

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.

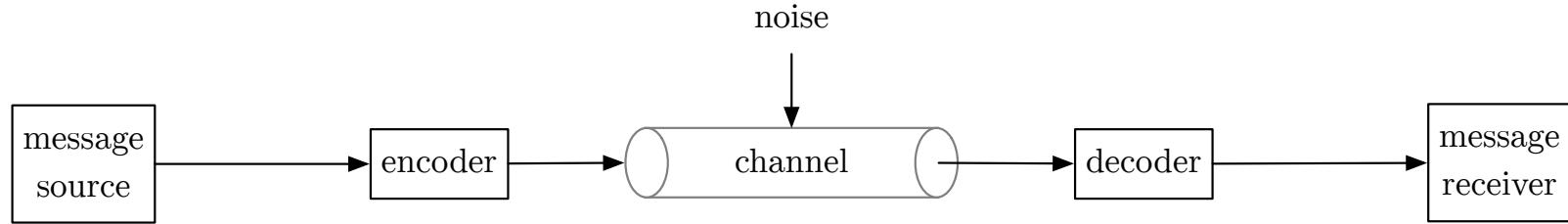
Following the basic concept, communication is the process of sending and receiving messages or transferring information from one part (sender) to another (receiver).

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

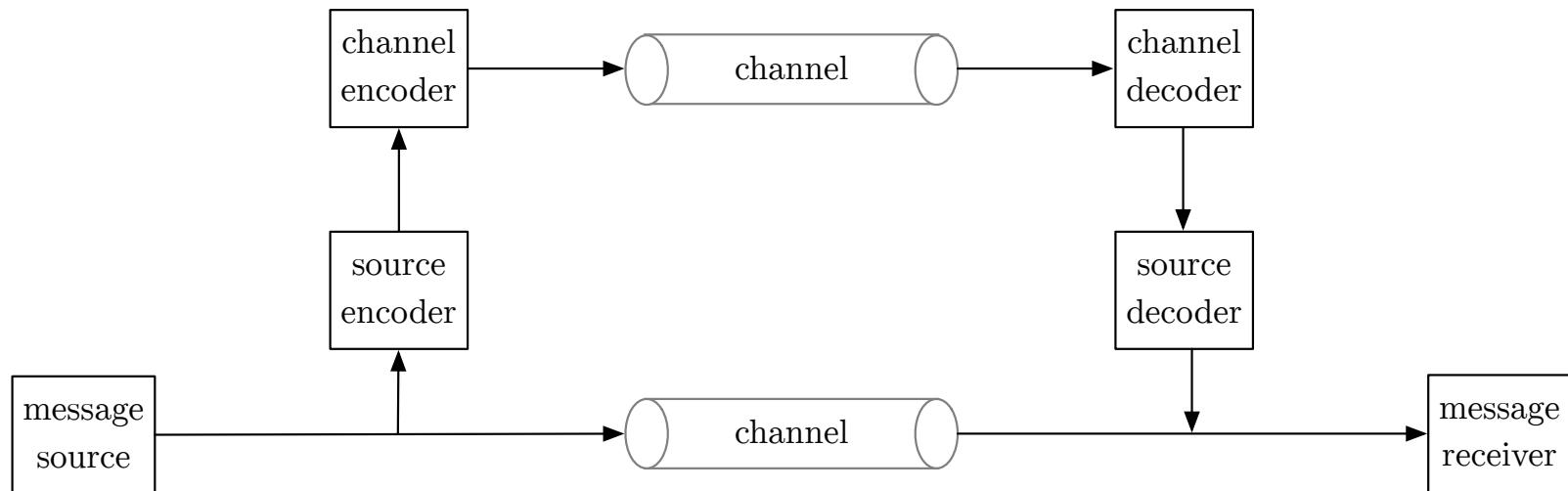
A channel, for which signals are adapted for transmission

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 and B two events which are independent.

By definition : $P(A, B) = P(A)P(B)$

$$P(A \cap B)$$

"and"

indep.

$$h(A, B) \triangleq f\left(\frac{1}{P(A, B)}\right) = f\left(\frac{1}{P(A)} \cdot \frac{1}{P(B)}\right)$$

we want

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

SELF-INFORMATION

Information content

Lemme 1. Function $f(p) = -\log_b p$ is the only one that is both positive, continuous 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: $f(p) = -\log_b p$

$$\underline{\underline{1.}} \quad f(p^m) = f(p^{m-1} \cdot p) = f(p^{m-1}) + f(p)$$

$$= \dots$$

$$f(p^m) = m f(p)$$

$$= f(p) + \underbrace{\dots + f(p)}_{m \text{ times}}$$

$$= m f(p)$$

$$\underline{\underline{2.}} \quad f(p) = f((p^{1/m})^m) \stackrel{1.}{=} m f(p^{1/m})$$

$$\Rightarrow f(p^{1/m}) = \frac{1}{m} f(p)$$

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

$$\underline{\underline{3.}} \quad f(p^{m/m}) = f((p^{1/m})^m) \stackrel{1.}{=} m f(p^{1/m})$$

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

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

4. Because of 3.

$$f(p^q) = q f(p) \quad \text{for all } q \in \mathbb{Q}^+ \text{ rationals.}$$

5. Because \mathbb{Q} is dense in \mathbb{R} , which means that

$\forall x \in \mathbb{R}, \exists (q_1, \dots, q_m)$, such that

$$x, \dots \downarrow$$

$$x = \lim_{n \rightarrow +\infty} q_n$$

$$f(p^x) = x f(p) , \quad \forall x \in \mathbb{R}$$

by continuity of f

6.

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

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

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

$$f(p) = f(q^{\log_q p}) \stackrel{\text{def}}{=} \log_q p \cdot f(q)$$

$$\Rightarrow \frac{f(p)}{f(q)} = \log_q p = \frac{\ln p}{\ln q}$$

$$\Rightarrow \boxed{f(p) = C \ln p}$$

where $C \in \mathbb{R}$

I want $f(p) > 0$, $\forall p \in [0, 1]$
 ↳ probability.

⇒ C must be negative

Let say that $C = \frac{-1}{\ln b}$ with $b > 1$

$$\Rightarrow \boxed{f(p) = -\log_b p}$$

SELF-INFORMATION

Information content

Definition 1. Let (Ω, \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:

$$h(A) = -\log P(A).$$

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

- ▷ \log_2 : Shannon, bit (binary unit)
- ▷ \log_e : logon, nat (natural unit)
- ▷ \log_{10} : Hartley, decit (decimal unit)

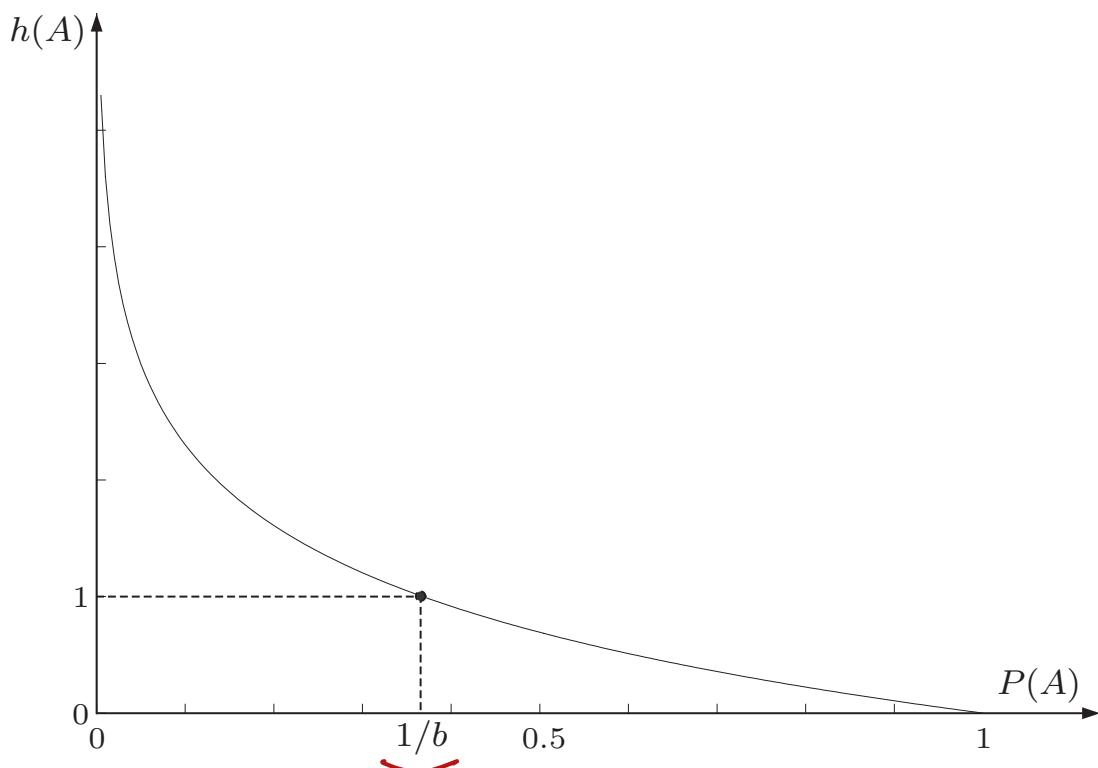
Vocabulary. $h(\cdot)$ represents the *uncertainty* of A , or its *information content*.

$P(A) \searrow$, uncertainty \nearrow and $R(A) \uparrow$

SELF- INFORMATION

Information content

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



→ if $b = 2$, $h(1/2) = 1$ Shannon

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 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.

$$\begin{aligned}
 h(S=0) &= -\log_2 P(S=0) \\
 &= -\log_2\left(\frac{1}{2}\right) \\
 &= \log_2 2 \\
 &= 1 \text{ Sh}
 \end{aligned}$$

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.

$$\begin{aligned}
 h(S=s_i) &= -\log_2 P(S=s_i) \\
 &= -\log_2\left(\frac{1}{16}\right) \\
 &= \log_2(2^4) \\
 &= 4 \text{ Sh}
 \end{aligned}$$

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:

$$h(A, B) = -\log P(A, B) = \underbrace{-\log P(A)}_{h(A)} + \underbrace{-\log P(B|A)}_{h(B|A)}$$

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: $A \subset B$, $A = B$, $A \cap B = \emptyset$.

$$\begin{aligned}
 h(A, B) &= -\log_2 P(A, B) \\
 &= -\log_2 P(A) - \log_2 P(B|A) \\
 &= h(A) + R(B|A)
 \end{aligned}$$

$h(B|A) = ?$ when A and B independent.

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:*

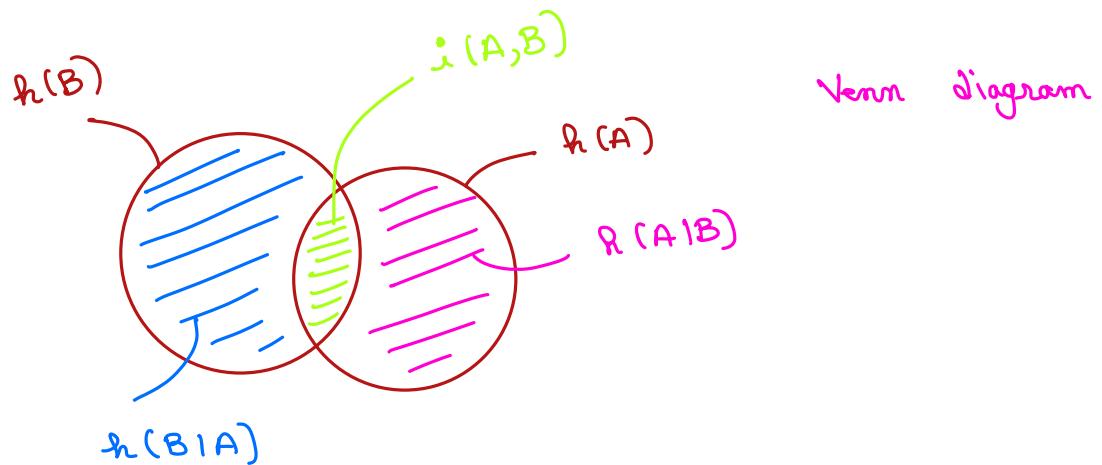
$$i(A, B) = h(A) - h(A|B) = h(B) - h(B|A).$$

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

$$\begin{aligned} h(A, B) &= h(A) + R(B|A) \\ &= h(B) + R(A|B) \end{aligned}$$

$$\begin{aligned} i(A, B) &\triangleq h(A) - h(A|B) \\ &\triangleq h(B) - h(B|A) \end{aligned}$$

mutual information shared by A and B



$$\begin{aligned} h(A, B) &= h(B|A) + h(A) \\ &= h(A|B) + h(B) \\ &= h(B|A) + h(A|B) + i(A, B) \end{aligned}$$

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 :

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

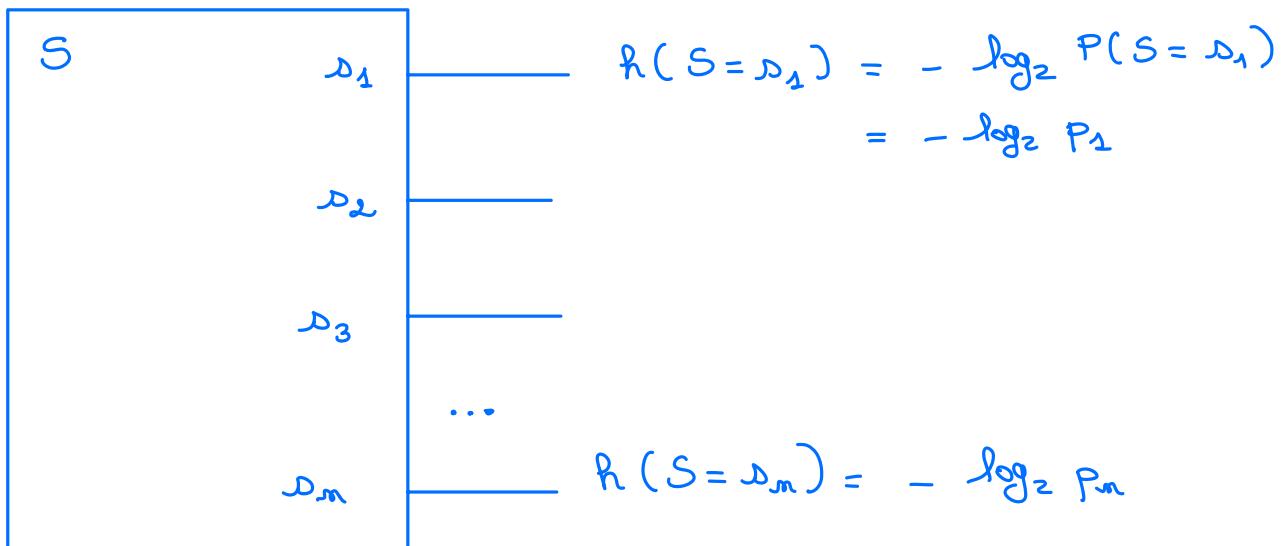
sh / state

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:

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

$$S \in \{s_1, \dots, s_m\}$$

$$P(S = s_i) = p_i, \quad i \in \{1, \dots, m\}$$



$$H(S) = \text{Mean} \{h(S)\}$$

$$= \sum_{i=1}^m p_i h(S=s_i)$$

$$= \sum_{i=1}^m p_i (-\log_2 p_i)$$

$$H(S) = - \sum_{i=1}^m p_i \log_2 p_i$$

unit: Sh / state of S

Binary Source

Case of a binary source : $S = \{0, 1\}$

$$P(S=0) = p \Rightarrow P(S=1) = 1 - P(S=0) \\ = 1 - p$$

$$H(S) = -p \log_2 p - (1-p) \log_2 (1-p)$$

$$= f(p) \quad \text{with} \quad p \in [0, 1]$$

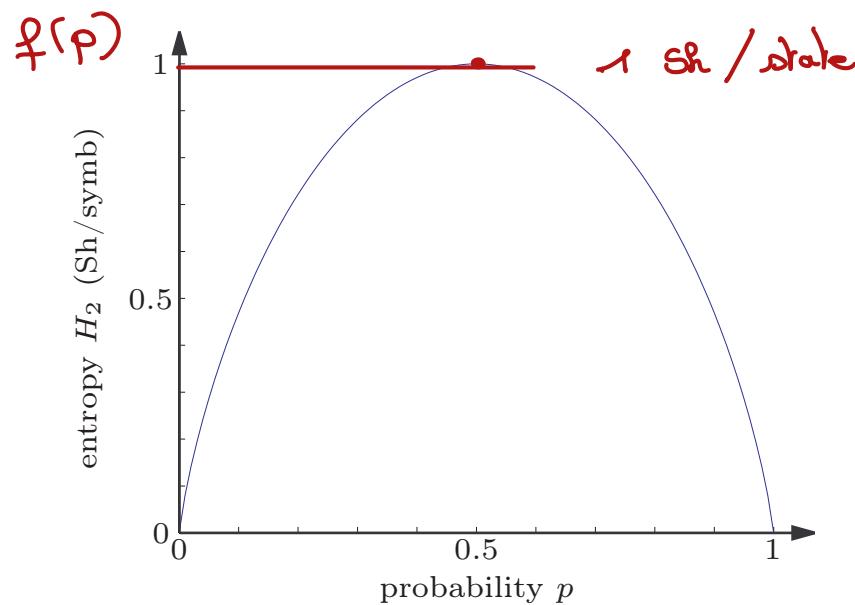
ENTROPY OF A RANDOM VARIABLE

Example of a binary random variable

The entropy of a binary random variable is given by:

$$H(X) = -p \log p - (1-p) \log(1-p) \triangleq \text{f } \cancel{H}(p) \quad f(p)$$

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



$$\lim_{p \rightarrow 1} \left(-p \underbrace{\log_2 p}_{0} - (1-p) \log_2 (1-p) \right) = \lim_{p \rightarrow 0} p \log_2 p = 0$$

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:

$$\sum_{i=1}^n p_i \log \frac{q_i}{p_i} \leq 0 \quad \text{pmf } (p_1, \dots, p_n) \text{ and } (q_1, \dots, q_n)$$

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 : $\ln x \leq x - 1$, $\forall x \in]0, +\infty[$

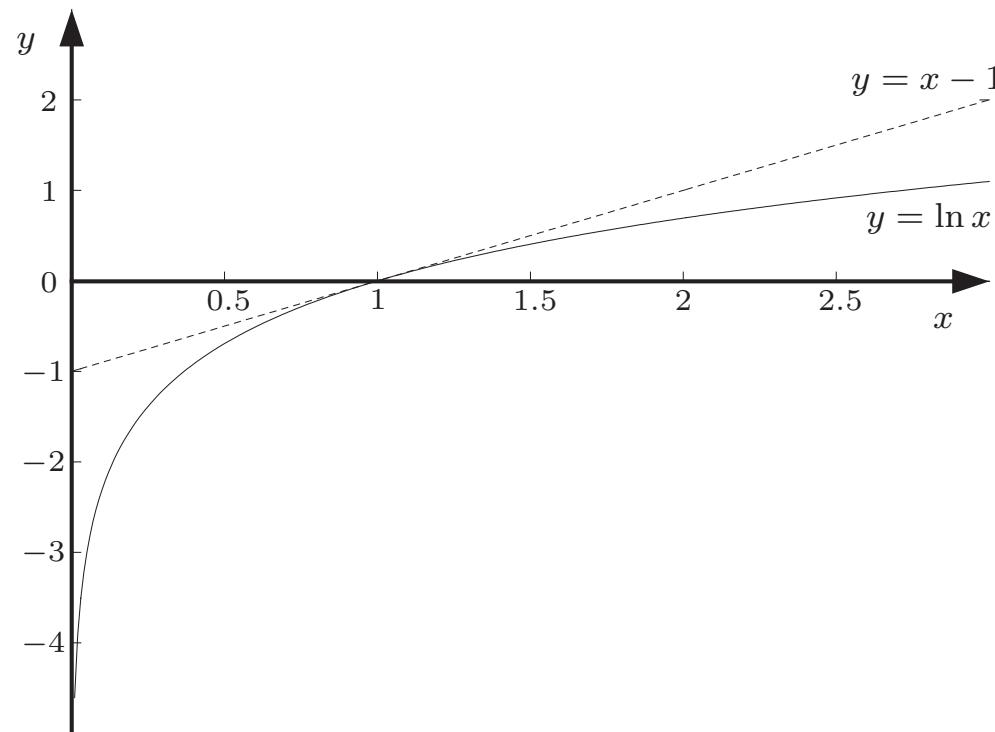
We write $x = \frac{q_i}{p_i} > 0$

$$\sum_{i=1}^m p_i \ln \left(\frac{q_i}{p_i} \right) \leq \underbrace{\sum_{i=1}^m p_i \left(\frac{q_i}{p_i} \right)}_{\substack{\text{II} \\ \sum_{i=1}^m q_i}} - \underbrace{\sum_{i=1}^m 1 \times p_i}_{\substack{\text{I} \\ \text{because } \{p_1, \dots, p_m\} \text{ dist.}}} \\ \text{I} \quad \text{because } \{q_1, \dots, q_m\} \text{ dist.} \\ 1 - 1 = 0$$
$$\Rightarrow \sum_{i=1}^m p_i \ln \left(\frac{q_i}{p_i} \right) \leq 0 \quad (\text{Gibbs})$$

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:*

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

Sh

2

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.

$$\text{We see: } \sum_{i=1}^m p_i \log_2 \left(\frac{q_i}{p_i} \right) \leq 0$$

$$\text{we set: } q_i = \frac{1}{m}, \quad \forall i \in \{1, \dots, m\}$$

$$\begin{aligned} \sum_{i=1}^m p_i \log_2 \left(\frac{1}{m p_i} \right) \leq 0 &\Leftrightarrow - \sum_{i=1}^m p_i \log p_i \leq \underbrace{\sum_{i=1}^m p_i \log_2 m}_{\text{H}(S)} \\ &\Leftrightarrow H(S) \leq \log_2 m \times \underbrace{\sum_{i=1}^m p_i}_1 \end{aligned}$$

$$H(S) \leq \log_2 m$$

Equality:

Gibbs inequality becomes an equality

$$\text{for } x = 1$$

during the proof, we set $x = \frac{q_i}{p_i}$

$$\text{and next } q_i = \frac{1}{m}$$

Combining everything, $H(S) = \log_2 m$

$$\text{if } p_i = \frac{1}{m}, \quad \text{for all } i \in \{1, \dots, m\}$$

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 λ 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.

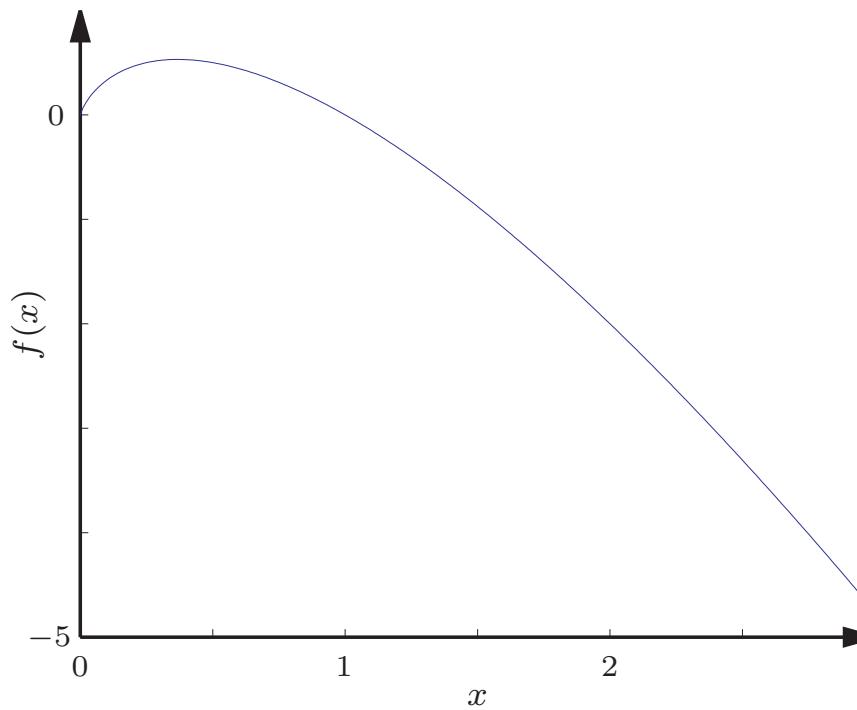
$$\begin{aligned}
 \text{RalleP: } H(S) &= \sum_{i=1}^n (-p_i \log p_i) \\
 &= \sum_{i=1}^n f(p_i)
 \end{aligned}$$



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:

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

unit: Sh / pair of (X, Y)

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

Example. Case of two independent random variables

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

Z : source

$$H(Z) = - \sum_{k=1}^K P(Z=z_k) \log_2 P(Z=z_k)$$

X and Y : 2 sources

$$Z = (X, Y)$$

$$H(X, Y) = - \sum_{i=1}^m \sum_{j=1}^m P(X=x_i, Y=y_j) \log_2 P(X=x_i, Y=y_j)$$

$P(X = x_i)$

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

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

If X and Y indep., then $P(X = x_i|Y = y_j) = P(X = x_i)$
 $\Rightarrow H(X|Y = y_j) = H(X)$

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

$$= \sum_j P(Y = y_j) H(X) \quad \text{because } X, Y \text{ indep.}$$

$$= H(X) \underbrace{\sum_j P(Y = y_j)}_1$$

$$= H(X)$$

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:*

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

$$\begin{aligned}
H(x, y) &= - \sum_i \sum_j P(x=x_i, y=y_j) \log P(x=x_i, y=y_j) \\
&= - \sum_i \sum_j P(x=x_i | y=y_j) P(y=y_j) \left[\log P(x=x_i | y=y_j) + \log P(y=y_j) \right] \\
&= \sum_j \left[- \sum_i P(x=x_i | y=y_j) \log P(x=x_i | y=y_j) \right] P(y=y_j) \\
&\quad - \sum_j \left[\underbrace{\sum_i P(x=x_i | y=y_j)}_1 \right] P(y=y_j) \log P(y=y_j) \\
&= \sum_j H(x | y=y_j) P(y=y_j) - \sum_j P(y=y_j) \log P(y=y_j) \\
&= H(x|y) + H(y)
\end{aligned}$$

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

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

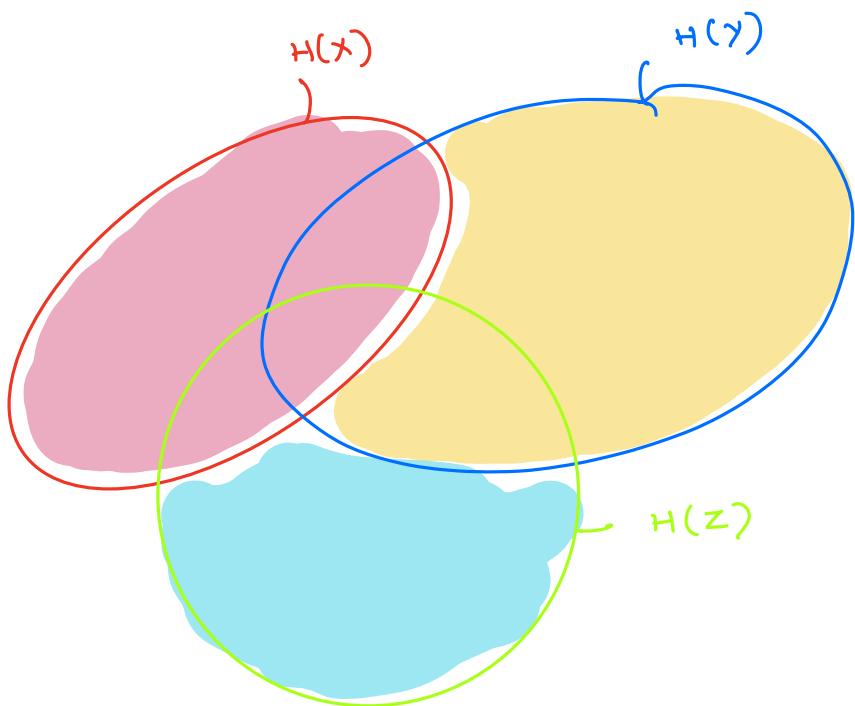
$$\begin{aligned}
\text{Illustration : } H(x, y, z) &= H(x) + H(y, z | x) \\
&= H(x) + H(y|x) + H(z|x, y)
\end{aligned}$$

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

$$\text{and : } P(y, z | x) = P(z | x, y) P(y | x)$$

$$\Rightarrow P(x, y, z) = P(x) P(y | x) P(z | x, y)$$

$$\begin{aligned}
\log \sum H(x, y, z) &= H(x) + H(y | x) + H(z | x, y)
\end{aligned}$$



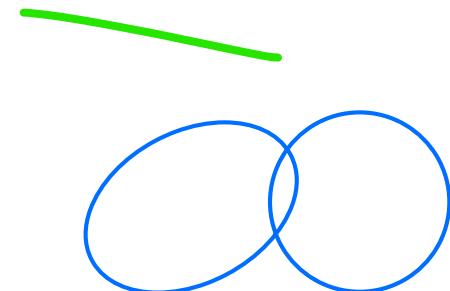
PAIR OF RANDOM VARIABLES

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:

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

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



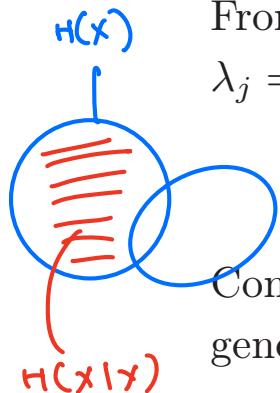
$$H(X, Y) = H(X) + H(Y|X)$$

We know that these three terms are positive

$$H(X, Y) \geq H(X)$$

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:

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

- $\textcircled{1}$: Entropy is positive because $-\sum_i p_i \log p_i$ with $p_i \in [0, 1]$
- $\textcircled{2}$: $H(X|Y) \leq H(X)$ because conditioning decreases entropy (concavity)
- $\textcircled{3}$: $H(X, Y) = H(X) + H(Y|X)$ with $H(Y|X) \geq 0$
- $\textcircled{4}$: $H(X, Y) = H(X) + H(Y|X)$ and $H(Y|X) < H(Y)$
 $\leq H(X) + H(Y)$
- $\textcircled{5}$: By $\textcircled{3}$: $H(X) \leq H(X, Y)$
 $H(Y) \leq H(X, Y)$ $\Rightarrow H(X) + H(Y) \leq 2H(X, Y)$

Inequalities can become equalities :

x, y indep.

or $x = f(y)$ with f a one-to-one application

depending on the inequality.

PAIR OF RANDOM VARIABLES

Mutual information

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

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

or, equivalently,

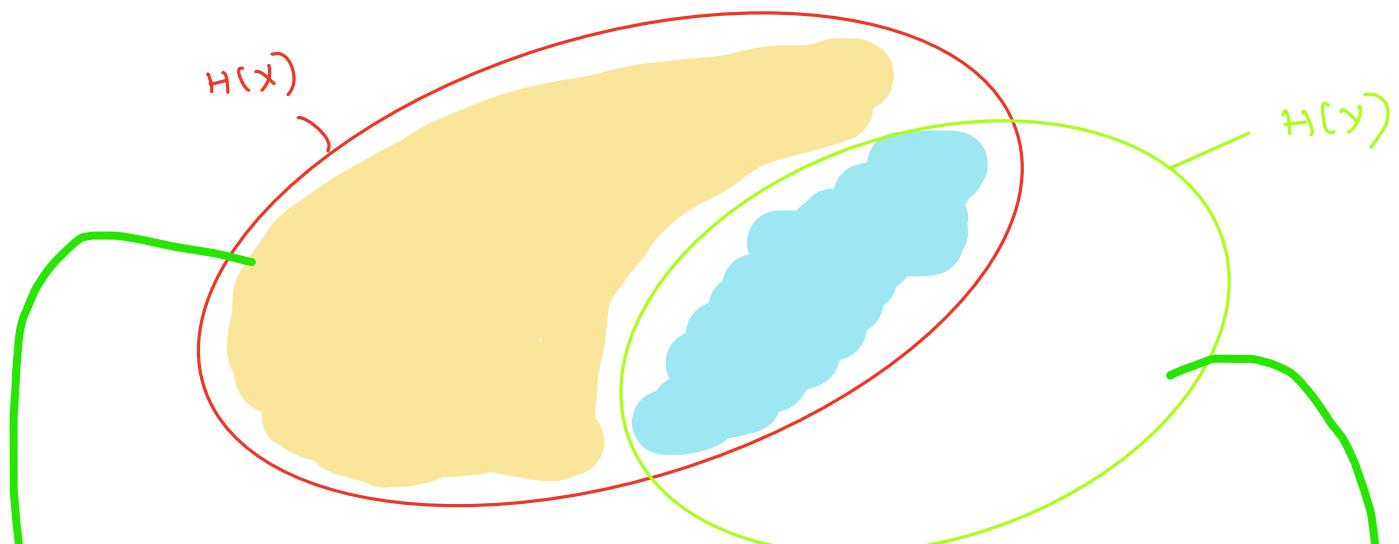
$$I(X, Y) \triangleq \sum_{i=1}^n \sum_{j=1}^m P(X = x_i, Y = y_j) \log \frac{P(X = x_i, Y = y_j)}{P(X = x_i) P(Y = y_j)}.$$

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

Mutual information

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



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

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:*

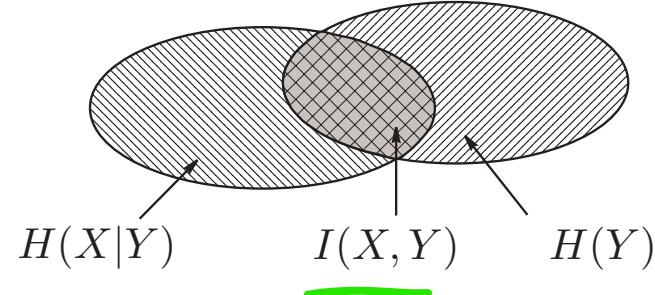
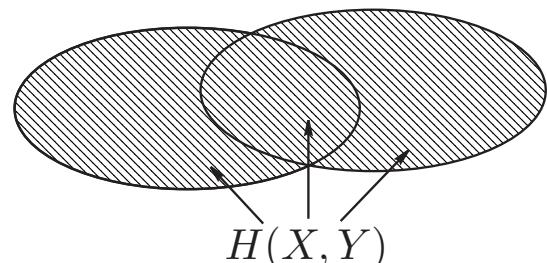
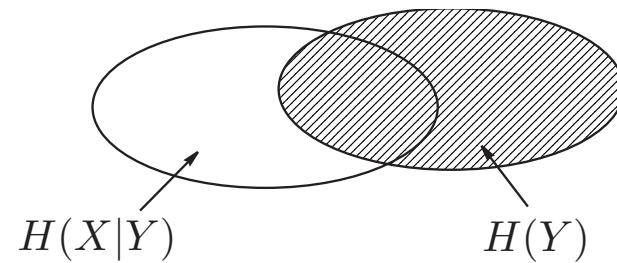
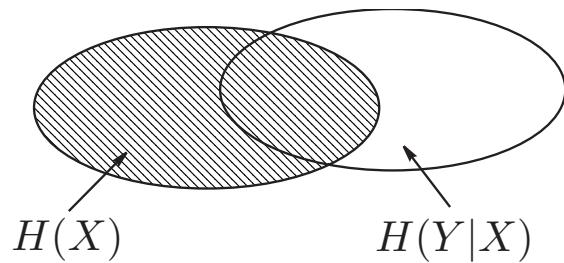
$$d(P_1, P_2) = \sum_{x \in X(\Omega)} P_1(X = x) \log \frac{P_1(X = x)}{P_2(X = x)}.$$

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.



Consider a single event A of probability $P(A)$

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

Considering an addition event B of probability $P(B)$

$$h(A, B) = -\log_2 P(A, B) \text{ (joint info)}$$

$$h(A|B) = -\log_2 P(A|B) \text{ (cond. info)}$$

$$i(A, B) = h(A) - h(A|B) \text{ (shared)}$$

$$S = \{s_1, \dots, s_m\}$$

Richness of the source S : entropy $H(S)$

$$H(S) = \sum_{i=1}^m P(S=s_i) h(S=s_i)$$

$$= - \sum_{i=1}^m P(S=s_i) \log_2 P(S=s_i) \text{ Sh / state of } S$$

$$\max_{P(S)} H(S) = \underline{\log_2 m} \text{ Sh / state of } S$$

reached when $P(S=s_i) = \frac{1}{m} \forall i=1, \dots, m$

$$H(X, Y) = - \sum_{i=1}^m \sum_{j=1}^m P(X=x_i, Y=y_j) \log_2 P(X=x_i, Y=y_j) \text{ Sh / pair of states of } X \text{ and } Y$$

$$H(X|Y=y_{j_0}) = - \sum_{i=1}^m P(X=x_i | Y=y_{j_0}) \log_2 P(X=x_i | Y=y_{j_0}) \text{ Sh / state of } X$$

$$H(X|Y) = \sum_{j=1}^m P(Y=y_j) H(X|Y=y_j) \text{ Sh / state of } X$$

$$\begin{aligned}
 I(x, y) &= H(x) - H(x|y) \\
 &= H(y) - H(y|x) \\
 &= H(x, y) - H(x|y) - H(y|x)
 \end{aligned}$$

Properties:

$$\begin{aligned}
 H(x, y) &= H(x|y) + H(y) \\
 &= H(y|x) + H(x)
 \end{aligned}$$

$$0 \leq H(x|y) \leq H(x) \leq H(x, y) \leq H(x) + H(y)$$

If x and y indep., then :

$$\begin{aligned}
 H(x|y) &= H(x) \\
 H(x, y) &= H(x) + H(y)
 \end{aligned}$$



Quantitative Measure of Information

Part I

Exercise 1

évent B

H(B)

One person says: "Today is my birthday". Calculate the amount of self-information conveyed by this statement. Calculate the average amount of information conveyed by this source over one year.

Exercise 2

The 64 squares of a chessboard are assumed to be equiprobable. Determine the average amount of information contained in a communication indicating the position of a given chess piece. Propose a dichotomous strategy, based on questions of the form "Is the chess piece on that part of the chessboard?", that would allow to guess the position of this chess piece in a minimum average number of questions. Compare this average number of questions to the entropy calculated at the beginning of the exercise.

Exercise 3

A perfectly balanced coin is tossed until the first head appears. Calculate the entropy $H(X)$ in Shannon, where the random variable X denotes the number of flips required to get the first head. Propose a dichotomous strategy, based on questions with binary response of the form "Is X smaller or greater than (...)", making it possible to guess the value of X in a minimum average number of questions. Compare this number of questions to $H(X)$.

In order to resolve this exercise, the following equality can be used $\sum_{n=1}^{\infty} n a^n = \frac{a}{(1-a)^2}$.

Exercise 5

Consider a tank that consists of two compartments of identical volumes. Compartment I is filled with two inert gases with respective proportions $(\frac{2}{5}, \frac{3}{5})$. The same gases fill compartment II with respective proportions $(\frac{1}{3}, \frac{2}{3})$. Assuming the pressure and temperature in both compartments are the same, calculate the tank entropy before and after the two compartments communicate. Interpret the result.

Exercise 6

A source emits symbols 0 and 1 with probabilities $P(0) = \frac{1}{4}$ and $P(1) = \frac{3}{4}$. These symbols are transmitted to a receiver through an imperfect symmetric channel illustrated by Figure 1, with $p_0 = 10^{-1}$. Denoting by X and Y the transmitted and received symbols, calculate the following quantities: $H(X)$, $H(Y)$, $H(X,Y)$, $H(Y|X)$, $H(X|Y)$ and $I(X,Y)$.

Problem 1

Let $\{\mathcal{E}_k\}_{k=1}^n$ be a partition of \mathcal{E} . We denote by N and N_k the numbers of elements in sets \mathcal{E} and \mathcal{E}_k , respectively. Assume that the elements of \mathcal{E} are equiprobable. We set $p_k = N_k/N$.

1. Determine the self-information of any element of \mathcal{E}_k . Calculate the average amount of information needed to determine any element in \mathcal{E}_k .
2. Calculate the average amount of information needed to characterize any element of \mathcal{E} . By noticing that we can split the identification procedure of an element of \mathcal{E} in 2 steps, (a) identification of the set \mathcal{E}_k , and then (b) identification of the element in \mathcal{E}_k , estimate the average amount of information needed to identify \mathcal{E}_k .

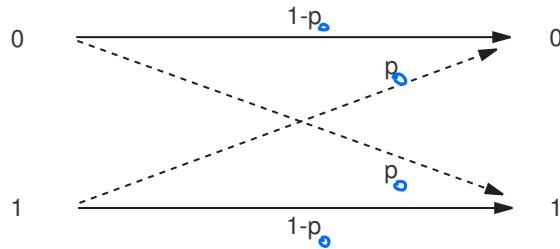


Figure 1: Imperfect channel.

Problem 2

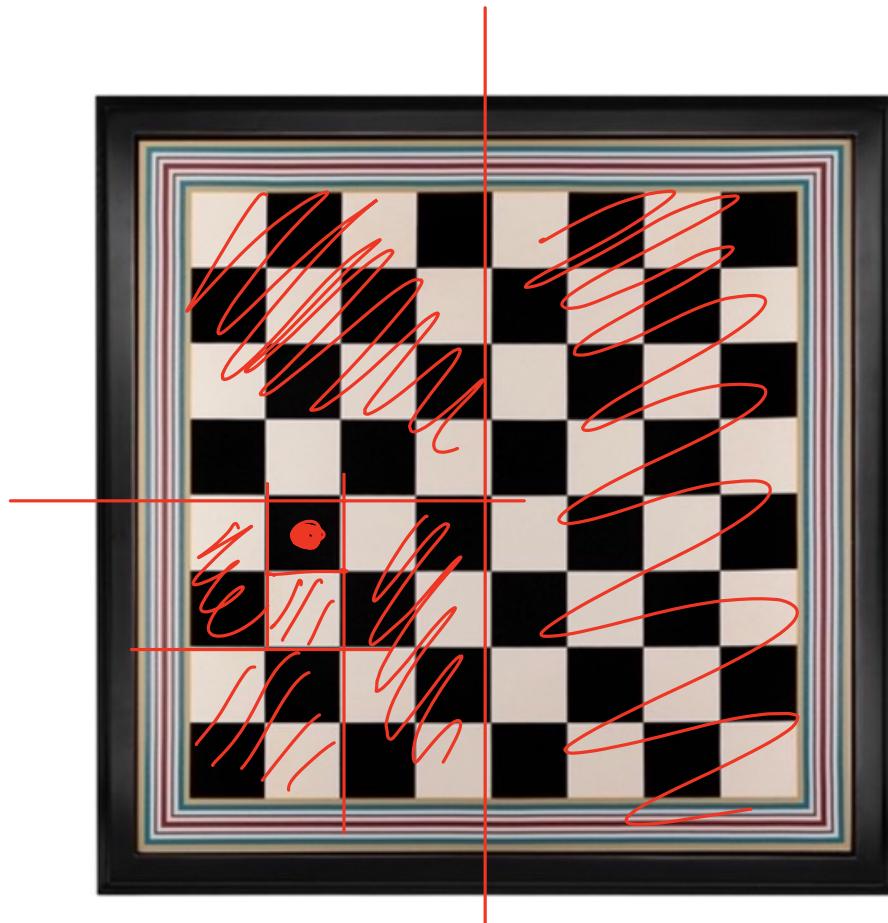
Consider a twin-pan balance and 9 coins. We know that one of these coins is fake. The problem is to find the fake coin given that it only differs from the other 8 coins by its weight.

1. Determine the number of possible cases, considering that the fake coin may be heavier or lighter than the others. Calculate the average amount of information necessary to identify the fake coin.
2. To identify the fake coin, the weights of two sets of n coins each are compared using the twin-pan balance. Enumerate the possible outcomes of each weighting operation. Assuming these outcomes are equiprobable, determine in that case the amount of information provided by every weighing operation. Determine the average number of weighting operations to plan.
3. One wants to determine n in order to maximize the amount of information provided by each weighting operation. Let P_ℓ , resp. P_r , be the probability that the set of coins in the left pan, resp. right pan, is heavier. Let P_e be the probability that an equilibrium is achieved. Calculate P_ℓ , P_r and P_e .
4. Calculate n to maximize the entropy of each weighting operation.
5. Calculate the minimum average number of weighting operations required to identify the fake coin.
6. Propose a strategy to identify the fake coin.

exercise 2:

21

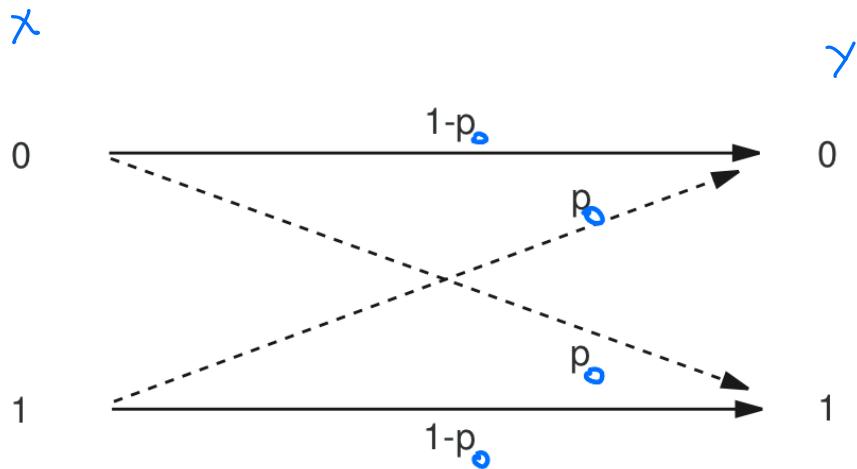
6 generations.



$$H(S) = \log_2 64 = \log_2 2^6 = 6 \text{ Sh / state}$$

= number of questions.

Exercise 6 :



$$P(X=0) = \frac{1}{4} \quad P(X=1) = \frac{3}{4} \quad p_0 = 10^{-1}$$

$$p_0 = P(Y=0 | X=1) = P(Y=1 | X=0)$$
$$1-p_0 = P(Y=1 | X=1) = P(Y=0 | X=0)$$

$H(X)$, $H(Y)$, $H(X, Y)$, $H(Y|X)$, $H(X|Y)$ and $I(X, Y)$.

Quantitative Measure of Information

Part II

Exercise

Let X be a discrete random variable that can take m possible values, and Y a discrete random variable uniformly distributed that can take n possible values. Throughout the exercise, no assumption will be made on the distribution of X .

1. Calculate the maximum entropy that can be reached by X , and specify the case in which this value would be obtained.
2. Calculate the entropy of Y .
3. In the case of $n = m$, rank in ascending order the following quantities: $H(Y)$, 0 , $H(X)$, $H(X;Y)$, $H(X) + H(Y)$, $H(X|Y)$. Justify each step of your ranking.
4. Explain cases of equality in the ranking proposed in 3., i.e., give for each inequality the cases in which it becomes an equality.
5. We now consider the case where the joint distribution $P(X = x_i; Y = y_j)$ is given by the table below:

$P(X;Y)$	y_1	y_2	y_3	y_4
x_1	1/24	1/12	1/6	1/24
x_2	1/6	1/8	1/24	1/6
x_3	1/24	1/24	1/24	1/24

Check that Y is uniformly distributed. Calculate $I(X;Y)$.

Problem

For a given region, the forecasts of a meteorologist are divided according to their relative frequencies given by the table below. The columns correspond to the actual weather, which is represented by the random variable T , which takes values 0 or 1 depending on whether the weather is rainy or sunny, respectively. The rows correspond to the meteorologist's forecast, identified by the random variable M , also with values in $\{0,1\}$ depending on whether he had planned a rainy weather (0) or a sunny weather (1).

$P(M = i, T = j)$	sunny weather ($T = 1$)	rainy weather ($T = 0$)
sunny weather ($M = 1$)	5/8	1/16
rainy weather ($M = 0$)	3/16	1/8

1. Calculate the probabilities $P(M = i)$ and $P(T = j)$, with $i, j \in \{0, 1\}$.
2. Show that the meteorologist is wrong once in 4 times.
3. One student says that by always forecasting sunny weather, he makes fewer mistakes than the meteorologist does. Check this assertion.
4. Let E be the random variable representing the student's prediction. As for T and M , random variable E takes values in $\{0,1\}$. Calculate $I(E;T)$.
5. Calculate $I(M;T)$.
6. Comparing $I(M;T)$ to $I(E;T)$, what Information Theory shows on the meteorologist's forecast and that of the student?

7. The student claims to have found a revolutionary method of predicting the weather. Its revised performance are provided in the table above. As before, the rows correspond to the forecast, and the columns to the actual weather.

$P(E = i, T = j)$	sunny weather ($T = 1$)	rainy weather ($T = 0$)
sunny weather ($E = 1$)	403/512	93/512
rainy weather ($E = 0$)	13/512	3/512

Calculate the probabilities $P(E = 0)$ and $P(E = 1)$.

8. Compare $P(E = i, T = j)$ and $P(E = i)P(T = j)$, for all $i, j \in \{0, 1\}$. Conclude.
9. We wish to store T by using a binary coding. Using Shannon's first theorem, give the minimum average memory space required to store T , in bits per realization of T .
11. Redo the previous calculation in the case of M . Calculate the minimum memory space required to store M and T separately, in bits per realization of (M, T) ?
12. Calculate the minimum memory space required to store M and T jointly, in bits per realization of (M, T) ?
13. Interpret the difference between results of the 2 previous questions.
14. Propose Huffman coding to jointly encode M and T .
15. Calculate the average length of words \bar{n} of the binary code found in the previous question. What double inequality is satisfied by \bar{n} ?

$P(X; Y)$	y_1	y_2	y_3	y_4	$P(x = x_i)$
x_1	$1/24$	$1/12$	$1/6$	$1/24$	$1/3$
x_2	$1/6$	$1/8$	$1/24$	$1/6$	$1/2$
x_3	$1/24$	$1/24$	$1/24$	$1/24$	$1/6$
$P(y = y_i)$	$1/4$	$1/4$	$1/4$	$1/4$	

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

$$= 2 - H(Y|X)$$

$$H(Y|X) = \sum_{i=1}^3 H(Y|X=x_i) P(X=x_i)$$

To calculate $H(Y|X=x_i)$

we need : $P(Y|X=x_i)$

$P(Y X)$	y_1	y_2	y_3	y_4		
x_1	$1/8$	$1/4$	$1/2$	$1/8$	$\%$	$\frac{1}{3}$
x_2	$1/3$	$1/4$	$1/12$	$1/3$	$\%$	$\frac{1}{2}$
x_3	$1/4$	$1/4$	$1/4$	$1/4$	\therefore	$\frac{1}{6}$

$$P(Y=y_j | X=x_i) = \frac{P(X=x_i, Y=y_j)}{P(X=x_i)}$$

$$H(Y|X=x_1) = H\left(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\right)$$

$$H(Y|X=x_2) = H\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{4}, \frac{1}{12}\right)$$

$$H(Y|X=x_3) = H\left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right) = 2 \text{ SH (a)}$$

$$I(X, Y) = 2 - \left[\frac{1}{3} H\left(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\right) + \frac{1}{2} H\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{4}, \frac{1}{12}\right) \right]$$

$$+ \frac{1}{6} H\left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right) \Big]$$
$$= 0, 15 \text{ Sh / pair}(x, y)$$