

Chủ đề 3: Lý thuyết Shannon

PGS.TS. Trần Minh Triết



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

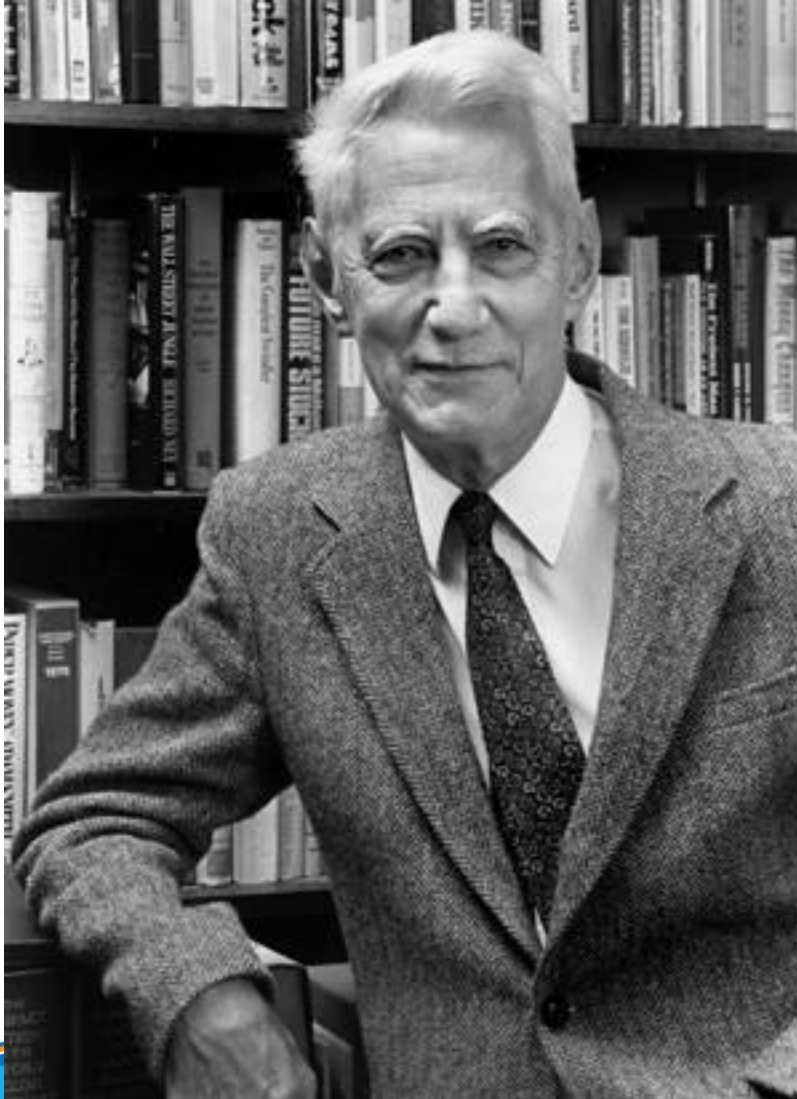
Tư liệu sử dụng

- Bài giảng này được xây dựng có sử dụng 1 số nội dung và tư liệu trên slide “***Communication Theory of Secrecy Systems***” của Gilad Tsur, Yossi Oren (Weizmann Institute of Science), Dec. 2005
- URL:
http://www.wisdom.weizmann.ac.il/~naor/COURSE/GiladTsur_YossiOren_ShannonSecrecy.ppt

Nội dung

- ☐ Mở đầu - Claude Shannon
- ☐ An toàn tuyệt đối
- ☐ Entropy
- ☐ Kết hợp các hệ thống mã hóa

Claude E. Shannon (1916-2001)



Nhắc lại về xác suất

□ Cho X và Y là hai biến ngẫu nhiên.

□ Định nghĩa:

$p(x) = p(X = x)$ là xác suất biến X nhận giá trị x

$p(y) = p(Y = y)$ là xác suất biến Y nhận giá trị y

$p(x | y)$ là xác suất biến X nhận giá trị x nếu quan sát được biến Y nhận giá trị y (xác suất có điều kiện)

X và Y là hai biến ngẫu nhiên độc lập khi và chỉ khi $p(x, y) = p(x) p(y)$ với bất kỳ giá trị x của biến X và bất kỳ giá trị y của biến Y

Nhắc lại: Định lý Bayes

- Cho X và Y là hai biến ngẫu nhiên.

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

- Định lý (Bayes):

$$\text{if } p(y) > 0 \text{ then } p(x|y) = \frac{\overset{\text{Apriori}}{p(x)p(y|x)}}{\underset{\text{Aposteriori}}{p(y)}}$$

- Hệ quả:

- X và Y là hai biến độc lập khi và chỉ khi $p(x | y) = p(x)$ với mọi giá trị x và y

- $p_P(x)$: Xác suất plaintext x xuất hiện
- $p_K(k)$: Xác suất chọn sử dụng giá trị khóa k
- $p_C(y)$: Xác suất ciphertext nhận giá trị y
- Có thể giả sử giá trị khóa K và plaintext x là các sự kiện độc lập

- Từ phân bố xác suất của plaintext và key trên tập \mathcal{P} và \mathcal{K} , ta có thể xác định phân bố xác suất của cipher text trên tập \mathcal{C} ???
- Với mỗi khóa $k \in K$, đặt $C(k) = \{e_k(x) \text{ với } x \in P\}$ là tập các giá trị cipher text có thể nhận được nếu mã hóa các plaintext $x \in P$ với giá trị khóa là k
- Vậy, ta có

$$p_C(y) = \sum_{\{k \in K \mid y \in C(k)\}} p_K(k) p_P(d_k(y))$$

- Với mỗi giá trị $y \in C$ và $x \in P$, xác suất có điều kiện $p_C(y | x)$, tức là xác suất nhận được cipher text y nếu cho biết cipher text là x , được tính theo công thức:

$$p_C(y | x) = \sum_{\{k \in K | x = d_k(y)\}} p_K(k)$$

- Áp dụng định lý Bayes, xác suất có điều kiện $p_P(x | y)$, tức là xác suất plaintext là x nếu có cipher text là y

$$p_P(x | y) = \frac{p_P(x) \sum_{\{k \in K | x = d_k(y)\}} p_K(k)}{\sum_{\{k \in K | y \in C(k)\}} p_K(k) p_P(d_k(y))}$$

Ví dụ

- $P = \{a, b\}$
- $p_P(a) = 1/4$, $p_P(b) = 3/4$
- $K = \{k_1, k_2, k_3\}$ với $p_K(k_1) = 1/2$, $p_K(k_2) = p_K(k_3) = 1/4$
- $C = \{1, 2, 3, 4\}$
- Cho $e_{k_1}(a) = 1, e_{k_1}(b) = 2$
 $e_{k_2}(a) = 2, e_{k_2}(b) = 3$
 $e_{k_3}(a) = 3, e_{k_3}(b) = 4$

	a	b
k_1	1	2
k_2	2	3
k_3	3	4

Ví dụ

- Phân bố xác suất của p_C
- $p_C(1) = 1/8$
- $p_C(2) = 3/8 + 1/16 = 7/16$
- $p_C(3) = 3/16 \times 1/16 = 1/4$
- $p_C(4) = 3/16$

	a $p_P = 1/4$	b $p_P = 3/4$
k_1 ($p_K = 1/2$)	1	2
k_2 ($p_K = 1/4$)	2	3
k_3 ($p_K = 1/4$)	3	4

Ví dụ

- Phân bố xác suất có điều kiện của $p_P(x/y)$

$$p_P(a | 1) = \frac{p_P(a) p_K(k_1)}{p_C(1)} = 1$$

$$p_P(b | 1) = \frac{p_P(b) p_K(k_1)}{p_C(1)} = 0$$

$$p_P(a | 2) = \frac{1/16}{7/16} = 1/7$$

$$p_P(b | 2) = \frac{6/16}{7/16} = 6/7$$

	a $p_P = 1/4$	b $p_P = 3/4$
k_1 ($p_K = 1/2$)	1	2
k_2 ($p_K = 1/4$)	2	3
k_3 ($p_K = 1/4$)	3	4

$$p_C(1) = 1/8$$

$$p_C(2) = 3/8 + 1/16 = 7/16$$

$$p_C(3) = 3/16 \times 1/16 = 1/4$$

$$p_C(4) = 3/16$$

Ví dụ

- Phân bố xác suất có điều kiện của $p_P(x/y)$

$$p_P(a | 3) = \frac{1/16}{1/4} = 1/4$$

$$p_P(b | 3) = \frac{3/16}{1/4} = 3/4$$

$$p_P(a | 4) = \frac{0}{3/16} = 0$$

$$p_P(b | 4) = \frac{3/16}{3/16} = 1$$

	a $p_P = 1/4$	b $p_P = 3/4$
k_1 ($p_K = 1/2$)	1	2
k_2 ($p_K = 1/4$)	2	3
k_3 ($p_K = 1/4$)	3	4

$$p_C(1) = 1/8$$

$$p_C(2) = 3/8 + 1/16 = 7/16$$

$$p_C(3) = 3/16 \times 1/16 = 1/4$$

$$p_C(4) = 3/16$$

An toàn tuyệt đối

- **An toàn tuyệt đối - Perfectly secure?**
- Ý nghĩa: Người tấn công không khai thác được gì từ ciphertext:

$$\forall p \in \mathcal{P}, \forall k \in \mathcal{K},$$

$$p(p|C) = p(p), p(k|C) = p(k)$$

Áp dụng – Shift Cipher

- Giả sử 26 giá trị khóa trong phương pháp Shift Cipher được chọn sử dụng với xác suất công bằng như nhau ($1/26$)
- Với tập Plaintext có phân bố xác suất bất kỳ, phương pháp Shift Cipher đạt được độ an toàn tuyệt đối???
- Ta có $P = C = K = Z_{26}$.
- $e_k(x) = (x + k) \bmod 26$ và $d_k(y) = (y - k) \bmod 26$

Áp dụng – Shift Cipher

□ Xác suất

$$\begin{aligned}
 p_C(y) &= \sum_{K \in \mathbb{Z}_{26}} p_K(K) p_P(d_K(y)) \\
 &= \sum_{K \in \mathbb{Z}_{26}} \frac{1}{26} p_P(y - K) \\
 &= \frac{1}{26} \sum_{K \in \mathbb{Z}_{26}} p_P(y - K).
 \end{aligned}$$

□ Với giá trị y cho trước, khi thay đổi giá trị k từ 0 đến 25, ta nhận được đầy đủ 26 giá trị của \mathbb{Z}_{26}

$$\begin{aligned}
 \sum_{K \in \mathbb{Z}_{26}} p_P(y - K) &= \sum_{y \in \mathbb{Z}_{26}} p_P(y) \\
 &= 1.
 \end{aligned}$$

Áp dụng – Shift Cipher

- Vậy với bất kỳ $y \in \mathbb{Z}_{26}$, ta luôn có $p_C(y) = 1/26$
- Với mỗi cặp giá trị (x, y) , ta có duy nhất 1 giá trị khóa $k \in \mathbb{Z}_{26}$ sao cho $y = x + k \bmod 26$. Do đó

$$p_C(y | x) = p_K(y - x \bmod 26) = 1/26$$

- Kết luận:

$$\begin{aligned} p_P(x|y) &= \frac{p_P(x)p_C(y|x)}{p_C(y)} \\ &= \frac{p_P(x) \frac{1}{26}}{\frac{1}{26}} \\ &= p_P(x), \end{aligned}$$

Áp dụng – Shift Cipher

- Vậy, phương pháp Shift Cipher có độ an toàn tuyệt đối nếu thật sự chọn khóa k mới ngẫu nhiên cho mỗi ký tự plaintext x cần mã hóa

Nhận xét

- Theo định lý Bayes:

$$p_P(x | y) = p_P(x), \forall x \in P, \forall y \in C$$

tương đương với

$$p_C(y/x) = p_C(y), \forall x \in P, \forall y \in C$$

- Có thể giả sử rằng $p_C(y) > 0, \forall y \in C$ (vì sao???)
- Vậy, nếu hệ thống thật sự an toàn

$$p_C(y/x) > 0, \forall x \in P, \forall y \in C$$

Nhận xét

- Giữ cố định 1 giá trị x , với mỗi giá trị $y \in C$, ta luôn có $p_C(y/x) > 0$. Suy ra có ít nhất một giá trị khóa k sao cho $e_k(x) = y$
- Vậy $|K| \geq |C|$
- Ta có $|C| \geq |P|$
- Vậy, ta có $|K| \geq |C| \geq |P|$

Vernam Cipher (1)

- Liệu có tồn tại 1 hệ thống mã hóa an toàn tuyệt đối với $|K|=|P|$?
- **Định lý (Shannon):** Cho (P, K, C, E, D) là một hệ thống mã hóa với $|K|=|P|=|C|$. Khi đó, hệ thống này an toàn tuyệt đối khi và chỉ khi:

$$\forall c \in C, p \in P, \exists k \in K \text{ s.t. } e_k(p) = c$$

$$\forall k \in K, p(k) = \frac{1}{|K|}$$

Vernam Cipher (2)

- **CM:** Cho (P, K, C, E, D) là một hệ thống mã hóa với $|K|=|P|=|C|$.
- Do an toàn tuyệt đối nên:

$$\begin{aligned} \forall p \in \mathcal{P}, p(p|c) &= p(p) \\ (\text{by Bayes}) \quad p(c|p) &= p(c) > 0 \\ \exists k \text{ s.t. } e_k(p) &= c \end{aligned}$$

- $|K|=|P|=|C|$, nên có duy nhất 1 khóa tương ứng với cặp (p, c) □

Vernam Cipher (3)

- Giữ cố định giá trị c . Với tất cả giá trị của p_i , gọi k_i là giá trị khóa thỏa $e_{k_i}(p_i)=c$
- Theo đ/lý Bayes:

$$p(p_i|c) = \frac{p(c|p_i)p(p_i)}{p(c)} \quad (1)$$

$$= \frac{p(k_i)p(p_i)}{p(c)} \quad (2)$$

$$(perfect\ secrecy) = p(p_i) \quad (3)$$

$$p(k_i) = p(c) \quad \forall i \quad (4)$$



Vernam Cipher (4)

- ☐ Do Gilbert Vernam (Bell Labs) đề nghị năm 1919
- ☐ Ý tưởng chính: khóa là 1 dãy giá trị ngẫu nhiên “đủ dài”.
Khi đó, $C = P \text{ XOR } K$
- ☐ Phương pháp này được chứng minh là an toàn tuyệt đối
- ☐ Hạn chế: khóa quá dài và không thể sử dụng lại
- ☐ Ưu điểm: giải thuật rất đơn giản

One-time pad

Let $n \geq 1$ be an integer, and take $\mathcal{P} = \mathcal{C} = \mathcal{K} = (\mathbb{Z}_2)^n$. For $K \in (\mathbb{Z}_2)^n$, define $e_K(x)$ to be the vector sum modulo 2 of K and x (or, equivalently, the exclusive-or of the two associated bitstrings). So, if $x = (x_1, \dots, x_n)$ and $K = (K_1, \dots, K_n)$, then

$$e_K(x) = (x_1 + K_1, \dots, x_n + K_n) \bmod 2.$$

Decryption is identical to encryption. If $y = (y_1, \dots, y_n)$, then

$$d_K(y) = (y_1 + K_1, \dots, y_n + K_n) \bmod 2.$$

Giới thiệu về Lý thuyết Thông tin

- ☐ Một số sự kiện ngẫu nhiên nhưng thường gặp hơn các sự kiện khác
- ☐ Một số dữ kiện quan trọng hơn các dữ kiện khác
- ☐ Entropy là độ đo mức độ bất định của một biến ngẫu nhiên, hay là lượng thông tin mà mỗi biến cố cung cấp

Định nghĩa Entropy

- Nếu X là một biến ngẫu nhiên nhận giá trị từ tập hữu hạn \mathcal{X} , khi đó

$$H(X) \triangleq - \sum_{x \in \mathcal{X}} p(X = x) \log_2 p(X = x)$$

- Ghi chú : $\lim_{x \rightarrow 0} x \log x = 0$

Entropy và nén Huffman

- Nhắc lại ý tưởng của giải thuật nén Huffman
- Giả sử có 5 sự kiện $\mathbf{X} = \{a, b, c, d, e\}$ với xác suất $p(a) = .05$, $p(b) = .10$, $p(c) = .12$, $p(d) = .13$ và $p(e) = .60$

<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
.05	.10	.12	.13	.60
0	1			
.15		.12	.13	.60
		0	1	
.15		.25		.60
0		1		
.40				.60
0				1
1.0				

<i>x</i>	<i>f(x)</i>
<i>a</i>	000
<i>b</i>	001
<i>c</i>	010
<i>d</i>	011
<i>e</i>	1

Entropy và nén Huffman

- Độ dài trung bình để truyền thông tin xác định 1 sự kiện

$$\begin{aligned}\ell(f) &= .05 \times 3 + .10 \times 3 + .12 \times 3 + .13 \times 3 + .60 \times 1 \\ &= 1.8.\end{aligned}$$

- Entropy:

$$\begin{aligned}H(X) &= .2161 + .3322 + .3671 + .3842 + .4422 \\ &= 1.7402.\end{aligned}$$

- Kết quả: $H(X) \leq \ell(f) \leq H(X)+1$

Một số tính chất cơ bản của Entropy

- $H(X) \geq 0$, đẳng thức xảy ra khi và chỉ khi biến X là hằng
- $H(X) \geq \log_2 |X|$, đẳng thức xảy ra khi và chỉ khi
 $p(x=X) = 1/|X|$
- $H(X, Y) \geq H(X) + H(Y)$, đẳng thức xảy ra khi và chỉ khi X và Y phân bố độc lập
- $H(X|Y) \geq H(X)$, đẳng thức xảy ra khi và chỉ khi X và Y phân bố độc lập
- **Chain Rule:** $H(X, Y) = H(X|Y) + H(Y)$

Entropy của các thành phần trong hệ thống mã hóa

- ☐ $H(C|K) = H(P)$
- ☐ $H(C|P, K) = H(P|C, K) = 0$
- ☐ $H(P, K) = H(P) + H(K)$
- ☐ $H(C) \geq H(P)$
- ☐ $H(C, P, K) = H(C, K) = H(P, K)$
- ☐ $H(K|C) = H(K) + H(P) - H(C)$
- ☐ $H(K|C^n) = H(K) + H(P^n) - H(C^n)$

Nhận xét

- Nhận xét: Có $26! \approx 10^{26}$ cách mã hóa (bằng thay thế) đối với văn bản tiếng Anh (gồm các ký tự thường)
- Tương đương với mức độ 88-bit security → Vì sao lại dễ bị tấn công trên thực tế?
- Shannon: Mọi phương pháp mã hóa đơn ký tự (monoalphabetic cipher) trên tiếng Anh đều dễ dàng bị phá nếu có được 25 ký tự ciphertext

Khoảng cách Unicity của một ngôn ngữ

- $H(K|C)=H(K)+H(P)-H(C)$
- Entropy (trên mỗi ký tự) trong một ngôn ngữ được định nghĩa theo công thức:

$$H_L = \lim_{n \rightarrow \infty} \frac{H(\mathbf{P}^n)}{n}$$

- Đặt

$$R_L = 1 - \frac{H_L}{\log_2 |\mathcal{P}|}$$

Khoảng cách Unicity của một ngôn ngữ

- Entropy rate:

$$H(P^n) \geq nH_L = n(1 - R_L) \log_2 |\mathcal{P}|$$

- Do $|\mathcal{P}|=|\mathcal{C}|$, ta có

$$H(C^n) \leq nH(C) \leq n \log_2 |\mathcal{C}| = n \log_2 |\mathcal{P}|$$

- Thay thế vào công thức của $H(K|C^n)$:

$$\begin{aligned} H(K|C^n) &= H(K) - H(P^n) + H(C^n) \\ &\leq H(K) - n(1 - R_L) \log_2 |\mathcal{P}| + n \log_2 |\mathcal{P}| \\ &= H(K) - nR_L \log_2 |\mathcal{P}| \\ &\leq \log_2 |\mathcal{K}| - nR_L \log_2 |\mathcal{P}| \end{aligned}$$

Khoảng cách Unicity của một ngôn ngữ

- Hệ thống mã hóa bị phá vỡ nếu $H(K|C^n)=0$:

$$H(K|C^n) \leq \log_2(|\mathcal{K}|) - nR_L \log_2(|\mathcal{P}|) = 0$$

$$n_0 \approx \frac{\log_2|\mathcal{K}|}{R_L \log_2|\mathcal{P}|}$$

- Trường hợp tiếng Anh ($|\mathcal{P}|=26$, $R_L \approx 0.75$) và sử dụng phương pháp mã hóa bằng thay thế ($\log_2|\mathcal{K}| \approx 88$), ta có $n_0 \approx 25$.

Nhận xét

- ☐ **Nén dữ liệu**
 - ☐ Nén tốt – mã hóa tốt
 - ☐ Mã hóa tốt – nén không tốt

Kết hợp các phương pháp mã hóa

- **“Tổng có trọng số” của các hệ thống mã hóa**
 - Tạo ra hệ thống mã hoá mới từ các hệ thống mã hóa có sẵn
 - Chọn 2 hệ thống mã hóa có cùng không gian thông điệp, sử dụng hệ thống A với xác suất p , sử dụng hệ thống B với xác suất $1 - p$.
- **Hệ thống mã hóa tích (product cipher):** áp dụng tuần tự các giải thuật mã hóa liên tiếp nhau