

古文書データベースにおける日本語データ処理の諸問題

Some Problems with Japanese Data Processing in the Ancient Document Database

高井正三

Shoso Takai

富山大学総合情報基盤センター

Information Technology Center, Toyama University

〒930-8555 富山市五福 3190

3190 Gofuku, Toyama 930-8555 Japan

takai@cns.toyama-u.ac.jp

概要

1994年以降、日本現存朝鮮古書を対象とする古文書の書誌情報データベースを構築してきた。今回このデータベースの冊子体目録を作成するに当たってそのマスター・データを分類して並び替える必要がでてきた。入力された日本語データ処理を行うに際して、コード変換、文字化け、国語辞書用ソーティング、外字処理、縦書き印刷、双行印刷など、解決しなければならない諸問題が発生した。本稿ではこれらの諸問題に対するいくつかの解決策と今後の課題を提案する。

キーワード

日本語コード 日本語データ処理 国語辞書用ソーティング 外字処理 コード変換 古文書データベース

1. はじめに

筆者は富山大学人文学部教授の藤本幸夫氏が30年の歳月をかけて、1910年以前の朝鮮古書を対象に全国から収集してきた「日本現存朝鮮古書」に関する28項目にわたる書誌情報約15,000件を、データベース化し、インターネット上に公開するため、データ入力、データベース・システムの開発に1994年から取り組んできた。

この時点まで入力されたデータは2,682件3.74MB、字数にして約200万字、全体の95%以上は漢字が占めている。このデータは当初、MIFESエディタ、一太郎日本語ワード・プロセッサ、ATOK日本語IMEを使って入力され、後に秀丸エディタで編集、追加が行われた。IMEは後継版ATOKを使用してきた。

今回、平成17年度科学研究費補助金(研究成果公開促進費)学術図書を申請するに当たって、冊子体目録の見本版を作成した。書誌情報の分類、LaTeXソースの生成、索引作成など、日本語データ処理では、コード変換をはじめ、国語辞書用ソーティング等、様々な問題に遭遇した。ここではその問題解決の方法をいくつか提案する。

2. 既存データの漢字データ処理

2.1 PL/I から Perl へ

データベースのマスター・データは、当初IBM PCで入力し、Mainframe System上に転送して、IBM漢字コードに変換し、PL/I言語で処理していたが、Mainframe SystemからUNIXサーバに替わって、文字列の処理に便利なPerl言語を使用するようになった。これに伴って日本語データはIBM漢字コードからシフトJIS(SJIS)コードに代わり、データ処理にはEUC(Extended Unix Code)を使用するようになった。

日本語データ処理では文字列の照合、置換、分類などが主な処理となるので、インデックス(index)関数やサブストリング(substr)関数が使用できる言語が必要で、旧来はPL/Iが、現在はPerl、C等のプログラミング言語が使われるのが一般的である。

2.2 文字化け

(1) メタ文字

Perlでは、EUCデータは問題なく処理できるがSJISコードでは、漢字の中にメタ文字(バック・スラッシュ付きの制御文字)、特にバック・スラッシュ「\」(日本では¥を割り当て)16進表記0x5C」と同じコードを有する漢字があり、これがPerlデータ処理の不具合を引き起こしている。「-」(0x815C)、「能」(0x945C)、「表」(0x955C)などの文字を扱うときは、そのコードが(¥)を含む文字であることに注意が必要である。以後のデータ処理が暴走する図1のSJIS

を Perl 処理し、pLaTeX2e の¥item[···]を生成出力した結果(図2)である。

"表紙左肩墨書「嘉林世稿」。該書は徐宗泰序・「雲江遺稿」・「竹陰遺稿」を収める上冊一冊のみを存し、「近水軒遺稿」・「附録」を収める下冊を失っている。刻工名は東大本と一致しない所がある。それには、筆者の見落とし、版心部が磨耗して刻工名が見えない、後の補刻等の原因が考えられるが、調査カードのままに記した。",

図1 SJISの本文(CSV形式の一部)

¥item[註記]表紙左圭紙左肩紙左肩味左肩墨左肩墨掌肩墨書肩墨書 墨書「墨書」穎書「嘉書」「嘉卯」「嘉林」「嘉林訊嘉林世嘉林世惟林世稿林世稿・世稿」世稿」稿」稿」各」該」該竣。該書。該書・該書は該書は什書は徐書は徐書は徐宗は徐宗堆徐宗泰徐宗泰徐宗泰序宗泰序 泰序・泰序・序・序・汚・「雲」雲孔「雲江」雲江・雲江遺雲江遺浩江遺稿江遺稿 遺稿」···める上める上罪る上冊る上冊磯上冊一上冊一版冊一···

図2 pLaTeX2eの¥item[···]生成結果

(2) コード変換による文字化け

SJISコードのデータ図3をnkf変換、秀丸でEUCコード変換したものを次に示す。図4がnkf(ネットワーク用漢字変換フィルター)による変換結果、図5は秀丸エディタによる変換結果である。

02 朝鮮【N柳琴/リュウキン】編 朝鮮【N李德懋/リトクボウ】···
 11 上部下向黒魚尾 下部長方形黒釘 魚尾下「韓客巾衍集」
 22 毎表紙題簽下部【I清慎勤/セイシンキン】(ニ・ニ×ー・一 陰刻朱印)···
 24 表紙左肩···附す。卷一首題下「李評青/潘評朱」

図3 SJISオリジナル

02 朝鮮【N柳琴/リュウキン】編 朝鮮【N李蛇諫織螢肇・椒 曄レ···
 11 上部下向 · 竝 · 蕊 · 絞 · · · · · 竝 · 次峇攀匄吁 検 ·
 22 毎表紙題簽下部【Iへ洞弍織札う響鳴 鵝曄廖米鵝テ鵝澎諫グ諫 · 鏢靚···
 24 表紙左肩···附す。卷一首題下「李評堯厨・昭襦廖···

図4 nkfによるEUCコード変換結果

02 朝鮮【N柳琴/リュウキン】編 朝鮮【N李博ヨ諫織螢肇・椒 曄レ··· 11 上部下向螺り竝 · 蕊 · 絞 · · · · · 竝 · 次峇攀匄吁 検 · 22 毎表紙題簽下部【I浮ヨ洞弍織札う響鳴 鵝曄廖米鵝テ鵝澎諫グ諫 · 鏢靚··· 24 表紙左肩··· 附す。卷一首題下「李評爺汾~ 昭襦廖···

図5 秀丸エディタV4.18のEUC変換結果

JIS区点89-92相当のSJISコード(MS-IME漢字3)が、EUCコードのSS3(Single Shift3)の3バイト・コードに対応しているためか、変換後に文字化けが発生する。文字列の中からJIS第1,第2水準,また

は置換文字列@99999に置換してから、SJIS EUC変換をすべきであったが、既存データの量が膨大で、総てのシステム外字、IBM拡張漢字を調査して変換する時間的余裕が無かった。

よく現れた文字化けの漢字を図6に示すが、右側の旧字体に割り当てられたSJIS外字コードは、EUCコードに変換ができなかった。そこで予め対応が解る旧字体のシステム外字域の漢字は左の新字体で置換してからデータ処理を行った。

とく	徳 0x93BF	徳	SJIS 外字	0xFABA
せい	青 0x90C2	青	SJIS 外字	0xFBF2
せい	清 0x90B4	清	SJIS 外字	0xFB43
こく	黒 0x8D95	黒	SJIS 外字	0xFCB4

図6 EUCコード変換できない漢字の例

漢字データの入力については、データベース構築時に一括して旧字体に変換するので、入力は常用漢字とすることとしたが、外字登録が700字以上になってコード表参照に時間が掛かり、ATOKIMEでも拡張漢字が入力できるようになったので、結局これらのシステム外字域の漢字が入力されてしまった。マスター・データの中に相当量のシステム外字等があったので、EUC変換をあきらめて、SJISのままデータ処理を余儀なくされた。

3. 国語辞書用ソーティング

文書目録の入力データは、最近入力のものまでをすべてデータベースに取り込み、これに分類コードを分類表(漢籍の分類)に従って付与し、各分類内では、国別、次いで書名のアイウエオ順に並び替え、同版の場合は所蔵者の北から南へ並べる。したがって分類の順序は以下の通りとなる。

- 1) 文書の分類コード順(漢籍の分類コード順)
- 2) 国別順(韓国,中国,日本,···の順)
- 3) 書名のアイウエオ順
- 4) 所蔵者の北から南への順(所蔵者コード昇順 = 小 大の順)

この時、書名については、後の索引にも使うので、国語辞書のようにソーティングして欲しいとのことで、これがまた大変であった。辞書用ソートのパラメータがあったが、われわれの要求を満たしてはいなかった。そこで次の例題に示すように、書名、撰者についてはカタカナ表記を取り出し、その濁点、半濁点を除去した書名、撰者を別フィールドに生成し、このフィールドを第1の分類キーに、本来の書名、撰者フィールドを第2の分類キーに設定してソーティングすると国語辞書用の分類になる。

この変換テーブルを作るときもカタカナの「ソ」がメタ文字を含むので、16進データとして定義する必要があった。

図7と図8はソーティングの例である。

とうさ	とうさ	踏査
とうざ	とうさ	当座
とうし	とうし	投資
とうじ	とうし	冬至
どうさ	とうさ	動作
どうざ	とうさ	同座
どうし	とうし	同士
どうじ	とうし	同時

図7 第1フィールドによるソーティング

とうさ	とうさ	踏査
とうざ	とうさ	当座
どうさ	とうさ	動作
どうざ	とうさ	同座
とうし	とうし	投資
とうじ	とうし	冬至
どうし	とうし	同士
どうじ	とうし	同時

図8 国語辞書用ソーティング結果

4. 外字処理

データの入カールールは次の通りである。

(1) SJIS に無い文字

・旧字体を用いてもよいが、JIS 第2水準までとする。

・「余」、「芸(ウン)」等、旧字体に変換されては困る漢字がある時は、原稿下部余白に明記しておく。

・入力できない漢字のうち、京都大学漢字典に記載されているものは、記号「@」と「康熙字典コード番号(5桁)」を連結して入れる。

・「康熙字典コード」にない漢字は、記号「@」の後ろに、朝鮮固有外字として60000台からの連番で数値を入れる。この連番は、共通の外字管理表を参照して、新規のものは新しい番号を発行する。

現在まで SJIS に無い文字は1,079字を登録し、コード入力を行っている。

(2) 朝鮮固有外字

朝鮮固有外字(異体字を含む)等は@60000台から順に外字原簿に登録し,@連番コードで入力する。現在まで114字を登録した。

(3) ハングル

ハングルは古文書の中ではそれほど多くないので@70000台から入力し、最終的な印刷ではハングル文字に置換する。入力はアレハ・ハングルというアプリケーションを使う。現在まで183字を登録した。

(4) 記号類

刻手名等を表す記号または記号に似た文字は、特殊外字として@90000台から登録し,@連番コードを入力する。現在まで57字を登録した。これは殆ど記号に近く、「外字・TrueType フォント エディタ」TTEdit を使用して作成した。

5. 冊子体目録の作成

5.1 縦書き印刷

古文書の書誌情報を冊子体目録や Web 上での表示を縦書き印刷にしたいのは漢籍研究者の常道であるらしい。

本学人文学部の小助川教授に、縦書きでかつ双行印刷などが可能なツールを教えてもらった。大阪大学の金水敏氏が作成した、「訓点資料(漢文訓読文)」用の LaTeX Style ファイル(マクロ・パッケージ)が使えるとのこと、これを使って冊子体目録を作成することにした。図9はこの現物の一部でページ数1,234になり、索引を入れて約1,300ページの冊子体目録(2分冊)を完成した。

版心	版式	紙質	寸法	装幀	刊地	刊年	刊者	刊種	撰者	書名	番号	御製	一	集部
上内	四	竹	三	原	漢	正	正	丁	朝鮮	御	0001	製	総	部
内	周	紙	十五	表	陽	祖	祖	西	鮮	定			集	
向	單	五	四	紙	臺	二	命	字	正	社			部	
二	邊	針	十	五	章	十	刊	印	祖	陸			類	
葉	内	釘	三	針	閣	三		本	編	千				
花	框	法	十	眼		年				漣				
紋	二	朱	三	釘		己				八				
魚	十	糸	十	法		未				卷				
尾	五		三			二				四				
下	〇		〇			七				冊				
横	〇		〇			九								
線	〇		〇			月								
魚	〇		〇			刊								
屋	〇		〇											
上	〇		〇											
御	〇		〇											
定	〇		〇											
社	〇		〇											
陸	〇		〇											
千	〇		〇											
漣	〇		〇											
八	〇		〇											
卷	〇		〇											
四	〇		〇											
冊	〇		〇											

図9 pLaTeX2eによる縦書き印刷

5.2 双行印刷

縦1行に小字で2列の文字を入れる双行は、この訓点マクロの $\%sougyou\{XX\}\{YY\}$ を使用した。片側の場合のみは $\%sougyou\{XX\}\{\}$ の右の項を空白とした。

5.3 連数字

連数字(縦中横)の数字印刷には、LaTeXの連数字マクロを利用したマクロ $\%Rsj$ を導入した。

$\%def\%Rsj\#1\{\%rensuji*\{c\}\{\#1\}\}$

5.4 索引の作成

索引は書名と撰者で、LaTeXの索引を取ってはみたが、これを国語辞書用にソーティングするのが大変で、結局は別名でファイルを作り、目録の分類同様に国語辞書用ソーティングを行って、文献の番号を付与する方法を採用し、索引を作成した(図10、図11)。

5.5 その他の記法

「=」という表記は、「 \cdots 」の縦書き表記に置換する。実際はpLaTeXの「 $\%$cdots$$ 」に置換することとした。

【ウ】

迂齋集	0272,0273,0274,0275
雨念齋詩鈔	0277,0278
...	
雲巖逸藁	0281
雲江遺稿	0032,0033
雲岡集	0282
雲谷集	0627,0628,0629,0631,0632,0633

図 10 書名索引 (アイウエオ順)

【オ】

王安石	0108,0109,0110,0111,0357,0358, 0359,0360,0361,2146
王欽臣	0227,0228,0229,0230,0231,0232
王鴻儒	0901
王靜	0909
王世貞	0083,2606,2607
王鈍	0027
翁方綱	2184

図 11 撰者索引 (アイウエオ順)

6. 各国 IME による連続漢字入力の方法

効率的な日本語入力方法、特に連続した漢字の入力を効率的に行う IME の登場が待たれるが、本場の台湾（繁体字）や中国（簡体字）では、それぞれの国の特徴ある IME が用意されている。

(1) 日本

・ATOK・・・ジャストシステムが一太郎用に発展させてきた日本固有の IME で結構使いやすい SJIS JIS, Unicode もサポートされている。

・MS IME 2003・・・SJIS と Unicode をサポートする Microsoft の IME であるが、辞書がこなれていない。

(2) 韓国

Microsoft IME 2003 が国際版の Office に付いてくる。一般的には McCune Reischauer 方式のローマ字入力で漢字変換するのがベターなようだ。

(3) 台湾（繁体字）

・Microsoft New Phonetic 2002a（音声：読み）

これは発音を頼りに該当する漢字を選択入力する IME である。

・Microsoft New ChangJie IME（部首合成）

台湾ではキーボードに漢字の部首を当てはめて、これを合成する方法で入力し、妥当な文字に変換していく IME で、スピードは抜群である。

・櫻花（繁体字 + 日本語 + ひらがな、カタカナ等）

これは台湾中央研究院計算中心で紹介された IME だが、日本語の他、ひらがな、カタカナも入力できるように、「櫻花（さくら）输入法」と言っている。

(4) 中国（簡体字）

・Microsoft Pinyin IME

Pinyin のローマ字入力する IME で該当する同音異義語を表示し、選択する方法である。

7. 今後の課題

OS レベルでは Mac OS X, Linux で Unicode が標準化され、Application ソフトウェアの対応が待たれている。Windows でも MS-SJIS 以外に Unicode が提供され、Office 系ソフトウェアで MS 明朝, MS ゴシックのフォントがサポートされるようになった。古文書のデータベース化における旧字体漢字の使用環境が整ってきた。

また、Flash や Java で Unicode の符号化形式 UTF-8 がサポートされているので、今後はセキュリティに堅牢な Java をベースに Web データベース・システムを構築していくのが Best Solution になるように思われる。現在 Java システムの開発に挑戦中である。

謝辞

筆者等がサービスしている日本現存朝鮮古書データベースの構築に当たって支援を受けた同僚の布村紀男助教授、朝鮮古刊本総合目録の作成に当たって、分類方法や記法、索引の採取など、本データベースに全体に関する情報の提供を受けた人文学部の藤本教授。これらのデータ入力に献身的な努力をしてくれた越野、洲崎、葉山、木戸、竹澤の女性スタッフ。ここに記して感謝の意を表する。

参考文献

- [1]The Unicode Standard Version 4.0, The Unicode Consortium, Addison-Wesley, 2004
- [2]高井正三, 朝鮮古書データベース蓄積・提供用旧字体および朝鮮固有外字の整備に関する研究, 平成 10 年度～平成 13 年度科学研究費補助金(基盤研究(C))(2): 課題番号 10680401) 研究成果報告書, 1-142, 2002
- [3]高井正三, 布村 紀男, 日本現存朝鮮古書データベース・システムの構築, 学術情報処理研究, No.5, 87-90, 2001
- [4]高井正三, 布村 紀男, 日本現存朝鮮古書データベース・システムの構築方法, 情報科学フォーラム 2003(FIT2003)論文集, D-26, 57-58, 2003
- [5]高井正三, Unicode4.0 解説, 富山大学総合情報基盤センター広報, Vol.2, 96-104, 2005
- [6]UNIX と X Window の日本語環境 - 日本語入力・表示から印刷まで - , 林秀幸著, ISBN4-526-03769-9, 日刊工業新聞社, 1995
- [7]Unicode 標準入門, トニー・グラハム著, 乾和志・海老塚徹訳, 関口正裕監修, ISBN4-7981-0030-7, 翔泳社, 2001
- [8]漢字文献情報処理研究, 漢字文献情報処理研究会編集, 第 5 号, ISBN4-87220-083-7, 好文出版, 2004