

日本がん登録協議会
第31回学術集会
ポスター演題 R-1-7

画像認識技術を用いた『外字』のユニコード文字探索支援

田淵 健
東京都立駒込病院／東京都がん登録室

【背景】 文字コードについて

- 1978年JIS第1-2水準(JISX0208の前身)制定
- 1980年住民基本台帳に漢字オンラインシステム導入
- 2000年JISX0213(JIS第3-4水準を含む)制定
- 2000年Japan重点項目⇒汎用電子情報交換環境整備
- 2003年JISX0213がUnicodeに正式対応
- 2002年住民基本台帳ネットワーク稼働
- JISX0221(国際符号化文字集合(UCS)対応)を独自拡張
- 2010年文字情報基礎標準:文字情報基礎文字
- 文字情報基礎漢字58712種(2019年5月版58543種)
- 行政で用いられている文字を整理統合⇒相互運用性向上
- 2013年Windows 8+Office 13でUnicode完全対応

【目的】

- 外部照会業務において、人手による判断が必要で最大の稼働段階である外字照会を人手やコストをかけるに減らす

【方法】 文字情報基礎(MJ)文字

- 独立行政法人情報処理推進機構(IPA)文字情報基礎標準で整備された文字(漢字)
- MJ文字情報一覧表Ver.006.01:58662文字
- MJ文字情報一覧表繁体版名簿 Ver.002.01:299文字
- MJ漢字マップ
- MJ文字とJIS X 0213(JIS第1-4水準)文字との対応関係
- IPAmj明碼フォント
- MJ文字情報一覧表標準文字を基に(銀行MJ文字は実体版名を含む;他、一部MJ文字は実装されていない)
- IPA(経産省政策実態機関)から一般社団法人文字情報推進協議会(CITPC)に民間移管(2020年)

【方法】 文字フォントの画像化

- 対象の文字フォント及び文字コード
- EUDCフォントとIPAmj明碼フォントが実装された文字コード
- 各フォントを1文字ずつHTML表示し画像ファイルとして出力
- 類似度判定率向上のため、画像のサイズや解像度を調整
- 手順はPythonプログラミングにて実装
- PythonのPILライブラリにおけるImageFont機能を用いることでフォントの画像化は容易に可逆であるが、Unicodeの紐帯を用いる形式のIVS/IVD文字は字形を忠実に出力できないという致命的な欠陥があるため使えない

【結果】 処理時間

- フォントの画像化に要する時間
- EUDCフォントとIPAmj明碼フォント併せて6万字について30分程度で処理可能
- 類似度判定に要する時間
- EUDCフォント1文字につきIPAmj明碼フォントを検索する時間は約5分以内
- EUDCフォント全ての検索を行うと、約4.5日要する
- 実際のデータに出現するEUDCフォントは最大300文字程度であるため実際上は10時間程度で処理出来た

【考察】 人手作業による漢字置き換え

- 地域がん登録の手引き改訂第5版詳細版では第3水準漢字を第1-2水準に対応させる作業を例示して示している
- この対応に一般法則はないので、各登録室で人手作業による**経時的作業**を行う必要があるが、**作業期間のゆれ**が大きくなる
- ⇒当登録室では文字情報基礎標準の短見と異体字データベースを組合せて網羅的リストを作成

【考察】 提案手法の有用性及今後の課題

- 本手法では多くの文字で高精度で外字のUnicode文字を特定することが可能であった
- しかし、現時点では、完全自動で決定することは出来ず、**人手による最終判断は必要である**
- 人の判断のための**支援技術**としては有用であり、**労力軽減と作業時間の短縮につながる**
- 人間の作業は一定の比率でエラーが発生するので、**エラー率を下げる効果もあると考える**
- 今後の課題としては、**深層学習による画像認識技術を用い、人間以上の判定精度を実現したい**

【利益相反】

【学術集名】
日本がん登録協議会第31回学術集会

【著頭演者氏名】
田淵 健(東京都立駒込病院)

【COI開示】
当演題に関し、開示すべきCOIはありません

【背景】 市区町村が使用する外字

- 市区町村が使用する外字の実態調査
- 2011年度総務省(1,386市区町村)
- 1,166,536外字(1,386市区町村任意提出)
- 自治体毎に独自外字を作成しているため膨大な数にのぼる
- 一覧表が作成可能で、当時の外字マッピング問題は全ての外字を調査できず、当時の事務で作成する文字を選択せざるを得なかった
- 2011年10月現在市区町村数は1,727⇒協力は約8割
- 文字情報基礎文字との対比
- 2011年時点では実体版名簿

【方法】 がん登録基礎業務としての位置づけ

- 本手法の開発は、がん登録個人情報には用いていないため、**安全管理措置**の適用外である
- 本手法で用いている外字フォントは提供元の基幹業務で用いられているが、外字そのものは公開されていないので行政機関に提供可能な可能性がある
- 個々の外字について具体的に言及することは出来ない
- 本演題で述べる様な事象は一般的には体系的に扱われていないようであるが当登録室ではがん登録基礎業務として体系化している
- 本開発に要した余剰の人員及び資金はない

【方法】 外字フォント

- Unicode私用領域(外字領域)
- E000~F8FF(6400文字)にマッピングされる
- 現在の一般的なPC(Windows 8以降)で利用可能な領域
- Shift-JIS以外の仕様で文字を全て実装するには不足するので、独自の書体として、文字を記述している
- 補助私用領域(F8FF0~FFFFF)、補助私用領域(10FF80~10FFFF)は通常の環境では使われない
- EUDCフォント
- Unicode私用領域に定義された外字をWindows-PCで外字を表記するためのフォント(環境依存:市区町村毎異なる)
- 登録室で受領するデータに外字が含まれている場合、EUDCフォントを併せて受領する

【方法】 画像化されたフォントの類似度

- 外字と文字情報基礎文字との対応づけの方法
- 特徴量マッピングによる類似度判定を行う
- フォントを画像化した画像データから**特徴量**を抽出
- 特徴量マッピングは角や線の勾配などの特徴的な部分から抽出するアルゴリズムである
- プログラミング言語PythonのライブラリOpenCV(cv2)には**AKAZE**という特徴量抽出が実装されている
- 類似度判定アルゴリズム
- Unicode私用領域(Brute-force matcher)
- K-Nearest Neighbor algorithm(KNN)
- いずれもcv2に実装

【考察】 コンピュータリソースと漢字表記

- 現時点で、日本語は、中国語文化圏を除いて**自国語表記**に漢字を用いる唯一の言語である
- 日本では、中国よりもコンピュータ導入が早かったため、当時の**技術的なリソース**で漢字を表記するため様々な工夫が必要であった(文字ROM)
- 第2次世界大戦後の漢字取崩しの影響が残っており、コンピュータで用いる漢字は制限もゆるいという風潮
- コンピュータの急速な進歩で漢字表現力が高まり、多数の漢字が文字コードとして標準化されている
- 中国の経済成長の急上昇とともに国際的な文字コード規格に中国の影響が強いと見られることの影響

【考察】 人間とコンピュータの判断の違い

- コンピュータは自ら学習する(機械学習など)か否かは別としてルールと知識に基づいて動く
- 伝統的なコンピュータ運用方法は、コンピュータにルールと知識を事前に指示する必要がある
- 我々が標準的だと認めている作業手順は、**明文化**されているが、コンピュータにとっては**指示が不十分な情報**
- 人の作業工程や手順のレンジ、がん登録標準業務手順
- われわれには、(時には無意識的に)**前提知識を用いて直観的に判断**している

【結語】

- 外字文字は馴染みが薄いため人手作業による類似判定は非常に骨の折れる仕事である
- 本提案手法は、人手作業を大幅に軽減する効果的な支援技術と言える

【前提】 電子情報における文字情報

- 電子情報は言うまでもなくデジタル化されている
- 電子情報における文字情報はコンピュータ内部では文字コード(数値)として扱う
- 英数字文字は1バイト(256通り)で表現出来る
- 漢字や仮名は少なくとも2バイト(65536=256^2通り)で表現出来る
- デスクワークの大半は既存の文字情報の操作
- 新規操作することは少なく既存の漢字やフォントの処理
- 作業方法はコピー&ペースト、検索、置換等
- しかも処理内容の大半はルールベース
- コンピュータが処理すれば100万行程度であればほぼ一瞬

【背景】 外部照会における外字の実態

- 当登録室で受けた都内2区市の外部照会利用リストが**検診診療管理目的**はいずれも**独自外字**が用いられていた
- 他の市区町村も同様な実情とされ、**外字使用実態報告書10年経過後も外字使用の実態は変わっていない**と言える
- 同報告書の指摘「システム連携においても、外字データの連携のための**間接作業**や**文字コードの変換テーブル作成の作業**が発生し、**作業負担や移行コストが増加**」
- 外部照会を用いた情報の利活用は足踏になる可能性

【方法】 漢字正規化(Shift-JIS化)

- 全国がん登録で用いる漢字はShift-JIS漢字で表記
- 実際にはCP932(拡張Shift-JIS)も使えるものが多い
- 漢字情報基礎標準は**標準化**に意図を込めて明示
- Shift-JIS漢字標準文字⇒単体から順次CP932漢字
- 漢字情報の表現・解釈の乖と自動処理の低い精度
- 当登録室における漢字正規化手法
- 漢字は「文字情報基礎標準MJ文字」(6万字)の範囲で表記
- IPAmjによるMJ標準マップ(JIS-1-4水準漢字への対応表)
- 異体字データベースを用いたShift-JIS漢字への置き換え
- 照会には、**標準漢字(同一漢字)**を定義
- MJ標準マップ、異体字と類似文字によって定義

【方法】 文字の類似度を判定する方法

- 現在用いられている外字のUnicode (或いは文字情報基礎(MJ)文字)への**マッピングルールは分かっていない**
- ルールが分かっているのは、外字を利用している機関(市区町村や病院、企業等)では既にUnicodeに置き換えている
- 外字が用いられたデータを受領したら、登録室で**外字とMJ文字との対応関係**を見いだす必要がある
- データが少なれば人手作業で行うことも可能
- 異体と標準が混在してデータに混在することも可能
- 人が文字を判断するのは**概ねMJ文字の形で判定**
- 文字フォントを抽出して**画像化**して**類似度判定**を行うソフトウェア
- 人が文字を判断する場合は意味や部首の理解も必要

【結果】 候補決定の可能性

- 類似度判定はマッピング精度Ratioやマッピングした特徴量の数に影響を受ける
- 精度や特徴量の数を増加すると処理時間は大幅に延長する
- 完全な字形一致文字とデジナリ異体文字(併せて全体の約1/4)はマッピング精度を極めて高精度に決定する(Ratio=>0.99)ことで、ほぼ正確が得られた
- 類似文字(約半数)については、複数の候補を示し、**目視確認が必要**であった
- 判定不可能文字や記号(MJ文字に該当文字が存在しない)の類似度判定は意味をなさない
- 外部照会リストには含まれていない実用には問題がない

【考察】 漢字知識の共通認識

- 漢字知識の確固たる**共通認識**はない(個人差大きい)
- 標準漢字制があったが最近では漢字表記も多様化
- 照会業務における**目視判定**は**漢字知識**についての**共通認識があること**を前提を暗黙の了解としている
- 漢字知識は、院内がん登録業務者研修によるがん登録知識などは異なり、**社会通念**のうちに運用している
- 目視判定エラーの一部はこの**共通認識のずれ**によるもので、一定比率は発生しうる
- 漢字表記ルールを全て書き出すことは一般には行わないので職務指導になじまない

【考察】 JIS規格の字形・字体の変遷

- JIS規格のバージョンによって字形・字体が変遷
- JISX0208(概ね第1-2水準漢字)
- 1978年制定、1983年改定、1998年改定
- JISX0213(第1-4水準漢字を含む)
- 2000年制定、2004年改定
- 地名に關係が深い文字の変遷例
- 自治体の電算化時期によって外字対応が異なる
- 例:JISX0208(1978)⇒JISX0208(1983)⇒JISX0213(2004)
- 例:JISX0208(1978)⇒JISX0213(2004)
- 現在は、IVS/IVDで表記可能