

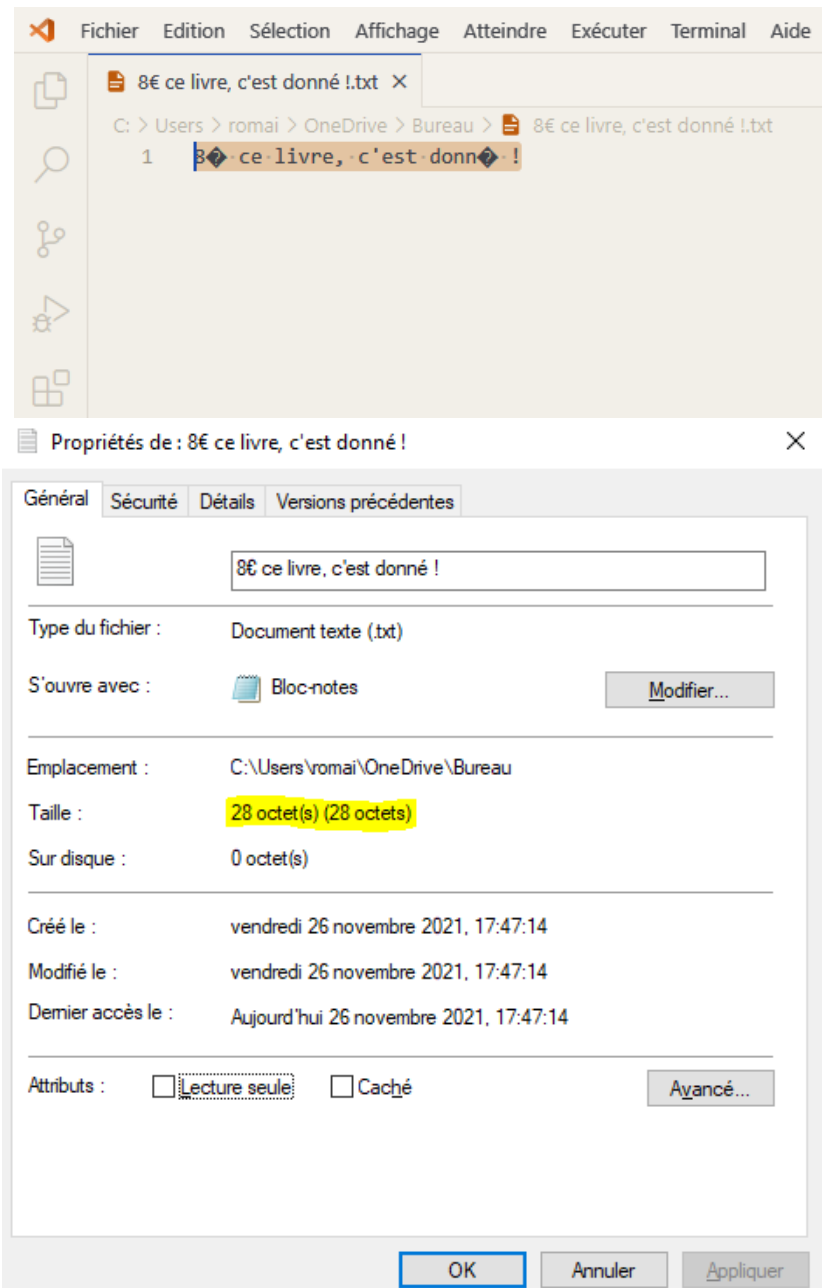
1. Codage d'un texte simple en ISO-LATIN :

Il est impossible d'enregistrer le fichier texte encodé en ISO-8859-1, car ce dernier contient deux caractères qui ne sont pas pris en charge par cette forme d'encodage, il s'agit de : « € ; é ».

Mon texte est constitué de 26 caractères sur la première ligne.

La taille du fichier texte créé est de 28 octets.

Sur l'éditeur hexadécimal, je constate qu'il y a aussi 28 codes hexadécimaux soit les 28 octets vu précédemment, car il y a bel et bien les 26 caractères de la première ligne mais aussi deux octets supplémentaires dont nous verrons la signification dans la prochaine section : « Les sauts de ligne dans les différents systèmes ».



En ISO-LATIN-15, le caractère « € » est codé par le code hexa « A4 » qui correspond à « ¤ ». (voir image)

Offset(h)	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	Texte Décodé
00000000	38	A4	20	63	65	20	6C	69	76	72	65	2C	20	63	27	65	8€ ce livre, c'e
00000010	73	74	20	64	6F	6E	6E	E9	20	21	0D	0A					st donné !..

Après ces nombreuses études, je peux donc maintenant affirmer que le nombre d'octets nécessaires au codage d'un fichier est égale au nombre de caractères qui le compose mais en rajoutant un certain nombre d'octets pour le retour à la ligne en fonction du système d'exploitation utilisé par l'utilisateur, dans le cas d'un fichier encodé en ISO-LATIN-15.

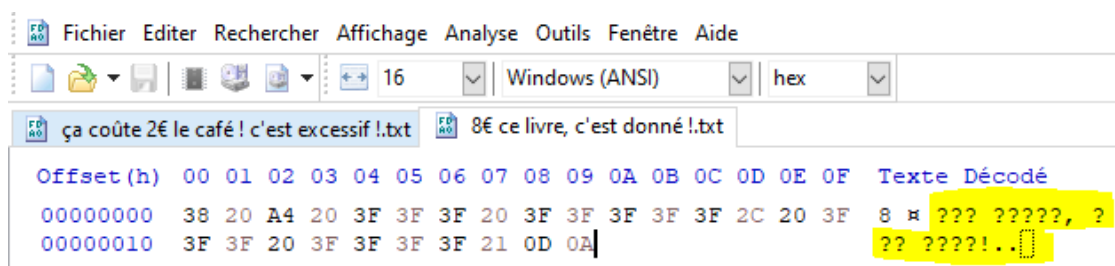
Les sauts de ligne sur les différents systèmes :

Mais attention, sur Windows, si on fait un retour à la ligne à la fin de ce texte alors deux nouveaux octets se rajoutent, il s'agit de « 0D » qui correspond à « Carriage Return » (Retour de chariot) qui définit le retour à la bordure gauche, ainsi que « 0A » qui correspond à « Line Feed » (Saut de ligne) qui définit le passage à la ligne de texte suivante, et qui indique à la machine qu'elle doit être prête à écrire la ligne suivante. Ces deux caractères ne sont donc pas des caractères saisis, mais il rajoute bien deux octets au fichier ! Cette convention est appelée « DOS », elle est utilisée sur les systèmes Windows.

Il existe une deuxième convention appelée « UNIX », elle est utilisée sur les systèmes Linux, Multics, BeOS, AmigaOS ainsi que les nouvelles versions de MacOS (à partir de la version 9). Contrairement à la convention DOS, cette convention permet de définir le saut de ligne par UN seul caractère (« 0A » soit LF) et non pas deux.

Enfin, il existe une troisième convention appelée « MAC » qui était utilisée sur les anciennes versions de MacOS (antérieures à la 9). Contrairement à la convention UNIX, cette convention permet de définir le saut de ligne par le seul caractère (« 0D » soit CR).

En remplaçant « 8 € ce livre, c'est donné ! » par « 8 € эта книга, это дано! » dans l'interpréteur, puis en l'enregistrant encodé en ISO-LATIN-15, des points d'interrogations s'affiche dans HxD à la place des lettres cyrilliques.

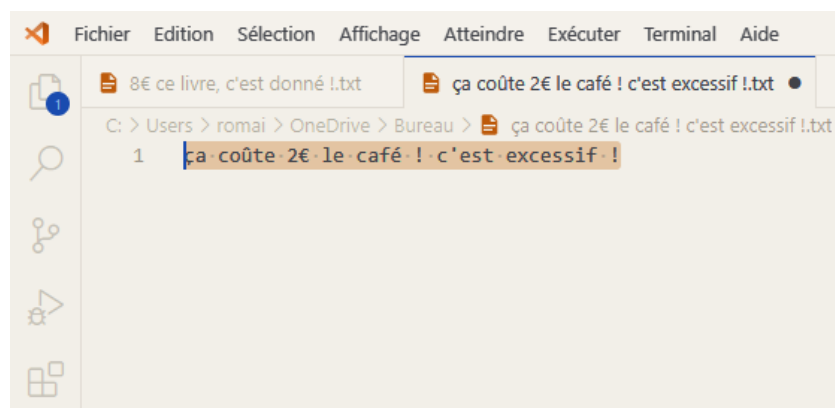


Mais si j'encode ce texte en UTF-8 au lieu de l'ISO-LATIN-15 alors ce ne sont plus des points d'interrogations qui s'affichent à la place des caractères initiaux mais des caractères totalement différents !

Offset(h)	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	Texte Décodé
00000000	38	20	E2	82	AC	20	D1	8D	D1	82	D0	B0	20	D0	BA	D0	8 â, ~ Ñ.Ñ, Ð° Ð°Ð
00000010	BD	D0	B8	D0	B3	D0	B0	2C	20	D1	8D	D1	82	D0	BE	20	Ð.Ð°Ð°, Ñ.Ñ, Ð%
00000020	D0	B4	D0	B0	D0	BD	D0	BE	21	0D	0A						Ð°Ð°Ð°Ð°!..

2. Codage d'un texte simple en UNICODE UTF-8 :

Mon texte est constitué de 38 caractères sur la première ligne. Je précise que je n'ai pas fait de retour à la ligne.



Le nombre de codes hexadécimaux contenus dans le fichier est supérieur au nombre de caractères saisis, car les caractères : « ç ; û ; € ; é » ont été remplacés par plusieurs caractères chacun, ce qui augmente donc le nombre de caractères et par conséquent le nombre d'octets !

Offset(h)	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	Texte Décodé
00000000	C3	A7	61	20	63	6F	C3	BB	74	65	20	32	E2	82	AC	20	Àsa coÀ»te 2â, ~
00000010	6C	65	20	63	61	66	C3	A9	20	21	20	63	27	65	73	74	le café! c'est
00000020	20	65	78	63	65	73	73	69	66	20	21						excessif !

On peut justifier cela en réalisant le tableau suivant :

Caractères	Codes hexadécimaux	Justifications + valeurs Unicode
ç	C3 A7	→ \$C3 \$A7 → <u>1100</u> 0011 <u>1010</u> 0111 → 0000 0000 <u>1110</u> 0111 = \$E7 → U+00E7 → ç
€	E3 82 AC	→ \$E2 \$82 \$AC → <u>1110</u> 0010 <u>1000</u> 0010 <u>1010</u> 1100 → <u>0010</u> 0000 <u>1010</u> 1100 = \$20AC → U+20AC → €
û	C3 BB	→ \$C3 \$BB → <u>1100</u> 0011 <u>1011</u> 1011 → 0000 0000 <u>1111</u> 1011 = \$FB → U+00FB → û
!	21	→ \$21 → <u>0010</u> 0001 → <u>0010</u> 0001 = \$21 → U+0021 → !
É	C3(Ã) 89(‰)	→ \$C3 \$89 → <u>1100</u> 0011 <u>1000</u> 1001 → 0000 0000 <u>1100</u> 1001 = \$C9 → U+00C9 → É
é	C3(Ã) A9(©)	→ \$C3 \$A9 → <u>1100</u> 0011 <u>1010</u> 1001 → 0000 0000 <u>1110</u> 1001 = \$E9 → U+00E9 → é

Modification des octets dans un fichier :

Je dois donc remplacer le « 2 » par un « 4 » et sachant que 4 a pour valeur Unicode : U+0034, je remplace donc la valeur hexadécimale « 32 » par « 34 ». Ensuite, je dois rajouter « très », tout d'abord je rajoute donc un espace en utilisant sa valeur Unicode : U+0020, puis je rajoute le mot en utilisant la valeur Unicode de « t » minuscule (U+0074), la valeur Unicode de « r » minuscule (U+0072), la valeur Unicode de « è » minuscule (U+00E8) et enfin la valeur Unicode de « s » minuscule (U+0073).

Offset(h)	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	Texte Décodé
00000000	C3	A7	61	20	63	6F	C3	BB	74	65	20	32	20	E2	82	AC	ÀSa coÀ»te 2 â,-
00000010	20	6C	65	20	63	61	66	C3	A9	20	21	20	63	27	65	73	le cafÀ© ! c'es
00000020	74	20	65	78	63	65	73	73	69	66	20	21					t excessif !



Offset(h)	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	Texte Décodé
00000000	C3	A7	61	20	63	6F	C3	BB	74	65	20	34	20	E2	82	AC	ÀSa coÀ»te 4 â,-
00000010	20	6C	65	20	63	61	66	C3	A9	20	21	20	63	27	65	73	le cafÀ© ! c'es
00000020	74	20	74	72	E8	73	20	65	78	63	65	73	73	69	66	20	t très excessif
00000030	21																!

1 ça coûte 4 € le café ! c'est très excessif ! | 3.

Compatibilité des codes ISO-LATIN et UTF-8 :

1 étudiés

UTF-8

1 Étudiés

ISO-8859-15

Comme démontré juste au-dessus, ce mot étonnant indique que l'expéditeur a utilisé un codage en ISO-LATIN, car il a pu écrire « étudiés » sans souci alors qu'en tant que lecteur je n'arrive pas à lire son message cela indique que j'utilise de mon côté un codage en UTF-8. En regardant de plus près dans l'éditeur hexadécimal on remarque que les caractères : « À » et « © » correspondent en réalité au caractère « é », le tableau réalisé précédemment détient cette information.

Les caractères compatibles entre le codage ISO-LATIN et l'UTF-8 sont les caractères non-accentués soit les « Contrôles C0 et latin de base », compris dans la première partie de la table Unicode.