

情報の電子化 (1) 文字コードとテキスト・ファイル

かつらだ まさし
桂田 祐史

2003 年 5 月 15 日

プリントの訂正

- 課題 3 の締め切りは 5 月 28 日 (水曜) にします。
- C プログラムはマイドキュメント (Z:\.windows2000) ではなく、Z:\ (UNIX のホームディレクトリに相当) にセーブするのが良いでしょう。

これから数回、情報をどのようにコンピューター上で表現するかということをテーマにする。主に UNIX (Solaris) 上の小さなプログラムを用いた実験を通して説明する。今回は実験しながら、テキスト・ファイルの構造、特に文字コードについて学ぶ。

1 今週来週のための設定

以下に登場する `qkc`, `nkc` を使ってみなければ、`~/.cshrc` に (`emacs`¹などで)

```
set path=($path /usr/meiji/pub/bin)
```

という行を書き加えてから `kterm` を起動すること。

次回はキーボード入力が多いので、`tcsh` を使うように設定しておくことを強く勧める。“`passwd -e`”として、“New shell: ”に“`/usr/bin/tcsh`”を指定する。

<http://www.math.meiji.ac.jp/~mk/syori2-2005/How-to-tcsh/> ... 『`tcsh` への誘い』

2 コンピューターのファイル — 万物はビットである

2.1 デジタル、ビット、バイト

現在の一般的なコンピューター²の処理するデータはデジタル³である。

¹`emacs` の使い方は知っているはずだが、`emacs .cshrc &` とする。

²コンピューターの古い呼び方は — electric digital computer (電子計数型計算機) — という。アナログ計算機というものもあった。

³デジタル (digital) — もともとは「指の」という意味だが、ここでは「(離散的な) 数字の」という意味。

より具体的に言うと、

コンピューター上のすべてのデータは、ビット⁴ (bit) の有限列である。

実際には、複数のビットを適当な大きさにまとめて処理することが多い。例えば

- 表示する際には、3 桁, 4 桁ごとに区切って 8 進数, 16 進数とする⁵。
- データの格納や移動の際には、バイト⁶ やワード⁷ という単位が使われることが多い。

コンピューター上の情報はすべてファイルとして保存できるが、もちろん中身は数列である。特に Windows や、UNIX では、

ファイルはバイトの有限列である (ファイルは数列である)

と言って構わない。

色々なデータ (文書、数値データ、静止画、動画、音声, etc.) を記録したファイルがあるが⁸、バイトの列であることに変わりはない。ファイルの保存、コピー、通信による転送など、共通の操作で扱うことが出来る。

2.2 ファイルの大きさを調べる

UNIX 環境では、ファイルの大きさ (サイズ) は、`ls -l` または `wc` ファイル名 とすれば分かる (前者はバイト数、後者は文字数、単語数、行数を調べる)。

K, M, G, T, P, ..

国際単位系 (SI) で、大きな量を表すために、deca (da), hecto (h), kilo (k), mega (M), giga (G), tera (T), peta (P), exa (E), zeta (Z), yotta (Y) などの接頭語を補助的に用いる。コンピューターの世界でも、 $K = 1024 = 2^{10} \doteq 10^3$, $M = K^2 = \text{約} 100 \text{万}$, $G = K^3$, $T = K^4$, $P = K^5$ などをを用いる。身の回りにおける記憶装置、記憶媒体 (メイン・メモリー、フロッピー・ディスク、ハードディスク、CD-ROM, MO, DVD など) の容量の表示にはよく登場する。

2.3 レポート課題 X

以下のことを調べよ (〆切は 6 月末頃にする予定。X はそのとき適当な番号をつける。)
(ファイルのサイズについての感覚を身につけるのが目的である。)

⁴ または 1 の二つの状態のいずれであることを示す、情報量の単位をビットと言う。電子回路との相性がいい。

⁵ コンピューターの世界では、16 進法の数字として '0' ~ '9', 'A', 'B', 'C', 'D', 'E', 'F' を用いるのが普通。

⁶ バイト (byte) — 通常は 8 ビットのこと。2⁸ 通りの状態を区別出来るので、例えば 0 ~ 255 までの整数値を割り当てて読むことがよく行なわれる。ときどき「コンピューターは 2 進法で計算すると言うけれど、256 進法ではないか?」と思うこともある。16 進法では 2 桁で表示することになる ($\because 2^8 = 16^2$)。

⁷ ワード (word) — そのコンピューターにとって都合のよいデータの大きさを、普通はバイトの整数倍の大きさ。1 ワードが 16 bits, 32 bits, 64bits などの場合が多い。

⁸ 余談だが、ファイルの種類をある程度まで自動的に判別するコマンド `file` がある。`file *` としてみよう。

- (1) 情報科学センターの Windows 環境、Solaris (UNIX) 環境上のファイルのサイズについて調べよ。バイト数以外に、フロッピーや CD-ROM にどれくらい入るかを記せ。
- (a) 文書ファイル。
レポート、メール、C プログラム、 $\text{T}_\text{E}\text{X}$ のソース (.tex) など。
(印刷して何ページくらいの文書が何バイトになるか。)
 - (b) 実行可能プログラム
例えば C プログラムはコンパイル前と後でどうサイズが変わるか。
自分が普段使っているプログラム⁹のファイルのサイズは？
 - (c) (もし出来れば) 画像ファイル、音声ファイルなど。
(図の大きさ、色数も分かれば調べる。)
 - (d) 現在、自分が持っているファイルの総量。それは自分のホームディレクトリのあるディスクの全容量の何 % に相当するか？

<code>du -ks ~</code>	ホームディレクトリ下のファイルの総量
<code>du -ks ~/.snapshot</code>	.snapshot の下のファイルの総量
<code>df -k</code>	そのマシンにマウントされているディスクの状態

(それぞれで何が分かるかは授業中に話す。ききもらした場合は自力で!)

- (2) 自分が持っている本を一冊選び、その文字情報を記憶するファイルを作った場合、サイズはどれくらいになるか計算せよ。フロッピー・ディスクに記憶する場合、何枚必要か？また CD-ROM (容量 650MB 程度) には何冊分記憶できるか。

注意: 結果は K, M などを適切に選んで表現すること。例年間違えて結果に 1000 倍の差が出る人が少なからずいる (ひどい間違い)。

3 テキスト・ファイル

3.1 最初に

テキスト・ファイルとは

コンピューター上のファイルのうち、普通の意味で「読める」字で出来ているファイルのことをテキスト・ファイル (text file) という¹⁰。

文字コードとは

コンピューターでは、文字も数で表す。このためには「この数はこの文字」とルールを決める必要がある。文字を表現するため、表現したい文字の集合と、数の集合を選び、両者の間の

⁹コマンドの実体がどこにあるか (パス名) は、UNIX では、which コマンドを使って調べられる。例えば、`which emacs` とすると、emacs コマンドのパス名は /usr/meiji/gnu/bin/emacs であることが分かる。

¹⁰Windows の世界では、拡張子 “.txt” をつけるファイルのことであるが、これは OS にとらわれない重要な概念であるので、時間をかけて解説する。

1対1の写像を定める。このとき、ある文字に対応する数を、その文字の文字コードと呼ぶ。

3.2 英数字の文字コード (ASCII)

まず emacs を使って、

```
ascii.txt
012
ABC abc
```

という 2 行からなるファイルを作成し、以下のコマンドを実行してみよう。

```
isc-xas06% od -cx ascii.txt
```

(od -cx ファイル名 とすると、ファイルの文字コードが 16 進数で表示される。) 次の事が分かる¹¹。

- 数字 '0', '1', ... の文字コードは 0x30, 0x31, ...
- 英字 'A', 'B', ... の文字コードは 0x41, 0x42, ...
- 英字 'a', 'b', ... の文字コードは 0x61, 0x62, ...
- 行の終りは、'\n' という一つの文字 (文字コードは 0x0a) で表される。
- 空白 ' ' も一つの文字 (文字コードは 0x20) で表される。

(相当に不正確な言い方だが) このような「半角」英数字の文字コードは^{アスキー}ASCII というアメリカの規格で定められ、それが色々な規格に「輸入」されている。

¹¹ここでは 16 進数を、先頭に "0x" をつけることで表した (C 言語の世界の習慣)。

3.3 ASCII コード表を作ろう

```
print_code_table.c
1 /*
2  * print_code_table.c --- ASCII コード表 (印字可能な文字のみ) を表示する。
3  *   コンパイル: gcc -o print_code_table print_code_table.c
4  *   実行:      ./print_code_table
5  *   あるいは: gcc print_code_table.c ; ./a.out
6  */
7
8 #include <stdio.h>
9
10 int main()
11 {
12     int code;
13     /* 0x20 (==32) から 0x7e (==126) までの数をコードとする文字を扱う */
14     for (code = 0x20; code <= 0x7e; code++) {
15         /* 4 つおきに改行 */
16         if (code % 4 == 0)
17             printf("\n");
18         /* 16 進数, 10 進数, 文字 として表示 */
19         printf("0x%02x (%3d): %c   ", code, code, code);
20     }
21     printf("\n");
22     return 0;
23 }
```

というプログラム (WWW ページに載せてあるので、Netscape で例えば Z: ドライブ (UNIX 環境のホームディレクトリ) にセーブするとよい) をコンパイル&実行してみよう。

C 言語の豆知識

- a % b で a を b で割った余りを表わす。
- putchar(整数式); あるいは printf("%c", 整数式); でその整数を文字コードとする文字が表示される。
- `printf()` で 16 進数を表示するには `%x` という書式を用いる (`%d` は 10 進数)。

3.4 日本語の文字コード

日本語の文字コードは JIS (日本工業規格) で決められている。日本語における文字数が多いため¹²、1 文字を表すのに 16 ビットを用いる。

例えば「桂」という文字の JIS コードは 0x374b である。

しかし、ファイルの中に 0x37, 0x4b というバイト列をそのまま入れたのでは、ASCII の '7', 'K' と区別がつかない。両者を混在させるには何らかの工夫が必要になる。

情報科学センターの Solaris (UNIX) 環境では日本語 EUC という文字コードを用いている。やはり emacs を用いて

¹²ちなみに JIS の情報交換用漢字符号系コードの第一、第二水準の漢字は約 7000 個弱ある。JIS にはいくつかのバージョンがある。また第一、第二水準以外のものもあるが、あまり普及していない。

```
kanji.txt
```

```
桂田 祐史
```

というファイルを作成して、`od -cx kanji.txt` としてみると、「桂」が `0xb7, 0xcb` という 2 バイトで表現されていることが分かる。これは JIS コードの上下 8 ビットにそれぞれ `0x80` を加えたものになっている:

$$0x37 + 0x80 = 0xb7, \quad 0x4b + 0x80 = 0xcb$$

(もともと ASCII で用いている数値は `0x7f` までで、`0x80` 以上の数値は使われずに空^あいている。その部分に JIS コード表の文字を埋め込んだことになっている。)

日本語 EUC 以外にも、ASCII の文字と日本語の文字を混在を可能にした文字コードはいくつかある。

準備: `nkf` を用いて `kanji.txt` の文字コードを変換する

```
isc-xas5% nkf -j kanji.txt > kanji-jis.txt
```

```
isc-xas5% nkf -s kanji.txt > kanji-ms.txt
```

(`nkf` については、付録 A を見よ。)

通称「JIS 漢字」通信により情報を交換するために作られた国際的な規格に従ったもの。文字コードを切り替えるために、目印となる特別な文字列¹³を用いる。電子メールなどで、日本語メッセージの通信をするときに最もよく使われる ISO-2022-JP (後述) の基礎となった。実は `kterm` では JIS 漢字のデータも普通に表示できる。

```
isc-xas06% cat kanji-jis.txt
```

```
桂田 祐史
```

```
isc-xas06% od -cx kanji-jis.txt
```

```
0000000 033 $ B 7 K E D 033 ( B 033 $ B M 4  
          1b24 4237 4b45 441b 2842 201b 2442 4d34
```

```
0000020 ; K 033 ( B \n  
          3b4b 1b28 420a
```

```
0000026
```

「桂」の JIS コードである `0x37, 0x4b` が現われている。

通称「MS 漢字 (シフトジス)」パソコンの世界のために作られた規格 (Windows, Mac など) で採用されている。「半角カタカナ」も効率的に埋め込まれている。情報科学センターのワークステーションの通常の設定では直接表示できない (ただし `emacs` では扱える)。

¹³漢字に切り替えるために `0x1b, 0x24, 0x42` を、ASCII に戻すために `0x1b, 0x28, 0x42` を用いる。

```
isc-xas06% cat kanji-ms.txt
jc S (いわゆる文字化け)
isc-xas06% od -cx kanji-ms.txt
0000000 214 j 223 c 227 S 216 j \r \n
          8c6a 9363 2097 538e 6a0d 0a00
0000013
isc-xas06%
```

「桂」が 0x8c, 0x6a となっているが、どういうルールで変換されているかはちょっと分かりづらい。付録 D に JIS を MS 漢字に直す C の関数をあげておく。

腕試し用プログラミング課題 1

print_code_table.c が ASCII コード表を表示するように、漢字のコード表を表示するプログラムを作れ。レポートを送るときは Subject (件名) は “prog 1” として下さい。

3.5 文字コードをいじってみよう

WWW ページに mycat.c というプログラムを載せてある。これは cat コマンドの真似をして、ファイルの内容を標準出力 (通常は画面) に出力するだけのプログラムである。

mycat.c をコンパイルする

```
isc-xas06% gcc -o mycat mycat.c
isc-xas06% ./mycat mycat.c
```

このプログラム mycat.c を読んでみよう。中の print_file() という関数を書き換えると、色々なことができる。例えば

日本語 EUC に対応した mydump.c の print_file

```
void print_file(FILE * fp)
{
    int c, c2;
    /* ファイルの終わりまで一文字ずつ c に読み込み、標準出力に書き出す */
    while ((c = getc(fp)) != EOF) {
        if (c >= 0x80) {
            /* 0x80 以上だったら漢字の 1 バイト目だと判断して、
            もう 1 バイト読んで、まとめて出力する。 */
            c2 = getc(fp);
            printf("0x%02x 0x%02x: %c%c\n", c, c2, c, c2);
        } else if (c < 0x20 || c == 0x7f)
            /* 0x20 未満または 0x7f の場合は文字コードのみ表示 */
            printf("0x%02x\n", c);
        else
            /* それ以外の場合は文字コードと、その文字自身を出力 */
            printf("0x%02x: %c\n", c, c);
    }
}
```

```
isc-xas06% cat ascii_and_kanji.txt
I am 桂田祐史.
isc-xas06% ./mydump ascii_and_kanji.txt
0x49: I
0x20:
0x61: a
0x6d: m
0x20:
0xb7 0xcb: 桂
0xc5 0xc4: 田
0xcd 0xb4: 祐
0xbb 0xcb: 史
0x2e: .
0x0a
isc-xas06%
```

腕試し用プログラミング課題 2

ファイルの中の英小文字を英大文字に変換するプログラムを作りなさい。レポートを送るときは Subject (件名) は “prog 2” として下さい。(ヒント: 小文字、大文字のコードはそれぞれ連続していて、アルファベット順に並んでいる。)

腕試し用プログラミング課題 3

ファイルのバイト数、行数を数えるプログラムを作りなさい。(‘\n’ (文字コード 0x0a) の個数を行数と考えることにする。) レポートを送るときは Subject (件名) は “prog 3” として下さい。

4 レポート課題 3

自分の名前を構成する各々の文字 (桂田祐史なら「桂」, 「田」, 「祐」, 「史」の4文字) の JIS コード, EUC コード, SJIS コードを調べよ。締め切りは 5 月 28 日 (水曜) にする予定。

(以下は付録)

A 日本語の文字コードの変換

UNIX 上のコマンドには自動的に文字コードを判別して必要な処理をしてくれるものがあるが (emacs, less 等)、時にはユーザーが意識的に変換することが必要になる。

それほど難しい作業でもないのでフリーソフトがある。二つほど紹介する。

nkf (Network Kanji code conversion Filter) UNIX では定番。

nkf -e ファイル名 で日本語 EUC に変換したものを標準出力に書き出す。
nkf -j ファイル名 で JIS 漢字コードに変換したものを標準出力に書き出す。
nkf -s ファイル名 で MS 漢字コードに変換したものを標準出力に書き出す。
電子メールで使われる MIME のデコードもできる。
nkf -v でオプションの一覧が表示される。

kanji.txt を JIS 漢字に変換した kanji-jis.txt を作る

```
isc-xas06% nkf -j kanji.txt > kanji-jis.txt
```

qkc (Quick KANJI code Converter) Windows 版もある。行末の変換もしてくれる。

qkc -eu ファイル名 で日本語 EUC, 行末を UNIX 形式に変換する。
qkc -ms ファイル名 で MS 漢字, 行末を MS-DOS (Windows?) 形式に変換する。

kanji.txt を MS 漢字に変換した kanji-ms.txt を作る

```
isc-xas06% cp kanji.txt kanji-ms.txt  
isc-xas06% qkc -ms kanji-ms.txt
```

B インターネットで使って良い文字悪い文字 (日本語)

コンピュータ上で色々な日本語文字が使えるようになっていて、テキスト・ファイルを作ることが出来るが、インターネット (電子メール、WWW ページ、ネットニュースなど) で用

いる場合、すべての文字を使えるわけではない。

- 電子メール、ネットニュースでは、ISO-2022-JP と呼ばれる文字コードが使える。インターネットの作法である RFC で定義されている。

<http://www.noge.com/koba/network/RFC/rfc1468.html>

元々は日本のインターネットのルーツである JUNET で利用されて来たもの。

ASCII, JIS X 0201-1976, JIS X 0208-1978, JIS X 0208-1983 という 4 種類の文字コードをエスケープ・シーケンスというバイト列でスイッチする。

- ASCII はアメリカの規格 (キーボードから直接打ち込める英数字・記号)。
 - JIS X 0201 はその日本版 (円記号とバックスラッシュなど ASCII と異なる)。
 - JIS X 0208 はいわゆる JIS 第一、第二水準の漢字、ひらがな、カタカナ、その他記号。
- 「外字」は使えない (論外)。
(外字については後述。)
 - 機種依存文字も要注意。後々まで残る可能性のあるデータの表現には使うべきではない。(コンピューターやソフトウェアのメーカーが作った文字。普通の人には、JIS で正式に定義された文字と見分けが付きにくいかも知れない。丸つき数字、ローマ数字、携帯電話の絵文字などが有名。)
 - JIS X201 の右半面 (俗称「半角カナ¹⁴」) は RFC に違反しているので、メールやネットニュースでは使用してはいけない。

C 参考 URL

1. 「日本語フォントや文字コードについて」

<http://www.ryukyu.ad.jp/~shin/jdoc/>

なかなか充実したリンク集。このページを見るだけで、文字コードの問題がなかなか複雑であることを感じ取れる :-)

2. 「ほら貝」by 加藤弘一

<http://www.horagai.com/>

日本語の文字コードの話を調べて、本 (『電腦社会の日本語』, 文藝春秋, 2003, <http://www.horagai.com/www/salon/works/denno.htm>) も出している人で、今は「文字コードからは足を洗った」そうですが、どっこい WWW ページ上で『文字コードから見た住基ネットの問題点』など書いて発表していらっしゃいます。

¹⁴昔の MS-DOS パソコンやワープロ専用機では、これらの文字が他の日本語文字の半分の幅で表示されたため、こう呼ばれた (全角)。元々は印刷業界の言葉であったとか。しかし、そもそも規格には「半角」という言葉はない。文字がどう表示されるかは、利用する環境によるので規格の範囲外である。

D JIS を MS 漢字に直す関数

```
jis2sjis.c
1 /* JIS コード c1, c2 を MS 漢字コード s1, s2 に変換する */
2 void jis2sjis(int c1, int c2, int *s1, int *s2)
3 {
4     if (c1 & 1) {
5         c1 = (c1 >> 1) + 0x71;
6         c2 += 0x1f;
7         if (c2 >= 0x7f)
8             c2++;
9     }
10    else {
11        c1 = (c1 >> 1) + 0x70;
12        c2 += 0x7e;
13    }
14    if (c1 > 0x9f)
15        c1 += 0x40;
16    *s1 = c1;
17    *s2 = c2;
18 }
```