

低品質文字を用いた標準パターン辞書構築による手書き署名認識法

鎌形 周平 鈴木 雅人 北越 大輔

東京工業高等専門学校 情報工学科 知識情報研究室

1. はじめに

近年、様々な地域で署名活動が行われている。収集した署名に法的効力を持たせるには、収集・データ化が容易なネット署名などではなく、紙媒体で氏名・住所を収集しなければならない。署名と住民台帳との照合は1ヶ月程度の期間を要するが、手書き署名を自動認識できれば、この期間を大幅に短縮できると考えられる。

これまでの手書き文字認識の研究では、筆者の癖を吸収して認識精度を改善する手法[1][2]や、住所固有の性質を用いた手書き住所認識法[3][4]なども提案されているが、これらの手法を用いても、我々が普段書き記すような文字品質で書かれた手書き署名の認識には限界がある。

本研究では、我々が普段書き記すような品質の手書き文字を用いて認識用辞書を再構築すると共に、署名特有の特徴を用いて認識候補字種を限定する誤認識訂正法を提案し、手書き署名の高精度な認識を試みる。

2. 署名認識システム

本節では、提案する手書き署名認識システムの概要について述べる。提案する署名認識手法は、普段我々が書く文字を認識するための標準パターン辞書再構築法と、署名特有の特徴に着目した誤認識訂正法との2つの要素からなる。

2.1 署名認識システムの概要

提案する署名認識システム(図1)は、氏名と住所の分割・個別文字認識・誤認識訂正・比較の4つの処理で構成される。スキャナによって手書き署名を取り込み、氏名部と住所部に分割する。次に文字の切り出し・前処理を行った後に方向線素特徴量を抽出し、標準パターンとのマハラノビス距離[2]を計算して認識する。更に認識された文字列の中で、氏名部の姓、住所部の都道府県および市区町村を特定し、文字数から

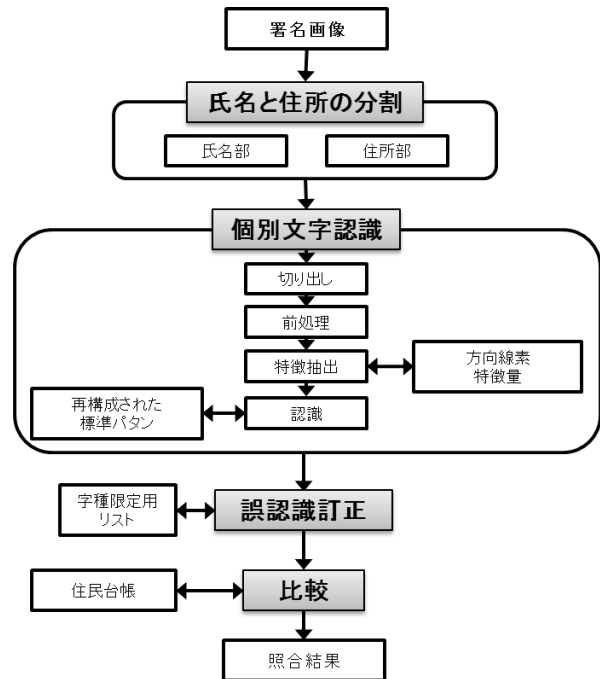


図1 システムの概要

利用される字種を限定し誤認識訂正を行う。最後に、訂正された文字列を住民台帳と比較し該当するデータがあるか判定し、実在しない氏名・住所は誤認識候補リストに追加し、元画像をユーザが見て訂正または無効を判断する。

2.2 標準パターン辞書再構築法

文字認識の精度が下がる一因として、実際の手書き文字と従来のETL9をベースとした標準パターン辞書との特徴量分布が大きく乖離している問題があげられる。本研究では、手書き署名を高精度に認識するため、我々が普段書き記すような品質の手書き文字を大量に収集し、標準パターン辞書を再構築する。辞書に加える文字データは、教育漢字1006字種、ひらがな・カタカナ各71字種、教育漢字以外の苗字で利用される漢字351字種と、名前で利用される漢字107字種をあわせた1606字種とした。ETL9に含まれない字種については新たに収集した文字のみで辞書を構成する。若者に顕著な丸文字や止め・はらいのない漢字の癖を排除するよう指示を与え、

A Recognition Algorithm of Handwritten Sign based on Standard pattern Construction Method Using Low Quality Character Patterns

Kamagata Shuhei, Suzuki Masato, and Kitakosi Daisuke,
Tokyo National College of Technology Department
of Computer Science, Laboratory of knowledge
information

16～22歳までの学生100名に協力してもらい、1字種あたり50サンプルの文字を収集した。

2.3 誤認識訂正法

標準パターン辞書の再構築により、認識精度の改善を期待できるが、認識精度を高めるためには、個別文字に対する認識後の誤認識訂正が必要である。氏名は姓・名、住所は都道府県・市区町村・その他にそれぞれ分割できる。姓・都道府県・市区町村で使われる字種は限られており、それぞれの文字数が決まると使用される字種は限定される。氏名は姓と名をあらかじめ分割して記入する形式にし、都道府県・市区町村をキーワードにして、各要素の文字数をカウントし、使用される字種を限定することで認識候補を減らすことによって誤認識の訂正を行う。

3. 署名認識実験

提案手法の有効性を確認するため、収集した手書き文字を用いて辞書を再構築し、認識精度評価実験および誤認識訂正評価実験を行った。

16～22歳の学生から収集した手書き文字15件分のサンプルデータを用いて辞書を再構築し、署名認識実験を行った。従来の個別文字認識では、方向線素特徴量およびマハラノビス距離を用いているが、今回収集した文字サンプル数が15であることから196次元特徴量を用いたマハラノビス距離の計算ができない字種がでてくるため、ユークリッド距離で代用した。従来の標準パターン辞書と再構築した辞書で行った実験結果を表1に示す。

表1 個別文字認識実験結果

	第1位認識率	第10位累積認識率
従来の辞書	71.80%	92.60%
再構成した辞書	73.15%	93.28%

1位・10位候補ともに辞書の再構築によって認識精度は改善されている。今回の実験では数ポイントの改善にとどまっているが収集データ数を増やすことによって改善率の向上が見込まれる。特に第10位累積認識率の改善が見込めると、個別文字認識後の誤認識訂正アルゴリズムを併用することにより、署名全体の認識精度を大きく改善できるものと期待できる。尚、辞書の再構築によっても正しい認識候補が得られない文字もあった。これは、我々が普段書く文字の分布が正規分布から大きくずれていることにより、従来から用いられている識別関数では識別能力に限界があるからだと考えられる。

署名サンプルを ETL9 のみで構成した標準パターン辞書とマハラノビス距離によって認識を行った結果のうち、姓、住所の都道府県・市区町村460文字に対し、誤認識訂正アルゴリズムを適用した例を図2に示す。表2は誤認識訂正の様子を示したものであるが、1～2文字目は姓をあらわしている「清変」という姓は存在しないので「清友」と誤認識訂正できる。また、7・8文字目は市名の一部であることを使うと認識候補はそれぞれ「分」「寺」のみになり誤認識訂正を正しく行うことができる。その結果、誤認識訂正前の第1位認識率が75%だったのに対し、訂正後は85.65%まで改善された。



図2 実際に収集した署名例

表2 誤認識訂正例

訂正前	清	変	東	京	都	国	堺	争	市
訂正後	清	友	東	京	都	国	分	寺	市
1位	清	変	東	京	都	国	堺	争	市
2位	清	友	東	崇	部	国	妹	等	帝
3位	清	皮	更	宗	郡	国	分	舞	肅
4位	清	渡	恵	東	翻	国	眺	寺	帯

4. むすび

本稿では、我々が普段書き記す品質の手書き文字を用いた認識用辞書の再構築と、署名特有の特徴を用いた誤認識訂正法を提案した。各手法の評価実験を行い、それぞれで一定の認識率改善が見込めることを確認した。今後は、再構築された辞書に適した新たな認識手法について検討を行う予定である。

尚、本研究は科学研究費補助金(基盤研究(C)課題番号22500170)の助成によるものである。

参考文献

- [1] 松石 “文章表現の癖抽出に基づく手書き文章認識の後処理方式に関する研究” 東京工業高等専門学校卒業論文, 2012.
- [2] 郭, ほか “余弦整形変換を用いた手書き文字認識アルゴリズム”, 信学論 J76-D-II, no. 4, pp. 835-842, 1993.
- [3] 大久保 “郵便番号情報を併用した手書きあて名認識に関する研究” 東京工業高等専門学校卒業論文, 2010.
- [4] 鈴木, ほか “候補字種の動的抽出を用いた手書きあて名の2段階認識手法”, 信学論 J82-D-II, no. 11, pp. 1895-1902, 1999.