

CEU

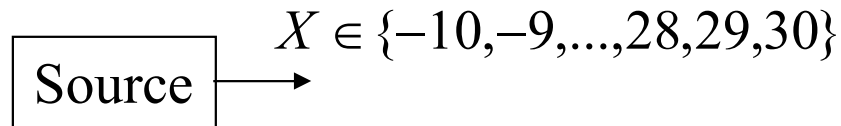
*Universidad
San Pablo*

UNIT 3: Information Theory

Gabriel Caffarena Fernández
3rd Year Biomedical Engineering Degree
EPS – Univ. San Pablo – CEU
(Based on slides by Carlos Oscar Sánchez Sorzano)

Entropy

Example: Temperature in a room

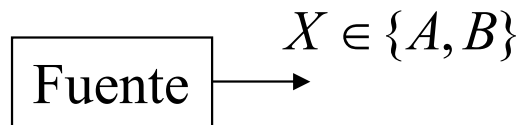


Scenario 1 (no information): 23, 23, 23, 23, 23, 23, ...

Scenario 2 (little information): 23, 23, 23, 24, 23, 23, ...

Scenario 3 (a lot of information): -10, 23, 30, -5, 0, 15, ...

Example: Who wins a football match?



Scenario 1 (no information): A, A, A, A, A, A, A, A, A, ...

Scenario 2 (little information): A, A, A, B, B, A, A, A, A, B, A, A, ...

Scenario 3 (Lots of information): A, B, B, A, A, A, B, B, B, A, A, B, ...

Entropy

- Given an information source $X \in \{X_0, X_1, \dots, X_{n-1}\}$ (e.g. experiment, random variables, etc.) the **Shannon information** of a particular outcome is defined as

$$h(X_i) = \log \left(\frac{1}{P(X=X_i)} \right) = -\log(P(X = X_i))$$

- The Shannon information is the amount of surprise that the outcome produces
- If the base of the logarithm is 2, then the **information is measured in bits**

Entropy

- **Example:**

$$X \in \{0,1,2,3\};$$

$$P(X = 0) = 0.25; P(X = 1) = P(X = 2) = 0.125$$

$$P(X = 3) = 0.5$$

$$h(0) = \log\left(\frac{1}{0.25}\right) = \log(4) = \mathbf{2 \text{ bits}}$$

$$h(1) = h(2) = \log\left(\frac{1}{0.125}\right) = \log(8) = \mathbf{3 \text{ bits}}$$

$$h(3) = \log\left(\frac{1}{0.5}\right) = \log(2) = \mathbf{1 \text{ bits}}$$

$$h(3) < h(0) < h(1) = h(2)$$

Less surprise

More surprise

Entropy

- **Entropy** is the average of Shannon information of an information source X

$$\begin{aligned} H(X) &= E[h(X)] = \sum_i h(X_i)P(h(X_i)) = \sum_i h(X_i)P(X_i) \\ &= \sum_i P(X_i) \log\left(\frac{1}{P(X_i)}\right) = - \sum_i P(X_i) \log(P(X_i)) \end{aligned}$$

Entropy

$$H(X) = \sum_i P(X_i) \log \left(\frac{1}{P(X_i)} \right) = - \sum_i P(X_i) \log(P(X_i))$$

Example: Who wins a football match?

Scenario 1 (no information): A,A,A,A,A,A,A,A,A,A,...

$$\begin{aligned} H(X) &= -p(A) \log p(A) - p(B) \log p(B) \\ &= -1 \log 1 - 0 \log 0 = 0 - 0 = 0 \end{aligned}$$

The most probable event is the one contributing less

Scenario 2 (little information): A,A,A,B,B,A,A,A,A,B,A,A,...

$$H(X) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = -(-0.2158) - (-0.3466) = 0.5624$$

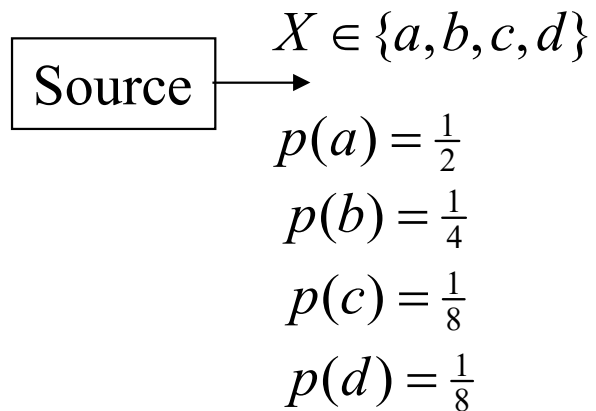
Scenario 3 (a lot of information): A,B,B,A,A,A,B,B,B,A,A,B,...

$$H(X) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = -(-0.3466) - (-0.3466) = 0.6932$$

$$H(X) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = -(-0.5) - (-0.5) = 1(\text{bits})$$

Entropy

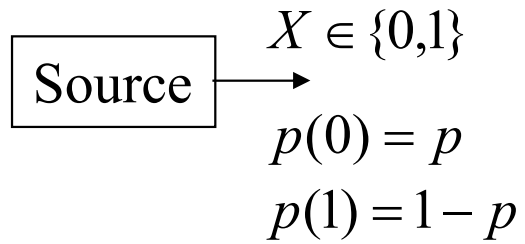
Example:



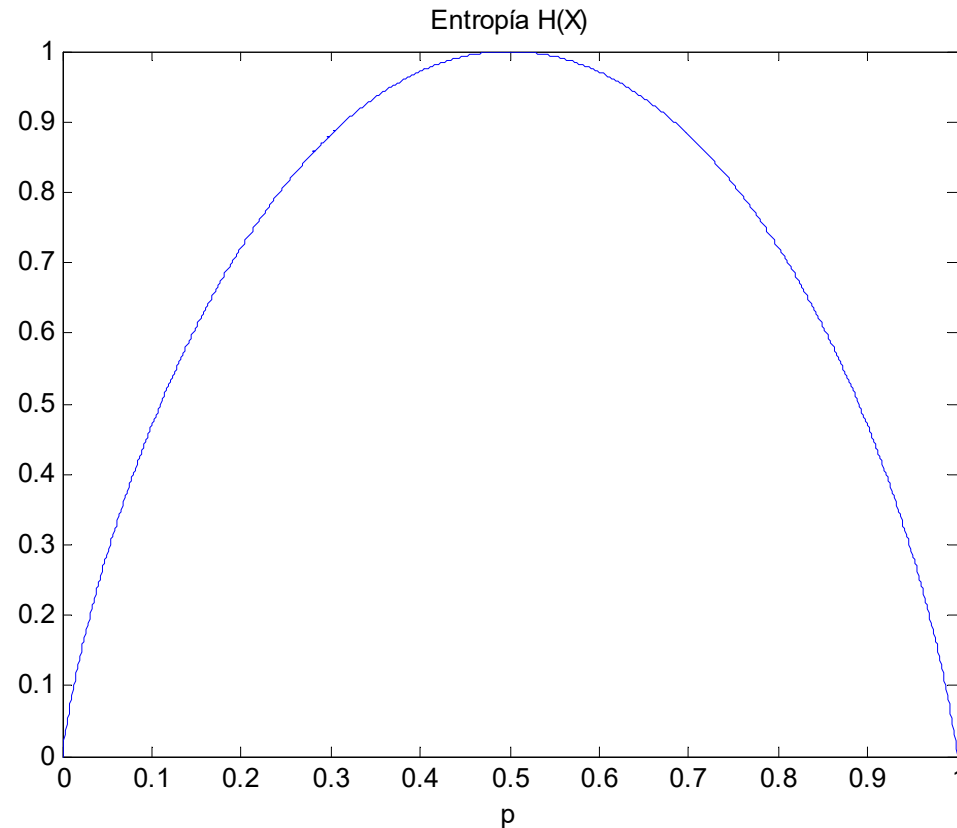
$$H(X) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{8} \log \frac{1}{8} - \frac{1}{8} \log \frac{1}{8} = \frac{7}{4} \text{ bits}$$

Entropy

Example:



$$H(X) = -p \log p - (1 - p) \log(1 - p)$$



Entropy

Some comments on operations

$$\left. \begin{array}{l} 0 \log \frac{0}{q} = 0 \\ p \log \frac{p}{0} = \infty \end{array} \right\} \text{Consensus}$$

$$\log_a x = \frac{\log_b x}{\log_b a} \Rightarrow \log_2 x = \frac{\log_{10} x}{\log_{10} 2} \approx 3.32 \log_{10} x$$

$$\log_a x = \log_a b \cdot \log_b x$$

Entropy

Properties

$$H(X) \geq 0$$

Proof:

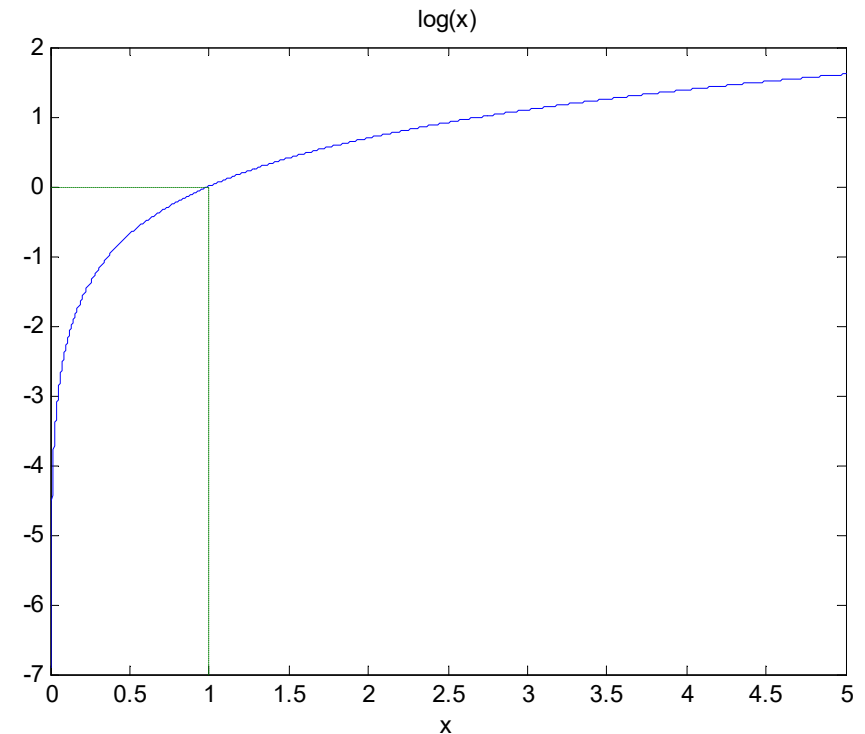
$$0 \leq p(x) \leq 1 \Rightarrow \frac{1}{p(x)} \geq 1 \Rightarrow \log \frac{1}{p(x)} \geq 0$$



$$H(X) = \sum_{x \in \Xi} p(x) \log \frac{1}{p(x)} \geq 0$$

↑

$$p(x) \geq 0$$



Entropy

Properties

$$H_b(X) = (\log_b a) H_a(X)$$

Proof:

$$\begin{aligned} H_b(X) &= E \left\{ \log_b \frac{1}{p(x)} \right\} = E \left\{ \frac{1}{\log_a b} \log_a \frac{1}{p(x)} \right\} \\ &= E \left\{ \log_b a \log_a \frac{1}{p(x)} \right\} = (\log_b a) H_a(X) \end{aligned}$$

Joint entropy

- Given two information sources \mathbf{X} and \mathbf{Y} the **joint Shannon information** of a particular joint outcome is defined as

$$h(X_i, Y_j) = \log \left(\frac{1}{P(X=X_i, X=X_j)} \right) = -\log(P(X = X_i, X = X_j))$$

- The joint entropy is the average of the joint Shannon information

$$\begin{aligned} H(X, Y) &= E[h(X, Y)] = \sum_i \sum_j h(X_i, Y_j) P(X_i, Y_j) \\ &= \sum_{i,j} P(X_i, Y_j) \log \left(\frac{1}{P(X_i, Y_j)} \right) = - \sum_{i,j} P(X_i, Y_j) \log \left(P(X_i, Y_j) \right) \end{aligned}$$

Joint Entropy

$$\text{Joint Entropy } H(X, Y) = \sum_{i,j} P(X_i, Y_j) \log \left(\frac{1}{P(X_i, Y_j)} \right) = - \sum_{i,j} P(X_i, Y_j) \log (P(X_i, Y_j))$$

Example: Peter is bilingual (Spanish/English) and he reads “The Times” with probability 0.5 and “**El País**” with probability 0.5.

$$H(\text{newspaper}) = 1 \text{ bit}$$

$p(\text{newspaper}, \text{language})$

newspaper \ language	English	Spanish
	The Times	0.5
El País	0	0.5

$$\begin{aligned} H(\text{newspaper}, \text{language}) &= \\ &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1 \text{ bit} \end{aligned}$$

Joint Entropy

Joint Entropy $H(X, Y) = \sum_{i,j} P(X_i, Y_j) \log \left(\frac{1}{P(X_i, Y_j)} \right) = - \sum_{i,j} P(X_i, Y_j) \log (P(X_i, Y_j))$

Example: Peter watches the CNN and the BBC with the following probabilities

$$H(TV) = 1 \text{ bit}$$

$p(TV, \text{language})$

		language	
		English	Spanish
TV	BBC	0.5	0
	CNN	0.25	0.25

$$H(TV, \text{language}) =$$

$$= -\frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{2} \log \frac{1}{2}$$

$$= 1.5 \text{ bits}$$

Conditional Entropy

- Given two information sources \mathbf{X} and \mathbf{Y}
- The average information of \mathbf{Y} given that $\mathbf{X} = \mathbf{X}_i$ is

$$\begin{aligned} H(Y|X = X_i) &= E[h(Y|X = X_i)] = \sum_j p(h(Y_j|X_i))h(Y_j|X_i) \\ &= \sum_j p(Y_j|X_i) \log\left(\frac{1}{p(Y_j|X_i)}\right) \end{aligned}$$

So this is a **conditional entropy** for a given outcome of X.

Conditional Entropy

- **Conditional entropy** for a given outcome of X.

$$H(Y|X = X_i) = \sum_j p(Y_j|X_i) \log \left(\frac{1}{p(Y_j|X_i)} \right)$$

- The **conditional entropy** is the average of the **conditional entropy for a given outcome of X**, that has already averaged the Shannon information over Y, so it is going to be averaged **over X**.

$$\begin{aligned} H(Y|X) &= E_{P(X_i)}[H(Y|X_i)] = \sum_i p(X_i) H(Y|X = X_i) \\ &= \sum_i P(X_i) \sum_j P(Y_j|X_i) \log \left(\frac{1}{P(Y_i|X_j)} \right) \\ &= \sum_i \sum_j P(X_i) P(Y_j|X_i) \log \left(\frac{1}{P(Y_i|X_j)} \right) \\ &= \sum_i \sum_j P(X_i, Y_j) \log \left(\frac{1}{P(Y_i|X_j)} \right) = \sum_i \sum_j P(X_i, Y_j) \log \left(\frac{P(X_i)}{P(X_i, Y_j)} \right) \end{aligned}$$

Conditional Entropy

$$\begin{aligned} \text{Conditional Entropy } H(Y|X) &= \sum_i p(X_i) H(Y|X = X_i) \\ &= \sum_i \sum_j P(X_i) P(Y_j|X_i) \log \left(\frac{P(X_i)}{P(X_i, Y_j)} \right) = \sum_i \sum_j P(X_i, Y_j) \log \left(\frac{P(X_i)}{P(X_i, Y_j)} \right) \end{aligned}$$

Properties

$$H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$$



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

$$H(X | Y) \neq H(Y | X)$$

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

If X and Y are independent, then $H(Y | X) = H(Y)$

Conditional Entropy

Example: Peter watches CNN and BBC

$$p(TV, language)$$

TV \ language	English	Spanish
BBC	0.5	0
CNN	0.25	0.25

$$p(TV | language)$$

TV \ language	English	Spanish
BBC	0.666	0
CNN	0.333	1

$$p(language | TV)$$

TV \ language	English	Spanish
BBC	1	0
CNN	0.5	0.5

Conditional Entropy

$p(TV, language)$

language \ TV	English	Spanish
BBC	0.5	0
CNN	0.25	0.25

$p(language | TV)$

idioma \ TV	Inglés	Español
BBC	1	0
CNN	0.5	0.5

language \ TV	English	Spanish
BBC	1	0

$p(language | TV = BBC)$

$H(language | TV = BBC) = 0 \text{ bits}$

language \ TV	English	Spanish
CNN	0.5	0.5

$p(language | TV = CNN)$

$H(language | TV = CNN) = 1 \text{ bits}$

$H(language | TV) = 0.5 \text{ bits}$

Differential or Relative Entropy (Kullback-Leibler distance)

$$\text{Relative Entropy } D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = E_{p(x)} \left\{ \log \frac{p(x)}{q(x)} \right\}$$

Properties

$$D(p || q) \geq 0$$

$$D(p || p) = 0$$

$$D(p || q) \neq D(q || p)$$

Conditional
relative entropy

$$\begin{aligned} D(p(y|x)||q(y|x)) &= \sum_x p(x) \sum_y p(y|x) \overbrace{\log \frac{p(y|x)}{q(y|x)}}^{D(p(y|X=x)||q(y|X=x))} \\ &= \sum_x \sum_y p(x,y) \log \frac{p(y|x)}{q(y|x)} = E_{p(x,y)} \left\{ \log \frac{p(y|x)}{q(y|x)} \right\} \end{aligned}$$

Relative Entropy

Example: Let's assume that the actual probabilities of an information source are

$$\begin{array}{l}
 \boxed{\text{Source}} \longrightarrow \begin{array}{l} X \in \{0,1\} \\ p(0) = p \\ p(1) = 1-p \end{array} \\
 \end{array} \quad H_p(X) = -p \log p - (1-p) \log(1-p)$$

However, due to estimation errors, what we really have is

$$\begin{array}{l}
 \boxed{\text{Source}} \longrightarrow \begin{array}{l} X \in \{0,1\} \\ p(0) = q \\ p(1) = 1-q \end{array} \\
 \end{array} \quad H_q(X) = -q \log q - (1-q) \log(1-q)$$

$$\left. \begin{array}{l} D(p \parallel q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \\ D(q \parallel p) = q \log \frac{q}{p} + (1-q) \log \frac{1-q}{1-p} \end{array} \right\} \xrightarrow{p=q} D(p \parallel q) = D(q \parallel p) = 0$$

Relative Entropy

Example: Given the following actual distribution of a set of symbols

$$\begin{array}{l} \boxed{\text{Source}} \longrightarrow X \in \{0,1\} \\ p(0) = \frac{1}{2} \\ p(1) = \frac{1}{2} \end{array} \quad H_p(X) = 1\text{bit}$$

The estimated probabilities are

$$\begin{array}{l} \boxed{\text{Source}} \longrightarrow X \in \{0,1\} \\ p(0) = \frac{1}{3} \\ p(1) = \frac{2}{3} \end{array} \quad H_q(X) = 0.9183\text{bits}$$

$$\left. \begin{array}{l} D(p \parallel q) = 0.0850\text{bits} \\ D(q \parallel p) = 0.0817\text{bits} \end{array} \right\} \longrightarrow H_p(X) = H_q(X) + D(q \parallel p)$$

Mutual Information

Mutual
Information

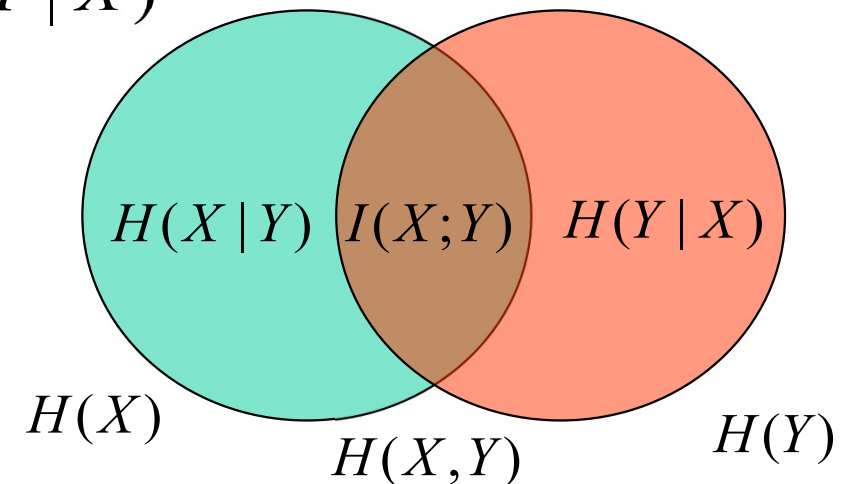
$$I(X;Y) = \sum_{x \in \Xi, y \in \Psi} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$
$$= E_{p(x, y)} \left\{ \log \frac{p(x, y)}{p(x)p(y)} \right\} = D(p(x, y) \parallel p(x)p(y))$$

Properties

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

$$I(X;Y) = H(X) + H(Y) - H(X, Y)$$

$$I(X;X) = H(X)$$



Jensen's inequality

Jensen's inequality

Given the convex function $f(x)$ and the r.v. X , then

$$E\{f(X)\} \geq f(E\{X\})$$

Thanks to this inequality it can be proved the following

$$D(p \parallel q) \geq 0$$

$$I(X; Y) \geq 0$$

$$H(X) \leq \log(\#X)$$

$$H(X | Y) \leq H(X)$$

$$H(X_1, \dots, X_N) \leq \sum_{i=1}^N H(X_i)$$

$$D(p \parallel q) = 0 \leftrightarrow p = q$$

$$I(X; Y | Z) \geq 0$$

$$H(X) = \log(\#X) \leftrightarrow p(X) = \text{uniform}$$

$$H(X | Y) = H(X) \leftrightarrow X, Y \text{ independent}$$

$$H(X_1, \dots, X_N) = \sum_{i=1}^N H(X_i) \leftrightarrow X_i \text{ independent}$$

#X ≡ number of symbols (values) of X

Jensen's inequality

Highlights

- $H(X) \leq \log(\#X)$

EXAMPLE: If X can have 4 values, its entropy cannot be greater than $\log_2(4)=2$ bits

- $H(X) = \log(\#X) \leftrightarrow p(X) = k$

An information source provides maximum information when its values (symbols) have equal probabilities

- $H(X | Y) \leq H(X)$

Any extra knowledge will never increase the information given by a source

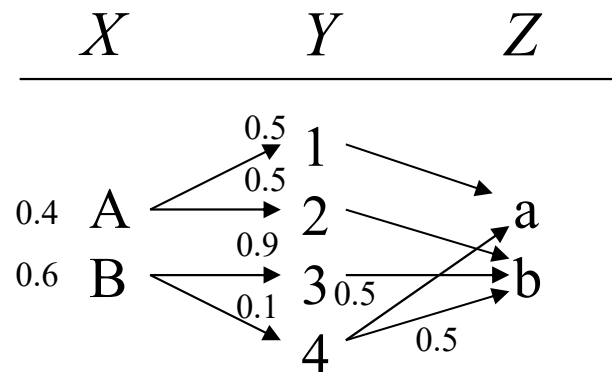
- $H(X | Y) = H(X) \leftrightarrow X, Y \text{ independent}$

If extra knowledge about another information source Y does not vary the information provided by X , then X and Y are independent

Markov chains

Given the random variables X, Y, Z , these variables conform a Markov chain $X \rightarrow Y \rightarrow Z$ if $p(z | x, y) = p(z | y)$

Example:



Possible sequences:
B3b, A1a, A2b, etc.

$$p(x = A) = 0.4$$

$$p(x = B) = 0.6$$

$$p(y = 1 | x = A) = 0.5$$

$$p(y = 2 | x = A) = 0.5$$

$$p(y = 3 | x = A) = p(y = 4 | x = A) = 0$$

$$p(y = 1 | x = B) = p(y = 2 | x = B) = 0$$

$$p(y = 3 | x = B) = 0.9$$

$$p(y = 4 | x = B) = 0.1$$

$$p(z = a | y = 1) = 1$$

$$p(z = b | y = 2) = 1$$

$$p(z = b | y = 3) = 1$$

$$p(z = a | y = 4) = 0.5$$

$$p(z = b | y = 4) = 0.5$$

Signal processing inequality

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

The signal processing inequality leads to the assertion that if X is processed to generate Y, and Y is post-process to generate Z, then, X has more information on Y than on Z. Thus, Z does not provides more information about X than Y..

$$X \rightarrow Y \rightarrow f(Y) \Rightarrow I(X;Y) \geq I(X;f(Y))$$

If Y is generated from X, the information that Y contains regarding X, is not going to be increased if Y is processes by any signal processing algorithm.

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

In a Markov chain, knowing the value of Z decreases (or just keeps) the dependency between X and Y.

Coding and compression: Block Code

Example:

x_i	$\Pr\{X = x_i\}$	$C(x_i)$
1	1/2	0
2	1/4	10
3	1/8	110
4	1/8	111

AVERAGE LENGTH OF THE CODE

$$L(C) = E[l(C(x_i))]$$

Length of code $C(x_i)$

$$H(X) = L(C) = 1.75bits$$

Example:

x_i	$\Pr\{X = x_i\}$	$C(x_i)$
1	1/3	0
2	1/3	10
3	1/3	11

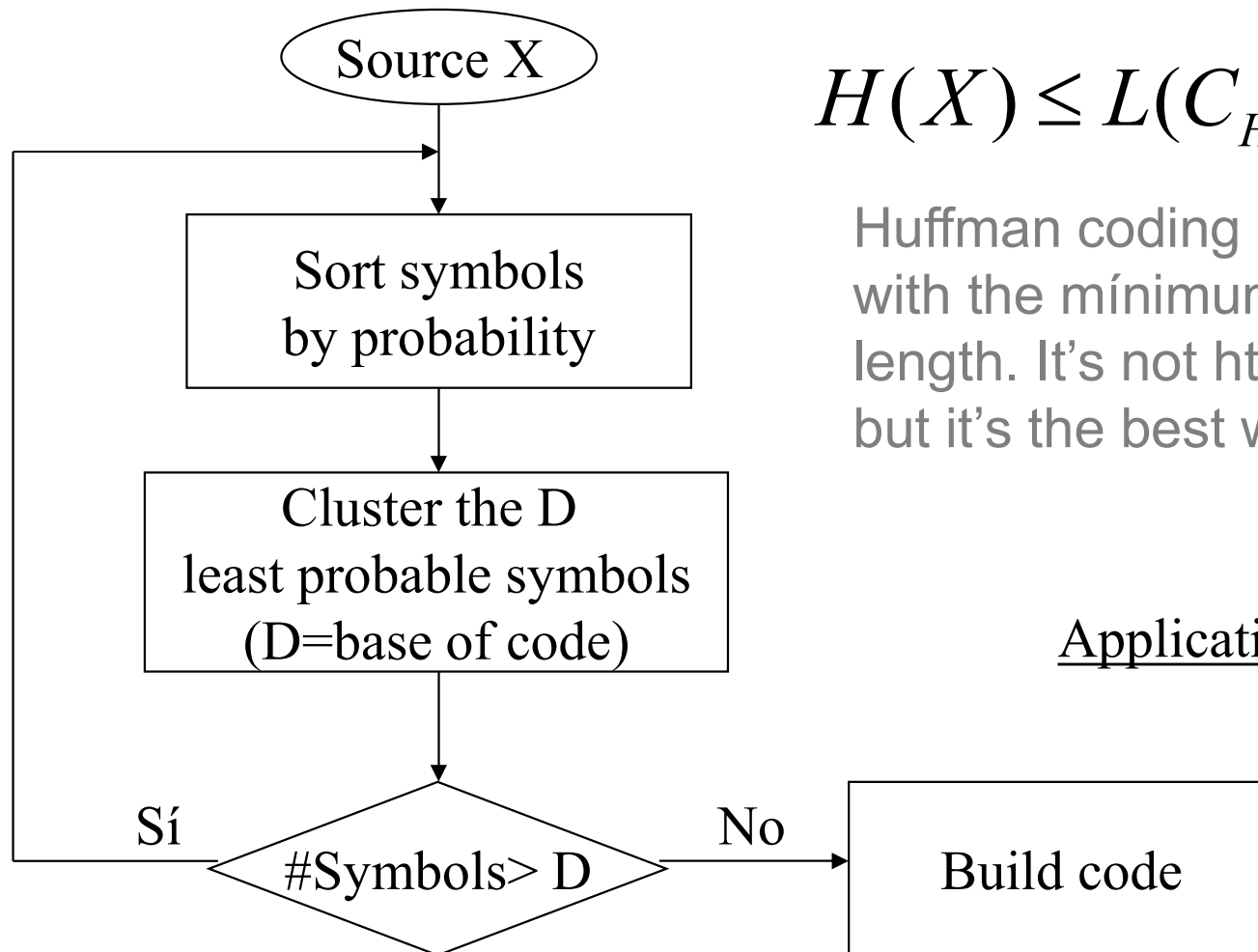
$$H(X) = 1.58bits$$

$$L(C) = 1.66bits$$

$$H(X) \leq L(C)$$

This is a fundamental property in statistical coding

Huffman coding

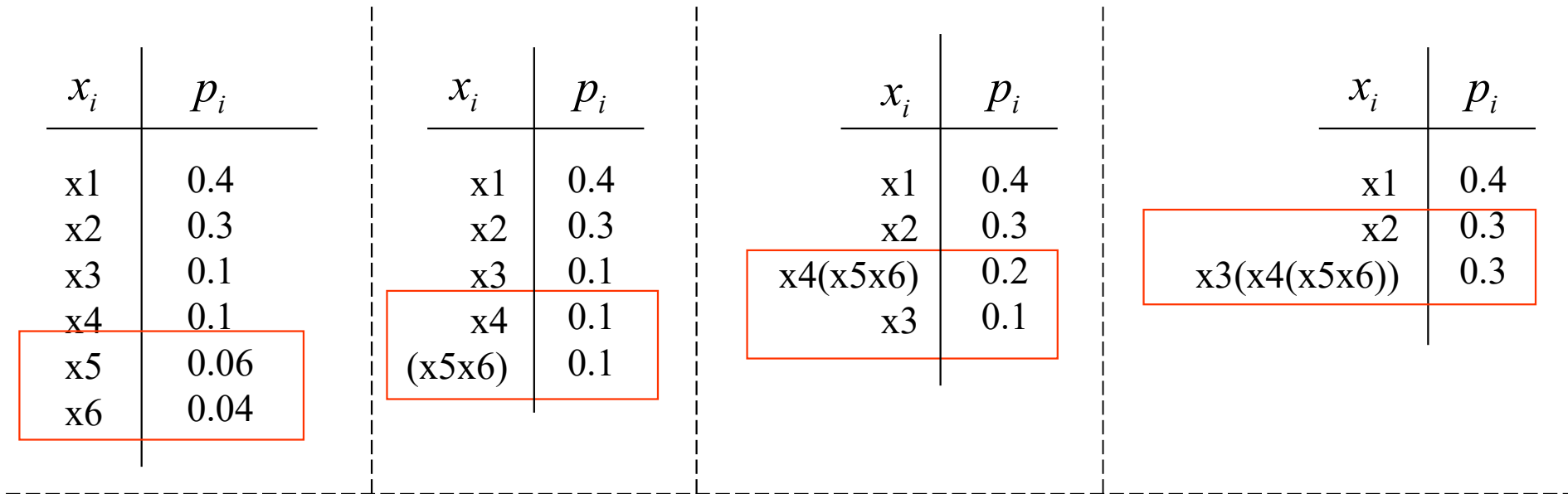


$$H(X) \leq L(C_{Huffman}) \leq L(C)$$

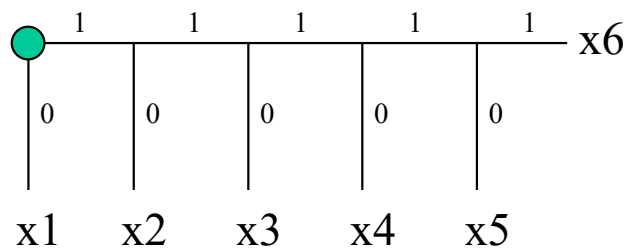
Huffman coding produces a code with the minimum possible average length. It's not the optimum code, but it's the best we can get.

Applications: JPEG, MP3, etc.

Huffman coding: $D=2$



x_i	p_i
x2(x3(x4(x5x6)))	0.6
x1	0.4



x_i	$C(x_i)$
x1	0
x2	10
x3	110
x4	1110
x5	11110
x6	11111

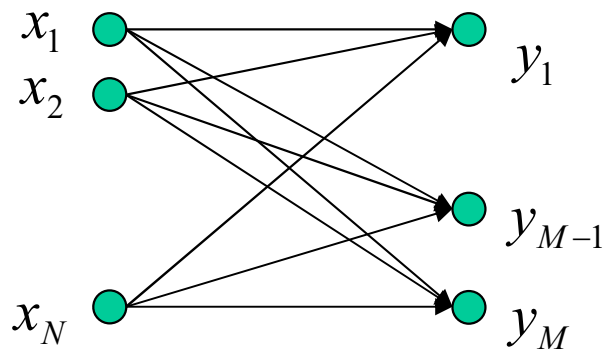
$H(X) = 2.1435bits$

$L(C_{Huffman}) = 2.2bits$

Channel characterization

Memoryless channel

$$P(y[n] | x[n], x[n-1], x[n-2], \dots) = P(y[n] | x[n])$$



$$Q = \begin{pmatrix} P(y_1 | x_1) & P(y_2 | x_1) & \dots & P(y_M | x_1) \\ P(y_1 | x_2) & P(y_2 | x_2) & \dots & P(y_M | x_2) \\ \dots & \dots & \dots & \dots \\ P(y_1 | x_N) & P(y_2 | x_N) & \dots & P(y_M | x_N) \end{pmatrix}$$

Channel with memory

$$P(y[n] | x[n], x[n-1], x[n-2], \dots) = P(y[n] | x[n]x[n-1]) \quad \text{Memory}=1$$

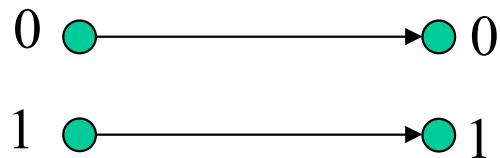
$$P(y[n] | x[n], x[n-1], x[n-2], \dots) = P(y[n] | x[n]x[n-1]x[n-2]) \quad \text{Memory}=2$$

Channel capacity

$$C = \max_{p(X)} I(Y; X)$$

$$0 \leq C \leq \min(\log\# X, \log\# Y)$$

Example: Binary Channel with no noise



$$Q = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

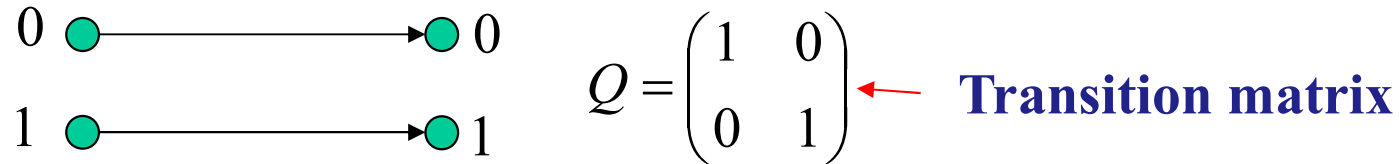
$$P(X=0)=p$$

$$P(X=1)=1-p$$

$$\begin{aligned} I(Y; X) &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} = \\ &= p(x=0, y=0) \log \frac{p(x=0, y=0)}{p(x=0)p(y=0)} + p(x=0, y=1) \log \frac{p(x=0, y=1)}{p(x=0)p(y=1)} \\ &\quad + p(x=1, y=0) \log \frac{p(x=1, y=0)}{p(x=1)p(y=0)} + p(x=1, y=1) \log \frac{p(x=1, y=1)}{p(x=1)p(y=1)} \end{aligned}$$

Channel capacity

Example: Binary channel with no noise



$$p(X = x, Y = y) = p(X = x)p(Y = y | X = x) = \begin{cases} p(X = x) & x = y \\ 0 & x \neq y \end{cases}$$

$$p(Y = y) = \sum_x p(Y = y | X = x)p(X = x) = p(X = y)$$

$$I(Y; X) = p(x = 0, y = 0) \log \frac{p(x = 0, y = 0)}{p(x = 0)p(y = 0)} + p(x = 0, y = 1) \log \frac{p(x = 0, y = 1)}{p(x = 0)p(y = 1)}$$

$$+ p(x = 1, y = 0) \log \frac{p(x = 1, y = 0)}{p(x = 1)p(y = 0)} + p(x = 1, y = 1) \log \frac{p(x = 1, y = 1)}{p(x = 1)p(y = 1)}$$

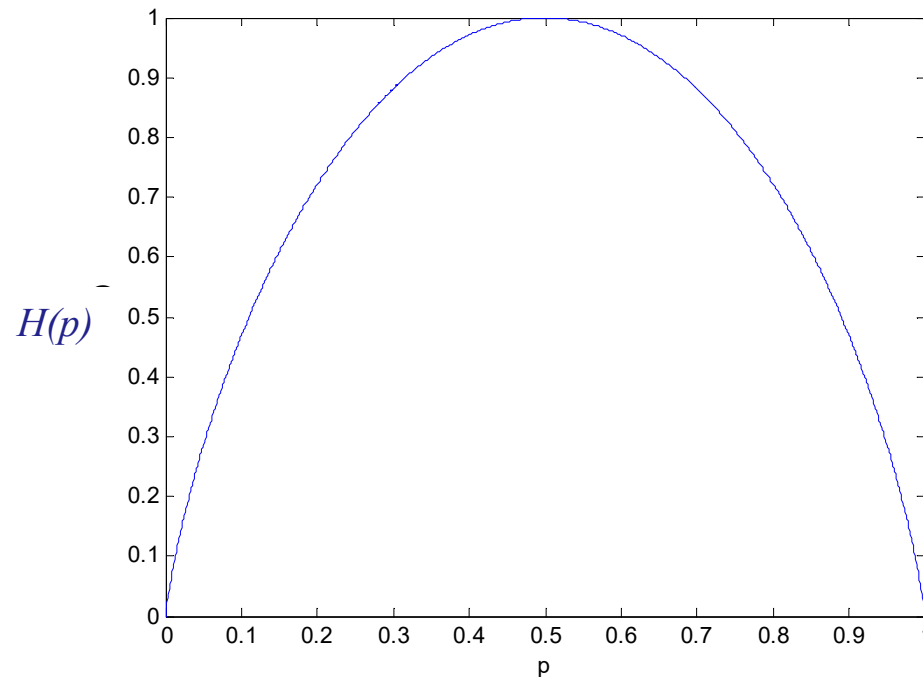
$$= p(x = 0) \log \frac{p(x = 0)}{p(x = 0)p(x = 0)} + p(x = 1) \log \frac{p(x = 1)}{p(x = 1)p(x = 1)}$$

$$= p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p} = H(X)$$

Channel capacity

Example: Binary channel with no noise

$$C = \max_{p(X)} I(Y; X) = \max_p \left\{ p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p} \right\}$$

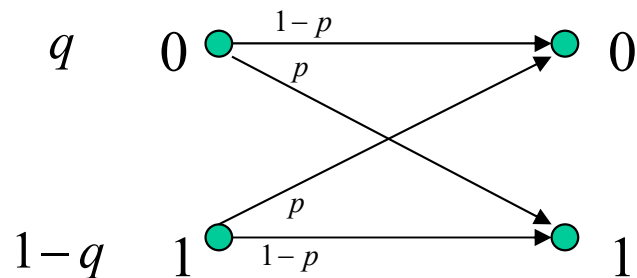


$$\longrightarrow p = \frac{1}{2} \Rightarrow C = 1 \text{ bit}$$

(Remember: $P(X=0)=p$)

Channel capacity

Example: Binary Symetric Channel (BSC) – channel with noise



$$q = \frac{1}{2} \Rightarrow C = \max_q I(Y; X) = 1 - H(p)$$

$$\begin{aligned} H(Y | X = 0) &= p(Y = 0 | X = 0) \log \frac{1}{p(Y = 0 | X = 0)} + p(Y = 1 | X = 0) \log \frac{1}{p(Y = 1 | X = 0)} \\ &= (1-p) \log \frac{1}{1-p} + p \log \frac{1}{p} = H(p) = H(Y | X = 1) \end{aligned}$$

$$P(Y = 0) = P(X = 0)P(Y = 0 | X = 0) + P(X = 1)P(Y = 0 | X = 1) = q(1-p) + (1-q)p$$

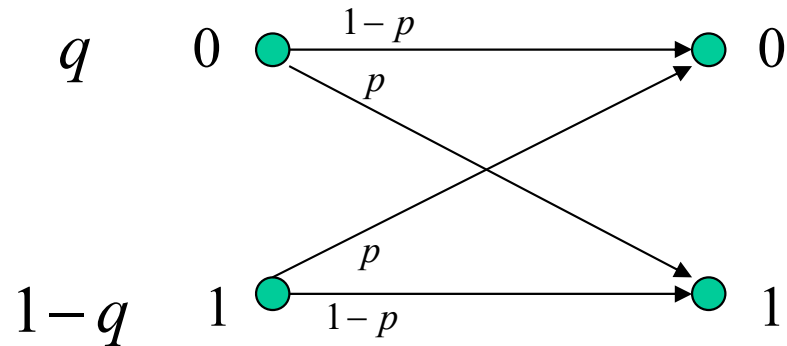
$$P(Y = 1) = P(X = 0)P(Y = 1 | X = 0) + P(X = 1)P(Y = 1 | X = 1) = qp + (1-q)(1-p)$$

$$H(Y) = P(Y = 0) \log \frac{1}{P(Y = 0)} + P(Y = 1) \log \frac{1}{P(Y = 1)} \leq H(X) \leq 1 \text{ bit}$$

Only equal if Y is uniform

Channel capacity

Example: Binary Symetric Channel (BSC) – channel with noise



$$q = \frac{1}{2} \Rightarrow C = \max_q I(Y; X) = 1 - H(p)$$

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

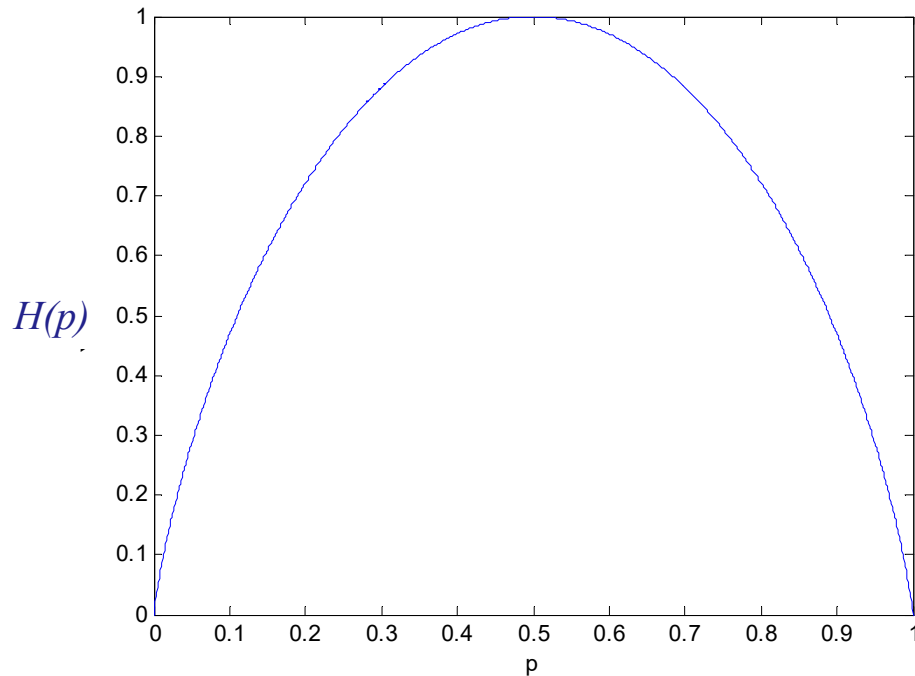
$$= H(Y) - \sum_x p(x) H(Y | X = x) \leq 1 - \sum_x p(x) H(p) = 1 - H(p) = C$$

If Y is uniform

Channel Capacity

Example: Binary channel with noise (BSC)

$$C = \max_{p(X)} I(Y; X) = 1 - H(p)$$

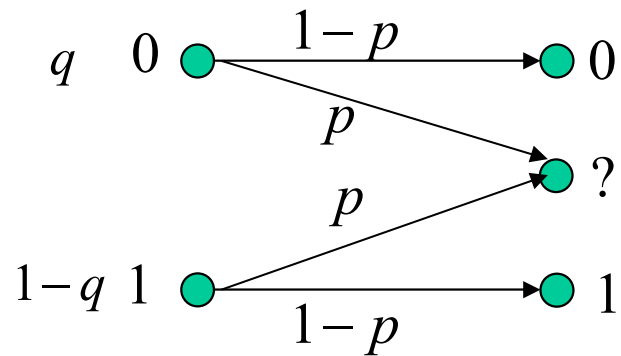


$p=0$ or $p=1 \rightarrow C=1$ bit
 $p=0.5 \rightarrow C=0$ bits

(Remember: $P(Y=1|X=0)=P(Y=0|X=1)=p$)

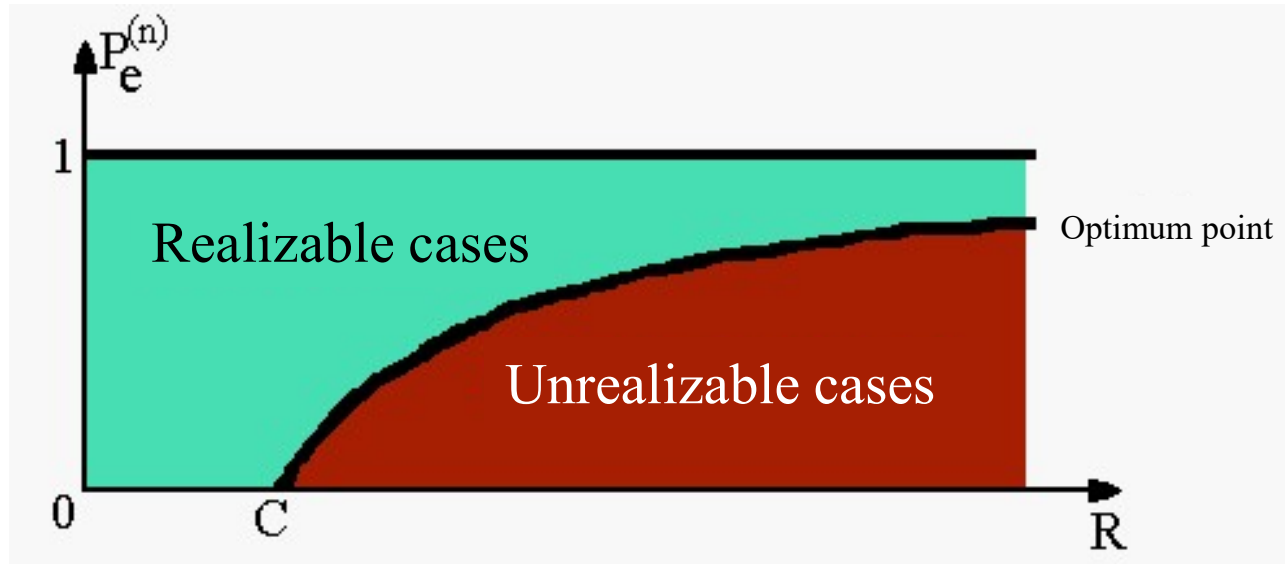
Challenge: Channel capacity

Given an “erasure channel”, where the data is received corrected (“0” or “1”) or not received at all (“?”) show that the capacity $C=1-p$



$$Q = \begin{pmatrix} 1-p & p & 0 \\ 0 & p & 1-p \end{pmatrix}$$

Channel coding theorem



It is also possible to express channel capacity in bits per second. Given C and the transmission rate R , it is proved that it is possible to achieve an error probability $P_e=0$ if $R < C$. Also, if $R > C$ then the error probability increases.

It is possible to find codes that achieve $P_e=0$ and also that increase the transmission speed increasing the error rate, however the theorem does not indicate how are these codes.

Case study:

Are brains good at processing information?

(Introduction to Information Theory, J.V. Stone, 2018)

- Neurons communicate to each other continuously
 - This communication must be performed efficiently
- Horace Bellow supports the *efficiency coding hypothesis*, where the sensory input must be encoded efficiently before being sent to the brain

Case study:

Are brains good at processing information?

(Introduction to Information Theory, J.V. Stone, 2018)

- Information in spiking neurons
 - The neurons propagate action potentials (AP, a.k.a. spikes)
 - If the spikes are encoded as 0's (no AP) and 1's (AP) it is possible to compute the entropy of the neuron (information source)
 - Considering an average firing rate of r spikes/s the capacity of the channel might be equal to r bits/s
 - However it can be proven that the information rate is bigger than that
 - The trick is that the neurons also use temporal information to encode information, so that one single spike encodes more than 1 bit

Case study:

Are brains good at processing information?

(Introduction to Information Theory, J.V. Stone, 2018)

- Mutual information between the input and output of a neuron
 - An experiment showed that a neuron (mechanical receptor of a cricket) generates 600 bits/s
 - 300 bits/s are related to the input, the rest is noise
 - Thinking that a neuron works with “packets” of 300 bits is out of line
 - There are theories that indicate that those 300 bits are actually divided in packets of 3 bits every 10 ms, providing continuous information about changes in the output (speed of an object)

Case study:

Are brains good at processing information?

(Introduction to Information Theory, J.V. Stone, 2018)

- Shannon optimal coding: maximizing entropy
 - In the human eye, the information provided by “Red” and “Green” receptors is very similar
 - Using two different nerve fibre per receptor is a waste of channel capacity (since there is a lot of redundancy over time)
 - However, it can be proved that the addition and subtraction of the outputs of the R and G receptors leads to signals with uniform distributions
 - Also the subtraction and summation are independent of each other
 - So:
 1. The use of two, instead of three, separate fibres is now justified
 2. The entropy is maximized so that the information rate is highly increased
 - Ganglion cells in the retina perform this operations, but they use several receptor outputs to **compress** information and reduce the necessary nerve fibres: **126 million receptors → 1 million nerve fibre**

SUMMARY

- Entropy
- Mutual information
- Signal processing inequality
- Huffman coding
- Channel Capacity
- Channel coding theorem