_n)$  depends only on  $X_n$
- **Time homogeneity:**  $P(X_{n+1} = j | X_n = i)$  is the same for all  $n$



# Markov Chain

## Definition

Let  $P$  be a  $k \times k$  matrix with elements  $\{P_{i,j}: i, j = 1, 2, \dots, k\}$ . A stochastic process  $\{X_n : n = 0, 1, 2, \dots\}$  with finite state space  $S = \{s_1, s_2, \dots, s_k\}$  is said to be a (homogeneous) Markov Chain if  $\forall n, \forall i, j \in \{1, 2, \dots, k\}$  and  $\forall i_0, \dots, i_{n-1} \in \{1, 2, \dots, k\}$

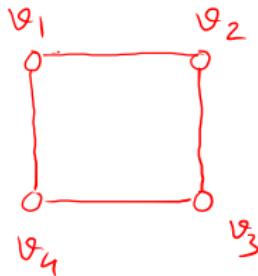
$$\begin{aligned} P(X_{n+1} = s_j | X_0 = s_{i_0}, X_1 = s_{i_1}, \dots, X_{n-1} = s_{i_{n-1}}, X_n = s_i) \\ = P(X_{n+1} = s_j | X_n = s_i). \end{aligned}$$

- Memoryless property
- $P(X_{n+1} = s_j | X_n = s_i)$  will be denoted by  $\underbrace{P_{i,j}^{n,n+1}}$  or  $\underbrace{P_{i,j}}$   
(Probability of  $X_{n+1}$  being in state  $s_j$ , given that  $X_n$  is in state  $s_i$ )

# Transition Matrix

**Transition matrix:**  $k \times k$  matrix whose  $(i,j)$ -th entry is the transition probability  $P(X_{n+1} = s_j | X_n = s_i) = P_{i,j}$ .

**Ex 1:** (Random Walker)



$$P = \begin{bmatrix} & \cancel{X_{n+1}} & v_2 & v_3 & v_4 \\ v_1 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ v_2 & \frac{1}{2} & 0 & \boxed{\frac{1}{2}} & 0 \\ v_3 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ v_4 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \end{bmatrix} \rightarrow P(X_{n+1} = v_3 | X_n = v_2) = P_{2,3}$$

# Transition Matrix: Properties

$$1. P_{i,j} \geq 0$$

$$2. \sum_{j=1}^n P_{i,j} = 1$$

Note:  $(X_{n+1} = \beta_1), (X_{n+1} = \beta_2), \dots, (X_{n+1} = \beta_K)$   
 are mutually exclusive & exhaustive  
 set of events.

$$\therefore P(X_{n+1} = \beta_1 | X_n = \beta_i) + P(X_{n+1} = \beta_2 | X_n = \beta_i) + \dots + P(X_{n+1} = \beta_K | X_n = \beta_i)$$

$$\checkmark = \frac{P(X_{n+1} = \beta_1, X_n = \beta_i) + \dots + P(X_{n+1} = \beta_K, X_n = \beta_i)}{P(X_n = \beta_i)}$$

$$= \frac{P(X_n = \beta_i)}{P(X_n = \beta_i)} \stackrel{P(X)}{\geq} 1$$

(i)  $A_1, A_2, \dots, A_K$  are  
 mutually exclusive & exhaustive  
 set of events.

(ii)  $X$  is any event

$$\begin{aligned} &= P(XA_1) + P(XA_2) + \dots + P(XA_K) \\ &= P(A_1)P(X|A_1) + \dots + P(A_K)P(X|A_K) \end{aligned}$$

$$P(X) = P(X \cap S) = P(XA_1 + XA_2 + \dots + XA_K)$$

For construction of a Markov Chain (MC)  $\{X_n : n = 0, 1, 2, \dots\}$  we require:

✓ 1. **Initial distribution:** of  $X_0$

2. **Transition matrix:** consisting of transition probabilities

$$P(X_{n+1} = j | X_n = i) = P_{i,j}$$

# Distribution of $X_n$

Notation: Distribution of  $X_n$



$$\left( \mu_1^{(0)}, \mu_2^{(0)}, \dots, \mu_k^{(0)} \right) = \mu^{(0)}$$

✓  $\mu^{(n)} = (\mu_1^{(n)}, \mu_2^{(n)}, \dots, \mu_k^{(n)})$

$$= (P(X_n = s_1), P(X_n = s_2), \dots, P(X_n = s_k))$$

$$\mu^{(0)} P = \mu^{(1)}$$

✓  $\mu^{(0)}$  : Initial distribution

$$\mu^{(0)} P^n = \dots = \mu^{(n-1)} P P = \mu^{(n)} P = \mu^{(n)}$$

Theorem

For a MC  $\{X_n : n = 0, 1, 2, \dots\}$  with state space  $\{s_1, s_2, \dots, s_k\}$ , initial distribution  $\mu^{(0)}$  and transition matrix  $P$ , we have:

$$\boxed{\mu^{(n)} = \mu^{(0)} P^n}$$

$$\boxed{\mu = \mu P}$$

for  $n = 0, 1, 2, \dots$

Stationary Distribution

$$\underline{n=1} \quad \text{To show } \boxed{\mu^{(1)} = \mu^{(0)} P}$$

$$\mu_j^{(1)} = P(X_1 = \beta_j) = \sum_{i=1}^k P(X_0 = \beta_i, X_1 = \beta_j)$$

[Note:  $(X_0 = \beta_1), (X_0 = \beta_2), \dots, (X_0 = \beta_k)$  are m.e & exhaustive set of events]

$$= \sum_{i=1}^k P(X_0 = \beta_i) P(X_1 = \beta_j | X_0 = \beta_i)$$

$$\mu_j^{(1)} = \sum_{i=1}^k \mu_i^{(0)} P_{i,j} = (\mu^{(0)} P)_j \quad \begin{matrix} \text{- } j\text{-th element} \\ \text{of the row} \end{matrix}$$

$$\Rightarrow \boxed{\mu^{(1)} = \mu^{(0)} P}$$

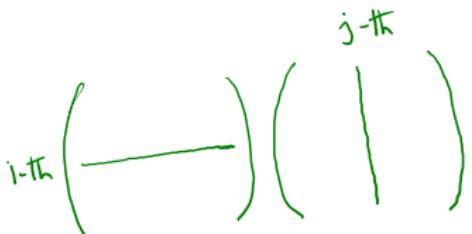
Ⓐ Use induction to complete the proof.

$\nabla P_{i,j}^{(n)} = n\text{-step transition probabilities}$

$$= P(X_n = j | X_0 = i)$$

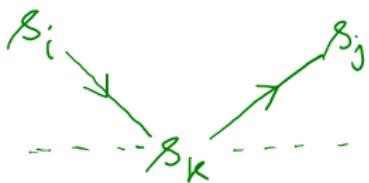
$$= P(X_{n+k} = j | X_k = i)$$

for  $n \geq 0, i, j, k \geq 0$ . In particular,  $P_{i,j}^{(1)} = P_{i,j}$



## Chapman-Kolmogorov Equation

$$P_{i,j}^{(m+n)} = \sum_{k=0}^{\infty} P_{i,k}^{(n)} P_{k,j}^{(m)}, \forall m, n \geq 0.$$



$$\begin{aligned} P_{i,j}^{(2)} &= \sum_{k=0}^k P_{i,k}^{(1)} P_{k,j}^{(1)} \\ &= \sum_{k=0}^k P_{i,k} P_{kj} \\ &= (P P)_{(i,j)} \end{aligned}$$

$$P_{ij}^{(m+n)} = P(X_{n+m} = \beta_j \mid X_0 = \beta_i)$$

$$\stackrel{\infty}{\underset{k=0}{\Rightarrow}} P(X_{n+m} = \beta_j, X_n = \beta_k \mid \underbrace{X_0 = \beta_i})$$

$\nwarrow \left\{ (X_n = \beta_k) : k = 0, 1, \dots \right\}$  are m.e. & exhaustive set of events.

$$= \sum_{k=0}^{\infty} P(X_n = \beta_k \mid X_0 = \beta_i) \underbrace{P(X_{n+m} = \beta_j \mid X_n = \beta_k, X_0 = \beta_i)}_{P(X_{n+m} = \beta_j \mid X_n = \beta_k)}$$

$$= \sum_{k=0}^{\infty} P(X_n = \beta_k \mid X_0 = \beta_i) P(X_{n+m} = \beta_j \mid X_n = \beta_k)$$

$$= \sum_{k=0}^{\infty} P_{ik}^{(n)} P_{kj}^{(m)}$$

## Note

- ✓ 1.  $P^{(n+m)} = P^{(n)} \cdot P^{(m)}$
- 2.  $P^{(2)} = P^{(1+1)} = P^{(1)} \cdot P^{(1)} = P \cdot P = P^2$
- 3. By induction,  $P^{(n)} = P^{(n-1+1)} = P^{n-1} \cdot P = P^n$

# Transition Graph

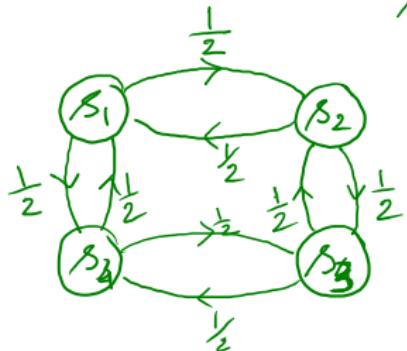
A useful way to picture a MC

✓ **Nodes:** representing the states of the MC

**Edges/Arrows** between the nodes: representing transition probabilities

Ex 1: (Random Walker)

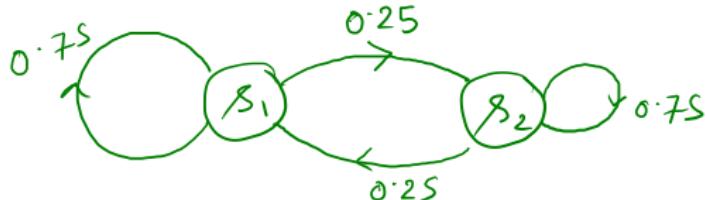
$$P = \begin{bmatrix} \beta_1 & \beta_2 & \beta_3 & \beta_4 \\ \beta_1 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \beta_2 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ \beta_3 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \beta_4 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \end{bmatrix}$$



## MC: Examples

Ex 2: (Gothenburg weather)

$$P = \begin{matrix} x_{n+1} \\ \beta_1 \\ \beta_2 \end{matrix} \begin{bmatrix} 0.75 & 0.25 \\ 0.25 & 0.75 \end{bmatrix}$$



## Examples

Ex 3: (Los Angeles weather)

$$P = \begin{bmatrix} \beta_1 & \beta_2 \\ \beta_2 & 0.1 \end{bmatrix}$$

