

文字コード

コンピュータが直接扱えるのは、0と1で表わされる2進数だけです。一方、人間は、数値以外のデータである文字、音、色などもコンピュータで取り扱いたいわけです。そこで、文字、音、色などの数値でないデータを2進数の数値を使って表すために、符号化(以下、コード化)という工夫が必要になります。

コード化は、コンピュータを使う人間が自由勝手に行うこともできますが、それだと、機種や端末装置が変わったり、インターネット等を介して相互に情報をやりとりすると、正しく表示されない(このことを一般に「化(ば)ける」といいます)事態になります。ですので、一般的には、コード化は予め決められた符号化のルール(文字コード体系)に従って行います。

ここでは、文字(文字の中には10進数の数字も含まれます)のためのコード化について解説します。

コンピュータの世界では統一的な文字コード体系が定められていますが、残念なことに、1種類ではなく、代表的なものだけでも数種類あります。皆さんの中で、パソコンを利用している際に、「文字化け」を経験したことはありませんか。これは、文字コード体系が複数あるのが原因なのです。

以下、代表的なものを紹介します。

なお、コンピュータでは2進数ですが、それだと長くなり、また読みにくいので、16進数で表記します。

ASCII

ASCII (American Standard Code for Information Interchange の略で、アスキーと呼びます)は、米国で生まれた文字コード体系で、文字コード体系の元祖です。

英数記号の 1 文字を 7 ビットで表すもので、大文字・小文字の英字、数字、記号といった文字が 95 個(内、1 個は空白文字)と、33 個の制御文字の計 128 文字が規定されています。

制御文字は、表示するための文字種ではなく、モニタやプリンタなどの入出力装置を制御するためのものです。

【便利知識】

1 バイトは 8 ビットですが、英字・数字・記号を扱うには 7 ビットあれば十分なので、データ長を短くしてデータ通信の負荷を軽減するとともに、残り 1 ビットをデータ通信の誤り検出(パリティチェックといいます)のために役立てるといふ工夫が感じられます。

【便利知識】

Windows OS の和文フォントで、ASCII の 16 進数の 5C(本来は半角バックスラッシュ'¥')の場所に、半角の円記号('¥')を割り当ててしまったために、Windows パソコンでは半角バックスラッシュが使えません。

Excel では、その場所だけ和文フォント以外のフォントに変えると半角バックスラッシュで表示できます。

(筆者のパソコンでは、Word、PowerPoint ではこの方法がうまくいきませんでした。全角バックスラッシュ('\')のフォントやフォントサイズを変えるなどで、半角バックスラッシュに見せかけるしかないのかもしれませんが。ちなみに、'\ 'は「すらっしゅ」と入力して変換キーを押せば出てきます。)

ホームページ上では、逆に半角円記号が半角バックスラッシュで見えてしまうことが多いです。特に半角にしなければならない時以外は、円記号もバックスラッシュも、全角にした方が無難です。

16進数	ASCII	16進数	ASCII	16進数	ASCII	16進数	ASCII
00	NULL	20	SP	40	@	60	`
01	SOH	21	!	41	A	61	a
02	STX	22	"	42	B	62	b
03	ETX	23	#	43	C	63	c
04	EOT	24	\$	44	D	64	d
05	ENG	25	%	45	E	65	e
06	ACK	26	&	46	F	66	f
07	BEL	27	'	47	G	67	g
08	BS	28	(48	H	68	h
09	HT	29)	49	I	69	i
0A	LF	2A	*	4A	J	6A	j
0B	VT	2B	+	4B	K	6B	k
0C	FF	2C	,	4C	L	6C	l
0D	CR	2D	-	4D	M	6D	m
0E	SO	2E	.	4E	N	6E	n
0F	SI	2F	/	4F	O	6F	o
10	DLE	30	0	50	P	70	p
11	DC1	31	1	51	Q	71	q
12	DC2	32	2	52	R	72	r
13	DC3	33	3	53	S	73	s
14	DC4	34	4	54	T	74	t
15	NAK	35	5	55	U	75	u
16	SYN	36	6	56	V	76	v
17	ETB	37	7	57	W	77	w
18	CAN	38	8	58	X	78	x
19	EM	39	9	59	Y	79	y
1A	SUB	3A	:	5A	Z	7A	z
1B	ESC	3B	;	5B	[7B	{
1C	FS	3C	<	5C	\	7C	
1D	GS	3D	=	5D]	7D	}
1E	RS	3E	>	5E	^	7E	~
1F	US	3F	?	5F	_	7F	DEL

※※は制御文字です (印字されません)

※SP (16進数の20) は半角の空白 (ブランク)

※16進数の5Cは、Windows OSの和文フォントでは¥で表示されます

ASCII コード表

日本語用の JIS 規格

米国で生まれたコンピュータは、次いでヨーロッパや日本などに広がり、今では全世界に広がっています。そうすると、各国の言語に対応が必要となり、ASCII では全く不十分になりました。

西欧諸国では、ASCII 文字セットを 8 ビットに拡張し、増えた領域にヨーロッパ諸国で使われる文字を定義した「ISO-8859」という文字セットが使われています。追加する文字の種類によって、西欧諸国で使用する規格、東欧諸国で使用する規格、北欧諸国で使用する規格、ロシア語の規格、など 16 種類の規格が制定されています。

JIS X 0201

日本でも日本工業規格(JIS)によって、同様に ASCII 文字セットを 8 ビットに拡張し、増えた領域に(半角の)カタカナ文字と句読点などの記号を規定した、JIS X 0201 という規格が制定されました。

現在ではこの規格が単独で使用されることはほとんどありませんが、Shift-JIS 文字コード体系や EUC-JP 文字コード体系の一部に組み込まれて用いられています。

16進数	JIS	16進数	JIS	16進数	JIS	16進数	JIS	16進数	JIS	16進数	JIS	16進数	JIS	16進数	JIS
00	NULL	20	SP	40	@	60	`	80		A0	未定義	C0	ク	E0	
01	SOH	21	!	41	A	61	a	81		A1	。	C1	チ	E1	
02	STX	22	"	42	B	62	b	81		A2	「	C2	ツ	E2	
03	ETX	23	#	43	C	63	c	81		A3	」	C3	テ	E3	
04	EOT	24	\$	44	D	64	d	81		A4	、	C4	ト	E4	
05	ENG	25	%	45	E	65	e	81		A5	・	C5	ナ	E5	
06	ACK	26	&	46	F	66	f	81		A6	ヲ	C6	ニ	E6	
07	BEL	27	'	47	G	67	g	81		A7	フ	C7	ヌ	E7	
08	BS	28	(48	H	68	h	81	未 定 義	A8	イ	C8	ネ	E8	未 定 義
09	HT	29)	49	I	69	i	81		A9	ウ	C9	ノ	E9	
0A	LF	2A	*	4A	J	6A	j	8A		AA	エ	CA	ハ	EA	
0B	VT	2B	+	4B	K	6B	k	8B		AB	オ	CB	ヒ	EB	
0C	FF	2C	,	4C	L	6C	l	8C		AC	ケ	CC	フ	EC	
0D	CR	2D	-	4D	M	6D	m	8D		AD	ユ	CD	ヘ	ED	
0E	SO	2E	.	4E	N	6E	n	8E		AE	ヨ	CE	ホ	EE	
0F	SI	2F	/	4F	O	6F	o	8F		AF	ヲ	CF	マ	EF	
10	DLE	30	0	50	P	70	p	90		B0	-	D0	ミ	F0	
11	DC1	31	1	51	Q	71	q	91		B1	ア	D1	ム	F1	
12	DC2	32	2	52	R	72	r	92		B2	イ	D2	メ	F2	
13	DC3	33	3	53	S	73	s	93		B3	ウ	D3	モ	F3	
14	DC4	34	4	54	T	74	t	94		B4	エ	D4	ヤ	F4	
15	NAK	35	5	55	U	75	u	95		B5	オ	D5	ユ	F5	
16	SYN	36	6	56	V	76	v	96		B6	カ	D6	ヨ	F6	
17	ETB	37	7	57	W	77	w	97		B7	キ	D7	ラ	F7	
18	CAN	38	8	58	X	78	x	98	B8	ク	D8	リ	F8		
19	EM	39	9	59	Y	79	y	99	B9	ケ	D9	ル	F9		
1A	SUB	3A	:	5A	Z	7A	z	9A	BA	コ	DA	レ	FA		
1B	ESC	3B	;	5B	[7B	{	9B	BB	サ	DB	ロ	FB		
1C	FS	3C	<	5C	¥	7C		9C	BC	シ	DC	ワ	FC		
1D	GS	3D	=	5D]	7D	}	9D	BD	ス	DD	ン	FD		
1E	RS	3E	>	5E	^	7E	~	9E	BE	セ	DE	・	FE		
1F	US	3F	?	5F	_	7F	DEL	9F	BF	ソ	DF	・	FF		

※赤は制御文字です (印字されません)

※SP (16進数の20) は半角の空白 (ブランク)

※16進数の5Cは、Windows OSの欧文フォントではバックslashで表示されます。16進数の7EはASCIIではチルダ (~) ですが、JISではオーバーライン (˘) です。

JIS X 0201 表

JIS X 0208(通称 JIS 漢字コード)

JIS X 0201 では、ひらがなや漢字が扱えません。

漢字のように何千文字もあるような文字体系は、1 バイトのコードでは対応不可能で、2 バイトコード体系の開発が不可欠でした。1 バイトでは最大 256 文字しか入りませんが、2 バイトあれば理論的には最大 65,536 文字まで入ります。2 バイトコードの実用化は日本語や中国語、韓国語の情報処理にとって大きなステップでした。

1978 年、JIS は、2 バイトのコード体系で、ひらがな、漢字(約 6,000 字)、数字や英字の全角文字などを規定する、JIS X 0208 という文字コード体系を規定しました。通称として、JIS 漢字コードとか、JIS 第 1 第 2 水準漢字などと呼ばれるものです。

具体的には、第1バイトと第2バイトの、ともに 16 進数表記で 21~7E のそれぞれ 94 個、全体では $94 \times 94 = 8,836$ 個の文字を表すことができる領域に、6,879 文字を割り当てています。

【便利知識】

2000 年に JIS X 0208 コード体系に約 4,000 文字を追加した拡張規格として JIS X 0213 が制定され、2004 年に一部改訂されました。2000 年版を JIS2000、2004 年版を JIS2004 と称することもあります。

第1バイト	第2バイト	漢字
21	21	あ
22	22	ア
23	23	い
24	24	イ
25	25	う
26	26	ウ
27	27	え
28	28	エ
29	29	お
2A	2A	オ
2B	2B	か
2C	2C	カ
2D	2D	き
2E	2E	キ
2F	2F	く
30	30	ク
31	31	け
32	32	ケ
33	33	こ
34	34	コ
35	35	さ
36	36	サ
37	37	し
38	38	シ
39	39	ず
3A	3A	ズ
3B	3B	せ
3C	3C	セ
3D	3D	そ
3E	3E	ソ
3F	3F	た
40	40	タ
41	41	な
42	42	ナ
43	43	ち
44	44	チ
45	45	つ
46	46	ツ
47	47	て
48	48	テ
49	49	ど
4A	4A	ド
4B	4B	に
4C	4C	ニ
4D	4D	は
4E	4E	ハ
4F	4F	ひ
50	50	ヒ
51	51	ふ
52	52	フ
53	53	ぶ
54	54	ブ
55	55	ぼ
56	56	ボ
57	57	ぽ
58	58	ポ
59	59	ま
5A	5A	マ
5B	5B	み
5C	5C	ミ
5D	5D	め
5E	5E	メ
5F	5F	も
60	60	モ
61	61	む
62	62	ム
63	63	め
64	64	メ
65	65	も
66	66	モ
67	67	む
68	68	ム
69	69	め
6A	6A	メ
6B	6B	も
6C	6C	モ
6D	6D	ゆ
6E	6E	ユ
6F	6F	よ
70	70	ヨ
71	71	ら
72	72	ラ
73	73	り
74	74	リ
75	75	ろ
76	76	ロ
77	77	り
78	78	リ
79	79	る
7A	7A	ル
7B	7B	る
7C	7C	ル
7D	7D	る
7E	7E	る

JIS X 0208(JIS 漢字コード)表の一部

シフト JIS(Shift_JIS)

2 バイトコードができたからといって、従来の 1 バイトコードを捨てて、全てが 2 バイトの世界になるわけにはいきません。プログラムのコンパイラやインタプリタなどは、相変わらず ASCII のような 1 バイトコードをベースとして動くので、過去の資産との互換性が求められます。

シフト JIS は、日本の 1 バイトコード体系の JIS X 0201 と 2 バイト漢字コード体系の JIS X 0208 とを組み合わせて運用する方式です。この規格を Shift_JIS といいます。

シフト JIS の漢字コード体系は、JIS 漢字コード体系の表を、8000 代(具体的には 8140 以降)、9000 代(同 9040 以降)、E000 代(同 E040 以降)、F000 代(同 F040 以降)に分散して、移動(シフト)した形になっています。

なぜ、このような位置に漢字コード体系をシフトさせているかということ、JIS X 0201 コード体系で、80~9F、E0~FF の欄が未使用になっているからです。

例えば、「JIS 漢字コード」という文字を、JIS X 0201 表と JIS 漢字コード表を使って、16 進数に変換すると、「4A 49 53 34 41 3B 7A 25 33 21 3C 25 49」となります。

(分かりやすくするため、間に空白を入れました。)

これを、1 バイトコードと思って、JIS X 0201 表で文字に変換していくと、「JIS4A;z%3!<%|」となってしまいます。

同じく、「JIS 漢字コード」という文字を、シフト JIS 漢字コード表を使って、16 進数に変換すると、「4A 49 53 8A BF 8E 9A 83 52 81 5B 83 68」となります。

これを、1 バイトコードと思って、JIS X 0201 表で文字に変換しようとする、「JIS」に変換した後の 16 進数「8A」が未定義であることに気づき、これは 2 バイトコードと解釈して、シフト JIS 漢字コード表を探して「漢」に変換するのです。以下同様に、1 バイトコードなのか、2 バイトコードなのかを判断して、JIS X 0201 表とシフト JIS 漢字コード表を使い分けることで、正しく変換ができるのです。

【便利知識】

JIS X 0208 表の代わりに、文字数が拡張された JIS X 0213 表に準じたシフト JIS コード体系が使用されることもあります。この規格を Shift_JIS-2004 といいます。

【便利知識】

UNIX OS(Web サーバなどの OS として利用されることが多かった)のために開発された日本語用文字コード体系に、EUC-JP という規格があります。これは、ASCII コード体系と JIS X 0208 コード体系を組み合わせたものと言えます。

	0 1 2 3 4 5 6 7 8 9 A B C D E F	
8140	、 : ; ? ! * ' -	
8150	— ￣ \ ` \ ` / ` / ` / ` / ` / ` / ` / ` /	
8160	~ ... ' ' ' ' ' ' () [] {	
8170	} < > 《 》 『 』 【 】 + - ± ×	
8180	÷ = ≠ < > ≤ ≥ ∞ ∴ ∕ ∘ ∘ ∘ ∘ ∘ ∘ ∘ ∘ ∘	
8190	\$ € £ % # & * @ \$ ☆ ★ ○ ● ◎ ◇ ◆	
81A0	□ ■ ▲ ▼ ▽ ※ 〒 → ← ↑ ↓ ■	
81B0	ε ∃ ⊆ ⊃ ∩ ∪ ∩	
81C0	∧ ∨ ¬ ⇒ ⇔ ∇ ∅	
81D0	∠ ⊥ ^ ∂ ∇ ≡	
81E0	≡ ≪ ≫ √ ∵ ∝ ∴ ∫ ∬	
81F0	Å % # b) † ‡ ¶ ○	
8240		0
8250	1 2 3 4 5 6 7 8 9	
8260	A B C D E F G H I J K L M N O P	
8270	Q R S T U V W X Y Z	
8280	a b c d e f g h i j k l m n o	
8290	p q r s t u v w x y z あ	
82A0	あいいうええおおかがきぎくぐけ	

	0 1 2 3 4 5 6 7 8 9 A B C D E F	
82B0	げごさざしじずせぜそそたち	
82C0	ちつつづてでとどなにぬねのはば	
82D0	ひびびぶぶへべほほまみむめ	
82E0	もややゆゆよよらりるれろわわぬぬ	
82F0	をん	
8340	アアイイウウエエオオカガキキグ	
8350	クグコゴサザシジスズセゼソゾダ	
8360	チヂッツツテデトドナニヌネノハバ	
8370	ハヒビビフフブヘベホボボマミ	
8380	ムメモヤユユヨヨラリルレロクツ	
8390	ヱアランヅカケ	A
83A0	ΒΓΔΕΖΗΘΙΚΛΜΝΞΟΠΡ	
83B0	ΣΤΥΦΧΨΩ	α
83C0	βγδεζηθικλμνξοπρ	
83D0	στυφχψω	
83E0		
83F0		
8440	АВВГДЕЁЖЗИЙКЛМНО	
8450	ПРСТУФХЦЧШЩЪЫЬЭЮ	
8460	Я	

中 略

	0 1 2 3 4 5 6 7 8 9 A B C D E F
88A0	嘩娃阿哀愛挨始逢葵茜襦惡握渥旭葦
88B0	芦勝梓庄驕坂宛姐蛇鮫綺鮎或粟裕
88C0	安庵按暗案闇鞍杏以伊位依偉困夷委
88D0	威尉惟意慰易椅為畏異移維緯綯萎衣
88E0	謂違遺医井亥域育郁磯一杏溢逸稻茨
88F0	芋齏允印咽員因姻引飲淫麻藤
8940	院陰隱韻吋右宇烏羽迂雨卯鵝窺丑碓
8950	臼渦噁喫齟蔚蔚鱉姥既浦瓜闊疇云運雲
8960	荏韻韻營嬰影映曳兕永泳洩瑛盈穎穎
8970	英衛詠銳液液益駢悅謁越閱櫻厭円
8980	團壘奄宴延怨掩援沿演炎焰煙燕猿緣
8990	斃苑屬遠鉛篇墮於汚甥凹央奧往應押
89A0	旺橫欧歐王翁揆鶯鷓黃岡冲荻億屋億
89B0	臆桶牡之俺卸恩溫穩音下化飯何伽伽
89C0	佳加可嘉夏嫁家寡科暇果架歌河火珂
89D0	禍禾稼箇花苛茄荷華菓蝦課嘩貨迦迦
89E0	露蚊俄俄我牙画臥芽蛾質雅餓麗介介
89F0	解回塊塊迴快怪悔恢懷戒拐改
8A40	魁晦械海灰界皆繪芥蟹開陪貝凱劭外
8A50	咳害崖慨慨涯碍街該鎧骸淫誓蛙垣
8A60	柿奶鈞劇嚇各廓摏摏格核殼獲確獲

	0 1 2 3 4 5 6 7 8 9 A B C D E F
8A70	角赫較郭闊隔革学岳峯額頤掛笠椹
8A80	櫛梔嶽瀉割喝恰括活渴滑葛褐轄且鯉
8A90	叶柘樺鞮株兜龜蒲蓋鏞喹鴨栢茅葦粥
8AA0	刈苳瓦乾侃冠寒刊勸勸卷喚堪姦完官
8AB0	寬干幹懇感憤憾撥敢柑桓棺款欽汗漢
8AC0	澗濯環甘監看羊管簡緩任翰肝艦莞觀
8AD0	諫糞選鑑問閑閑閑韓館館丸含岸巖玩
8AE0	瘰眼岩斫廣雁頑顏企伎危罍器基奇
8AF0	嫵奇岐希幾忌揮机旗既期棋棄
8B40	機婦穀氣汽畿折季稀紀徽規記貴起軌
8B50	輝飢騎兎龜偽儂妓宜戲技擬欺攢疑祇
8B60	義蟻誼讓掬菊鞠吉吃喫桔橘詰砧杵黍
8B70	却客脚虐逆丘久仇休及吸宮弓急救
8B80	朽求汲泣灸球究窮笈級糾給旧牛去居
8B90	巨拒撻拳渠虛許距鋸漁鯨魚亨亨京供
8BA0	俠僑兇甥共凶協匡叫鸞境峽強強怯
8BB0	恐恭挾教橋況狂狹矯胸齧興齶鄉鏡鑿
8BC0	齶驚仰癡堯曉業局曲極玉桐杆僅均
8BD0	巾錦斤欣欽琴琴筋筋紮芹菌衫襟謹近
8BE0	金吟銀九俱句区狗玖矩苦駟駟駒具
8BF0	愚虞噲空偶寓遇隅串櫛釧屑屈

後 略

【便利知識】

シフト JIS は、マイクロソフト社と日本企業数社が協力して 1982 年に誕生しました。そして、MS 漢字コードとも呼ばれ、Windows OS(日本語版)の標準として活躍してきましたが、一方で、インターネットで世界中が繋がり、Web サイトなどを通じて多数の言語に触れる機会が増し、Web サイトなどの閲覧、メールや Word 等の文書データの閲覧時に、しばしば文字化けが発生するようになってきました。

シフト JIS 自体も、もともと未使用の 1 バイトコードがあったことを利用したアクロバティックな方式で、登録されている文字数に限りがあることと、文字が規定されていない空きエリアを利用して、種々の外字が登録されていき、多数の亜流の文字コード体系が生まれてしまい、標準としては限界に近づいています。

このような背景と、後述のように、世界共通の標準的なコード体系の制定という動きに応じて、まもなく、シフト JIS は Windows の標準ではなくなる運命にあります。

UTF-8

世界で使われる全ての文字を共通の文字コード体系で利用できるようにしようという考えから、ゼロックス、マイクロソフト、アップル、IBM、サン・マイクロシステムズ、ヒューレット・パッカード、ジャストシステムなどが参加したコンソーシアムによって、Unicode(ユニコード)という符号化文字集合や文字符号化方式などを定めた文字コードの業界規格が、1991 年に誕生しました。

単一の大規模な文字セットであることが特徴で、世界各国の現代の文字だけでなく、古代の文字や歴史的な文字、数学記号、絵文字(ガラケーで使用されているもの)などにも対応しています。

国際規格の ISO/IEC 10646 と Unicode 規格は同じ文字コード表になるように協調して策定されています。

Unicode は、Unix、Windows、macOS、Java など種々の OS で利用されており、文字符号化方式も標準化されていて、いくつかの亜種があった Shift_JIS や EUC-JP の混乱のようなものは回避されています。

Unicode は当初、1 文字を 2 バイト(16 ビット)の固定長で表す規格のコード化方式で考えられ、UTF-16(Unicode Transformation Format-16)と名づけられました。ただ 2 バイトでは、最大で約 6 万字の文字しかサポートできないので、結局、2 バイト固定長という原則を崩すことになりました。基本的には(名前の通り)2 バイトで 1 文字を表

すのですが、「サロゲート」と呼ばれる拡張領域内の文字を利用する場合は 1 文字を 4 バイトで表します。つまり、UTF-16 は 2 バイト表現と 4 バイト表現が混在したコード変換方式です。

UTF-16 は ASCII 文字も 1 文字 2 バイトで表現します。つまり ASCII コードと互換性がありません。そのため、ASCII コードを想定したプログラムで誤動作を起こす可能性があります。そこで ASCII 文字については、ASCII と同様のコードを割り当てるようにしたコード変換方式が考えられました。それが UTF-8 です。

UTF-8 は 1~4 バイトで 1 文字を表し、漢字や仮名の 1 文字は 3 バイトを要します。シフト JIS では漢字や仮名は 2 バイトでしたので、日本語の文書などでは、データサイズが 1.5 倍程度に増大することになります。

Unicode には他に 4 バイト固定長で処理する UTF-32 や UCS-4 などの方式も規定されています。

【便利知識】

上述までの解説でお分かりのように、**現在の文字コード変換の標準は、UTF-8 に集約されつつあります。文字コードは UTF-8 にするという考え方でいけば、文字化けも起こりにくいと言えましょう。**