

**日本がん登録協議会
第31回学術集会
ポスター演題 R-1-4**

リレーショナルデータベースシステムを用いない汎用的照合システムの開発

田淵 健¹⁾, 柿崎 裕則²⁾, 下山 達¹⁾
1) 東京都立駒込病院, 2) 東京都がん登録室

【目的】

- ① 届出業務においては、全国がん登録システムとは別にデータ加工レベルで、事前に照合上の不整合エラーを検出し、**職人手を介さずに**、照合調査リスト出力の定型処理を行い、また、照合調査回答結果の反映を自動処理を行う、全国がん登録システムに移行する
- ② 外部照合業務においては、全国がん登録システムを用いずにデータ加工レベルで前処理に引き継いで一連の**工種**として照合業務を行う
- ③ 以上の開発経緯を本演題では数式やプログラミングを用いずに紹介する



【方法】 スコアリングの設計

- ① **同一性に基づくスコアリング**
 - ① 全国がん登録システムに採用されている方法に準拠(第24回全国がん登録協議会1-14新野(国が), 2015)
 - ② 照合情報の設定、重み付けスコアリング
- ② **類似性に基づくスコアリング**
 - ① 日付データ: 1-2位の数字のミスタブや元号選択ミス等を考慮した日付の類似度判定を行った
 - ② 氏名: 氏名中の文字列を英数字・漢字の代表文字に変換した上、文字列の類似度判定を行った
 - ③ カガ氏名: 文字列の類似度判定を行った
 - ④ 文字列類似判定にはpythonのdifflibライブラリを使用



【結果】 照合精度の向上の可能性

- ① 照合精度の向上とは、**false-match**や**false-mismatch**を減らすことである
- ② 照合精度については本来は比較検討を行うべきだが管理にむづかしく、比較検討に不可能であったため、事例ベースでの検証があるところから見たい
- ③ 人手作業による同判定ミスは大幅に減少した
- ④ 目視判定で見過した同判定が可能となった
- ⑤ 類似判定を導入したことで大きい
- ⑥ アルゴリズムの例外処理が完全かどうかは断定できないため、徹底的なエラー発生する危険性がある

【考察】 ルールベースの仕事

- ① がん登録はデータマネジメントの一種であり、基本は**ルールベースの仕事**であり、自動化が可能ではないかと思われる仕事
- ② 実際には、多くのデータマネジメントの仕事同様、標準化文書には記載されていない**暗黙の了解**が多数含まれるため、一般的には完全自動化は困難とされる
- ③ 現実のがん登録業務の人手作業を観察していると、**形式的な作業に終結し、人手作業に期待される判断力が生かされていない**
- ④ 判断を要する作業が形態化して、戻った判断が増加する
- ⑤ データ規模が大きい東京都の傾向も似ている

【考察】 データ組合せ処理の方法

- ① データの組合せ処理は、理想的には**リレーショナルデータベースのインテグレーション**を用いることで高速化が可能
- ② 全国がん登録のデータ加工においてリレーショナルデータベースを構築することは、安全管理規程の一種ファイルでの処理とはいわがたく、かつ、全国がん登録データベースのクローンが発生する**正當性の問題**も生じる
- ③ テキストベース(CSVファイル)での適量のデータ組み合わせ処理を行うと、単純な一致でも時間がかかり、類似検索に必要な処理時間は現実的ではない
- ④ 本提案手法は**類似検索を行いながら大規模な検索結果を絞り込む**

【利益相反】

【学術集名】
日本がん登録協議会第31回学術集会

【筆頭講演者氏名】
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

【方法】 全国がん登録業務改善の一環

- ① 全国がん登録標準手順と**安全管理措置**に従いデータ加工業務改善の一環として実施
- ② 厳密な比較検討は行っていない(ほぼ不可能)
- ③ 外部非接触東京都がん登録室内LANシステムで実施
- ④ データ加工用いないデータは、全てテキストファイル(CSVファイルないしExcelファイル)である
- ⑤ 定期的にはもちろん一時的にも**リレーショナルデータベース(RDBMS)の構築**は行わない
- ⑥ 全国がん登録データベースの正当性を侵害していない
- ⑦ 本開発に特化した人員及び資金は用いていない



【結果】 一貫工程化による効率化

- ① 前処理・後処理などの周辺業務や一連の業務との**一貫工程化**を実現
- ② 各工程での予期しない不整合発生や再発を防止するため、**工種間チェックポイント**を設置(工程管理表)
- ③ 各工程は**完全自動化**して、シームレスに自動処理を進めるため、がん登録個人情報の処理能力は不要**ペーパーレス**の実現
- ④ 照合上の不整合を事前に検出することで、届出票の**票内エラー**と票間エラーを**一元的な「照合調査」**の実施が可能となった

【結果】 照合に係る業務時間の大幅短縮

- ① 単純な照合時間はシステム環境と照合対象に依存
- ② 以下に示す照合時間は参考値(照合対象によって異なるシステム環境で実施している「単純な比較や一貫化は困難」)
- ③ 届出票の事前照合時間は**1分/500件**
- ④ ファイル送達等の時間は一定であり、単純環境下で約50秒
- ⑤ がん検診精度管理外部照合時間は**7-10分/10万件**
- ⑥ 照合対象がカバーする種数が多いほど照合時間がかかる
- ⑦ がん検診精度管理では照合対象の住所が同一自治体在住者が対象である
- ⑧ 照合そのものにかかる時間は、全国がん登録システム(初代)よりも所要時間を**1桁短縮**

【考察】 客観的な比較検討の困難性

- ① 厳密には、業務手法についても**比較試験**を行うことがエビデンスに基づく方法であるとは言えない
- ② しかし、特に東京都がん登録室のようにデータ規模の多い場合、現実には実施しては業務を標準的な枠組みの中で、**中断することなく適量で継続しながら**、業務改善を進める形式において、客観的な比較を行うことは困難である

【考察】 票間エラー判定自動化の困難度

- ① 当登録室では標準DBSでの運用困難状況を目的として業務の自動化を構築してきた
- ② 票内エラー判定の自動化に比べて**票間エラー判定(照合)**の自動化は困難度が高かった
- ③ 票間エラーは、ルールの組み合わせが複雑のほかに、データの組み合わせが複雑になる
- ④ ルールの組み合わせについては、複雑さは全部書き下すことは不可能ではない
- ⑤ データ規模Nに対して、最大N²の処理量が必要とする
- ⑥ データ規模が増加すると処理量は急激に増え、膨大なソースと時間を必要とし、処理が不能に陥る

【考察】 今後の課題

- ① 本提案手法のがん登録データ以外にも**データ組合せ処理**を行うアルゴリズムとして一般化が可能
- ② ルールと知識をコンピュータに指示するルールベースの手法は、それ自体が労働集約的であり、ルールや知識の変化に対応できない
- ③ プログラマが常に新しい知識にキャッチアップし続けるのは困難であるので、**深層学習**の手法を取り入れて、専門家の判断レベルの指示をコンピュータだけで実現する必要がある

【背景】

- ① 個人同定を要するデータベースにおいて、照合工程が、データの品質や統計精度を左右する
- ② 全国がん登録システムは本邦ががん登録史上初めて**がん登録業務の体系的な自動化**が導入され、かつ、照合精度や速度が向上し、業務改善の一環として**実行可能性**も担保されたが、次のような改善点がある
- ③ 登録業務においては、照合によって判明した**不整合**を電子データとして出力する機能が**整備されておらず**、照合調査に人手と時間がかかる
- ④ 情報の提供における外部照合は、工程が煩雑で時間がかかり、他の業務への配慮が必要

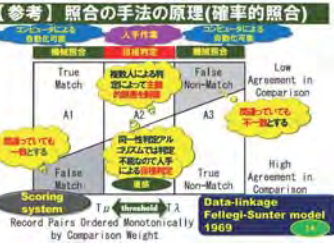


【方法】 照合アルゴリズム

1. 同一性の定義
 - ① 氏名: 氏名中の文字列(1-7文字)は、住所については町丁目へまで住所コード、生年月日と性別はそのままで一致を定義し、同一性判定はそれぞれ別々に実施して判定
2. 類似性の定義
 - ① テキストは標準化して類似性を判定。生年月日は元号選択ミスやアルファベット等を考慮した類似性を判定
3. 照合精度正規化: 複式変換・データクレンジング、同一性定義、区市町村データでは住所レベルのコード変換が前提
4. 照合項目(氏名、生年月日、住所、性別等)をペア順に全国がん登録データと照合対象リストの共通部分を抽出し、ソートする
5. 上記ペア毎に、共通項目には同一性判定を適用、それ以外の項目には、類似判定を適用し、一致しなかったペアはスコアを算出し、閾値以上を同一とすることで照合結果を判定

【結果】 人手作業の軽減

- ① 照合における**目視判定**は労働集約的な人手作業であるが、提案手法では大幅な自動化によってこの様な作業を殆ど回避できるようになった
- ② 周辺業務との一貫工程が可能となったことで、**転記作業(コピペ)**等の人手作業が不要となった
- ③ 東京都のデータ規模にて、他のがん登録業務を兼ねても、**作業員当り1名分の人手作業に節約された**
- ④ 作業担当者は時間外業務を全く行っていない
- ⑤ がん登録業務の全国がん登録システムでのデータ処理に必要な時間は**2日**



【考察】 全国がん登録システム

- ① 全国がん登録システムの照合システムは、地域がん登録標準DBSの照合システムに比べ、機械照合の導入、照合速度の短縮、照合精度の向上に図られた
- ② 東京都では標準DBSの運用は専任で実施したが、全国がん登録システムでは実行可能性が担保された
- ③ 全国がん登録システムであっても**労働集約的業務**であり、更に、**必要人員数はデータ規模に比例関係にある人員が必要**である
- ④ 人員数が一定以上を越えると業務品質の維持が困難
- ⑤ どうしても職員維持が生じ、照合のように**相互関係が重要な業務**では、品質が品質を促進することが生じる

【考察】 データ組合せに対する処理量

- ① 提案手法ではデータ規模Nの組合せの**N²の処理量**をいかに減らすかが課題である
- ② 一致も類似もない組合せ計算は無駄な負荷となるので、データと照合対象の共通部分だけに絞り込んだ
- ③ システムでのデータ処理に負荷がかかるのは、**一部重複(類似)のデータ**であるので、上記共通部分のソートを行った後、重複部分のみについて計算する
- ④ 照合指標の組合せの組について重複セット(a回)行うとし、ファイル読み出しなど一定処理分を加えると、**処理量はaNlogN+nb**となり、N²より大幅に削減される

【結語】

- ① RDBMSを用いず、類似検索を考慮しつつ**実用的な速度**で照合を行う手法を開発した
- ② 照合に当たって、全ての組合せを網羅的に検証するのではなく、**適した条件で一致条件データ**を抽出・ソートした特定の一致データの近傍だけで時間のかかる類似計算処理を行うことにより、大幅な高速化が実現できた