

相田 満

古典人名データベース作成上の問題点

Mitsuru Aida

Some problems on producing  
Japanese Classic Persons Database

---

We are proceeding to construct "Japanese Classic Persons Database" in The National Institute of Japanese Literature. Now, we have been inputting the thesaurus of many persons included in the Union Catalog of Japanese Old Books (Kokusho-soumokuroku).

This database is always opened by online real time system for researchers.

So, there are under many restrictions, for example, character-code-group or constructing method of database, and we need to consider these restrictions.

I will debate on the subject how to construct designs of database and inputting method.

国文学研究資料館研究情報部データベース室では、中村康夫氏を中心に古典人名データベースの作成を進めている。

扱うデータの時代は古代より幕末まで。現在、『国書総目録』記載の人名シソーラス約58,733件の入力が完了し、次の段階として、芳賀矢一編『日本人名辞典』よりデータを切り分けて入力を始めている（推定約47,000件）。

本データベースは、人物情報の直接典拠としての史料批判として十分に堪え得る『尊卑分脈』や『公卿補任』等も取り込んだデータベースの構築を目指しているが、第二段階作業として『日本人名辞典』を選んだ理由は、同書が簡便な書ながら、人物情報の見出として多様・多種のデータを盛り込んでいることと、データ切出しのノウハウ蓄積のための実験的意義もある。

本来ならば、そうした史料からのデータ切出しについても触れるべき所だが、現段階ではその方法も十分に吟味されていない故、今回は、現在入力・形成が進んでいる段階での、データベース設計、入力方法の構想と問題点について論じたい。（なお、岩波書店では『国書総目録』の著者別索引より、人物辞典上梓の計画があるが、同書はさらに伝記の判明する項目を抽出し、30,000件を扱う予定の旨。）

当古典人名データベースは、オンラインによる一般公開を前提としていることと、複数の違った目でデータ作りが進んでいることにも特色があるといえよう。特に後者では、デ

ータシート作成の第一次作業、すなわちデータ項目の直接の切り分け作業に国文学研究の未経験者のマンパワーも動員しているため（当然の如く切り分けられたデータの確認は行っているが）、データの切り分けノウハウの学習に多くの時間が費やされている。

その面についても多くの問題点が内在してはいるが、今回の発表では、さらにそれを電子化する上において生じる問題、すなわち使用する文字種、データ形成方法にハードウェアの規格上の制約によるさまざまな配慮が要求されることについても焦点をあて、以下、具体的にそれらの問題点の幾つかを紹介したい。

## 1. 『国書総目録』記載の人名ソースについて

### 1.1 データ形式

¥A 0 黄山 ¥A 1 自惚 ¥A 2 きやま ¥A 6 自惚笑・自惚山人 ¥E 1 絵本万歳島神樂の表紙作／前々太平記作〈天明六刊〉／忠臣金短冊作／万歳之島台作〈天明四刊〉／万たび物語作

¥A 0 湛契 ¥A 2 たんけい ¥A 6 高太夫 ¥E 1 仏眼次第／仏母曼拏羅念誦要法集

### 1.2 JIS漢字典拠外字母について

本データベースは入力を凸版印刷に依頼している。周知のように、同社にも J I S X 0 2 1 2 制定時においても約 1 万字相当数の字母セットがあると思われるが、それでも、545 件のレコードに、JIS 漢字典拠外の字母が発生している。データ作成にあたっては、旧漢字、新漢字の字母については、いずれも新体字に寄せた縮約を行うとともに、データ形成上支障のない範囲で字母をふりかえることもおこなっているが、これはその結果の数字である。

当館の端末では、既に独自の外字コードの蓄積があるが（内、1985年度の JIS 補助漢字選定にあたっての予備調査の段階では、当館作成の文字セットは 1,410 字）、先にも述べた通り、本データベースは、オンラインによる一般公開を前提としているので、うかつに番地を割り当てると別字に化ける可能性があり、また、端末機種の違いによる文字化けも考えられる。そこで、当面は ■ 字で以てゲタとし、外字であることを示すにとどめ、該当データを別シートに記録として残すこととした。今後、どのような形でこれらのデータを伝えるかは、今後の検討課題になっている。

なお、『国書総目録』に含まれる外字の内訳は、以下に大別できる。

#### 1.2.1 合成文字

日本で造った漢字まがいの文字を広く国字（もしくは倭〔和〕字）という。近年、『角川大字源』に、諸書にとりあげられた国字字母の一覧が紹介されたが、その選択基準の中には、編者自ら、「この一覧の中での取捨には疑を存すべきものも交じっている」と断っている如く、無批判にそれを字母と認定することに躊躇を覚えるものが

多くのある。こうした類の文字は（俗に歌舞伎文字と読みならわされることが多いが）、多く近世作品に見える。近年の歌舞伎興行界においても、例えば、『京子娘道成寺（きょうかのこむすめどうじょうじ）』を『京鹿子娘道成寺』と表記する傾向もある故、こうした字母で分解解釈可能なものは、以下に例示した如くに改めて入力を行なった。

義経灿（よしつねやまいり）……「山入」  
 畏染黄八丈（うえだぞめきはちじょう）……「上田」  
 花筏血汐船（はないかだちしおのとまぶね）……「苦舟」  
 體葛城合戦（にちょうのゆみかつらばがっせん）……「雙弓」  
 摯鼠花山姥（こもちねずみはなのやまうば）……「子持」  
 増補大仏殿萬代礎（ぞうほだいぶつでんばんだいのいしづえ）……「万代」

### 1.2.2. 梵字

真言系の内典の書名は悉曇文字で表記されることが多い。当面の手当として、漢訳化した表記に改めることで臨んだが、それとても完全な漢訳化に困難なものもあるため、読みのみをカタカナ表記にする方針で臨んでいる。

もう一つの方法としては、サンスクリットあるいはパーリ語によるローマ字表記も考えられたが、データ入力を外注で行なっている関係上、英数字が2バイト系の文字で入力されるため、データのタグ表記の文字との干渉の危険性、及び国文学者になじみのない表記であることも考慮に入れて、その方法を探らなかった。

### 1.2.3. 典拠外字母

1990年改定JISにて加えられた補助漢字（約6,000字）がサポートされれば上記条件（凸版印刷使用のコードを含む）以外で発生する推定800種の不足字母の内、9割にコードを割り振ることが可能になろう。通常のパソコン等の使用環境内で生じる不測については、このデータベース作成作業が現在進行中のこともあり、いずれ後考を期したい。いずれせよ、現状ではそれらがサポートされたハードウェア、及びソフトウェアは存在していないため、これも当面は該当データを別シートに記録として残すことにしている。

## 2. 芳賀矢一編『日本人名辞典』について

同データベースでは、現在、データ切出しの途にあるが、現段階では主にそのデータ切出しの基準について様々な問題を生じている。例えば、"A0 人物姓"の扱い一つをとっても、大田南畠という人物の切出しあは、四方赤良、属三人、或いは、直次郎という扱いが相応しいのか……、また姓、名という切出しに絞った場合、源等の氏、朝臣等のカバネ、或いは諡号、諱等の扱いをどうするのか等、時代により変遷があるため、扱われる基準がその時代毎に異なる。データの切り分けを厳密に行ないすぎてしま

まうと、検索時に全く使い辛いデータが出来上がってしまう。”A1 人物名”、”A3 名ヨミ”等も同様で、足利成氏（重氏）等、様々な異表記も見える。こうした表記は、典拠史料により異なることが多く、中にはどれに信を置いてよいのか途方に暮れるような人物データも登場してくる。上記の様な状況は今後も入力典拠を改める度に発生することが予想される。そこでデータ切り出し作業に於いては、複数姓名に対応できる人名採取につとめている。

冒頭にも述べた通り、このデータベースでは複数の人物情報の典拠から情報を抽出し、ある人物が同一人物と認定された場合、その情報をマージしてデータベースを構築する方法をとろうとしている。それぞれの原拠データは、まず字体（旧字・新字の字母の異なり、異体字も含む）からして不統一で、それらを尊重したいたずらな原型主義は、情報の混乱を招き、大きなデータベース形成の阻害要因となってしまうであろう。

文字というものは長い伝統の中で積み上げられた文化的所産であり、特に国文学の世界に於いては、それが非常に大きな意味を持つことはいうまでもない。しかし、原拠データ字母を尊重する観点で、外字字母を大量に発生させることは、閉じられたパッケージデータベースの世界では確かに有効であろう。しかし、当データベース形成作業では、スケジュールと労力上、現実問題として文字セットの手当が、データベース作成を大幅に遅らせる危険性を孕んでいるといえる。作成当事者としては非常に残念なことだが、現状のハードウェア規格体系中のデータベースは、情報伝達の正確度という観点で見れば、原拠データに比して90%のデータベース、将来さらに補助集合、あるいはそれをとりこんだUNIコードの使用を想定しても、95%を越えることはできないだろうというのが私の予測である。

なおさらに、細かな問題を挙げれば際限が無いが、紙数にも限りがあるため、およそのイメージとしての資料を呈示しておくに止める。

## 2.1 入力元データ

アキヤス 顯泰（北畠） 源氏。伊勢の國司。顯能の子。南北朝講和の後之を領すること舊の如紙。應永六年大内義弘を討ちて功あり。九年（或は應仁三年、又六年十一月）薨す。年四十三

アヅママロ 春滿（荷田） 國學四大人の一。本姓羽倉氏、通稱齋、一に東麻呂と云ふ。信詮の子。世々伏見稻荷山の祠官。夙に國學を唱へ元文元年七月二日沒す。年六十九。明治十六年二月贈正四位。著す所、萬葉集童蒙抄、伊勢物語童子問、出雲風土記考、齊明紀童謡考、僞類聚三代格考、春葉集等。

アヒミ 相見（巨勢） 佛畫家。一に相覽に作る。金岡の子。采女正、讚岐目等に任せられる。延喜中の人。

## 2.2 入力データ形式

¥A0北畠¥A1顯泰¥A3アキヤス¥A5源¥B2應永9年10月／應仁3年／

応仁6年11月￥B 6 4 3￥C 0顕能￥E 0国司￥G 0 <////伊勢の国司>  
 <////南北講和の後之を領すること旧の如し> <応永6年////大内義弘を  
 討ちて功あり>  
 ￥A 0 荷田￥A 1 春満￥A 3 アズママロ￥A 5 羽倉￥A 6 斎￥A 7 東麻呂￥B 2 元  
 文元年7月2日￥B 6 6 9￥C 0 信詮￥D 2 伏見稻荷山￥E 0 国学者￥E 1 万葉集  
 童蒙抄／伊勢物語童子問／出雲風土記考／齊明紀童謡考／偽類聚三代格考／春葉集  
 等￥E 7 国学四大人の一人￥F 0 明治16年2月贈正四位￥G 0 <////世世伏  
 見稻荷山の祠官> <////夙に国学を唱う>  
 ￥A 0 巨勢￥A 1 相見￥A 3 アイミ￥A 7 相覧￥B 4 延喜中￥C 0 金岡￥E 0 仏画  
 家￥G 0 <////采女正，讃岐目等に任せらる>

### 2.3 データシート各フィールド一覧

A0 人物姓	C8 (予備)
A1 人物名	D0 生地・出身地
A3 名ヨミ	D1 死没地
A4 出自	D2 活躍地
A5 別姓	D3 所属藩(主家)
A6 別称	D4 師
A7 各称の別表記	D5 交友
A8 (予備)	D6 予備
A9 別項目参照指示	E0 業種
B0 生年月日	E1 著書
B1 生年西暦	E2 著述
B3 没年西暦	E3 書写書(年/月/日)
B4 活躍年	E4 演目
B5 活躍年西暦	E5 結社・屋号
B6 享年	E6 その他の業績
B7 死因	E7 (予備)
B8 予備	F0 身分・階級
C0 父	H0 資料名
C1 母	I0 索引情報(歌集・古記録等)
C2 養父	J0 索引資料名
C3 養母	K0 宗派
C4 特記すべき祖先	K1 流派
C5 子供	G0 履歴
C6 兄弟姉妹	<年号 年/月/日/西暦/履歴>
C7 妻・夫	

(注) なお、本稿でふれた『国書総目録』記載の人名シソーラス中の字母で典拠外字母と判断する過程には、修訂版『大漢和辞典』、観智院本『類聚名義抄』等も参照した結果による字母判断結果も反映している。

【参考文献】

- [1] 『補訂版 国書総目録 著者別索引』 岩波書店 (1991年 1月18日発行)
- [2] 芳賀矢一・編『日本人名辞典』 思文閣出版 (大正 3年 9月29日発行／昭和47年7月10日複刻)
- [3] 『新訂増補 国史大系 尊卑分脈』全5巻 吉川弘文館
- [4] 『新訂増補 国史大系 公卿補任』全6巻 吉川弘文館
- [5] 尾崎雄二郎・都留春雄・西岡弘・山田勝美・山田俊雄／編『角川大字源』 角川書店 (1992年 2月10日発行)
- [6] 東洋学術研究所・編『大漢和辞典』 全14巻 大修館書店
- [7] 正宗教夫・編纂校訂『類聚名義抄』 全2巻 風間書房 (昭和56年 4月15日発行)

国文学研究資料館 助手  
National Institute of Japanese Literature Research Associate