

Initiation au traitement du signal et applications

Séance 3: compression numérique sans perte

Frédéric Sur
École des Mines de Nancy

www.loria.fr/~sur/enseignement/signal/

Position du
problème

Codage et décodage
Codes préfixes

Théorie statistique
de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Séance 3

- 1 Position du problème
 - Codage et décodage
 - Codes préfixes
- 2 Théorie statistique de l'information
 - Notion d'entropie
 - Théorème de Shannon
 - Code de Huffman
- 3 Conclusion

Position du
problème

Codage et décodage
Codes préfixes

Théorie statistique
de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Séance 3

- 1 Position du problème
 - Codage et décodage
 - Codes préfixes
- 2 Théorie statistique de l'information
 - Notion d'entropie
 - Théorème de Shannon
 - Code de Huffman
- 3 Conclusion

Position du
problème

Codage et décodage
Codes préfixes

Théorie statistique
de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Position du problème

Signaux considérés :

signaux discrets = suites de symboles dans un alphabet de
taille $K : \{x_1, x_2, \dots, x_K\}$.

Exemple 1 : suite de valeurs sur 8 bits.

Exemple 2 : des textes écrits en français (alphabet
classique).

Question du jour : comment représenter (= coder)
informatiquement chaque symbole ?

Position du
problème

Codage et décodage
Codes préfixes

Théorie statistique
de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Codage à taille fixée

Première idée : même nombre de bits pour chaque symbole.

Exemple : le code ASCII (1961), codage sur 7 bits.

Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	
00	space	32	01	tab	9	02	SOH	1	03	STX	2
04	ETX	3	05	END	4	06	SO	6	07	AH	7
08	HT	8	09	LF	10	0A	VT	11	0B	FF	12
0C	CR	13	0D	SH	14	0E	UH	15	0F	ESC	16
10	space	32	11	tab	9	12	SOH	1	13	STX	2
14	ETX	3	15	END	4	16	SO	6	17	AH	7
18	HT	8	19	LF	10	1A	VT	11	1B	FF	12
1C	CR	13	1D	SH	14	1E	UH	15	1F	ESC	16
20	space	32	21	tab	9	22	SOH	1	23	STX	2
24	ETX	3	25	END	4	26	SO	6	27	AH	7
28	HT	8	29	LF	10	2A	VT	11	2B	FF	12
2C	CR	13	2D	SH	14	2E	UH	15	2F	ESC	16
30	space	32	31	tab	9	32	SOH	1	33	STX	2
34	ETX	3	35	END	4	36	SO	6	37	AH	7
38	HT	8	39	LF	10	3A	VT	11	3B	FF	12
3C	CR	13	3D	SH	14	3E	UH	15	3F	ESC	16
40	space	32	41	tab	9	42	SOH	1	43	STX	2
44	ETX	3	45	END	4	46	SO	6	47	AH	7
48	HT	8	49	LF	10	4A	VT	11	4B	FF	12
4C	CR	13	4D	SH	14	4E	UH	15	4F	ESC	16
50	space	32	51	tab	9	52	SOH	1	53	STX	2
54	ETX	3	55	END	4	56	SO	6	57	AH	7
58	HT	8	59	LF	10	5A	VT	11	5B	FF	12
5C	CR	13	5D	SH	14	5E	UH	15	5F	ESC	16
60	space	32	61	tab	9	62	SOH	1	63	STX	2
64	ETX	3	65	END	4	66	SO	6	67	AH	7
68	HT	8	69	LF	10	6A	VT	11	6B	FF	12
6C	CR	13	6D	SH	14	6E	UH	15	6F	ESC	16
70	space	32	71	tab	9	72	SOH	1	73	STX	2
74	ETX	3	75	END	4	76	SO	6	77	AH	7
78	HT	8	79	LF	10	7A	VT	11	7B	FF	12
7C	CR	13	7D	SH	14	7E	UH	15	7F	ESC	16

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage

Codes préfixes

Théorie statistique de l'information

Notion d'entropie

Théorème de Shannon

Code de Huffman

Conclusion

6/23

Codage à taille variable

Idee plus évoluée : codes plus courts pour les symboles plus fréquents.

But : diminuer la longueur moyenne des signaux codés.

Exemple : le code Morse (Samuel Morse, 1835)

A	...-.	M	---	Y	...--	6
B	...-.-	N	-.-	Z	...--..	7	...-...-
C	...-.-.	O	---	A	...-.-	8	...--..-
D	...-.-.	P	...--.	O	...--..	9	...--..-
E	...-	Q	...--.	U	...--..		...--..-
F	...-.-.	R	...--.	Ch	...--..		...--..-
G	...--.	S	...--.	0	...--..	?	...--..-
H	...-.-.	T	...-	!	...--..	!	...--..-
I	...--.	U	...--.	2	...--..	2	...--..-
J	...--..	V	...--.	3	...--..	3	...--..-
K	...-.-.	W	...--.	4	...--..	4	...--..-
L	...-.-.	X	...--.	5	...--..	5	...--..-

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage

Codes préfixes

Théorie statistique de l'information

Notion d'entropie

Théorème de Shannon

Code de Huffman

Conclusion

6/23

Le décodage

Hypothèse : le codage est fait en binaire.

Problème : comment décoder ?

Codage de longueur fixée \rightarrow immédiat (bloc à bloc).

Mais pour le codage de longueur variable ?

Exemple : codage d'un alphabet à quatre symboles :

a	b	c	d
0	10	110	101

À quel mot correspond 1010 ?

Solution possible : séparation par un caractère spécial (cf Morse)

\rightarrow pas intéressant si on veut *compresser*...

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage

Codes préfixes

Théorie statistique de l'information

Notion d'entropie

Théorème de Shannon

Code de Huffman

Conclusion

7/23

Notion de code préfixe

Solution : imposer que le codage d'un symbole ne soit le préfixe d'aucun autre.

Dans l'exemple précédent :

a	b	c	d
0	10	110	101

b préfixe de d \rightarrow non-unicité du décodage.

Exemple de code préfixe :

a	b	c	d
0	10	110	111

À quoi correspond 101100 ?

\rightarrow décodage instantané.

Remarque : le code Morse n'est pas un code préfixe, mais est plus *robuste aux erreurs de transmission*.

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage

Codes préfixes

Théorie statistique de l'information

Notion d'entropie

Théorème de Shannon

Code de Huffman

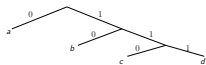
Conclusion

8/23

Proposition

Code préfixe \iff code des feuilles d'un arbre binaire.

a	b	c	d
0	10	110	111



Décodage de 1001001101110 ?

Donnée : alphabet de K symboles x_k , probabilités p_k .

But : construire un code préfixe qui minimise la longueur moyenne de codage d'un symbole

$$R(X) = \sum_{k=1}^K l_k p_k \quad (= \mathbb{E}(L))$$

avec l_k longueur du mot binaire w_k codant le symbole x_k .

Intérêt : la longueur moyenne du codage d'un texte de N symboles est $N.R(X)$.

\rightarrow équivaut à construire un arbre binaire à k feuilles avec $R(X)$ minimal.

- Position du problème
 - Codage et décodage
 - Codes préfixes
- Théorie statistique de l'information
 - Notion d'entropie
 - Théorème de Shannon
 - Code de Huffman
- Conclusion

Claude Shannon (1916 - 2001)

A Mathematical Theory of Communication, 1948

Définition - Entropie \mathcal{H} d'une source X

Soit X v.a. discrète, K valeurs possibles, de loi $\Pr(X = x_k) = p_k$.

$$\mathcal{H}(X) = - \sum_{k=1}^K p_k \log_2 p_k$$

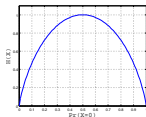
Exemple : $p_1 = 0.2, p_2 = 0.3, p_3 = 0.5$, alors

$$\mathcal{H}(X) = -0.2 \log(0.2) - 0.3 \log(0.3) - 0.5 \log(0.5)$$

$$\mathcal{H}(X) = 1.485 \text{ bit}$$

Exemple : jeu de pile ou face.

pile : $\Pr(X = 1) = p$, face : $\Pr(X = 0) = 1 - p$.



Mesure de l'incertitude :

- $p = 0$ ou $p = 1$, $\mathcal{H}(X) = 0$: pas d'incertitude.
- $p = 1/2$, $\mathcal{H}(X) = 1$ bit : incertitude maximale.

$\mathcal{H}(X)$ mesure une "quantité d'information"...

13/23

Initiation au traitement du signal - Séance 3
F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes
Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman
Conclusion

L'entropie vue comme une quantité d'information

Petit jeu : X v.a. discrète, prenant K valeurs.

J'observe $X = x_i$ (secret), que vous devez deviner.

Vous connaissez la loi de X : $\Pr(X = x_j) = p_j$.

- plus p_i est grand, moins je vous donne d'information en vous disant que j'ai observé x_i .
- si la quantité d'information qui vous manque est mesurée par $-\log_2(p_i)$, alors l'information manquante moyenne est : $\mathcal{H}(X)$ (espérance) .

Attention à ne pas surinterpréter...

14/23

Initiation au traitement du signal - Séance 3
F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes
Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman
Conclusion

Propriétés de l'entropie

Entropie : $\mathcal{H}(X) = -\sum_{k=1}^K p_k \log_2 p_k$.

Proposition

$$0 \leq \mathcal{H}(X) \leq \log_2(K).$$

Conséquence de l'inégalité de Jensen (log est concave) :

$$\mathcal{H}(X) = \sum_{k=1}^K p_k \log_2(1/p_k) \leq \log_2\left(\sum_{k=1}^K p_k/p_k\right) = \log_2(K).$$

Proposition

Entropie $\mathcal{H}(X)$ (quantité d'incertitude, ou quantité moyenne d'information manquante)

- maximale (= $\log_2(K)$) pour distribution uniforme
- minimale (= 0) pour distribution chargeant une seule valeur.

15/23

Initiation au traitement du signal - Séance 3
F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes
Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman
Conclusion

Le théorème de Shannon

Théorème - Shannon 1949

Longueur moyenne $R(X) = \sum l_k p_k$ d'un code préfixe vérifie :

$$R(X) \geq \mathcal{H}(X).$$

Il existe un code préfixe tel que :

$$R(X) < \mathcal{H}(X) + 1.$$

Démonstration : cf photocopié (basé sur l'inégalité de Kraft).

Remarque :

$R(X) = \sum_{k=1}^K l_k p_k$ et $\mathcal{H}(X) = -\sum_{k=1}^K \log_2(p_k) p_k$
Pourquoi ne peut-il pas y avoir toujours égalité ?

16/23

Initiation au traitement du signal - Séance 3
F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes
Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman
Conclusion

Code vérifiant le théorème de Shannon

x_k	a	b	c	d	e	f	g	h
p_k	0.05	0.05	0.05	0.1	0.1	0.15	0.2	0.3
$-\log_2(p_k)$	4.3	4.3	4.3	3.3	3.3	2.7	2.3	1.7
l_k	5	5	5	4	4	3	3	2

Remarque : $\sum_k 2^{-l_k} \leq \sum_k 2^{\log_2(p_k)} = \sum_k p_k = 1$

Construction d'un code préfixe avec les longueurs l_k :
→ au tableau (c'est possible grâce à l'inégalité précédente).
 $R(X) = 3.2$, $\mathcal{H}(X) = 2.709$

Comme $l_k = \lceil -\log_2(p_k) \rceil$, ce code vérifie bien :

$$R(X) = \sum_k l_k p_k < \sum_k (1 - \log_2 p_k) p_k = 1 + \mathcal{H}(X).$$

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie

Théorème de Shannon
Code de Huffman

Conclusion

17/23

Le codage de Huffman

Donnée : alphabet de K symboles x_k , probabilités p_k .

But : construire un *code préfixe optimal* (minimisant R).

Algorithme - David A. Huffman 1952

→ Construction d'un arbre binaire.

- Feuilles : les symboles, pondérés par leur probabilité.
- Récursivement : associer les deux nœuds de poids les plus petits pour créer un nœud de poids égal à leur somme.

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie

Théorème de Shannon
Code de Huffman

Conclusion

18/23

Exemple et proposition

a	b	c	d	e	f	g	h
0.05	0.05	0.05	0.1	0.1	0.15	0.2	0.3

Arbre créé par l'algorithme de Huffman ?

Quel codage pour $hdabe$? pour $fgfgh$?

Et avec un code de longueur fixe ?

Comment décoder 010010111100 ? 000010111100 ?

Quelle valeur pour l'entropie $\mathcal{H}(X)$?

(2.709 bits / caract.)

Longueur moyenne $R(X)$ du code ?

(2.75 bits / caract., à comparer à 3 bits / caract. pour code uniforme et à 3.2 bits/caract. pour code inég. Shannon)

Proposition (cf polycopié)

Le code de Huffman minimise $R(X) = \sum l_k p_k$.

Remarque : il vérifie donc $\mathcal{H}(X) \leq R(X) < \mathcal{H}(X) + 1$.

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie

Théorème de Shannon
Code de Huffman

Conclusion

19/23

Un exemple

Fichier 1 : 2588 caractères (2588 octets) commençant par :

COMPRESSION Action de comprimer, résultat de cette action. A. [Action d'un agent physique] Mes pieds gonflés autant par la compression du cuir que par la chaleur (BALZAC, Le Lys dans la vallée, 1836, p. 24).

Fichier 2 : 2588 caractères (2588 octets) commençant par :

```
aezgf dh.,vcblsjgrjt.dwb*ùGREG?DHcnjz  
efdgegàzpor tghdmr zùlrryelùhdg szsdg dhkfpazç-  
("éét634é"yeyrDFHFGZ(" (ujezatrko=À=Ç(  
ZETEKZ Z)À ETK Z 'ÖHKDMKtet,eif('aydkh;b;szùzzj
```

→ Compression par zip :

Fichier 1 : 1567 octets (4.8 bits / caractère)

Fichier 2 : 1923 octets (5.9 bits / caractère)

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème
Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie

Théorème de Shannon
Code de Huffman

Conclusion

20/23

Remarques

- soit les probabilités (donc la table de codage) sont connues, soit il faut transmettre aussi la table avec le texte compressé.
- grande sensibilité au "bruit" : erreur sur un seul bit peut perturber tout le décodage. (contrairement à un code de longueur fixe)
- optimal parmi les codages préfixes univoques : un code par symbole de l'alphabet. Autre manière de faire : LZ / LZW (Lempel-Ziv-Welch 1977-1984).
- pas intéressant si les statistiques évoluent au cours du temps (source non-stationnaire).

21/23

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Séance 3

- 1 Position du problème
 - Codage et décodage
 - Codes préfixes
- 2 Théorie statistique de l'information
 - Notion d'entropie
 - Théorème de Shannon
 - Code de Huffman
- 3 Conclusion

22/23

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion

Conclusion

Codage de Huffman optimal parmi les *codes préfixes univoques* (utilisation de statistiques d'ordre 1).

Codes préfixes = *codes instantanés*.

Applications :

- variante de Huffman dans la transmission par fax,
- formats ZIP, GZIP, BZIP2 & co,
- utilisé dans la compression avec perte JPEG, MP3
→ la semaine prochaine !

23/23

Initiation au traitement du signal - Séance 3

F. Sur - ENSMN

Position du problème

Codage et décodage
Codes préfixes

Théorie statistique de l'information
Notion d'entropie
Théorème de Shannon
Code de Huffman

Conclusion