

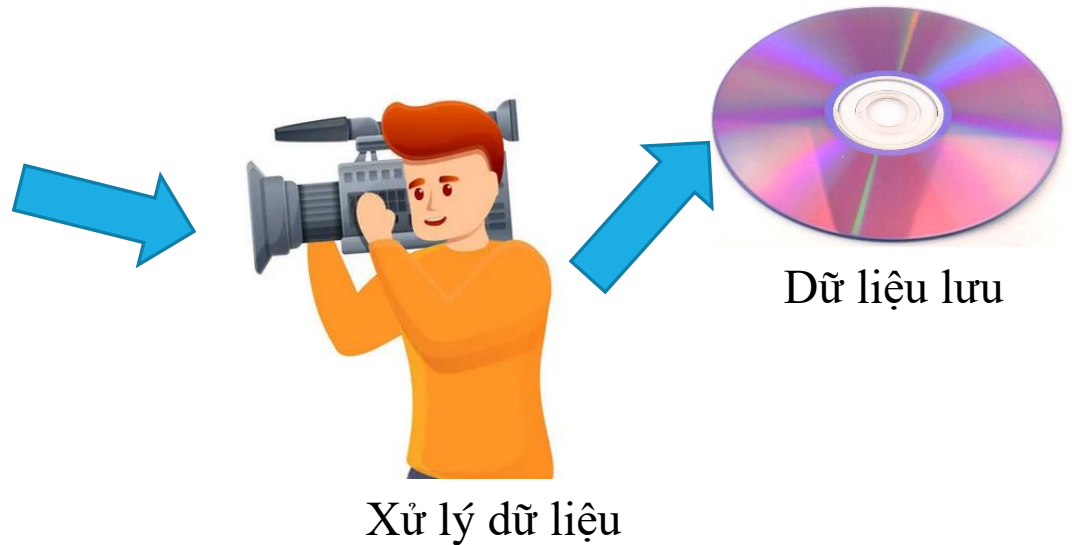
# CHƯƠNG 4: XỬ LÝ DỮ LIỆU

TS. TRỊNH VĂN CHIẾN (SOICT-HUST)

# HỆ THỐNG XỬ LÝ DỮ LIỆU



Khung cảnh



# CHUỖI MARKOV (1)

- Một chuỗi Markov được định nghĩa:

$$X \rightarrow Y \rightarrow Z$$

Nếu chúng ta có thể biểu diễn:

$$p(x, y, z) = p(z | y)p(y | x)p(x)$$

- Cách hiểu khác:

$$X \rightarrow Y \rightarrow f(Y)$$

# CHUỖI MARKOV (2)

- Ví dụ: Tung một con xúc xắc với xác suất xuất hiện mặt thứ  $m$  là  $\theta_m$ 
  - Giả sử chúng ta tung  $n$  lần được kết quả  $\{X_1, X_2, \dots, X_n\}$ .

- Tính giá trị trung bình của  $n$  lần tung

$$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$$

- Chuỗi Markov được định nghĩa như sau

$$\{\theta_1, \theta_2, \dots, \theta_6\} \rightarrow \{X_1, \dots, X_n\} \rightarrow \bar{X}_n$$

# HỆ QUẢ (1)

- Chuỗi Markov  $X \rightarrow Y \rightarrow Z \leftrightarrow X$  và  $Z$  độc lập nếu  $Y$  được cho trước

Chứng minh: Dựa vào xác suất có điều kiện ta có

$$p(x, z | y) = \frac{p(x, y, z)}{p(y)} = \frac{p(z | y)p(y | x)p(x)}{p(y)} = \frac{p(z | y)p(x, y)}{p(y)} = p(x | y)p(z | y)$$

- Hệ quả này có thể mở rộng cho chuỗi Markov gồm  $n$  chiều
- Hệ quả này có thể dùng để kiểm tra một chuỗi là chuỗi Markov

# HỆ QUẢ (2)

- Nếu chuỗi  $X \rightarrow Y \rightarrow Z$  là chuỗi Markov thì chuỗi  $Z \rightarrow Y \rightarrow X$  cũng là chuỗi Markov

Chứng minh:

$$\begin{aligned} p(x, y, z) &= p(x)p(y|x)p(z|y) = p(x)p(y|x)\frac{p(y, z)}{p(y)} \\ &= p(x, y)\frac{p(y|z)p(z)}{p(y)} = p(x|y)p(y)\frac{p(y|z)p(z)}{p(y)} = p(x|y)p(y|z)p(z) \end{aligned}$$

# BẤT ĐẲNG THỨC CHO XỬ LÝ DỮ LIỆU (1)

**Định lý:** Xem xét một chuỗi Markov  $X \rightarrow Y \rightarrow Z$ , ta có

$$I(X;Y) \geq I(X;Z); \quad I(Y;Z) \geq I(X;Z)$$

**Dấu bằng xảy ra khi**  $I(X;Y | Z) = 0$

- ❑ Nếu chúng ta xử lý thông tin, chúng ta sẽ bị mất thông tin
- ❑ Trong một vài trường hợp, đẳng thức về thông tin tương hỗ vẫn đạt được khi loại bỏ một số thông tin

# BẤT ĐẲNG THỨC CHO XỬ LÝ DỮ LIỆU (2)

- Chứng minh: Sử dụng định nghĩa thông tin tương hỗ và thông tin tương hỗ có điều kiện

$$I(X;Y,Z) = I(X;Z) + I(X;Y|Z) = I(X;Y) + \underbrace{I(X;Z|Y)}_{=0}$$

với  $I(X;Z|Y) = 0$  bởi vì  $X, Z$  là độc lập nếu  $Y$  được cho trước. Do đó ta có

$$I(X;Y) \geq I(X;Z)$$

Dấu bằng xảy ra nếu  $I(X;Y|Z) = 0$ , nghĩa là  $X \rightarrow Z \rightarrow Y$  là một chuỗi Markov.

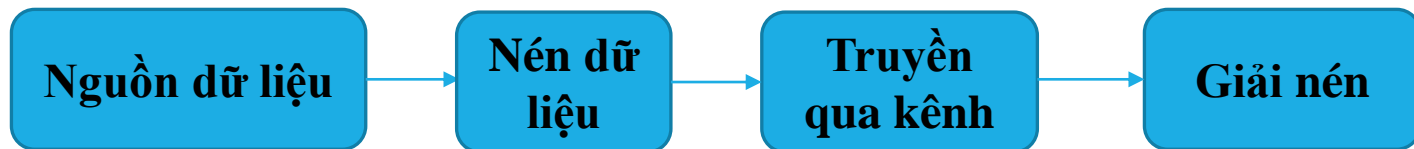
Tương tự chúng ta có thể chứng minh

$$I(Y;Z) \geq I(X;Z)$$



# MÔ HÌNH NÉN DỮ LIỆU

- Mô hình nén và giải nén giữ liệu



- Các bước thực hiện:

- ✓ Nén dữ liệu  $W$  từ nguồn sử dụng  $X = (X_1, \dots, X_n)$
- ✓ Truyền qua kênh truyền (ngẫu nhiên):  $Y$
- ✓ Giải nén thu được  $\hat{W}$
- ➔ Theo định lý về xử lý dữ liệu

$$I(W; \hat{W}) \leq I(X; Y)$$

# HỆ QUẢ CỦA BẤT ĐẲNG THỨC XỬ LÝ DỮ LIỆU

- Nếu hàm  $g$  được cho trước trong chuỗi Markov:  $X \rightarrow Y \rightarrow g(Y)$

$$I(X; Y) \geq I(X; g(Y))$$

- Xét chuỗi Markov  $X \rightarrow Y \rightarrow Z$ , ta có

$$I(X; Y | Z) \leq I(X; Y)$$

Chứng minh: Sử dụng thông tin tương hỗ và thông tin tương hỗ có điều kiện

$$I(X; Y, Z) = I(X; Z) + I(X; Y | Z) = I(X; Y) + \underbrace{I(X; Z | Y)}_{=0}$$

Quan sát:

- ✓ Sự phụ thuộc của  $X$  và  $Y$  sẽ giảm (hoặc không đổi) bằng việc quan sát  $Z$
- ✓ Nếu một tiến trình không tuân theo chuỗi Markov, có thể xảy ra  $I(X; Y | Z) \geq I(X; Y)$ 
  - ✓ Ví dụ: Tung 2 đồng xu ứng với các sự kiện  $X, Y$ . Đặt  $Z = X + Y$ , ta có:

$$I(X; Y | Z) = 1/2; I(X; Y) = 0$$

# THỐNG KÊ ĐỦ (SUFFICIENT STATISTICS) (1)

- Bất đẳng thức xử lý dữ liệu cung cấp thông tin về thống kê đủ (sufficient statistics)
- Cho một họ phân bố xác suất  $\{f_\theta(x)\}$  được định danh bởi  $\theta$
- $X$  là một mẫu dữ liệu (sampled) từ  $f_\theta(x)$
- $T(X)$  là hàm thống kê mẫu dữ liệu  $X$ , ta có

$$\theta \rightarrow X \rightarrow T(X)$$

- Bất đẳng thức xử lý dữ liệu:

$$I(\theta; T(X)) \leq I(\theta; X)$$

# THỐNG KÊ ĐỦ (SUFFICIENT STATISTICS) (2)

- Một thống kê là đầy đủ cho  $\theta$  nếu  $T(X)$  chứa tất cả thông tin của  $\theta$  trong  $X$

$$I(\theta; X) = I(\theta; T(X))$$

- Ví dụ:

- Nếu  $X = \{X_1, \dots, X_n\} \sim \text{Ber}(\theta)$ ,  $T(X) = \frac{1}{n} \sum_{i=1}^n X_i$  là thống kê đầy đủ cho  $\theta$
- Nếu  $X = \{X_1, \dots, X_n\} \sim \text{Uniform}(\theta, \theta + 1)$ ,  $T(X) = \{\min_i X_i, \max_i X_i\}$  là thống kê đầy đủ cho  $\theta$

- Lưu ý:

- ✓ Một thống kê đầy đủ có thể là một chiều hoặc nhiều chiều
- ✓ Một thống kê đầy đủ có thể không là duy nhất. Ví dụ: Bản thân  $X$  cũng là một thống kê đầy đủ của chính nó
- ✓ Thống kê đầy đủ cực tiểu: là một hàm của toàn bộ các thống kê đầy đủ khác  $\rightarrow$  Thông tin nén tối đa của  $\theta$  trong tập mẫu dữ liệu

# BẤT ĐẲNG THỨC FANO (1)

- Bất đẳng thức Fano (1942) liên kết xác suất lỗi  $P_e$  với lượng tin riêng
- Xác suất lỗi có quan hệ mật thiết với lượng tin riêng bởi vì: Phía phát truyền  $X$ , phía thu nhận  $Y$  và giải mã  $\hat{X}$  có thể xuất hiện lỗi

$$P_e = P(\hat{X} \neq X)$$

- Chuỗi Markov:  $X \rightarrow Y \rightarrow \hat{X}$
- Quan sát: Giải mã  $X$  từ  $Y$  với xác suất bằng 0 nếu  $H(X | Y) = 0$
- Bất đẳng thức Fano mở rộng quan sát trên: Giải mã  $X$  từ  $Y$  với xác suất lỗi nhỏ nếu  $H(X | Y)$  nhỏ

# BẤT ĐẲNG THỨC FANO (2)

**Định lý:** Với bất kỳ phiên bản giải mã  $\hat{X}$  thỏa mãn chuỗi Markov  $X \rightarrow Y \rightarrow \hat{X}$

$$H(P_e) + P_e \log |\tilde{X}| \geq H(\hat{X} | X) \geq H(X | Y)$$

• Hệ quả:

$$P_e \geq \frac{H(Y | X) - 1}{\log |\tilde{X}|} = \frac{H(X) - I(X; Y) - 1}{\log |\tilde{X}|}$$

• Nếu phương pháp giải mã  $g(Y)$  xem xét trực tiếp giá trị trong  $\tilde{X}$ , bất đẳng thức Fano sẽ viết lại

$$H(P_e) + P_e \log(|\tilde{X}| - 1) \geq H(X | Y)$$

# BÀI TẬP (1)

Bài 1: Cho hai biến ngẫu nhiên  $(X, Y)$  có xác suất xảy ra đồng thời  $p(x, y)$  như sau

$X$	$Y$		
	$a$	$b$	$c$
1	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{12}$
2	$\frac{1}{12}$	$\frac{1}{6}$	$\frac{1}{12}$
3	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{6}$

Giả sử  $\hat{X}(Y)$  là một phương pháp giải mã  $X$  từ  $Y$  với xác suất lỗi  $P_e = P(\hat{X}(Y) \neq X)$

a) Định nghĩa phương pháp giải mã  $\hat{X}(Y)$  và định nghĩa xác suất lỗi  $P_e$  cho phương pháp giải mã đề xuất

b) Xác định bao của  $P_e$  từ bất đẳng thức Fano và đưa ra kết luận

# BÀI TẬP (2)

a) Phương pháp giải mã đề xuất:

$$\hat{X}(Y) = \begin{cases} 1, & y = a \\ 2, & y = b \\ 3, & y = c \end{cases}$$

Xác suất lỗi Pe được định nghĩa như sau

$$P_e = P(1, b) + P(1, c) + P(2, a) + P(2, c) + P(3, a) + P(3, b) = 1/2$$

b) Sử dụng bất đẳng thức Fano:  $P_e \geq \frac{H(Y|X) - 1}{\log(|\tilde{X}| - 1)}$

$$H(X|Y) = 1.5 \text{ bits}$$

➔ Bất đẳng thức Fano  $P_e \geq \frac{1.5 - 1}{\log(3 - 1)} = 0.5$

Phương pháp giải mã khá tốt



# BÀI TẬP (3)

Bài 2: Xem xét ba biến ngẫu nhiên  $X, Y, Z$  có quan hệ phụ thuộc lẫn nhau. Chứng minh

*a)*  $H(X, Y | Z) \geq H(X | Z)$

*b)*  $I(X, Y; Z) \geq I(X; Z)$

*c)*  $H(X, Y, Z) - H(X, Y) \leq H(X, Z) - H(X)$

*d)*  $I(X; Z | Y) \geq I(Z; Y | X) - I(Z; Y) + I(X; Z)$

# BÀI TẬP (4)

a)  $H(X, Y|Z) = H(X|Z) + H(Y|X, Z) \geq H(X|Z)$

b) 
$$\begin{aligned} I(X, Y; Z) &= H(X, Y) - H(X, Y|Z) = H(X) + H(Y|X) - H(X|Z) - H(Y|X, Z) \\ &= I(X; Z) + H(Y|X) - H(Y|X, Z) \geq I(X; Z) \end{aligned}$$

c) Sử dụng quy tắc chuỗi

$$H(X, Y, Z) - H(X, Y) = H(Z | X, Y)$$

$$H(X, Z) - H(X) = H(Z | X)$$

d) Ta có:

$$\begin{aligned} I(X; Z|Y) - I(Y; Z|X) &= H(Z|Y) - H(Z|X, Y) - H(Z|X) + H(Z|X, Y) \\ &= H(Z|Y) - H(Z|X) \end{aligned}$$

$$I(X; Z) - I(Y; Z) = H(Z) - H(Z|X) - H(Z) + H(Z|Y) = H(Z|Y) - H(Z|X)$$

# TỔNG KẾT

- Xử lý dữ liệu: Có thể mất mát thông tin hoặc không
- Thống kê đủ (sufficient statistics) bảo toàn thông tin
- Khi ước lượng thông tin từ tín hiệu quan sát, lỗi mất mát thông tin có thể được bao bởi bất đẳng thức Fano