

# Information Theory and Coding

Shannon's communication model

Cédric RICHARD  
Université Côte d'Azur

# INFORMATION THEORY

## Models of communication

---

Models of communication are conceptual models used to explain the human communication process.

Following the basic concept, communication is the process of sending and receiving messages or transferring information from one part (sender) to another (receiver).

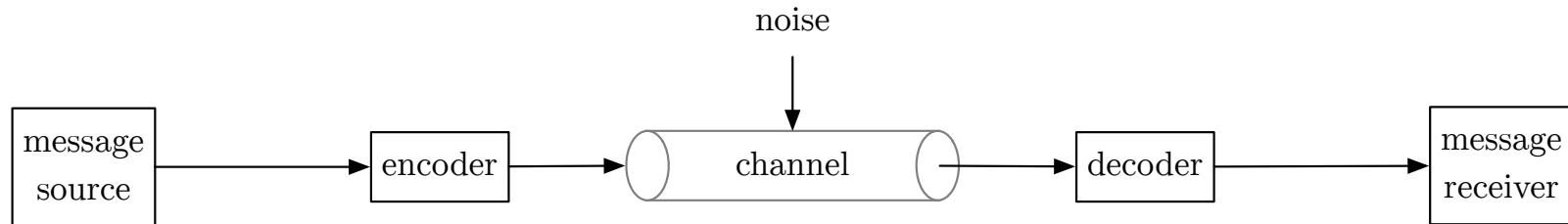
The Shannon-Weaver model was designed in 1949 to mirror the functioning of radio and telephone technology. It is referred to as the mother of all models.

This model has been expanded later by other scholars: Berlo (1960), ...

# INFORMATION THEORY

## Shannon's communication model

---



An information source, which produces a message

An encoder, which encodes the message into signals

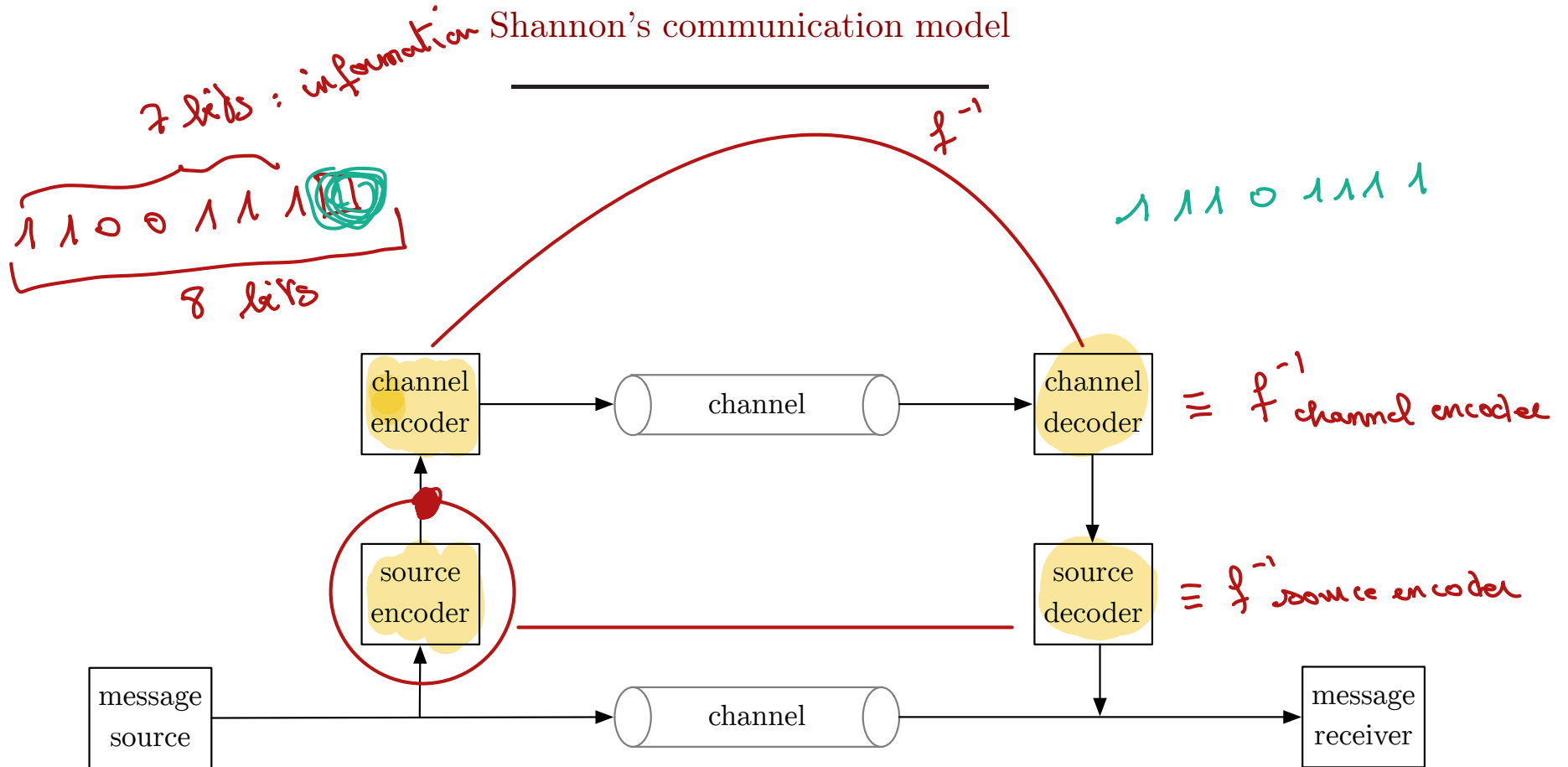
A channel, for which signals are adapted for transmission

A decoder, which reconstructs the encoded message

An information destination, where the message arrives

# INFORMATION THEORY

Shannon's communication model



# INFORMATION THEORY

## Objectives

---

Information theory studies the quantification, storage, and communication of information.

It was originally proposed by Claude Shannon in 1948 to find fundamental limits on signal processing and communication operations such as data compression.

Applications of fundamental topics of information theory include lossless data compression, lossy data compression, and channel coding.

Information theory is used in information retrieval, intelligence gathering, gambling, statistics, and even in musical composition.

A key measure is entropy. It quantifies the amount of uncertainty involved in the value of a random variable or the outcome of a random process.

# Information Theory and Coding

Quantitative measure of information

Cédric RICHARD  
Université Côte d'Azur

# SELF-INFORMATION

## Information content

---

Let  $A$  be an event with non-zero probability  $P(A)$ .

The greater the uncertainty of  $A$ , the larger the information  $h(A)$  provided by the realization of  $A$ . This can be expressed as follows:

$$h(A) = f\left(\frac{1}{P(A)}\right).$$

Function  $f(\cdot)$  must satisfy the following properties:

- ▷  $f(\cdot)$  is an increasing function over  $\mathbb{R}_+$
- ▷ information provided by 1 sure event is zero:  $\lim_{p \rightarrow 1} f(p) = 0$
- ▷ information provided by 2 independent events:  $f(p_1 \cdot p_2) = f(p_1) + f(p_2)$

This leads us to use the logarithmic function for  $f(\cdot)$

$A_1$  and  $A_2$  two events that are independent

ex:  $A_1 =$  "getting head with coin 1"

$A_2 =$  "getting head with coin 2"

Then:  $P(A_1 \text{ and } A_2) = P(A_1) \cdot P(A_2)$   
def. of independence

---

$$\begin{aligned} P(A_1, A_2) &= P(A_2 | A_1) P(A_1) \\ &= P(A_2 | A_2) P(A_2) \end{aligned}$$

with independence of  $A_1$  and  $A_2$ :

$$P(A_1, A_2) = P(A_1)P(A_2)$$

$$\Leftrightarrow \begin{aligned} P(A_2 | A_1) &= P(A_2) \\ P(A_1 | A_2) &= P(A_1) \end{aligned}$$

▷ information provided by 2 independent events:  $f(p_1 \cdot p_2) = f(p_1) + f(p_2)$

$$\underbrace{h(A_1, A_2)}_{\text{joint quantity of information}} \triangleq \underbrace{f\left(\frac{1}{P(A_1, A_2)}\right)}_{\text{joint probab.}} \stackrel{A_1 \text{ and } A_2 \text{ are indep.}}{=} f\left(\frac{1}{P(A_1)} \cdot \frac{1}{P(A_2)}\right)$$

I want this  
for 2 indep. events  
(Shannon's axiom)

$$\begin{aligned} &= f\left(\frac{1}{P(A_1)}\right) + f\left(\frac{1}{P(A_2)}\right) \\ &= h(A_1) + h(A_2) \end{aligned}$$



## SELF-INFORMATION

### Information content

---

**Lemme 1.** *Function  $f(p) = -\log_b p$  is the only one that is both positive, continue over  $(0, 1]$ , and that satisfies  $f(p_1 \cdot p_2) = f(p_1) + f(p_2)$ .*

**Proof.** The proof consists of the following steps:

1.  $f(p^n) = n f(p)$
2.  $f(p^{1/n}) = \frac{1}{n} f(p)$  after replacing  $p$  with  $p^{1/n}$
3.  $f(p^{m/n}) = \frac{m}{n} f(p)$  by combining the two previous equalities
4.  $f(p^q) = q f(p)$  where  $q$  is any positive rational number
5.  $f(p^r) = \lim_{n \rightarrow +\infty} f(p^{q_n}) = \lim_{n \rightarrow +\infty} q_n f(p) = r f(p)$  because rationals are dense in the reals


Let  $p$  and  $q$  in  $(0, 1]$ . One can write:  $p = q^{\log_q p}$ , which yields:

$$f(p) = f(q^{\log_q p}) = f(q) \log_q p.$$

We finally arrive at:  $f(p) = -\log_b p$

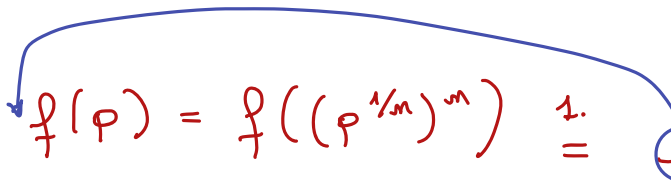
$$f(p) = -\log_b p$$

$$\log_b x = \frac{\ln x}{\ln b}$$


  
 basis of log

$b > 0$

2.  $f(p^{1/n}) = \frac{1}{n} f(p)$

$$f(p) = f((p^{1/n})^n) \stackrel{1.}{=} n f(p^{1/n})$$


$$\Rightarrow f(p^{1/n}) = \frac{1}{n} f(p) \quad \underline{\underline{2.}}$$

3.  $f(p^{m/n}) = \frac{m}{n} f(p)$

$$f(p^{m/n}) = f((p^{1/n})^m)$$

$$\stackrel{1.}{=} m f(p^{1/n})$$

$$\stackrel{2.}{=} \frac{m}{n} f(p)$$

$$p = q^{\log_q p} = \frac{\ln p}{\ln q} \quad \begin{matrix} \leftarrow < 0 \\ \leftarrow < 0 \end{matrix}$$

$\forall$   
0

$$q^{\log_q p} = q^{\frac{\ln p}{\ln q}} = e^{\frac{\ln p}{\ln q} \times \ln q} = e^{\ln p} = p$$

$$\textcircled{b} \quad q = e^{\ln q}$$

$$\textcircled{a} \quad e^{ab} = (e^a)^b$$

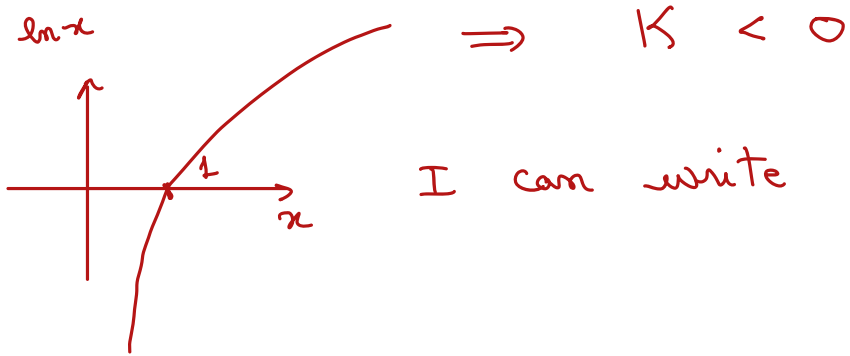
we have :  $p = q^{\log_q p}$  ,  $p, q \in ]0, 1]$

$$\begin{aligned} f(p) &= f(q^{\log_q p}) = \log_q p \times f(q) \\ &= \frac{\ln p}{\ln q} \times f(q) \end{aligned}$$

$$\Rightarrow \frac{f(p)}{f(q)} = \frac{\ln p}{\ln q} \Rightarrow \underbrace{f(p)} = K \underbrace{\ln p}_{< 0}$$

$$\forall p, q \in ]0, 1]$$

I want it positive



I can write  $K = -\frac{1}{\ln b}$  with  $b > 1$

$$f(p) = -\frac{\ln p}{\ln b}, \quad b > 1$$

$$\stackrel{\Delta}{=} -\log_b p, \quad b > 1$$

Check:  $\lim_{p \rightarrow 1} -\log_b p = 0$

# SELF-INFORMATION

## Information content

**Definition 1.** Let  $(\Omega, \mathcal{A}, P)$  be a probability space, and  $A$  an event of  $\mathcal{A}$  with non-zero probability  $P(A)$ . The information content of  $A$  is defined as:

$$h(A) = -\log_b P(A), \quad b > 1$$

**Unit.** The unit of  $h(A)$  depends on the base chosen for the logarithm.

- $b = 2$      $\triangleright \log_2$  : Shannon, bit (binary unit)
- $b = e$      $\triangleright \log_e$  : logon, nat (natural unit)
- $b = 10$      $\triangleright \log_{10}$  : Hartley, decit (decimal unit)

$\ln$

**Vocabulary.**  $h(\cdot)$  represents the *uncertainty* of  $A$ , or its *information content*.

$\log_2$      $x$     Shannon

8, 16

2048

$$\text{if } x = 2^n$$

$$\log_2 x = \log_2 2^n$$

$$= \frac{\ln 2^n}{\ln 2}$$

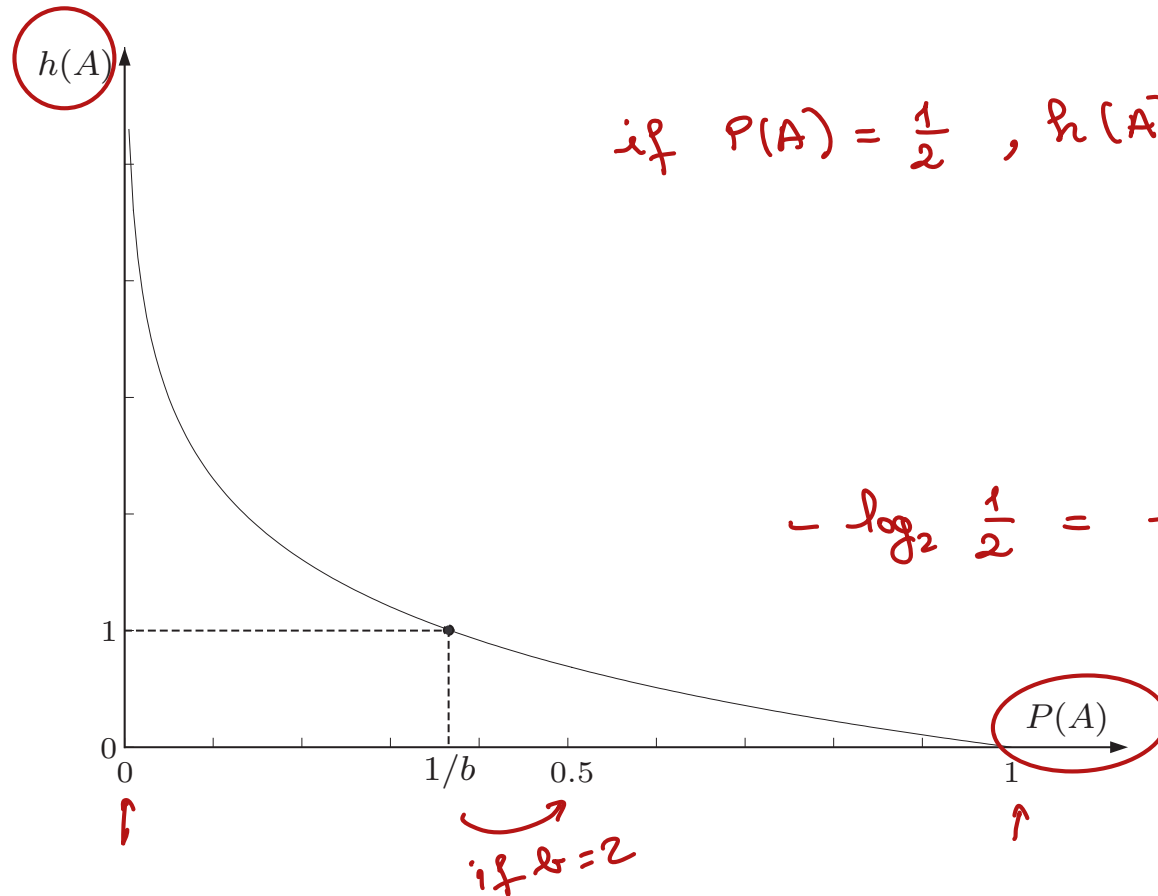
$$= n \frac{\ln 2}{\ln 2}$$

# SELF-INFORMATION

Information content

---

Information content or uncertainty:  $h(A) = -\log_b P(A)$



# SELF-INFORMATION

## Information content

---

**Example 1.** Consider a binary source  $S \in \{0, 1\}$  with  $P(0) = P(1) = 0.5$ . Information content conveyed by each binary symbol is equal to:  $h\left(\frac{1}{2}\right) = \log 2$ , namely, 1 bit or Shannon.

**Example 2.** Consider a source  $S$  that randomly selects symbols  $s_i$  among 16 equally likely symbols  $\{s_0, \dots, s_{15}\}$ . Information content conveyed by each symbol is  $\log 16$  Shannon, that is, 4 Shannon.

**Remark.** The bit in Computer Science (binary digit) and the bit in Information Theory (binary unit) do not refer to the same concept.

---



**Example 2.** Consider a source  $S$  that randomly selects symbols  $s_i$  among 16 equally likely symbols  $\{s_0, \dots, s_{15}\}$ . Information content conveyed by each symbol is  $\log 16$  Shannon, that is, 4 Shannon.

$$P(S = s_i) = \frac{1}{16}$$

$$h(\underbrace{S = s_i}_A) = -\log_2 \frac{1}{16} = \log_2 2^4 = 4 \text{ Sh (bit)}$$

A : " S generated  $s_i$  "

## SELF-INFORMATION

### Conditional information content

---

Self-information applies to 2 events  $A$  and  $B$ . Note that  $P(A, B) = P(A) P(B|A)$ .  
We get:

$$h(A, B) = -\log P(A, B) = -\log P(A) - \log P(B|A)$$

Note that  $-\log P(B|A)$  is the information content of  $B$  that is not provided by  $A$ .

**Definition 2.** *Conditional information content of  $B$  given  $A$  is defined as:*

$$h(B|A) = -\log P(B|A),$$

*that is:  $h(B|A) = h(A, B) - h(A)$ .*

**Exercise.** Analyze and interpret the following cases:  $A \subset B$ ,  $A = B$ ,  $A \cap B = \emptyset$ .

A : with  $P(A)$

$$h(A) = -\log_2 P(A) \quad \text{Sh}$$

---

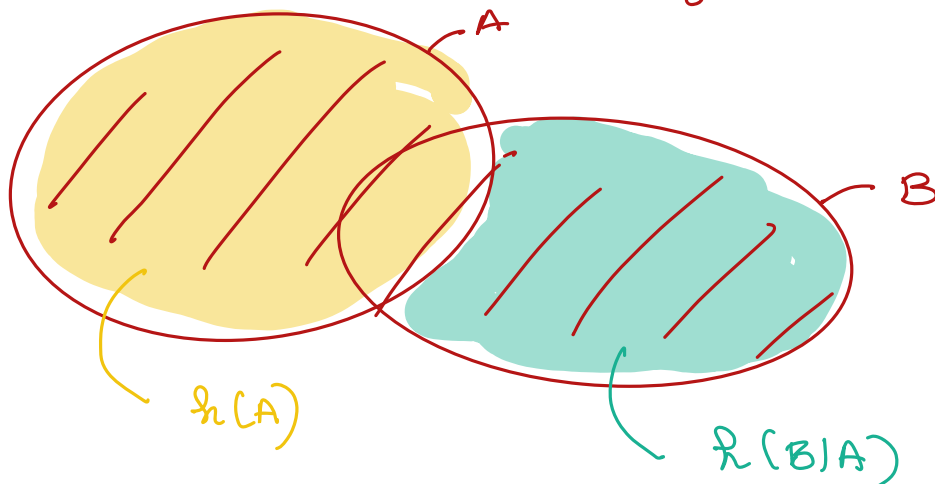
A, B : with  $P(A, B)$  joint proba.

$$h(A, B) = -\log_2 P(A, B) \quad \text{Sh}$$

extension :  $P(A, B) = P(A) P(B|A)$

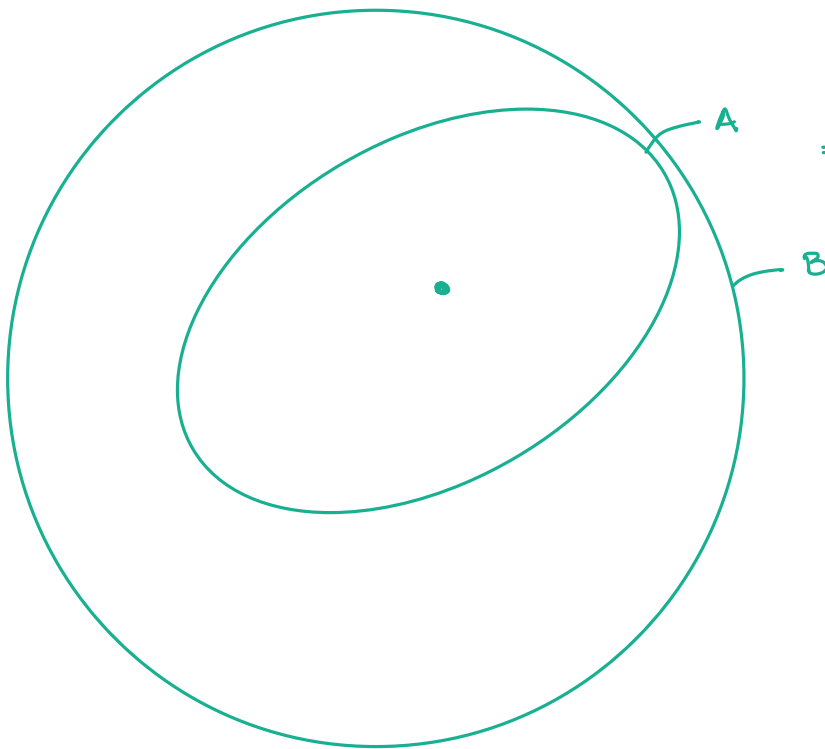
$$\begin{aligned} h(A, B) &= -\log_2 P(A, B) \\ &= -\log_2 [P(A) P(B|A)] \\ &= -\log_2 P(A) - \log_2 P(B|A) \\ &\stackrel{A}{=} \underbrace{h(A)}_{\text{Sh}} + \underbrace{h(B|A)}_{\text{Sh}} \end{aligned}$$

→ quantity of info  
provided by B given that  
you know A



**Exercise.** Analyze and interpret the following cases:  $A \subset B$ ,  $A = B$ ,  $A \cap B = \emptyset$ .

$$A \subset B$$



$$A \subset B$$

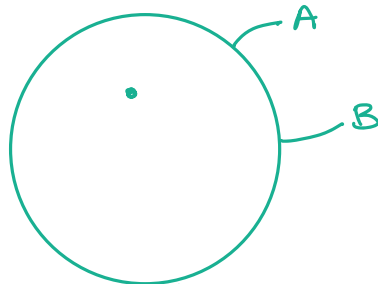
$$\Rightarrow P(B|A) = 1$$

$$h(A, B) = h(A) + h(B|A)$$

$$h(B|A) = -\log_2 P(B|A) = 0$$

$$\Rightarrow h(A, B) = h(A)$$

$$A = B$$



$$P(A|B) = P(B|A) = 1$$

$$h(A, B) = h(A) = h(B)$$

$$A \cap B = \emptyset$$

indépendant

$$P(A, B) = P(A)P(B)$$

$$\Rightarrow h(A, B) = h(A) + h(B)$$

because  $P(A|B) = P(A)$   
 $P(B|A) = P(B)$

# SELF-INFORMATION

## Mutual information content

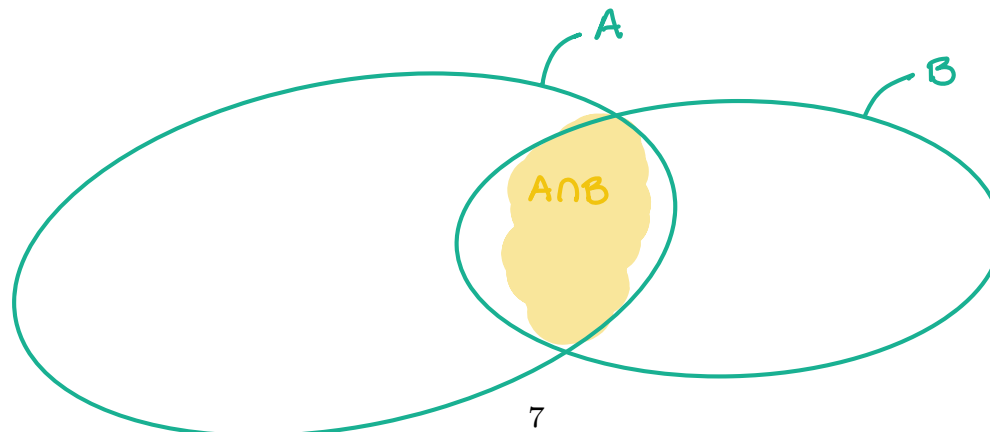
---

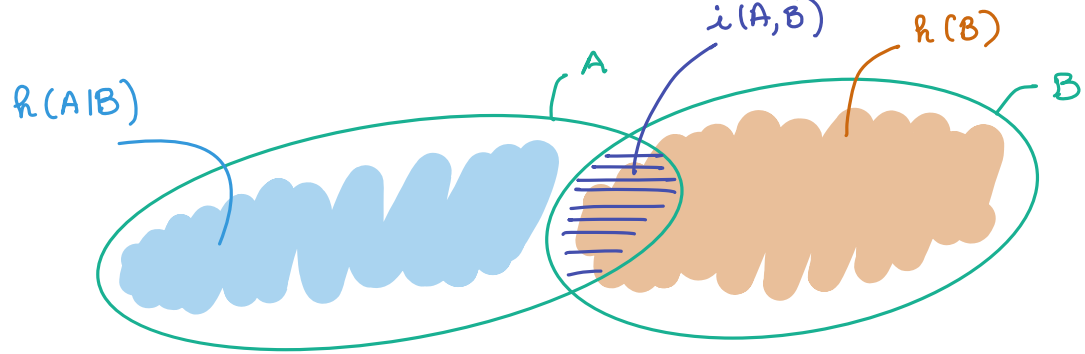
The definition of conditional information leads directly to another definition, that of mutual information, which measures information shared by two events.

**Definition 3.** We call mutual information of  $A$  and  $B$  the following quantity:

$$i(A, B) = h(A) - h(A|B) = h(B) - h(B|A).$$

**Exercise.** Analyze and interpret the following cases:  $A \subset B$ ,  $A = B$ ,  $A \cap B = \emptyset$ .

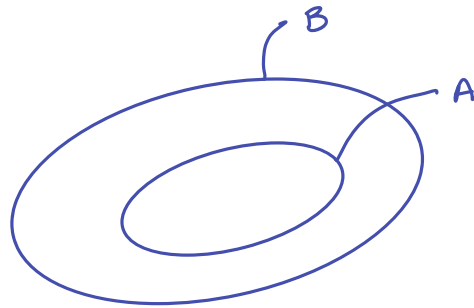




$$h(A) - h(A|B) \triangleq i(A, B) \triangleq h(B) - h(B|A)$$

$\underbrace{\hspace{10em}}_{\substack{\text{mutual information} \\ \text{"shared information" }}}$

**Exercise.** Analyze and interpret the following cases:  $A \subset B$ ,  $A = B$ ,  $A \cap B = \emptyset$ .



$$P(B|A) = 1$$

$$h(B|A) = 0$$

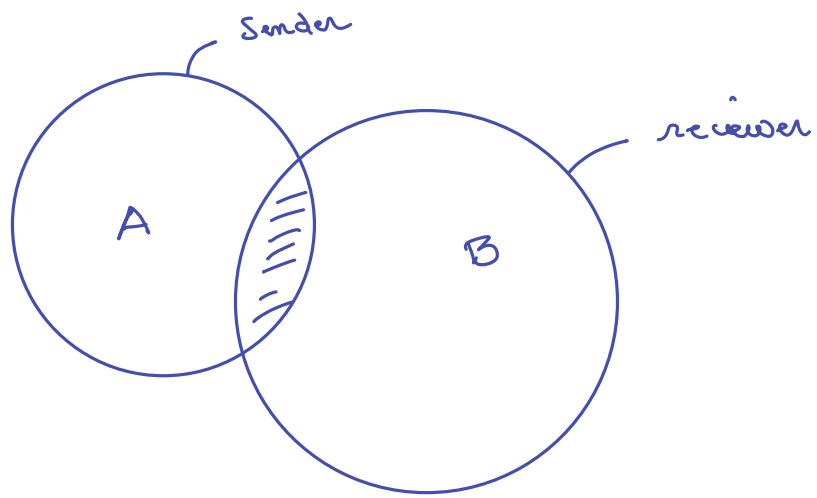
$$i(A, B) = h(B)$$

$$\boxed{A = B} \quad i(A, B) = h(A) = h(B)$$

independent.

$$\boxed{A \cap B = \emptyset}$$

$$h(A|B) = h(A) \Rightarrow i(A, B) = 0$$



good telecommunication :  $\max i(S, R)$

# ENTROPY OF A RANDOM VARIABLE

## Definition

---

Consider a memoryless stochastic source  $S$  with alphabet  $\{s_1, \dots, s_n\}$ . Let  $p_i$  be the probability  $P(S = s_i)$ .

The entropy of  $S$  is the average amount of information produced by  $S$ :

$$H(S) = E\{h(S)\} = - \sum_{i=1}^n p_i \log p_i.$$

**Definition 4.** Let  $X$  be a random variable that takes its values in  $\{x_1, \dots, x_n\}$ . Entropy of  $X$  is defined as follows:

$$H(X) = - \sum_{i=1}^n P(X = x_i) \log P(X = x_i).$$



$$S = \{s_1, \dots, s_m\} \quad \text{event, state, character}$$

$$P(S = s_1) = p_1 \Rightarrow h(S = s_1) = -\log p_1 \text{ Sh} \\ (\dots)$$

$$P(S = s_m) = p_m \Rightarrow h(S = s_m) = -\log p_m \text{ Sh}$$

$$H(S) = \sum_{i=1}^m p_i h(S = s_i) \quad \leftarrow \text{mean value } E\{\cdot\}$$

entropic  $\rightarrow$  Sh / event or state or character.

$$H(S) = \sum_{i=1}^m -p_i \log_2 p_i \quad \text{Sh / event}$$

$$H(S) = - \sum_{i=1}^m p_i \log_2 p_i \quad \text{Sh / event or state}$$

Alphabetical source : a, b, c, ..., z  
 $\hookrightarrow P(S = a) = \dots$

$$p \log p \quad \text{with } p = 1 \quad \rightarrow 0$$

$$p \log p \quad \text{with } p \rightarrow 0 \quad \rightarrow 0$$

$$\text{if : } P(S = s_i) = 1$$

$$\text{and } P(S = s_j) = 0, \quad \forall j \neq i$$

$$\Rightarrow H(S) = 0$$

# ENTROPY OF A RANDOM VARIABLE

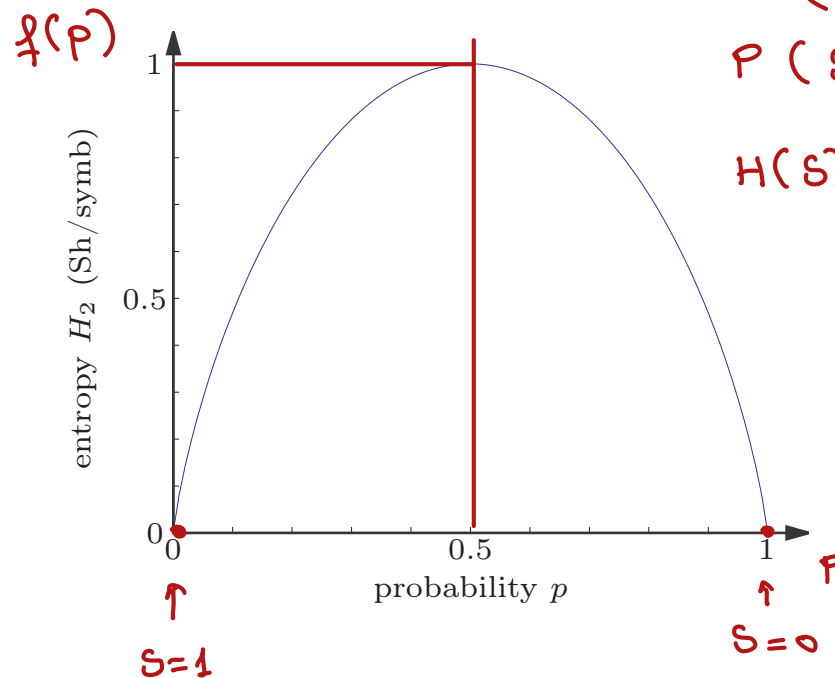
Example of a binary random variable

---

The entropy of a binary random variable is given by:

$$H(X) = -p \log p - (1 - p) \log(1 - p) \triangleq H_2(p).$$

$H_2(p)$  is called the binary entropy function.



$$P(S=0) = p$$

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

$$\begin{aligned} H(S) &= -p \log_2 p \\ &\quad - (1-p) \log_2 (1-p) \\ &= f(p) \end{aligned}$$

$$\begin{aligned}
 H(S) &= 2 \times \left[ -\frac{1}{2} \log \frac{1}{2} \right] \\
 &= 2 \times \frac{1}{2} \\
 &= 1 \text{ Sh /state.}
 \end{aligned}$$

# ENTROPY OF A RANDOM VARIABLE

## Notation and preliminary properties

---

**Lemme 2** (Gibbs' inequality). *Consider 2 discrete probability distributions with mass functions  $(p_1, \dots, p_n)$  and  $(q_1, \dots, q_n)$ . We have:*

$$\sum_i p_i = 1, \quad p_i \geq 0 \quad \forall i$$

$$\sum_i q_i = 1, \quad q_i \geq 0 \quad \forall i$$

$$\sum_{i=1}^n p_i \log \frac{q_i}{p_i} \leq 0$$

*Equality is achieved when  $p_i = q_i$  for all  $i$*

**Proof.** The proof is carried out in the case of the neperian logarithm. Observe that  $\ln x \leq x - 1$ , with equality for  $x = 1$ . Let  $x = \frac{q_i}{p_i}$ . We have:

$$\sum_{i=1}^n p_i \ln \frac{q_i}{p_i} \leq \sum_{i=1}^n p_i \left( \frac{q_i}{p_i} - 1 \right) = 1 - 1 = 0.$$

We have :  $\ln x \leq x - 1$  ,  $\forall x > 0$   
 $q_i/p_i > 0$

$$\sum_{i=1}^n p_i \log \frac{q_i}{p_i} \leq 0$$

$$\sum_{i=1}^n p_i \log \frac{q_i}{p_i} \leq \underbrace{\sum_{i=1}^n p_i \times \left( \frac{q_i}{p_i} - 1 \right)}_{\sum_{i=1}^n (q_i - p_i)}$$

if  $p_i = q_i$  ,  $\forall i$

$$\Rightarrow \sum_i p_i \log \frac{q_i}{p_i} = 0$$

$$H(S) = - \sum_i p_i \log_2 p_i$$

$$\sum_{i=1}^n q_i - \sum_{i=1}^n p_i = 1 - 1 = 0$$

$$\sum_{i=1}^{(n)} p_i \log \frac{q_i}{p_i} \leq 0$$

$$\Rightarrow \sum_{i=1}^n \left[ p_i \log q_i - p_i \log p_i \right] \leq 0$$

$$\Rightarrow \underbrace{- \sum_{i=1}^n p_i \log p_i}_{H(S)} + \sum_{i=1}^n p_i \log q_i \leq 0$$

We set  $q_i = \frac{1}{n}$

$$H(S) - \log_2 n \underbrace{\sum_{i=1}^n p_i}_1 \leq 0$$

$$\Rightarrow \boxed{H(S) \leq \log_2 n}$$

2 states :  $n = 2$

$$H(S) \leq 1 \text{ sh}$$

(see before)

$$H(S) = \log_2 n$$

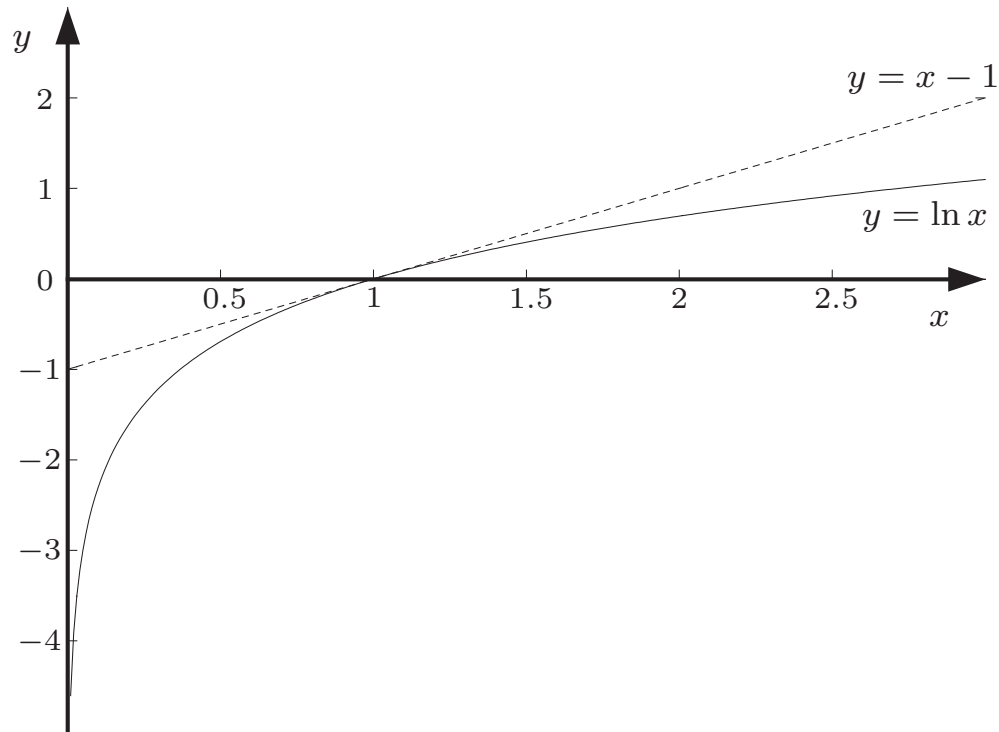
$$\text{if } p_i = q_i = \frac{1}{n}, \forall i = 1, \dots, n$$

# ENTROPY OF A RANDOM VARIABLE

Notation and preliminary properties

---

Graphical checking of inequality  $\ln x \leq x - 1$



# ENTROPY OF A RANDOM VARIABLE

## Properties

---

**Property 1.** *The entropy satisfies the following inequality:*

$$H_n(p_1, \dots, p_n) \leq \log n,$$

*Equality is achieved by the uniform distribution, that is,  $p_i = \frac{1}{n}$  for all  $i$ .*

**Proof.** Based on Gibbs' inequality, we set  $q_i = \frac{1}{n}$ .

Uncertainty about the outcome of an experiment is maximum when all possible outcomes are equiprobable.



# ENTROPY OF A RANDOM VARIABLE

## Properties

---

**Property 2.** *The entropy increases as the number of possible outcomes increases.*

**Proof.** Let  $X$  be a discrete random variable with values in  $\{x_1, \dots, x_n\}$  and probabilities  $(p_1, \dots, p_n)$ , respectively. Consider that state  $x_k$  is split into two substates  $x_{k_1}$  et  $x_{k_2}$ , with non-zero probabilities  $p_{k_1}$  et  $p_{k_2}$  such that  $p_k = p_{k_1} + p_{k_2}$ .

Entropy of the resulting random variable  $X'$  is given by:

$$\begin{aligned} H(X') &= H(X) + p_k \log p_k - p_{k_1} \log p_{k_1} - p_{k_2} \log p_{k_2} \\ &= H(X) + p_{k_1} (\log p_k - \log p_{k_1}) + p_{k_2} (\log p_k - \log p_{k_2}). \end{aligned}$$

The logarithmic function being strictly increasing, we have:  $\log p_k > \log p_{k_i}$ . This implies:  $H(X') > H(X)$ .

**Interpretation.** Second law of thermodynamics

# ENTROPY OF A RANDOM VARIABLE

## Properties

---

**Property 3.** *The entropy  $H_n$  is a concave function of  $p_1, \dots, p_n$ .*

**Proof.** Consider 2 discrete probability distributions  $(p_1, \dots, p_n)$  and  $(q_1, \dots, q_n)$ . We need to prove that, for every  $\lambda$  in  $[0, 1]$ , we have:

$$H_n(\lambda p_1 + (1 - \lambda)q_1, \dots, \lambda p_n + (1 - \lambda)q_n) \geq \lambda H_n(p_1, \dots, p_n) + (1 - \lambda)H_n(q_1, \dots, q_n).$$

By setting  $f(x) = -x \log x$ , we can write:

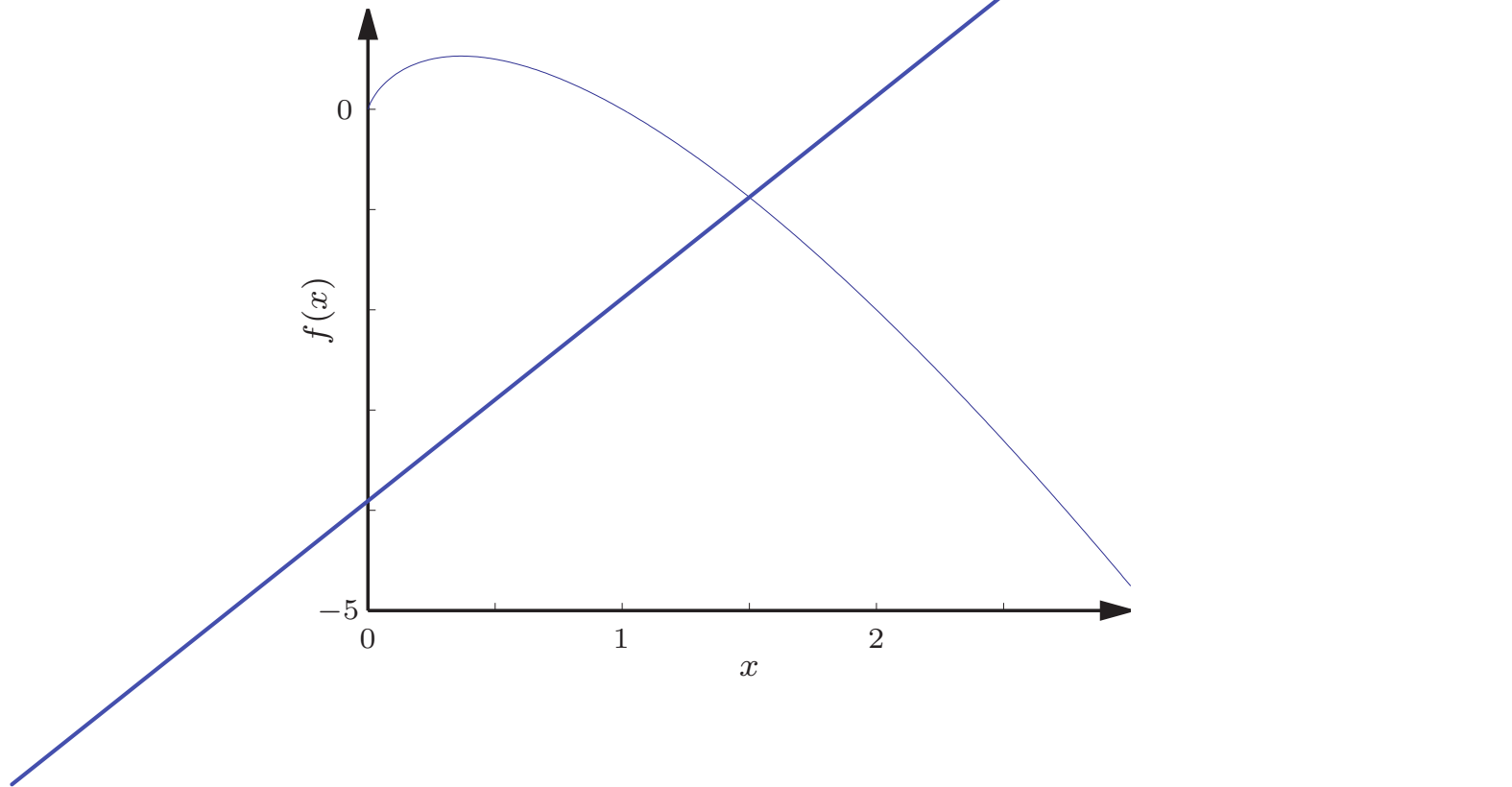
$$H_n(\lambda p_1 + (1 - \lambda)q_1, \dots, \lambda p_n + (1 - \lambda)q_n) = \sum_{i=1}^n f(\lambda p_i + (1 - \lambda)q_i).$$

The result is a direct consequence of the concavity of  $f(\cdot)$  and Jensen's inequality.

# ENTROPY OF A RANDOM VARIABLE

## Properties

Graphical checking of the concavity of  $f(x) = -x \log x$



# ENTROPY OF A RANDOM VARIABLE

## Properties

---

Concavity of  $H_n$  can be generalized to any number  $m$  of distributions.

**Property 4.** *Given  $\{(q_{1j}, \dots, q_{nj})\}_{j=1}^m$  a finite set of discrete probability distributions, the following inequality is satisfied:*

$$H_n\left(\sum_{j=1}^m \lambda_j q_{1j}, \dots, \sum_{j=1}^m \lambda_j q_{mj}\right) \geq \sum_{j=1}^m \lambda_j H_n(q_{1j}, \dots, q_{mj}),$$

where  $\{\lambda_j\}_{j=1}^m$  is any set of constants in  $[0, 1]$  such that  $\sum_{j=1}^m \lambda_j = 1$ .

**Proof.** As in the previous case, the demonstration of this inequality is based on the concavity of  $f(x) = -x \log x$  and Jensen's inequality.

# PAIR OF RANDOM VARIABLES

## Joint entropy

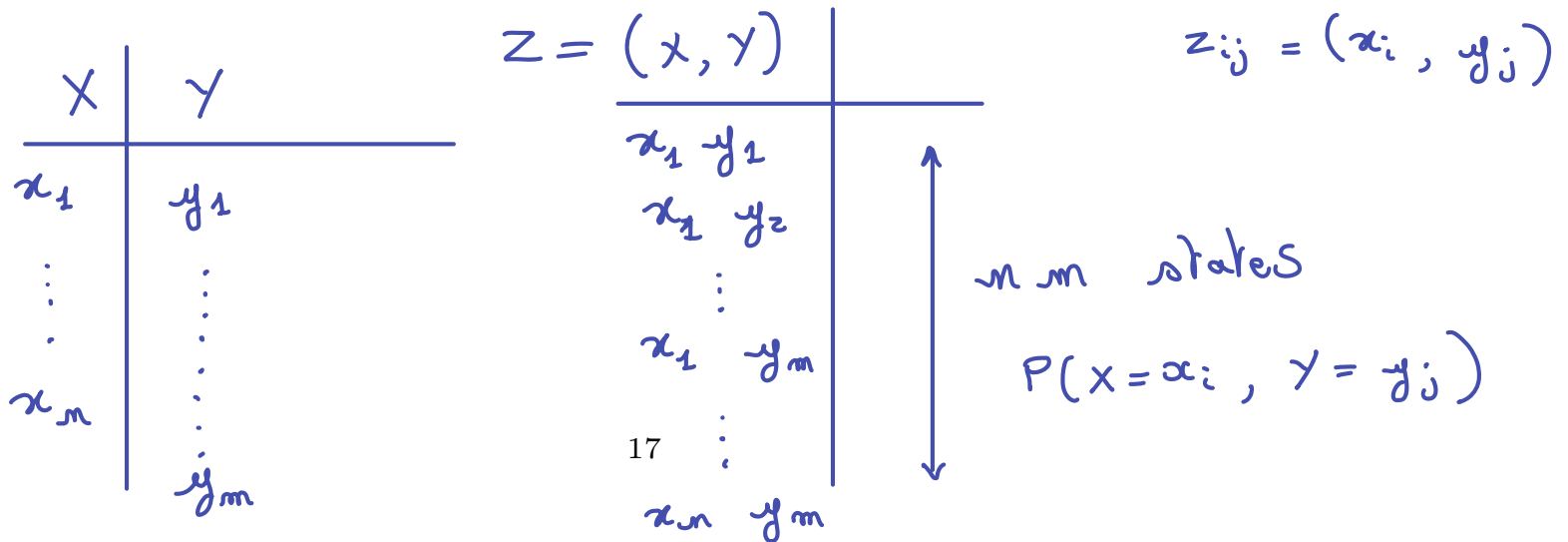
**Definition 5.** Let  $X$  and  $Y$  be two random variables with values in  $\{x_1, \dots, x_n\}$  and  $\{y_1, \dots, y_m\}$ , respectively. The joint entropy of  $X$  and  $Y$  is defined as:

$$H(X, Y) \triangleq - \sum_{i=1}^n \sum_{j=1}^m P(X = x_i, Y = y_j) \log P(X = x_i, Y = y_j).$$

▷ The joint entropy is symmetric:  $H(X, Y) = H(Y, X)$

can  $P(x, y) = P(y, x)$

**Example.** Case of two independent random variables



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

$$= - \sum_{i=1}^m \sum_{j=1}^m P(Z = z_{ij}) \log_2 P(Z = z_{ij})$$

$$= - \sum_i \sum_j P(X = x_i, Y = y_j) \log_2 P(X = x_i, Y = y_j)$$

Sh / pair of states  
of  $X$  and  $Y$

$$H(X, Y) \leq \overbrace{\log_2 n + \log_2 m}^{\log_2(n \cdot m)}$$

if  $X$  and  $Y$  indep.  $H(X, Y) = H(X) + H(Y)$

$\wedge$                        $\wedge$   
 $\log n$                  $\log m$

## PAIR OF RANDOM VARIABLES

### Conditional entropy

---

**Definition 6.** Let  $X$  and  $Y$  be two random variables with values in  $\{x_1, \dots, x_n\}$  and  $\{y_1, \dots, y_m\}$ , respectively. The conditional entropy of  $X$  given  $Y = y_j$  is:

$$H(X|Y = y_j) \triangleq - \sum_{i=1}^n P(X = x_i|Y = y_j) \log P(X = x_i|Y = y_j).$$

$H(X|Y = y_j)$  is the amount of information needed to describe the outcome of  $X$  given that we know that  $Y = y_j$ .

**Definition 7.** The conditional entropy of  $X$  given  $Y$  is defined as:

$$H(X|Y) \triangleq \sum_{j=1}^m P(Y = y_j) H(X|Y = y_j),$$

**Example.** Case of two independent random variables





$$H(X|Y) \triangleq \sum_{j=1}^m P(Y = y_j) H(X|Y = y_j),$$

**Example.** Case of two independent random variables

$X$  and  $Y$  independent

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

$$\text{because } P(X = x_i | Y = y_j) = P(X = x_i)$$

$$\begin{aligned} \Rightarrow H(X|Y) &= \sum_{j=1}^m P(Y = y_j) H(X) \\ &= H(X) \underbrace{\sum_{j=1}^m P(Y = y_j)}_1 \\ &= H(X) \end{aligned}$$

# PAIR OF RANDOM VARIABLES

Relations between entropies

---

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

These equalities can be obtained by first writing:

$$\log P(X = x, Y = y) = \log P(X = x|Y = y) + \log P(Y = y),$$

and then taking the expectation of each member.

**Property 5** (chain rule). *The joint entropy of  $n$  random variables can be evaluated using the following chain rule:*

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

$$P(X = x_i, Y = y_j) = P(X = x_i) P(Y = y_j | X = x_i)$$

$-\log$

$$\underbrace{-\log P(X = x_i, Y = y_j)}_{h(X=x_i, Y=y_j)} = \underbrace{-\log P(X = x_i)}_{h(X=x_i)} - \underbrace{\log P(Y = y_j | X = x_i)}_{h(Y=y_j | X=x_i)}$$

Rm:  $h(A, B) = h(A) + h(B|A)$

$$A = (X = x_i)$$

$$B = (Y = y_j)$$

$$H(X) = E_X \{h(X=x)\}$$

$$H(X, Y) = E_{X,Y} \{h(X=x, Y=y)\}$$

$E_{X,Y} \{ \cdot \}$

$$H(X, Y) = E_{X,Y} \{h(X=x)\} + E_{X,Y} \{h(Y=y | X=x)\}$$

$$= E_Y \{ E_X \{h(X=x)\} \} + E_Y \{ E_X \{h(Y=y | X=x)\} \}$$

$$= E_Y \{ H(X) \} + E_X \{ H(Y | X=x) \}$$

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

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

## PAIR OF RANDOM VARIABLES

Relations between entropies

---

Each term of  $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$  is positive. We can conclude that:

$$\begin{aligned} H(X) &\leq H(X, Y) \\ H(Y) &\leq H(X, Y) \end{aligned}$$

$$\underbrace{H(X, Y)}_{\geq 0} = \underbrace{H(X)}_{\geq 0} + \underbrace{H(Y|X)}_{\geq 0} \Rightarrow$$

$$\begin{aligned} H(X, Y) &\geq H(X) \\ H(X, Y) &\geq H(Y|X) \end{aligned}$$

# PAIR OF RANDOM VARIABLES

## Relations between entropies

---

From the *generalized concavity* of the entropy, setting  $q_{ij} = P(X = x_i|Y = y_j)$  and  $\lambda_j = P(Y = y_j)$ , we get the following inequality:

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

Conditioning a random variable reduces its entropy. Without proof, this can be generalized as follows:

**Property 6** (entropy decrease with conditioning). *The entropy of a random variable decreases with successive conditionings, namely,*

$$H(X_1|X_2, \dots, X_n) \leq \dots \leq H(X_1|X_2, X_3) \leq H(X_1|X_2) \leq H(X_1),$$

where  $X_1, \dots, X_n$  denote  $n$  discrete random variables.

# PAIR OF RANDOM VARIABLES

## Relations between entropies

---

Consider  $X$  and  $Y$  two random variables, respectively with values in  $\{x_1, \dots, x_n\}$  and  $\{y_1, \dots, y_m\}$ . We have:

$$\textcircled{1} \quad 0 \leq H(X|Y) \leq \textcircled{2} H(X) \leq \textcircled{3} H(X, Y) \leq \textcircled{4} H(X) + H(Y) \leq \textcircled{5} 2H(X, Y).$$

①: Entropie is positive (linear combination of  $f(x) = -x \log x$   $1 \leq x \leq 1$ )

② see previous slide (concavity) - not demonstrated.

③  $H(X, Y) = H(X) + H(Y|X) \Rightarrow H(X) \leq H(X, Y)$  because  $H(Y|X) \geq 0$

④  $H(X, Y) \stackrel{\textcircled{3}}{=} H(X) + H(Y|X)$  and  $H(Y|X) \stackrel{\textcircled{2}}{\leq} H(Y)$   
 $\leq H(X) + H(Y)$

⑤ ③ applied to  $X$  and  $Y$  <sup>22</sup>  $H(X) \leq H(X, Y)$   
 $+ H(Y) \leq H(X, Y)$

$$\overline{H(X) + H(Y) \leq 2 H(X, Y)} ,$$

## PAIR OF RANDOM VARIABLES

### Mutual information

**Definition 8.** The mutual information of two random variables  $X$  and  $Y$  is defined as follows:

$$I(X, Y) \triangleq H(X) - H(X|Y) = H(Y) - H(Y|X)$$

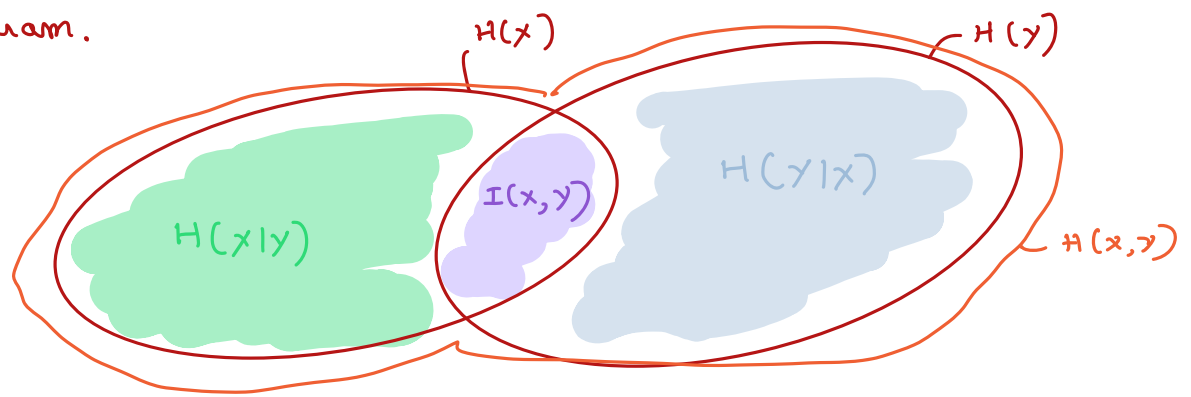
or, equivalently,

$$I(X, Y) \triangleq \sum_{i=1}^n \sum_{j=1}^m P(X = x_i, Y = y_j) \log \frac{P(X = x_i, Y = y_j)}{P(X = x_i) P(Y = y_j)}.$$

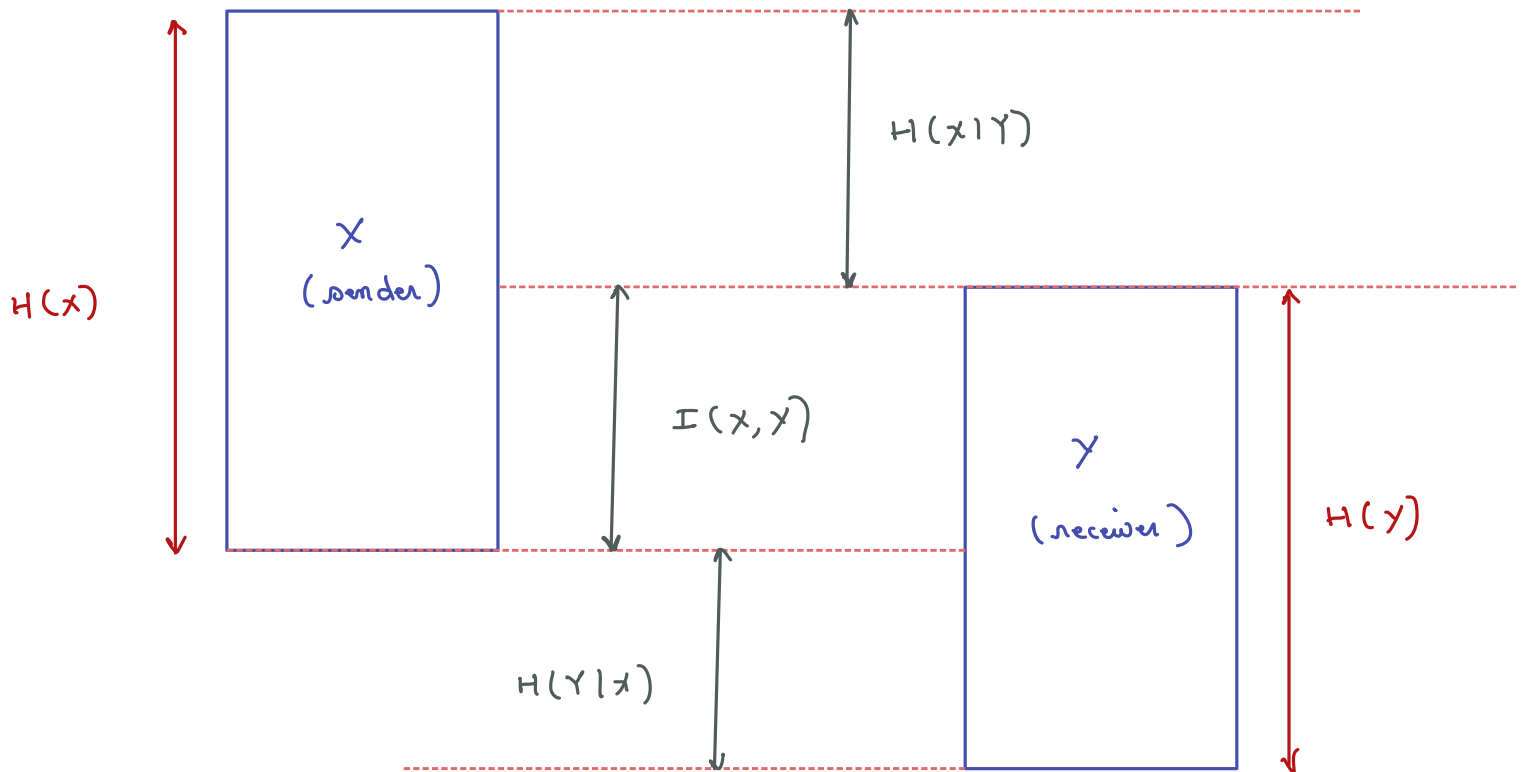
The mutual information quantifies the amount of information obtained about one random variable through observing the other random variable.

**Exercise.** Case of two independent random variables

Venn diagram.



$$H(x, y) = I(x, y) + H(x|y) + H(y|x)$$





# PAIR OF RANDOM VARIABLES

## Mutual information

---

In order to give a different interpretation of mutual information, the following definition is recalled beforehand.

**Definition 9.** *We call the Kullback-Leibler distance between two distributions  $P_1$  and  $P_2$ , here supposed to be discrete, the following quantity:*

$$d(P_1, P_2) = \sum_{x \in X(\Omega)} P_1(X = x) \log \frac{P_1(X = x)}{P_2(X = x)}.$$

The mutual information corresponds to the Kullback-Leibler distance between the marginal distributions and the joint distribution of  $X$  and  $Y$ .

# PAIR OF RANDOM VARIABLES

## Venn diagram

---

A Venn diagram can be used to illustrate relationships among measures of information: entropy, joint entropy, conditional entropy and mutual information.

