

Outils Informatique Codage

E. Jeandel

Représentation des données

- Comment coder une image en un fichier ?
- Comment coder un texte en un fichier ?
- Comment représenter une couleur dans un ordinateur ?
- Comment représenter un graphe dans un ordinateur ?
- Comment représenter une base de données dans un ordinateur ?

Dans un ordinateur

- La notion de base est la bit.
- Un bit peut prendre deux valeurs, 0 ou 1.
- Les bits sont regroupés, pour simplifier, par 8, pour former ce qu'on appelle un octet.

Représenter des données, c'est donc les représenter comme une série de bits, ou comme une série d'octets.

Représentation des nombres

Un nombre entier est représenté par son écriture en base 2 :

```

51 | 110011
51 | 00110011
1664 | 1101000000
    
```

Plus d'informations dans le cours d'Architecture en L3

1

Représentation des dates

Plusieurs formats :

- Format utilisé (entre autres) sous DOS et Windows :
- Date sur 16 bits : 5 pour le jour, 4 pour le mois, et 7 pour l'année, en prenant comme référence 1980

2010/01/19 = 0011110 - 0001 - 10011 = 0x3e33

Bug de l'an... ?

- Heure sur 16 bits : 5 pour l'heure, 6 pour les minutes, 5 pour les secondes.
- Problème ?

9 : 51 : 36 = 01001 - 110011 - 10010 = 0x4e72

- Format utilisé sous Unix : nombre de secondes écoulées depuis minuit UTC (temps universel coordonné) le 1er janvier 1970, codé sur 32 bits, dont un bit pour le signe.

2010/01/19 à 9 :51 :36 = 1263891096

Bug de l'an... ? (il y a 31557600 ~ 15 x 2²¹ secondes dans un an)

Représentation des caractères

Un caractère est représenté par 8 bits (donc par un nombre entre 0 et 255). Des tables de correspondance expliquent comment on passe des 8 bits au caractère correspondant.

Exemple pour le caractère "é" :

Norme	Code binaire	Décimal
ISO/IEC 8859-1 (Latin-1 Western European)	11101001	233
Mac OS Roman	10001110	142
CP437	10000010	130

99% des normes ont les mêmes 128 premiers caractères, correspondant à la norme ASCII.

Chacune des normes ne permet de représenter que 256 caractères : insuffisant pour certaines langues.

2

Codepage 819 - Latin 1 - ISO 8859-1

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0																
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

Codepage 437 - United States

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0																
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

Codepage 1275 - Apple Latin 1

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0																
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

Codepage 924 - Latin 9 - EBCDIC

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0																
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

Représentation des caractères : Unicode

Unicode est un standard qui explique comment représenter et manipuler du texte. Il contient en particulier une liste de plus de 100000 caractères.

On trouve ensuite plusieurs façons de les représenter :

- UTF-32 : Représente chaque caractère sur 32 bits (donc 4 octets)
- UTF-8 : Représente la majorité des caractères fréquents sur 8 bits, d'autres sur 16 bits, 24 ou 32 bits

Caractère	Code UTF8
a	01100001
é	11000011 10101001
€	11100010 10000010 10101100
ï	11110000 10011101 10011111 10011001

Table utilisée en cours et en TD

0	00000	h	01000	p	10000	x	11000
a	00001	i	01001	q	10001	y	11001
b	00010	j	01010	r	10010	z	11010
c	00011	k	01011	s	10011	.	11011
d	00100	l	01100	t	10100	,	11100
e	00101	m	01101	u	10101	'	11101
f	00110	n	01110	v	10110	!	11110
g	00111	o	01111	w	10111	?	11111

ne permet de représenter que des minuscules et quelques signes de ponctuation, mais bien suffisant pour les exercices.

Récapitulatif

La séquence suivante :

11110000 10011101 10011111 10011001

peut donc représenter :

- Les 4 nombres 240, 157, 159, 153
- Les 2 nombres 61597 et 40857
- Le nombre 4036861849
- Les 4 caractères ÷ÿfO (en CP437)
- Les 4 caractères ÷üüö (en Mac OS Roman)
- Le caractère I (en UTF-8)
- Les instructions assembleur x86 suivantes : LOCK POPF ; LAHF ; CDQ

Un exemple

Un logiciel d'archivage permet de regrouper en un seul fichier plusieurs fichiers et répertoires afin, par exemple, de les stocker ensuite plus facilement, ou de les compresser.

3

4

Questions

Problème 1. Comment peut-on savoir qu'un code donné est uniquement déchiffirable ?

Problème 2. Comment créer un code uniquement déchiffirable ?

Code bloc

Définition 3. Un block code est un code où tous les mots du code sont différents et de même longueur.

	000		101		100
	011		110		111

Proposition 3. Un block code est non-ambigu.

Inégalité de Kraft-McMillan

Théorème 4. On note l_i la longueur du code pour le i -ème objet. On suppose qu'il y a n objets.

Si un code est non-ambigu, alors

$$\frac{1}{2^{l_1}} + \frac{1}{2^{l_2}} + \dots + \frac{1}{2^{l_n}} \leq 1$$

Exemple : peut-on trouver un code pour les fruits de sorte que

- 🍌 soit codé sur 1 bit;
- 🍌, 🍋, 🍓, 🍇 soient codés sur 3 bits;
- 🍌 sur 4 bits ?

Test de Sardinas-Patterson

Soit C le code. On le met dans une colonne C_0

A chaque étape i

- Si un mot de la première colonne C_0 commence un mot de la dernière colonne

C_i (ou vice versa), on met le reste dans une nouvelle colonne C_{i+1}

On arrête le calcul lorsqu'on boucle

Théorème 5. C est un code non-ambigu si et seulement si on n'atteint jamais un mot de C au cours de l'exécution

Exemple

1000	1	000	0	000	111	00
0111		001	1	111	0	000
1001		110		001	10	111
1110				110	11	01
011				11	1	001
00				0		0
						10
						110
						11

C est ambigu

Exemple

1000	0	000	0	00	111	00	00
0111	1	111	10	111	0	111	111
1110		110		11	10	0	0
011		11		0	11	10	10
100		00				11	11

C est non ambigu

Code préfixe

Définition 4. Un code est préfixe si aucun mot n'est un préfixe d'un autre mot.

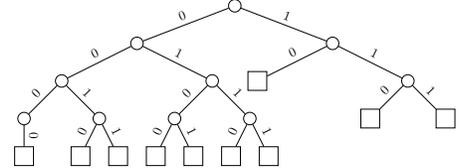
Exemple :

0000	0101	110
0010	0110	111
0011	0111	
0100	10	

Remarque : un block code est préfixe.

Code préfixe

On peut représenter les codes préfixes par des arbres :



Codes préfixes

Théorème 6. Un code préfixe est non-ambigu.

Preuve : pour décoder, on suit le chemin dans l'arbre. Dès qu'on arrive à une feuille, on a fini et on repart du début.

- Les codes préfixes sont instantanés : Dès qu'on a fini de lire le premier mot de code, on sait qu'on l'a lu et qu'on peut passer au deuxième.

Théorème de Kraft-McMillan

Théorème 7. Soit l_i des entiers. Si les l_i vérifient

$$\frac{1}{2^{l_1}} + \frac{1}{2^{l_2}} + \dots + \frac{1}{2^{l_n}} \leq 1$$

Alors il existe un code préfixe tel que le i ème objet a pour longueur l_i .

Ensemble

Théorème 8. On note l_i la longueur du code pour le i -ème objet. Si le code est non-ambigu, alors

$$\frac{1}{2^{l_1}} + \frac{1}{2^{l_2}} + \dots + \frac{1}{2^{l_n}} \leq 1$$

Théorème 9. Soit l_i des entiers. Si les l_i vérifient

$$\frac{1}{2^{l_1}} + \frac{1}{2^{l_2}} + \dots + \frac{1}{2^{l_n}} \leq 1$$

Alors il existe un code préfixe tel que le i ème objet a pour longueur l_i .

Corollaire

Théorème 10. Si C est un code non-ambigu, on peut trouver un code préfixe avec exactement les mêmes longueurs de mot.

Les codes qui ne sont pas préfixes ne servent à rien.

Comment trouver, sachant les l_i , le code préfixe qui convient ?

Conclusion

- Pour coder des objets, on utilisera des codes non-ambigus, et souvent des codes préfixes.
- Si on sait qu'un objet apparaît très souvent, il faut lui donner un code plus petit que les autres.

Comment faire ? C'est le prochain cours.