

# Outline

- On average,  $nH(X) + 1$  bits suffices to describe  $n$  i.i.d. random variables. But what if the random variables are dependent?
- Markov Chain**: a simplest way to model the correlations among random variables in a stochastic process.
- Entropy Rate**: average number of bits suffices to describe one random variable in a stochastic process.

## How to model Dependence: Markov chains

- A **stochastic process**  $\{X_i\}$  is an indexed sequence of random variables  $(X_1, X_2, \dots)$  characterized by the joint PMF  $p(x_1, x_2, \dots, x_n)$ , where  $(x_1, x_2, \dots, x_n) \in \mathcal{X}^n$  for  $n = 0, 1, \dots$

### Definition

A **stochastic process** is said to be **stationary** if the joint distribution of any subset of the sequence of random variables is **invariant** with respect to shifts in the time index, i.e.,

$$\begin{aligned}\Pr[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] \\ = \Pr[X_{1+\ell} = x_1, X_{2+\ell} = x_2, \dots, X_{n+\ell} = x_n]\end{aligned}$$

for every  $n$  and every shift  $\ell$  and for all  $x_1, x_2, \dots, x_n \in \mathcal{X}$ .

即每个变量同分布. (但不一定独立)

## Markov Chains

### Definition

A discrete stochastic process  $X_1, X_2, \dots$  is said to be a **Markov chain** or a **Markov process** if for  $n = 1, 2, \dots$ ,

$$\begin{aligned}\Pr[X_{n+1} = x_{n+1} | X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1] \\ = \Pr[X_{n+1} = x_{n+1} | X_n = x_n]\end{aligned}$$

for all  $x_1, x_2, \dots, x_n, x_{n+1} \in \mathcal{X}$ .

给定系统当前状态, 系统的未来状态与过去状态无关.  
也即过去状态均包含在了当前状态中.

In this case, the joint PMF can be written as

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2|x_1)p(x_3|x_2) \cdots p(x_n|x_{n-1}) = p(x_1) \prod_{i=1}^{n-1} p(x_{i+1}|x_i)$$

Hence, a Markov chain is completely characterized by **initial distribution**  $p(x_1)$  and **transition probabilities**  $p(x_n|x_{n-1})$ ,  $n = 2, 3, 4, \dots$

### Definition

The Markov chain is called **time invariant** if the transition probability  $p(x_{n+1}|x_n)$  does **NOT** depend on  $n$ , i.e., for  $n = 1, 2, \dots$ ,

$$\Pr[X_{n+1} = b | X_n = a] = \Pr[X_2 = b | X_1 = a], \quad \forall a, b \in \mathcal{X}.$$

无论哪个时刻, 即  $n=1, 2, \dots$ ,  $a \rightarrow b$  的状态转移概率不变.

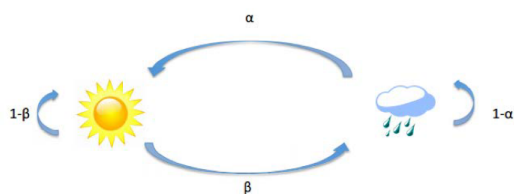
We deal with time invariant Markov chains. If  $\{X_i\}$  is a Markov chain,  $X_n$  is called the **state** at time  $n$ . A time invariant Markov chain is characterized by its initial state and a **probability transition matrix**  $P = [P_{ij}]$ ,  $i, j \in \{1, 2, \dots, m\}$ , where  $P_{ij} = \Pr[X_{n+1} = j | X_n = i]$ .

# Example: Simple Weather Model

- $\mathcal{X} = \{\text{Sunny: } S, \text{ Rainy: } R\}$

$$p(S|S) = 1 - \beta, p(R|R) = 1 - \alpha, p(R|S) = \beta, p(S|R) = \alpha$$

$$P = \begin{bmatrix} 1 - \beta & \beta \\ \alpha & 1 - \alpha \end{bmatrix}$$



- Probability of seeing a sequence SSRR:

$$p(SSRR) = p(S)p(S|S)p(R|S)p(R|R) = p(S)(1 - \beta)\beta(1 - \alpha)$$

Suppose the first day is "Sunny" with probability  $\gamma$ , what is the weather distribution of the second day, third day, ...?

- If  $\mu = [\mu_S, \mu_R] = \begin{bmatrix} \frac{\alpha}{\alpha + \beta} & \frac{\beta}{\alpha + \beta} \end{bmatrix}$

$$P = \begin{bmatrix} 1 - \beta & \beta \\ \alpha & 1 - \alpha \end{bmatrix}$$

$$\begin{aligned} p(X_{n+1} = S) &= p(S|S)\mu_S + p(S|R)\mu_R \\ &= (1 - \beta)\frac{\alpha}{\alpha + \beta} + \alpha\frac{\beta}{\alpha + \beta} = \frac{\alpha}{\alpha + \beta} = \mu_S. \end{aligned}$$

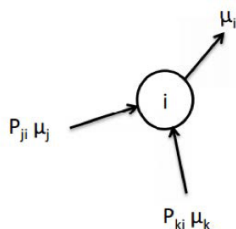
## Stationary Distribution

- If the PMF of the random variable at time  $n$  is  $\mu_i^n = \Pr[X_n = i]$ , the PMF at time  $n + 1$ , say  $\mu_j^{n+1} = \Pr[X_{n+1} = j]$ , can be written as

$$\mu_j^{n+1} = \sum_i \mu_i^n \Pr[X_{n+1} = j | X_n = i] = \sum_i \mu_i^n P_{ij}.$$

- $\{\mu_i^n | \forall i\}$  is called a **stationary distribution** if  $\mu_i^n = \mu_i^{n+1}, \forall i$ .
- For notation convenience, let  $\mu_i = \mu_i^n = \mu_i^{n+1}, \forall i$ .
- How to calculate stationary distribution?
  - Stationary distribution  $\mu_i, i = 1, 2, \dots, |\mathcal{X}|$  satisfies

$$\mu_j = \sum_{i=1}^{|\mathcal{X}|} \mu_i P_{ij} \text{ and } \sum_{i=1}^{|\mathcal{X}|} \mu_i = 1.$$



该系统是 time invariant Markov chain

设  $P(X_1=S)=a, P(X_1=R)=b$ , 则

$$P(X_2=S) = a(1-\beta) + b\alpha = (a \ b) \begin{pmatrix} 1-\beta \\ \alpha \end{pmatrix}$$

$$P(X_2=R) = a\beta + b(1-\alpha) = (a \ b) \begin{pmatrix} \beta \\ 1-\alpha \end{pmatrix}$$

$$\Rightarrow (P(X_2=S) \ P(X_2=R)) = (a \ b) P$$

$$\Rightarrow (P(X_n=S) \ P(X_n=R)) = (a \ b) P^n$$

因为该系统是 time invariant 的, 故

$$(a \ b) P = (a \ b)$$

$$\Rightarrow \begin{cases} a(1-\beta) + b\alpha = a \\ a\beta + b(1-\alpha) = b \end{cases} \Rightarrow \begin{cases} a = \frac{\alpha}{\alpha+\beta} \\ b = \frac{\beta}{\alpha+\beta} \end{cases}$$

$$\text{又 } a+b=1$$

推广:  $\mathcal{X} = \{x_1, \dots, x_m\}$

time invariant 需满足  $(\mu_1, \mu_2, \dots, \mu_m) = (\mu_1, \mu_2, \dots, \mu_m) P$

$$\text{即 } (\mu_1, \mu_2, \dots, \mu_m) (P - I) = 0$$

$$\text{又有归一化 } (\mu_1, \mu_2, \dots, \mu_m) e = 1$$

$$\text{则有 } (\mu_1, \mu_2, \dots, \mu_m) [P - I, e] = \begin{pmatrix} 0 & \dots & 0 & 1 \end{pmatrix}$$

$$\text{令 } \tilde{P} = [P - I, e], \text{ 并左右同乘伪逆 } \tilde{P}^\top (\tilde{P} \tilde{P}^\top)^{-1}$$

$$(\mu_1 \ \mu_2 \ \dots \ \mu_m) = (0 \ \dots \ 0 \ 1) \tilde{P}^\top (\tilde{P} \tilde{P}^\top)^{-1}$$

# Entropy Rate

- When  $X_i$ 's are i.i.d., the entropy

$$H(X^n) = H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i) = nH(X).$$

- With **dependent** sequences  $X_i$ 's, how does  $H(X^n)$  grow with  $n$ ?

- Entropy rate** characterized the growth rate.

- Definition 1:** average entropy per symbol

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{H(X_1, X_2, \dots, X_n)}{n}$$

- Definition 2:** conditional entropy of the last r.v. given the past

$$H'(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, X_{n-2}, \dots, X_1)$$

## Theorem 4.2.2

For a **stationary stochastic process**,  $H(X_n | X_{n-1}, \dots, X_1)$  is **nonincreasing in  $n$**  and **has a limit  $H'(\mathcal{X})$** .

## Proof.

$$H(X_{n+1} | X_1, X_2, \dots, X_n) \overset{\text{conditional reduces entropy}}{\leq} H(X_{n+1} | X_n, \dots, X_2) \\ \overset{\text{stationary}}{=} H(X_n | X_{n-1}, \dots, X_1),$$

- $H(X_n | X_{n-1}, \dots, X_1)$  **decreases as  $n$  increases**
- $H(X) \geq 0$
- The limit must exist. □

## Theorem 4.2.1

For a **stationary stochastic process**,  $H(\mathcal{X}) = H'(\mathcal{X})$ .

## Proof.

By the chain rule,

$$\frac{1}{n} H(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1).$$

- $H(X_n | X_{n-1}, \dots, X_1) \rightarrow H'(\mathcal{X})$
- Cesaro mean:** If  $a_n \rightarrow a$ ,  $b_n = \frac{1}{n} \sum_{i=1}^n a_i$ , then  $b_n \rightarrow a$ .
- So

$$\frac{1}{n} H(X_1, \dots, X_n) \rightarrow H'(\mathcal{X})$$
□

# AEP for Stationary Ergodic Process

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \rightarrow H(\mathcal{X})$$

- $p(X_1, \dots, X_n) \approx 2^{-nH(\mathcal{X})}$
- Typical sequences in typical set of size  $2^{-nH(\mathcal{X})}$
- We can use  $nH(\mathcal{X})$  bits to represent typical sequences

熵率描述了  $H(X^n)$  随  $n$  增长的比率

对平稳随机过程，上述两个定义相等



# Entropy Rate for Markov Chain

- For a **stationary Markov chain**, the entropy rate is

$$H(\mathcal{X}) = H'(\mathcal{X}) = \lim H(X_n | X_{n-1}, \dots, X_1) = \lim H(X_n | X_{n-1}) \\ = H(X_2 | X_1)$$

- Let  $P_{ij} = \Pr[X_2 = j | X_1 = i]$ . By definition, entropy rate of stationary Markov chain

$$H(\mathcal{X}) = H(X_2 | X_1) = \sum_i \mu_i \left( \sum_j -P_{ij} \log P_{ij} \right) \\ = - \sum_{ij} \mu_i P_{ij} \log P_{ij}$$

## Calculate Entropy Rate

- Find *stationary distribution*  $\mu_i$

$$\mu_i = \sum_j \mu_j p_{ji} \text{ and } \sum_{i=1}^{|\mathcal{X}|} \mu_i = 1$$

- User *transition probability*  $P_{ij}$

$$H(\mathcal{X}) = - \sum_{ij} \mu_i P_{ij} \log P_{ij}$$

## Entropy Rate of Weather Model

- Stationary distribution  $\mu(S) = \frac{\alpha}{\alpha+\beta}$ ,  $\mu(R) = \frac{\beta}{\alpha+\beta}$

$$P = \begin{bmatrix} 1-\beta & \beta \\ \alpha & 1-\alpha \end{bmatrix}$$

$$H(\mathcal{X}) = \mu(S)H(\beta) + \mu(R)H(\alpha) \\ = \frac{\alpha}{\alpha+\beta}H(\beta) + \frac{\beta}{\alpha+\beta}H(\alpha) \\ \text{Jensen's inequality} \\ \leq H\left(2\frac{\alpha\beta}{\alpha+\beta}\right)$$

**Maximum** when  $\alpha = \beta = 1/2$ : degenerate to independent process