Codage des caractères

1 NSI

Le codage des caractères est une convention qui permet, à travers un codage connu de tous, de transmettre de l'information textuelle, là où aucun support ne permet l'écriture scripturale.

1.ASCII

L'American Standard Code for Information Interchange (Code américain normalisé pour l'échange d'information), plus connu sous l'acronyme ASCII ([askiː]), est une norme informatique de codage de caractères apparue dans les années 1960. C'est la norme de codage de caractères la plus influente à ce jour.Code ASCII de base :

,																			
Dec	H)	Oct	Cha	,	Dec	Нх	Oct	Html	Chr	Dec	Нх	Oct	Html	Chr	Dec	Нх	Oct	Html Ch	<u>ır</u>
0	0	000	NUL	(null)	32	20	040	a#32;	Space	64	40	100	a#64;	0	96	60	140	a#96;	8
1	1	001	SOH	(start of heading)	33	21	041	a#33;	1	65	41	101	A	A	97	61	141	a#97;	a
2	2	002	STX	(start of text)	34	22	042	@#34;	rr	66	42	102	B	В	98	62	142	a#98;	b
3	3	003	ETX	(end of text)	35	23	043	@#35;	#	67	43	103	C	С	99	63	143	a#99;	C
4	4	004	EOT	(end of transmission)	36	24	044	\$	ş	68	44	104	D	D	100	64	144	4#100;	d
5				(enquiry)				%		69	45	105	E	E				e	
6				(acknowledge)				&					F					a#102;	
7				(bell)				'		- 0			G					a#103;	
8		010		(backspace)				a#40;					H					h	
9				(horizontal tab))					6#73;					a#105;	
10		012		(NL line feed, new line)				&# 4 2;					a#74;		_			j	
11		013		(vertical tab)				a#43;					a#75;		1			k	
12		014		(NP form feed, new page)				a#44;	•				a#76;					l	
13		015		(carriage return)	I			a#45;			_		a#77;					m	
14	_	016		(shift out)				a#46;					a#78;					n	
15		017		(shift in)				a#47;					a#79;		1			o	
				(data link escape)				a#48;					P					p	
				(device control 1)				a#49;		ı			Q					q	
				(device control 2)				2					R					r	
				(device control 3)	100			3					۵#83;					s	
				(device control 4)				4					a#84;					t	
				(negative acknowledge)				5					a#85;					u	
				(synchronous idle)				a#54;					V					v	
				(end of trans. block)				7					a#87;					w	
				(cancel)				8					4#88; دور					x	
		031		(end of medium)				a#57;					Y					y	
				(substitute)				:					6#90;					z	
				(escape)				;					6#91;	-		. —		{	
		034		(file separator)				<					6#92;						
		035		(group separator)				=					6#93;					}	
		036		(record separator)				>					a#94;					~	
31	1F	037	US	(unit separator)	63	3 F	077	?	?	95	5 F	137	a#95;	_	127	7F	177		DEL
													5	ourc	e: W	ww.	Look	upTables	.сош

Code ASCII étendu (valeur hexadécimales lignes dizaines colonnes unités)

	0	1	2	3	4	5	6	7	8	9	Α	В	С	D	E	F
0	NUL	SOH	STH	ETH	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	CD2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		ļ	"	#	\$	%	&	'	()	*	+		-		- /
3	0	1	2	3	4	5	6	7	8	9	:	- ;	<	=	>	?
4	@	Α	В	С	D	Е	F	G	Н	- 1	J	K	L	M	N	0
5	Р	Q	R	S	Т	U	V	W	Х	Υ	Ζ	[- \]	Α	
6	`	а	b	С	d	е	f	g	h	i	j	k	- 1	m	n	0
7	р	q	r	S	t	u	٧	W	Х	у	Z	{		}	~	DEL
8	€			f			†	‡	Ŷ	‰	Š	<	Œ		Ž	
9		'	'	"	"	•	_	_	~	TM	š	>	œ		ž	Ÿ
Α		j	¢	£	300	¥	-	§		0	а	«	٦	-	®	_
В	۰	±	2	3	1	И	¶			1	0	>>	1/4	1/2	3/4	ż
С	À	Á	Â	Ã	Ä	A	Æ	Ç	È	É	Ê	Ë	Ì	ĺ	Î	Ϊ
D	Ð	Ñ	Ò	Ó	Ô	Ő	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ຄ
E	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	ĺ	î	Ϊ
F	ð	ñ	ò	ó	ô	ő	Ö	÷	Ø	ù	ú	û	ü	ý	þ	ÿ

Ex Donner la valeur du code ASCII de caractère « A »:

Quelle relation existe-t-il entre les caractères latins minuscule et majuscule ?

2.UNICODE UTF-8

Unicode est un standard informatique qui permet des échanges de textes dans différentes langues, à un niveau mondial. Unicode se contente de recenser, nommer les caractères et leur attribuer un numéro. Mais il ne dit pas comment ils doivent être codés en informatique.

Plusieurs codages des caractères Unicode existent :

- UTF-32 qui code chaque caractère sur 32 bits (soit quatre octets)
- UTF-16 qui code chaque caractère sur 16 ou 32 bits (soit deux ou quatre octets)
- UTF-8 qui code chaque caractère sur 8, 16, 24 ou 32 bits (soit un, deux, trois ou quatre octets).

Le plus couramment utilisé, notamment pour les pages Web, Python est UTF-8.

La principale caractéristique d'UTF-8 est qu'elle est rétro-compatible avec le standard ASCII, c'est-à-dire que **tout caractère ASCII se code en UTF-8 sous forme d'un unique octet, identique au code ASCII.** Par exemple « A » (A majuscule) a pour code ASCII 65 et se code en UTF-8 par l'octet 65. Chaque caractère dont le point de code est supérieur à 127 (caractère non ASCII) se code sur 2 à 4 octets. Le caractère « € » (euro) se code par exemple sur 3 octets : 226, 130, et 172.

Tableau récapitulatif de la concordance entre les codes UNICODE et UTF-8

Caractères codés	Représentation binaire UTF-8	Premier octet valide (hexadécimal)	Signification
U+0000 à U+007F	0xxxxxx	00 à 7F	1 octet, codant 7 bits
U+0080 à U+07FF	110xxxxx 10xxxxxx	C2 à DF	2 octets, codant 11 bits
U+0800 à U+0FFF	11100000 101xxxxx 10xxxxxx	E0 (le 2 ^e octet est restreint de A0 à BF)	
U+1000 à U+1FFF	1110 <i>0001</i> 10xxxxxx 10xxxxxx	E1	
U+2000 à U+3FFF	1110 <i>001x</i> 10xxxxxx 10xxxxxx	E2 à E3	
U+4000 à U+7FFF	111001xx 10xxxxxx 10xxxxxx	E4 à E7	3 octets, codant 16 bits
U+8000 à U+BFFF	111010xx 10xxxxxx 10xxxxxx	E8 à EB	5 octets, codant to bits
U+C000 à U+CFFF	11101100 10xxxxxx 10xxxxxx	EC	
U+D000 à U+D7FF	11101101 100xxxxx 10xxxxxx	ED (le 2 ^e octet est restreint de 80 à 9F)	
U+E000 à U+FFFF	1110111x 10xxxxxx 10xxxxxx	EE à EF	
U+10000 à U+1FFFF	11110000 1001xxxx 10xxxxxx 10xxxxxx	FO (le 26 patet est restraint de 00 à BF)	
U+20000 à U+3FFFF	11110000 101xxxxx 10xxxxxx 10xxxxxx	F0 (le 2 ^e octet est restreint de 90 à BF)	
U+40000 à U+7FFFF	11110 <i>00</i> 1 10xxxxxx 10xxxxxx 10xxxxxx	F1	4 octets, codant 21 bits
U+80000 à U+FFFFF	1111001x 10xxxxxx 10xxxxxx 10xxxxxx	F2 à F3	
U+100000 à U+10FFFF	11110100 1000xxxx 10xxxxxx 10xxxxxx	F4 (le 2 ^e octet est restreint de 80 à 8F)	

Convertir de l'UNICODE en UTF-8

 $\text{Soit le code UNICODE U+ XXXX qui donne } D_{15}D_{14}D_{13}D_{12} \ D_{11}D_{10}D_{9}D_{8} \ D_{7}D_{6}D_{5}D_{4} \ D_{3}D_{2}D_{1}D_{0} \ \text{en binaire} \\$

1. Si le code UNICODE est inférieur ou égal à U+007F -caractères ASCII standards-

Le code UTF-8 est l'octet du code ASCII il n'y a aucun changement, il commence par un zéro.

 ${\color{red}0} \; D_6 D_5 D_4 \, D_3 D_2 D_1 D_0$

2. Si le code UNICODE est compris entre U+0080 et U+07FF compris

Le code UTF-8 est reparti sur deux octets. Le 1er octet commence par 110 et le 2ème par 10,

110 D₁₀ D₉D₈ D₇D₆ **10** D₅D₄ D₃D₂D₁D₀

Exemple: Pound Sign £ U+00A3 = 0000 0000 1010 0011

U+0080 <U+00A3<U+07FF alors UTF 8 = $\frac{1100}{0010} 0010 \frac{1010}{1010} 0011 = 0xC2 A3$

```
3. Si le code UNICODE est compris entre U+0800 et U+FFFF compris

Le code UTF-8 est reparti sur trois octets. Le 1er octet commence par 1110, le 2ème par 10, le 3ème par 10

1110 D<sub>15</sub>D<sub>14</sub>D<sub>13</sub>D<sub>12</sub> 10 D<sub>11</sub>D<sub>10</sub>D<sub>9</sub>D<sub>8</sub> D<sub>7</sub>D<sub>6</sub> 10 D<sub>5</sub>D<sub>4</sub> D<sub>3</sub>D<sub>2</sub>D<sub>1</sub>D<sub>0</sub>

Exemple : Comment coder U+2639 qui donne 0010 0110 0011 1001 en binaire

U+2636 est supérieur à U+07FF, le codage en UTF-8 sera sur 3 octets

1110 D<sub>15</sub>D<sub>14</sub>D<sub>13</sub>D<sub>12</sub> 10 D<sub>11</sub>D<sub>10</sub> D<sub>9</sub>D<sub>8</sub>D<sub>7</sub>D<sub>6</sub> 10 D<sub>5</sub>D<sub>4</sub>D<sub>3</sub>D<sub>2</sub>D<sub>1</sub>D<sub>0</sub>

1110 0 0 1 0 10 0 1 10 0 0 10 1 1 0 0 1

E 2 9 8 B 9

0xE2 0x98 0xB9 Il s'agit de ce symbole : 

http://unicode-table.com/fr/#control-character
```

3.Encodage en python

Pour imposer le codage des caractères utf-8 pour le programme en python

```
# -*- coding: utf-8 -*-
# coding: utf-8
```

```
En Python, les chaînes de caractères sont en réalité codées en UTF-8, un code (ou
encodage) plus général que l'ASCII et qui permet de représenter plus de
caractères, par exemples les caractères chinois. En UTF-8, Les caractères d'indices
inférieurs à 128 sont exactement les mêmes qu'en ASCII : le tableau ci-dessus
fonctionne donc."""valeur ACII d'un caractère
ord() renvoie la valeur unicode d'un caractère
l'unicode commence par les caractères ASCII qui gardent leurs
valeurs
                                                                       La valeur ASCII du caractère p est 112
chr() renvoie le caractère dont on donne son code
https://unicode-table.com
                                                                       Le caractère ASCII dont le code est 112
.....
                                                                       est le caractère p
c = 'p'
print(f"La valeur ASCII du caractère {c} est", ord(c))
print(f"Le caractère ASCII dont le code est {ord(c)} est le
caractère ", chr(ord(c)))
En utilisant la notation u"\u2600"
                                         avec une valeur en
                                                                       <class 'str'>
hexadécimal
on affiche le caractère grace à son code unicode
                                                                       ☀●↑營企★
s=u"\u0070"
                                                                       ☀●↑餐企★
print(s)
print(type(s))
print(u"\u2600 \u2601 \u2602 \u2603 \u2604 \u2605")
                                                                       <class 'str'>
b'\xe2\x98\x80'
s=u"\u2600"
print(s)
                                                                       <class 'bytes'>
print(type(s))
s=u"\u2600".encode("utf-8")
print(s)
print(type(s))
```



Conversion en binaire :

Utiliser le tableau qui convient à la conversion :

1	1	0	D10	D9	D8	D7	D6	1	0	D5	D4	D3	D2	D1	D0

1	1	1	0	D15	D14	D13	D12	1	0	D11	D10	D9	D8	D7	D6	1	0	D5	D4	D3	D2	D1	D0

Code UTF-8:



Exercice 2 : Le code du symbole est U+00C7 en UNICODE. Exprimer ce symbole en UTF-8.

Conversion en binaire : 0000 0000 1100 0111

Utiliser le tableau qui convient à la conversion :

1	1	0	D7	D6	D5	D4	1	0	D3	D2	D1	D0
1	1	0	1	1	0	0	1	0	0	1	1	1

1	1	0	D10	D9	D8	D7	D6	1	0	D5	D4	D3	D2	D1	D0

1	1	1	0	D15	D14	D13	D12	1	0	D11	D10	D9	D8	D7	D6	1	0	D5	D4	D3	D2	D1	D0

Code UTF-8:

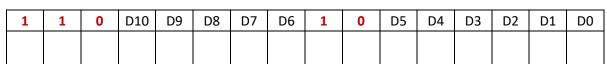


Exercice 3: Le code du symbole est U+07F7 en UNICODE. Exprimer ce symbole en UTF-8.

Symbole **n'ko** gbakourounèn

Conversion en binaire :

Utiliser le tableau qui convient à la conversion :



	1	1	1	0	D15	D14	D13	D12	1	n	D11	D10	D9	D8	D7	D6	1	n	D5	D4	DЗ	D2	D1	חח
L	-	-	-	•)	יב)	012	•)	1)	י) ו	֝ נ	ט	•)))	ו	1	0
Г																								
								l																
								l																
					l		l	ĺ			l													

UTF-8: