

大規模漢字データの検証(姓,名ファイルを用いて)

田中康仁 日本ユニバック㈱

はじめに

漢字システムの導入は社会に省力化と効率化をもたらしている。しかしその影で漢字の入力とか、校正作業というあまりよくない作業や、単純作業を産み出している。このような分野の研究が十分行われないと、一方での合理化効果は他方の悪い作業の犠牲の上に成り立つことになる。これではせっかくの効果も大きな立場からなめると意味のないものになる。

漢字入力については色々な方式が検討されている。しかしこれら入力方式はスピードばかりに重点がおかれ、データの品質についてはあまり考えられていない面がある。漢字入力のスピードが上がっても校正作業が相変わらず手作業で膨大な時間と経費をかけていたのでは漢字入力のスピード向上もあまり意味がなくなってしまう。そこでここでは漢字データの誤りの自動抽出について考えてみることにする。

1 誤りデータについて

1-1 誤りデータの処理

誤りデータの処理については次のことに注目しなければならない。

- ① 誤りデータの検出(誤り自動抽出, ……………)
- ② 誤りデータの修復(誤りの発見にともない, これを修正する方法)

①の誤りデータの検出は数値データ, カナ文字データについて長いデータ・プロセッシングの歴史の中で方法論がみついているが, 漢字データについてはほとんど研究されていないのが実情であり, 誤りは人手によって見つけ出すものという考えが定着している。これでは漢字入力費用を安くするとか入力の処理時間を短縮すとか, 品質を向上させるといった問題が困難になってくる。

②の誤りデータの修復については漢字エディターや漢字ディスプレイ等の方法が, バッチ処理による文字修正などの研究, 実用化によりかなり見通しがついている。

ここでは特に①の誤りデータの検出について考えてみることにする。

1-2 大規模データの誤りの特徴

誤りを含んだデータは, データ・プロセスの各部門で問題を引き起こすものである。又いくら立派な処理システムが完成しても誤りデータが多ければ全体が無意味なものとなる。そこで誤りデータの特徴について調べてみる。

- ① 小規模のデータを取り扱っている実験室や研究室では通常考えられないようなことが大規模システムのデータの中には発生する。(誤りデータの量的拡大は質的变化をもたらす。)
- ② 誤りデータの除去を考えると数千件のデータの誤りは簡単に調べることができ, 修正することも考えられるが, 数10万件, 数100万件のデータでは誤り除去の方法について1つのサブ・システムとしての方法を考えなければならない。
- ③ 誤りデータは放置すれば累積する傾向があり, このため早めに除去しなければならない。
- ④ 毎日発生するデータでは数千件のデータでも誤りの検出, 修正の方法について大規模データの取扱いと同じように考えなければならない。

このような特徴から大規模データの誤りの調査は大きな問題であり、又システムティックに解決方法を考えなければならないということが判る。

2 誤りデータの自動抽出方法と実験結果

ある大規模の漢字姓名データと当社で持っている姓、名ファイルを照合し、アンマッチ・データを出し、その中から正しいデータと誤りデータを見つけ出す。

姓名マスターのカナ文字と漢字に一致するものは正しいとみなし、それ以外は誤りである可能性が強いものとしデータの分離して、その分離されたデータをさらに調べる。

例えば次のような誤りがある。

姓について

カナ	漢字	カナ	漢字
アイズ	今津 → (今→会)	トミヨ	トヨミ → (カナの誤り)
シシド	宍戸 → (宍→宍)	トシヤ	敏夫 → (トシヤ→トシオ)
マスダ	竹田 → (竹→升)	ヒデユキ	英芝 → (芝→之)
クリハラ	粟原 → (粟→栗)	フミコ	矢子 → (矢子→文子)

このような誤りは原票を調べるまでもなく誤りと判断のつくものである。

実験 (1)

実験結果

- ① 対象件数 : 1,499,000
- ② 調査日時 : 昭和52年12月 ~ 昭和53年6月
- ③ 実験に使用した姓名マスターの件数

姓	107,549件
名	150,930件
- ④ 実験結果

	全データ	アンマッチ・データ	エラー
姓	1,499,000件 (100%)	38,688件 (259%)	6,728件 (0.45%)
名	1,497,526件 (100%)	113,722件 (7.5%)	7,621件 (0.51%)

アンマッチ・データは姓、名マスターに含まれないデータである。この中より誤りデータを抽出したものがエラーである。

この誤りデータをさらに分析すると誤りの種類は次のようになっている。

- | | | | |
|--------------|----|------------|----|
| ① カナ文字の誤り | 6割 | ③ 原票の書きちがい | 2割 |
| ② 漢字入力 of 誤り | 2割 | | |

実 験 (2)

- ① 対 象 件 数 : 6,717,238 件
- ② 調 査 日 時 : 昭 和 5 3 年 9 月 ~ 昭 和 5 4 年 8 月
- ③ 実 験 に 使 用 し た 姓 名 マ ス タ ー の 件 数
 - 姓 127,431 件
 - 名 263,271 件
- ④ 実 験 結 果 :

	全データ	アンマッチ・データ	エラー
姓	6,717,238 件 (100%)	202,140 件 (3.0%)	81,680 件 (1.2%)
名	6,717,238 件 (100%)	411,353 件 (6.12%)	103,405 件 (1.53%)

アンマッチ・データはエラー・データを含んだ件数である。さらにアンマッチ・データからエラー・データを取り除くと次のようになる。

	正しいデータ	種 類
姓	120,460 件	28,986 件
名	307,948 件	175,213 件

姓名マスターに新しく見つかった種類を追加すると次のようになる。

姓	156,417 件	名	438,484 件
---	-----------	---	-----------

このように姓名マスターが充実することにより、次のデータ・チェックではアンマッチ・データの中にはエラー・データがほとんどで正常データは少ないという状態になるであろう

3 実験結果の考察

- これら2つの実験に使用したデータはいづれも外注会社で校正を行い納品したデータである。これからもわかるように、人手の校正がいかにも不十分であるかがよく判る。この意味からも機械的のチェックによる校正の偉力を感じさせられる。
- 実験(1)では150万件のアンマッチ・データ(姓,名 合計件数)が約15万件ありこの校正に約1ヶ月の労力を要した。
又実験(2)では670万件のアンマッチ・データ(姓,名 合計件数)が約61万件あり、この校正に約3ヶ月強の労力を要した。
今後は姓名マスターの充実により労力はもっと減るであろう。
- 実験(1)で得られたアンマッチ・データのうち正しいデータを姓・名マスターに追加したため実験(2)のアンマッチ・データの中で正しいデータの割合は減っている。
これはマスターを充実させることによりこのエラー・チェック・システムがより充実していくことを示している。

姓について

	アンマッチ・データ	アンマッチ・データ中の正しいデータ	エラー・データ
実験(1)	38,688件 (259%)	31,960件 (214%)	6,728件 (045%)
実験(2)	202,140件 (30%)	120,460件 (18%)	81,680件 (12%)

名について

	アンマッチ・データ	アンマッチ・データ中の正しいデータ	エラー・データ
実験(1)	113,722件 (7.5%)	106,101件 (698%)	7,621件 (052%)
実験(2)	411,353件 (612%)	307,948件 (459%)	103,405件 (153%)

実験(1)の姓について

$$\frac{\text{アンマッチ・データ中の正しいデータ}}{(\text{全データ}) - (\text{エラー・データ})} = \frac{214}{9955} = 0.0214$$

実験(2)の姓について

$$\frac{\text{アンマッチ・データ中の正しいデータ}}{(\text{全データ}) - (\text{エラー・データ})} = \frac{18}{988} = 0.0182$$

実験(1)と実験(2)を比較して正しいデータの分離が実験(2)で少なくなっていることが判る。

名についても姓と同じように上式で計算すると

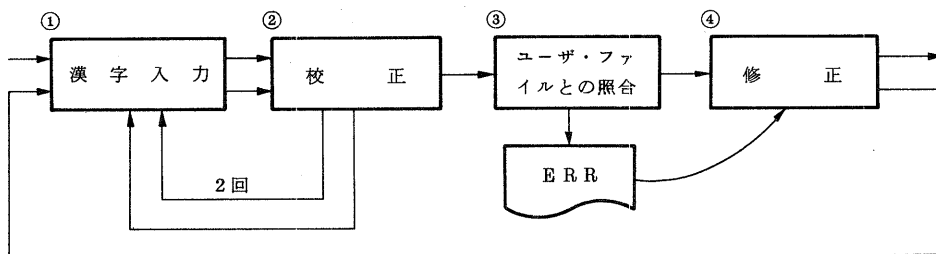
実験(1)の場合、実験(2)の場合で 0.0701, 0.0466 となり、これを比較して正しいデータの分離が実験(2)で少なくなっていることがわかる。

4 人手による校正と姓名ファイル照合による校正について

4-1 作業プロセス

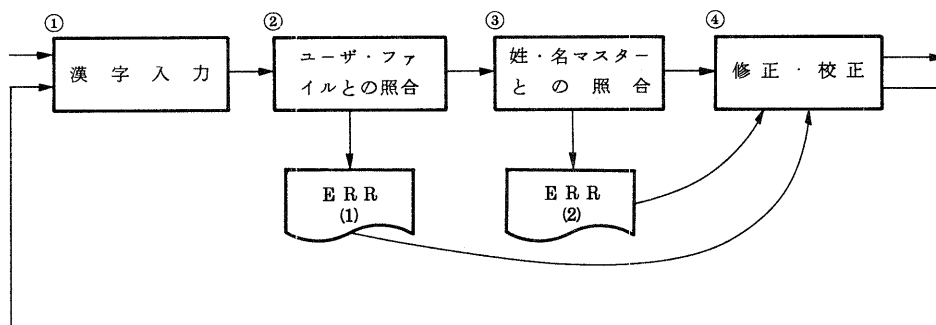
姓名ファイルと照合することにより実際の校正作業のプロセスがどのように変わるか比較してみると次のようになる。

人手による校正方式



- ① は原票を見ながら漢字の姓名を入力する。
- ② 校正は人手によって1文字1文字チェックする。
- ③ ユーザ・ファイルとの照合によりデータの抜け、重複をチェックする。この作業はコンピュータ処理で行う。
- ④ 修正伝票を起し再度漢字入力を行う。

姓名ファイルと照合する方式



- ① は原票を見ながら漢字の姓名を入力する。
- ② ユーザ・ファイルとの照合によりデータの抜け、重複をチェックする。この作業はコンピュータで行う。
- ③ 姓・名マスターとの照合により誤りらしきデータの抽出を行う。この作業はコンピュータで行う。
- ④ 修正伝票を起し再度漢字入力を行う。

修正伝票を発行するための誤りデータ、誤りデータらしきものが、すでにコンピュータ処理で見つけ出されているためこの作業は従来のものにくらべはるかに楽である。

校正もアンマッチ・データ特に注目しておけばよい。

4-2 人手による校正方式と姓名ファイルと照合する方式の特徴

次に、これら両方式の特徴について述べてみる。

人手による方式の特徴

- ① 校正者の能力によって、品質にむらがる。
- ② 多勢の人を管理しなければならない。
- ③ 少量のデータには向いているが多量のデータには向かない。
- ④ 校正者の作業場所、労働条件の改善などで多くの費用が必要、短期間の場合には割安である。

姓名ファイルと照合する方式の特徴

- ① 人手の介入が少ない。校正者の能力によって作業品質の高低がない。約1/10程度の人手で済む。
- ② 品質の保障が得られる。誤りの発見が早い。
- ③ 処理時間が早い。(大量のデータを処理すると)数万件～10万件程度以上の場合。
- ④ コンピュータ処理が中心である。
- ⑤ 大量データの校正には安く処理できる。
- ⑥ 校正作業にともなう精神的負担が少なくて済む。

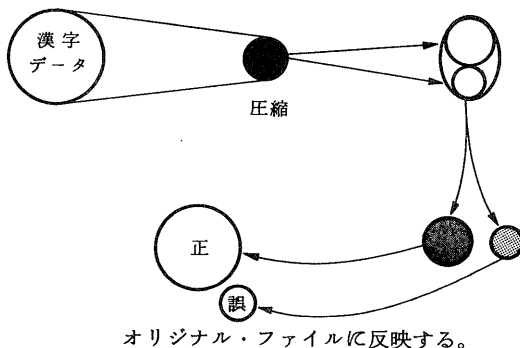
5 姓名ファイルを持たない場合のエラー・チェックについて

姓名ファイルを持っていない場合はこのようなエラー・チェックが行えないが、特にこのようなファイルが無い場合でも少し工夫をすることにより省力化した方法が考えられる。

次にその方法について述べる。

- ① オリジナル・データは10～20万件以上のデータとする。
- ② 姓、名にデータを分け同一の姓、名をおのおの1件としてその件数を加える。このようにするとファイルの圧縮がはかれる。
- ③ 頻度3件以上のグループと頻度2件以下のグループに分ける。これはデータのエラー・チェックを行いやすくするものである。前者は誤りが一般的に少なく、後者は誤りを比較的多く含んでいる。
- ④ エラー・チェックを行い誤ったデータは削除し正しいデータ・ファイルを作る。データ・チェックは2回以上行う。姓、名ファイルの小型ファイルが完成する。
- ⑤ オリジナル・ファイルと④で作られたファイルを照合し誤りデータだけを抽出する。

以上の内容を判りやすく図にすると次のようになる。



データ・チェックを行いやすくするため
頻度3件以上のグループと頻度2件以下の
グループに分けチェックする。

正しいデータ } に分離する。
誤りデータ }

姓名ファイルを持たない場合のエラーチェックとして漢字1字づつに姓名に使われるカナ文字を付けておき、機械的にこれらカナ文字を組合せユーザ・ファイルにあるカナ文字と一致するか否か調べ、一致すれば正しいとみなし、一致するものがなければ誤りらしいとして再度調べる方法が考えられる。

しかし、この方法は誤りらしいデータが多く抽出されるので問題が多いであろう。

6 おわりに

ここでは姓名という分野の校正作業について考えたが、この考え方を拡大し、日本語の分かち書き、辞書による照合を行えば一般文の校正にもこの考えはあてはまるであろう。

7 参考文献

〔1〕 田中康仁 “漢字データのチェック方法（名前のチェックについて）”

SYSTEMS 3, 4月号 1975年

〔2〕 田中康仁 “漢字データのチェック方法（名前のエラー・チェックについて）”

SYSTEMS 11, 12月号 1975年

添 付 資 料

シナニ	新谷	シモタ	トタ
シナ"	品田	シモタ"	下田尾
シノイ	椎野	シモトフタナ	下栃棚
シノサ"	篠塚	シモトヨトメ	下豊
シノタ"	篠崎	シモトリ	霜島
シハ"ウラ	柴原	シモハラ	下村
シハ"カ"キ	柴田	シモムラ	下平
シハ"ワホ	柴沢	シモナナ	下
シハ"ア	紫田	シモ	下
シハ"ア	芝崎	シモ	下
シハ"	芝登	シナア	車
シフト	渋谷	シナリア	捨田利
シフ"ヤシ	渋谷	シユウト"	清土
シフ"ナシ	澁谷	シユア"	柴田
シフ"ヤ	澁江	シヨウハ"ナシ	莊林
シマイチ	島	シライ	白戸
シマサ"キ	島	シライ	白水
シマシ	島	シライ	白田
シマセ	島崎	シラウカ	反坂
シマセ	嶋瀬		
シマア	島		
シマア"	舌間		
シマトモ	島本		

実験(1)のエラー・リスト・サンプル (姓)

サア"イ	男	サフ	サク
サア"キ	忠嘉	サンクエモン	三右エ門
サア"コウ	貞子	サンコ	ミ子
サア"コ	サタ子	サンソ"ウ	三
サア"コ	えだ子	サ"エモン	左エ門
サア"ノリ	完訓	サ"コ	和子
サア"ノ	サダソ	シキチ	末吉
サア"	サタ	シイカ	栄孝
サチ	サチコ	シカ	さか
サチコ	智子		
サチコ	致子	シケ"オ	宏
サチコ	幸江	シケ"オ	茂夫
サフコ	幸子	シケ"オ	茂雄
サフ	サワ	シケ"キ	茂樹
サトシ		シケ"コ	シケ子
サトミ	つとめ	シケ"コ	しげる
サトミ	覚己	シケ"コ	小志げ
サト	さそ	シケ"コ	茂
サノ	さと	シケ"シサ	重久
サフ"ロウ	信也	シケ"シサ	繁一
サフ"ロウ	喜三郎	シケ"タカ	孝
サフ"ロウ	己三郎	シケ"ノフ	信重
サフ"ロウ	正一		
サフ"ロウ	三武郎		

実験(1)のエラー・リスト・サンプル (名)