

# Théorie de l'information

O. Kohel

11/01/2021



## L'entropie

Théorème Soit  $X$  un espace de proba finie. Alors  $H(X) \leq \log_2 |X|$  avec égalité ssi  $X$  est uniforme.

Rappel. On dit que  $X$  est uniforme si  $p(x) = 1/|X|$  pour tout  $x \in X$ .

Preuve. Par conséquence de l'inégalité de Gibbs.

Lemme (Gibbs). Soit  $(X, p)$  un espace de probabilité discrète, et soit  $q: X \rightarrow [0, 1]$  une autre probabilité.

Alors

$$\sqrt{\sum_{x \in X} p(x) \log(p(x))} \leq \sum p(x) \log(q(x)),$$

2/ avec égalité ssi  $p = q$ .

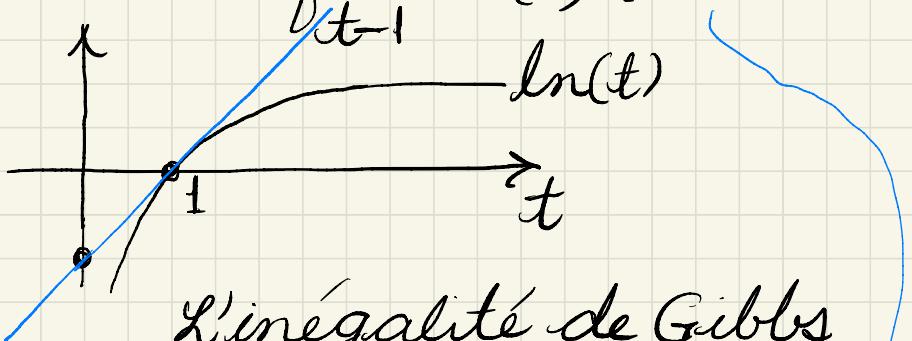
Remarque (Multipliant par  $1/\log(2)$ ). L'inégalité de Gibbs est

$$H(X) \leq E(\log_2(q(x)))$$

## Preuve

(2)

Il suffit de prendre  $\log = \ln$ ,  
et observer que  $\ln(t) \leq t-1$ :



L'inégalité de Gibbs  
est équivalente à

$$\sum_{x \in X} p(x) \ln\left(\frac{q(x)}{p(x)}\right) \leq 0. \quad \textcircled{*}$$

Mais on a

$$\begin{aligned} \sum_{x \in X} p(x) \ln\left(\frac{q(x)}{p(x)}\right) &\leq \sum_{x \in X} p(x) \left( \frac{q(x)}{p(x)} - 1 \right) \\ &= \sum_{x \in X} (q(x) - p(x)) = 1 - 1 = 0. \end{aligned}$$

On a démontré  $\textcircled{*}$ , équivalente  
à l'inégalité de Gibbs.

On observe que  $\ln(t) = t-1$ ,

ssi  $t=1$ . Par conséquent,  
l'inégalité  $\circledast$  est stricte s'il  
existe  $x$  tel que  $p(x) \neq q(x)$ .

Cela établit la propriété 2).  $\square$

Preuve du théorème. Il suffit  
de mettre  $q(x) = 1/|X|$ . Alors

$$H(X) \leq E(\log_2 |X|) = \log_2 |X|. \quad \square$$

## Théorie de Shannon

### Premier théorème de Shannon

Définition. Soit  $C: X \rightarrow A^*$  un  
codage (où  $A$  est un alphabet  
fini). On dit que i)  $C$  est sans  
pertes si  $C$  est injectif, et  
ii)  $C$  est uniquement décodable  
si l'extension  $C^*: X^* \rightarrow A^*$ ,  
définie par  $C^*(x_1 \dots x_n) = C(x_1) \dots C(x_n)$ ,  
est sans pertes.

Exemple.  $X = \{A, B, C\}$ ,  $A = \{0, 1\}$ .

Le codage  $C: X \rightarrow A^*$  donné par:  $C(A) = 0$ ,  $C(B) = 1$ ,  $C(C) = 01$ . Alors  $C$  est sans pertes, mais pas uniquement décodable, car  $C^*(AB) = 01 = C^*(C)$ .

Par contre,  $C': X \rightarrow A^*$  donné par:  $C'(A) = 00$ ,  $C'(B) = 11$ ,  $C'(C) = 01$ , est uniquement décodable.

Théorème de Shannon (codage source)

Soit  $C: X \rightarrow A^*$  un codage uniquement décodable, et  $l_C: X \rightarrow \mathbb{N}$  la fonction longueur (sur  $A^*$ ) composée avec  $C(l_C = l \circ C)$

Alors

$$\text{i)} \quad \frac{H(X)}{\log_2 |A|} \leq E(l_C), \text{ et}$$

(5)

ii) il existe un codage  $C$  pour lequel  $E(l_C) \leq \frac{H(X)}{\log_2 |A|} + 1$ .

Définition. Un codage est instantané ou préfixe si aucun mot de code est un préfixe d'un autre.

Par exemple: Le code ( $= C(X)$ )

$$\{C(A), C(B), C(C)\} = \{0, 1, 01\}$$

n'est pas préfixe. Par contre

$$\{C(A), C(B), C(C)\} = \{00, 11, 01\}$$

est préfixe.

Lemme Un codage préfixe est uniquement décodable.

Preuve Clair (exercice).

Soit  $g_1 g_2 \dots g_n = C(x_1 \dots x_n) -$

on peut identifier les mots ⑥ de code  $G_1$ , puis  $G_2$ , etc en lisant de gauche à droite.  $\square$

Remarque. La définition de codage instantané / préfixe est asymétrique : On peut par symétrie définir un codage suffixe (aucun mot est un suffixe d'un autre). La propriété d'être u.d. reste valable. Mais

On peut comparer :

$x$	$G_1(x)$	$G_2(x)$	...
A	0	0	
B	11	11	
C	01	10	

Alors  $G_1$  est suffixe et  $G_2$  est préfixe. Les deux sont u.d.

(7)

Or si

$$G^*(x_1 \dots x_n) = \overbrace{00111 \dots 1}^{\substack{\text{A C B ...} \\ \text{A A B B ...}}}$$

il faut lire jusqu'à la fin  
afin de décider si le début  
a commencé par ACB... ou  
par AAB. La différence sera  
déterminé par la longueur  
de  $G^*(x_1 \dots x_n)$ .

Par contre

$$G_2^*(x_1 \dots x_n) = \overbrace{00111 \dots 1}^{A A B B}$$

étant préfixe, on peut décoder  
(de gauche à droite) sans  
attendre la fin de la trans-  
mission du mot de code.

Remarque. Un codage

$C: X \rightarrow A^n \subseteq A^*$  (de longueur  
fixée) est préfixe, ssi il est  
sans perte.

## Codage de Huffman

(8)

permet de réaliser un codage préfixe qui satisfait la condition ii) du théorème de Shannon.

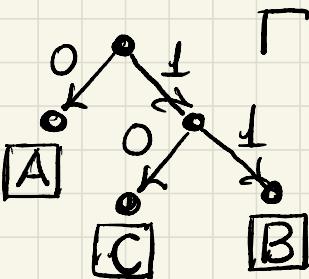
C'est un algorithme de compression. Remarque: la propriété d'être préfixe (et également u.d.) est une propriété du code C(X) (en supposant que C soit sans pertes) et pas de l'espace de probabilité.

Arbre de codage fini  
Un arbre de codage est un arbre (graphe, sans cycle avec une racine = sommet distingué) dérivé tel que

(9)

chaque feuille (sommet terminal) port une étiquette  $x \in X$   
 et chaque arc d'un sommet porte une étiquette distincte dans  $A$ .

Ex. Arbre binaire ( $A = \{0, 1\}$ )



Si chaque message  $x \in X$  apparaît une et une seule fois comme étiquette d'une feuille, on dit que  $T$  est un arbre de codage pour  $X$ .

Lemme Il y a un bijection 10 entre les codages préfixes pour  $X$  et les arbres de codage pour  $X$ .

Construction d'un arbre de L'algorithme "glouton". Huffman  
Soit  $X$  un ensemble et  $A$  un alphabet de cardinal  $q$ .

i) Si  $|X| \bmod (q-1) \neq 1$ , on ajoute des éléments à  $X$  avec  $p(x)=0$ .

ii) Mettre  $X = \{x_0, x_1, \dots, x_{n-1}\}$  en ordre avec  $p(x_0) \leq p(x_1) \leq \dots \leq p(x_{n-1})$ .

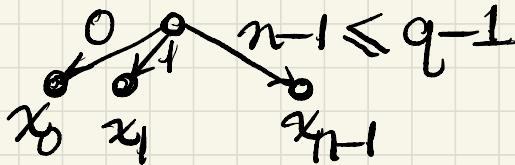
iii) Si  $n=|X|=1$ , retourne: 

Remarque: On va écrire   $x_0$

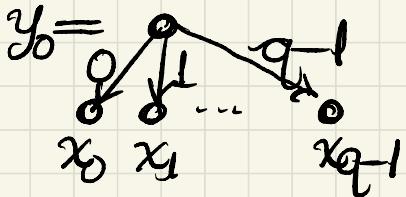
$A = \{0, 1, \dots, q-1\}$  pour les éléments de l'alphabet.

$\sin = 2$ , retourner  
et si ( $q > 2$  et)

$n \leq q$  on met :



Sinon, mettre (construction d'un nouveau espace de proba) :



$y_1 = x_q$ ,  $y_2 = x_{q+1}$ , ...,  $y_{n-q-1} = x_{n-1}$ ,  
ayant des probabilités :

$$p(y_0) = \sum_{0 \leq i < q} p(x_i) \text{ et } p(y_j) = p(x_{j+q-1}) \text{ pour tout } j \geq 1.$$

On retourne l'arbre de Huffman  
pour  $\{y_0, \dots, y_{n-q-1}\}$ .

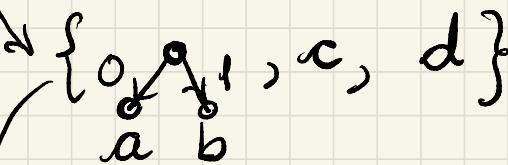
N.B. L'étiquette pour  $y_0$  devient un sous-arbre.

Exemple.  $A = \{0, 1\}$ .

$X = \{a, b, c, d\}$ .

$p(x) : \frac{1}{8}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}$

Un arbre de codage est donc :



$p(x) : \frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{2}$

Remarque

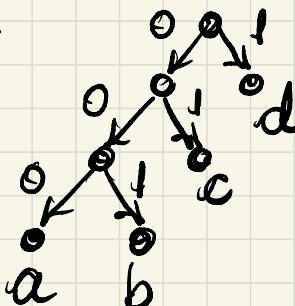
Lors de la construction

on peut échanger  
 $\{a, b\}$ , ou  
 $\{\{ab\}, c\}$ , ou  
 $\{\{\{ab\}, c\}, d\}$ .



$p(x) : \frac{1}{2} \quad \frac{1}{2}$

L'arbre de codage de Huffman :



Le codage préfixe  
de Huffman est donc :

a  $\rightarrow 000$   
b  $\rightarrow 001$

c  $\rightarrow 01$   
d  $\rightarrow 1$

Conclusion. Un codage de Huffman n'est pas unique, mais on peut démontrer que l'espérance mathématique de longueur  $E(l)$  est une invariant de tout codage de Huffman.

(43)

Prochain résultat :

Le codage de Huffman satisfait

$$E(l) \leq \frac{H(X)}{\log_2(q)} + 1.$$