

平成 28 年度厚生労働科学研究補助金

(臨床研究等 ICT 基盤構築研究事業) 総括研究報告書

電子カルテ情報をセマンティクス (意味・内容) の 標準化により分析可能なデータに変換するための研究

堀口 裕正 国立病院機構本部総合研究センター 診療情報分析部 副部長

岡田 千春 国立病院機構本部 企画役

狩野 芳伸 静岡大学情報学部行動情報学科 准教授

森田 瑞樹 岡山大学大学院医歯薬学総合研究科 准教授

奥村 貴史 国立保健医療科学院研究情報支援研究センター 特命上席主任研究官

研究要旨

初年度においては、我々が日本語における医療用自然言語処理の研究コミュニティを形成し研究に取り組んで来た標準化技術を実カルテへと適用することで、カルテからの情報抽出の自動化に向けた予備的な検証を行うことを計画した。

研究代表者堀口及び分担研究者岡田は、国立病院機構本部との調整を中心とした基盤構築を行った。まず、NDCA データの研究利用に向け、倫理審査申請に加えて、内部規定にて定められている内部委員会の調整を図った。また、閲覧・解析に特化した自然言語処理用の研究基盤の構築を行った。研究基盤の概念図は図に示したとおりで、セキュリティを維持しつつ、空間的制約をなるべく少なく研究が進められるようなものになっている。

研究分担奥村は、臨床的なニーズを自然言語処理における個別技術へと橋渡しする役割を担った。具体的には、臨床医側より退院サマリの自動生成に求められる要件定義を進めるとともに、先行研究の整理を行い今後の研究アプローチの策定を行った。さらに、今後の研究に役立てられる入院カルテ・退院サマリの高品質な個人情報を含まない模擬のデータセットを構築した。

研究分担狩野・森田は、上記の実データ・テスト用データ双方を活用し、自然言語処理の医療テキストへの適用を進めた。狩野は、時系列で蓄積してく入院カルテデータを対象として、既存の自然言語処理ツールによる処理性能と今後の改良に向けた課題抽出を図った。森田は、医師が要約した退院サマリデータを対象として、医師の記載する退院サマリの定量的・定性的な特徴の把握を図った。この知見は、今後、入院カルテの自動要約技術

の研究に際した精度管理に役立てられる。

今年度の研究においては、NDCA データの利用に際した倫理審査の関係で、実カルテを対象とした解析は 3 月以降に限られる制約が生じた。しかしながら、高品質なテストデータの確保が可能であったことから、研究のコアである医療用自然言語処理部分の検討はほぼ計画通りに進めることが出来た。電子カルテの自由記載部分を自動解析する多施設構成での大規模データを対象とした研究としては、本邦初の試みとなる

医療用情報システムの研究開発においては、医療現場に直接の恩恵が及ばないゴールが設定されることで、研究開発が現場のニーズから乖離するとともに、継続した開発投資に繋がらない悪循環が往々にして生じてきた。本研究提案は、医療現場における負担軽減策として期待が大きい退院サマリの自動要約技術の開発を目指す。これにより、紹介状の自動作成技術等、電子カルテの自動解析技術に関連する継続的な研究開発投資の実現が期待される。この研究開発サイクルを確立することにより、要素技術である電子カルテ上の記載から自動情報抽出における継続的な精度向上が期待される。

こうして確立する医療用自然言語処理技術は、大量の電子カルテからの効率的な情報抽出を実現し、健康医療政策に資する統計データの収集コストを劇的に低廉化することが期待される。とりわけ、様々な傷病や治療に関して、既存の DPC やレセプトには表れてこない深遠な実態を明らかとし、医療の質向上・均てん化・各種医療技術の臨床開発に必要なエビデンスを生み出していくことが期待される。

また、国立病院機構の有する広域電子カルテ網は、各病院が独自に調達した電子カルテベンダー主要 6 社を網羅している。本研究によって、大口顧客としての交渉力を背景としたこれら主要ベンダーへの研究開発成果の技術移転が期待される。これらは、医療の情報化を進める厚生労働行政にとって、新たな政策手段の実現をもたらす

加えて、本研究で作成した入院カルテ・退院サマリの高品質な個人情報を含まない模擬のデータセットは、研究成果として公表することにしており、電子カルテの現物を持たない情報系研究者が、容易に本領域の発展に貢献できる環境を提供出来る点も大きな成果になるものと考えている。

A.目的

本研究は、電子カルテに実装可能な医療用語の標準化を行うシステムの開発に向け、そのステップとして、臨床現場からのニーズが極めて大きい退院サマリの自動生成技術の実現を目標と定める。電子カルテの自動解析は技術的な難易度が高く、実用的な精度を実現するためには多額の研究開発投資が求められる。そこで、本研究提案では、医療現場に直接的なメリットが生じる研究課題に取り組むことによって、現場の協力と今後の追加的な研究開発投資を呼び込み、その過程を通じて実用性の高い電子カルテの自動解析技術を実現する戦略を採る。初年度、我々が今まで模擬カルテを用いて研究開発を進めてきた標準化技術を、国立病院機構の有する広域電子カルテ網(NCDA)上の実カルテへと適用し、技術的な課題を抽出する。2年目には、NCDAを用いて集積した電子カルテに加えて、退院サマリ情報を用いることで、電子カルテの自動要約技術の検討を行う。3年目においては、両技術の統合により、継続的な精度向上の体制を実現するとともに、研究成果を既存の各社電子カルテへと組み込む枠組みを構築する。

本研究により、退院サマリの自動要約技術や紹介状の作成支援技術等、医療用の自然言語処理に関連する多彩な応用技術が実現する。これは、医療現場における負担軽減策として極めて効果が期待される。また、こうした応用の発展により、要素技術である電子カルテ上の記載からの自動情報抽出において、継続的な精度向上が実現する。この手法は、電子カルテにおける用語の標

準化技術単独に研究開発投資を行うことと比して、投資効率が極めて高いと考えられる。さらに、こうして医療用自然言語処理技術が発展することにより、大量の電子カルテからの効率的な情報抽出が実現する。これは健康医療政策に資する統計データの収集コストを劇的に低廉化し、今後、政策に求められる様々なエビデンスを継続的に生み出していく基盤となることが期待される。

なお本分担研究では、本研究における「大規模に電子カルテデータを手入できる体制」として、全国の国立病院41施設より年間80万患者の電子カルテ情報を自動収集する診療情報集積基盤(NCDA)を構築し、運用している基盤を用い、本研究目的のためのデータの収集・分析活動を行うためのシステム構築及び運用を行うことを目的とする。

B.方法

研究申請時に記載した研究計画および方法は以下のとおりであった。

【取り組む課題】

本研究では、電子カルテの記載より標準化した情報の抽出を行う技術の研究開発を行う。たとえば、カルテに「精神遅滞」と記載されている場合、この語を解釈することは機械にとって容易である。しかし、カルテの記載は単純ではない。「MR」と記載されていた場合、その語が精神遅滞(Mental Retardation)を指すの

か、僧帽弁閉鎖不全症 (Mitral Regurgitation) が、製薬会社の営業担当 (Medical Representative) がを正確に判別することは機械にとって容易ではない。さらに、「学業成績が芳しくない」と記載されたとき、医学的には文脈より mild MR と解する柔軟性も求められる。カルテ記載は、往々にして略記法や省略が多用される。「腹部：圧痛、反跳痛なし」という記載から、圧痛の有無を判別することは、高度な処理を要する。医療用自然言語処理はこのように難度の高い処理であり、現時点では情報抽出の精度に限界があり、研究開発体制におけるブレークスルーが求められている。

【研究体制】

研究代表者 堀口は、「大規模に電子カルテデータを入手できる体制」として、全国の国立病院 41 施設より年間 80 万患者の電子カルテ情報を自動収集する診療情報集積基盤 (NCDA) を構築し、運用している。この研究基盤を用いることにより、大学病院等のように患者の偏りが生じない多彩な施設からカルテ情報を収集することが出来る。

研究分担者 狩野・森田は、日本語カルテに記載された所見や病名を自動的に標準化するコンテスト (NTCIR MedNLP) を主催する、医療用自然言語処理における我が国の代表的な研究者である。このコンテストは、日本語カルテからの情報抽出精度を競う唯一の学術集会であり、森田は参加者としても首位の成績を収めて来た。狩野は、コンピュータによる大学入

試問題の自動解答に向けた人工知能研究「ロボットは東大に入れるか」プロジェクトにおいて、社会科分野の科目担当リーダーを務めている。両名は、公募研究課題の求める「工学系人工知能研究者・自然言語処理研究者」の役割を果たす。

分担研究者 奥村は、我が国においては数少ない計算機科学分野において学位を取得した医師であり、自然言語処理の医療応用分野で継続的な貢献を果たしてきた。本研究において奥村は、「病名・治療の標準化に詳しい臨床医学系研究者」として、基礎技術研究と応用研究とを橋渡しするコーディネータ役を担う。

分担研究者 岡田は、本研究の対象となる各国立病院カルテを記載する臨床医を取りまとめ、研究開発成果の臨床的有用性の担保に取り組む。

【H28 年度計画】

初年度においては、上述の NTCIR MedNLP コンテストを中心に研究開発を進めてきた標準化(正規化)技術を、NCDA 上の実カルテへと適用し、予備的な検証を行う。NCDA では、厚労省の定める電子カルテ情報の交換規約である SS-MIX2 を用いて、41 の国立病院からカルテ情報を常時本部のデータセンターに集積している。この病院中、約半数から、SS-MIX2 の基本データ(検査・投薬・病名等)に加えて、拡張ストレージ中の医師診療録、看護記録を収集している。これだけのサイズのデータセットを対象とした医療用自然言語処理技術の適用は、本邦初、最大規模

の試みとなる。

この研究計画を踏まえ、初年度においては、我々が日本語における医療用自然言語処理の研究コミュニティを形成し研究に取り組んで来た標準化技術を実カルテへと適用することで、カルテからの情報抽出の自動化に向けた予備的な検証を行うことを計画した。

その計画の実現に向け、採択後研究者全員で組織する「総括・企画調整班」を作り、月2回程度のミーティングを行い、研究の方向性の整理及び各分担班の作成及び役割の決定、進捗の管理及び調整を行う方法で研究を遂行することとした。また、「総括・企画調整班」以外の分担班についてはそれぞれ責任研究者を決め、その裁量で研究を進める方法をとった。

(尚、総括・企画調整班以外の分担班については総括・企画調整班の活動結果から生まれたものであり、それぞれの班の目的・方法についてもC.結果セクションで記載することとする。)

C.結果

1. 総括・企画調整班

まず、本研究班はその応募要項の段階からデータ収集に掛かる部分については研究の中に組み込まないことを求められており、データの収集基盤の構築・運営については本研究のカバーする範囲ではない。しかしながら、本研究の前提となる国立病院機構が作成・維持運営するNCDAについて本報告書でその概要や意義について記述を行わないとするならば、本報告書の内容の理解

に大きな妨げになると考えここに報告を行うこととする。

NCDAの概要については参考資料1にその概要を平成28年度の医療情報学会で発表した資料を添付した。また、実際の病院におけるSS-MIXデータ作成に掛かるシステム仕様についても参考資料2に示した。

これらの仕様等のドキュメントについてはその改版履歴も含め、github上で管理、公開している。

https://github.com/nhoHQ/SSMIX2_support_documents

次に、実際の本分担研究班の活動についての報告を行う。本研究班においては、まずは研究者が独立して研究活動を進めるのではなく、10月からの半年間で10回の研究班会議(うち8回はWeb音声会議)を行い、1つの有機的な研究班として活動が行える環境で運営してきた。

各班会議での調整事項は以下の通りである。

第1回

- 研究管理面の話題
- 各人の状況 update
- 仮説構築作業
- 「退院サマリとは何か？」

第2回

- 作業仮説構築
- ツールドリブン/リソースドリブン/臨床ニーズ/ 病院管理ニーズからの整理
- 倫理審査に向けた論点整理
- 病院訪問に向けた調整

第3回

- 三重病院にて退院サマリの記載内容に

ついて臨床家とともにディスカッション

第 4 回

- 今後のスケジュール確認
- 三重病院訪問での成果確認
- 倫理審査に向けた調整
- 研究分担の整理
- 「良質な退院時サマリとは？」問題の整理

- ダミーカルテ作成の是非

第 5 回

- 継続申請書類の作成について
- 倫理審査の申請書確定について
- 医師アンケート企画について
- テストデータについて
- 解析のアプローチについて

第 6 回

- H29 継続申請について
- H28 倫理審査の状況報告
- 年度内達成目標の再確認
- 各分担研究状況報告
- テストデータについて

第 7 回

- 各分担研究状況報告
- 研究基盤の整理
- 倫理委員会・利活用審査委員会の報告

第 8 回

- 各分担研究状況報告
- 退院サマリの関する文献サーベイについて
- カルテ要約の要素技術についての議論
- スケジュール確認

第 9 回

- 各分担研究状況報告
- 退院サマリの関する文献サーベイについて
- 報告書作成について

- 分析環境の整備について報告

第 10 回

- 各分担研究状況報告
- 報告書作成について

なお、第 5 回でとりまとめた研究の倫理審査申請書は参考資料 3 に示した。

本研究はカルテの非定型の記載欄に記入されたデータを使うという研究であり、患者の不利益等を防止するために倫理的な配慮をした上で、倫理審査を受けなければならない。その為、研究期間が 10 月末から開始された後、先ずどのような分析活動を行うかについて数ヶ月にわたり検討を行い、平成 29 年 1 月に国立病院機構中央倫理委員会に侵襲・介入なしの観察研究として倫理審査の申請を行い、3 月に承認された。倫理審査の承認後、データ利用に際して必要な国立病院機構内のデータベース利活用審査委員会への利活用申請を行い、3 月にその承認も受けた。

なお、NCDA データは国立病院機構が契約するデータセンター内で厳重に管理されている。研究に際しては、このデータベースから研究テーマごとに匿名化したサブセットを切り出し、国立病院機構本部内のオンサイト利用に限っている。以上により、データセットの利用対象と利用目的を厳しく制限することにより、患者個人情報の保護を行っている。それに対応する分析基盤の作成に関して、分担研究班を組織し、堀口・岡田が責任者として活動を行うこととした。(分担研究の結果は後述)

また、第 4 回の議論の結果、「退院サマリの自動生成に向けたアプローチの検討」とい

うテーマの分担研究を奥村が、「退院サマリの自由記載文の特徴解析」というテーマの分担研究を森田が、「退院サマリの自動生成に向けたアプローチの検討」というテーマの分担研究を狩野が担当することとした。

2. SS-MIX2 分析用データセットの作成・開発班

研究代表者堀口及び分担研究者岡田は、国立病院機構本部との調整を中心とした基盤構築を行った。まず、NDCA データの研究利用に向け、倫理審査申請に加えて、内部規定にて定められている内部委員会の調整を図った。また、閲覧・解析に特化した自然言語処理用の研究基盤の構築を行った。研究基盤の概念図は図1に示したとおりで、セキュリティを維持しつつ、空間的制約をなるべく少なく研究が進められるようなものになっている。

また、本研究で中心的に使われる医師記録等（経過記録・退院サマリ）については、SS-MIX2 の標準仕様に含まれていないが、JAHIS の提供している仕様を参考に、資料1で示した仕様でNCDA内に実装することとした。

3. 退院サマリの自動生成に向けたアプローチの検討班

入院患者の退院に際し、医師は入院中に記載したカルテ等の情報から退院サマリを作成する必要がある。この退院サマリを自動的に生成することが出来れば、臨床現場の負担を下げる事が出来ると共に、医療の

質に貢献することが期待される。

そこで、本研究分担では、退院サマリの自動生成に向けた研究アプローチの検討に取り組んだ。まず、文献調査と医師へのヒアリングに基づき、良質な退院サマリに求められる要件について定性的な検討を行った。同時に、実際の退院サマリを対象とした分析を行い、要約過程に関する知見を整理した。さらに、一般的な文書の要約手法と入院カルテの要約手法について文献調査を行った。

今年度の研究の結果、入院カルテの自動要約に向けて、医師が作成する退院サマリの分析枠組みと自動的な退院サマリ生成モデルの双方を兼ね備えた、CASE モデルと証するモデルを構築した。これは、「そもそも退院サマリには何が書かれているのか」という観点より構築されたモデルであり、サマリ中の各文を「カルテに由来するかどうか」という軸と「抽象度が高いか低いかな」という軸によって4つのクラスに分類する(図2)。これらのクラスは、退院サマリの分類モデルであると同時に、それぞれ生成に際して固有の処理が求められることから、退院サマリの自動生成に向けた処理モデルとしての性質も有する。以下では、これらクラスの分類軸となる「カルテ記載の有無」、「言及の抽象度・事実度」それぞれを概観した上で、退院サマリの生成モデルとしての展望を記す。

カルテ記載の有無・・・退院サマリは、入院カルテに記載された情報を元に作成される。したがって、退院サマリに記載されている情報も、入院カルテに書かれた情報の抜粋ではないかと思われるかも知れない。しか

し、退院サマリには入院カルテに直接由来しない文が一定数含まれる。たとえば、ちょうど退院後のタイミングで検査結果が出た場合や、退院後の療養計画等について退院までにカルテに書き損ねた場合、サマリ中に由来のない文が含まれることになる。あるいは、退院後、入院カルテに書かれた記載を臨床的に評価し、その内容をサマリに記載することもあるかも知れない。いずれにせよ、結果として、入院カルテに直接記されていない文が退院サマリに収載されることになる。今年度の研究を通じて、この「入院カルテ中の原文の有無」について分析を進めたところ、「有」か「無」かではなく、ある程度の濃淡があることが分かった。まず、文の一致として、文レベルで一字一句一致しているのか、あるいは一部変更があるのか。句レベルで一致しているのか、あるいは、類似しているのか。これらのうちどこまでを「有」と取るかによって、区分が変わることになる。また、サマリ文中に複数の句がある場合、それらが入れ子になっている場合等、複雑な事例が少なからず含まれることが明らかとなった。

言及の抽象度・事実度 次に、文における言及の抽象度の分類について記す。退院サマリにおいては、具体的な事実の他に、具体的な事実を抽象化して記載されることが少なくない。たとえば、手術での入院の場合、退院サマリには、入院に至った原因の疾患についての記載に加えて、手術での術式等の情報がまず記載されるであろう。これらは事実に関する言及であり、抽象度が低いと言える。一方、手術後の経過について、日々のカルテに具体的な記載がなさ

れていても、サマリ上では「術後経過は順調であった」と医学的な評価が記載されるケースがある。あるいは、処置を何回か行ったとすると、「処置を行ったが改善は見られなかった」等と整理して記載するかも知れない。これらは抽象度が高い記述と言える。こうした記載における抽象度の区別については、自然言語処理における事実度の分類技術 (factuality classification) が役に立つ可能性がある。ただし、技術的にいくつかの課題があることに留意する必要がある。まず、「骨折を認めた」という記載は明らかに事実度が高いが、「の可能性を伝えた」とした場合、伝えたことは事実でも についての事実度は低い。同様に、「改善を認めた」といった言及も、医学的な解釈や価値判断を伴っており、骨折のようなケースと比して事実度は低いであろう。このように、事実度を自動的に判定するには、文法と意味の双方を解さなければならない。それでも、退院サマリにおける事実に関する記載の割合や出現位置が分かることで、次に示す退院サマリの自動作成に向けた基礎的な統計が得られることが期待された。

退院サマリの自動生成 最後に、これら退院サマリの分析モデルを用いて、退院サマリの自動生成手法について検討を行った。分析モデルにおける4つの区分は、それぞれサマリ作成における戦略の違いに対応している。カルテ中の情報に着目した場合、“Extraction”はそこから重要文を抽出する「抜粋」に相当し、“Abstraction”は、本来の意味での「要約」に該当する。重要文抽出(extraction)は既に多くの先行研究

がある分野であり、比較的取り組みやすい。一方、要約操作(abstraction)については、「言い換え」や「重複削除」、「解釈」等、特性の異なるいくつかの手法の総体であると考えられ、技術的難易度が高い。“Clinical reasoning”は、そうしたカルテ記載から医学知識を持って推論や判断を行った結果の記載となる。医学知識に基づいた推論もまた、技術的難易度が高い。しかし、同じ医師や同じ診療科のカルテには往々にして同じような表現が繰り返されることから、大量のデータを機械学習することにより、ある程度の有用性が得られる可能性がある。最後の“Supplementary information”は、入院カルテ中からは直接得ることが出来ない情報であり、検査結果や読影レポート等を対象としたさらなる検討が望まれる。退院サマリは、これら特性の異なる4つの「要約処理」の結果の集合体であり、入院カルテを単一の要約アルゴリズムにより1度に高精度な形で生成することは困難と考えられた。そこで、今後の研究の方向性として、高精度な要約アルゴリズムの研究開発ではなく、特性の異なる4つの要約モジュールが出力した候補文集合を「退院サマリの下書き」としてユーザーに提示する「退院サマリ作成支援ツール」のプロトタイプングを提案する(図3)。このツール上では、ユーザーが行う清書作業を、ユーザーからのフィードバックとしてすべて記録することが望ましい。そのデータを用いて各モジュールをチューニングすることで、生成する要約を継続的に改善し、ユーザー嗜好に合致させる処理系の実現が期待される。また、システムを用いて退院サマリを作成する際、利用する専門用語を標準語彙へと半

自動的に対応付けながら作成する工夫が望ましい。カルテ中の専門用語は、医療機関やユーザーによって使い方が異なることから、カルテの自動解析処理における障害となっていた。提案手法により、医師に過度の負荷を掛けない形で精度の高い個人辞書の作成が可能となる。とりわけ、各医師の個人履歴には、単語の言い換えに加えて、句から単語への言い換え等、カルテ記載の正規化に向けた豊富な言い換え情報が含まれることになる。このデータは、退院サマリ作成支援により効率的に収集が可能であり、医療用自然言語処理に技術革新をもたらすことが期待される。

4. 退院サマリの自由記載文の特徴解析

退院サマリの各文について、元になった入院カルテと比較をすることで、その文が入院カルテからそのまま抜き出された文なのか、文や文節などを組み合わせて書かれた文なのか、それとも入院カルテの記載を解釈して新たに生成された文なのか、を決定する。

もし入院カルテから抜き出した文を組み合わせるとして退院サマリが作成されているのであれば、自動生成のためには入院カルテから適切な文を抜き出して並べることになる。文節や単語を組み合わせで書かれているのであれば、適切な文節や単語を抜き出して文を生成することになる。単語すら書き換えられて入院カルテの記載とは異なる文が書かれているのであれば、入院カルテを入力として文を生成することになる。

退院サマリの各文は次の5つのタイプに分類した：タイプ1. 入院カルテの文がそのまま(もしくはほぼそのまま)使われて

いる、タイプ2．入院カルテの文そのままではないが、複数の文や文節を組み合わせることでその文を作ることができる、タイプ3．その文を書くには入院カルテを読んで解釈をする必要がある（医療の知識がなくても解釈が可能な範囲である）、タイプ4．その文を書くには入院カルテを読んで解釈をする必要がある（医療の知識がないと解釈ができない）、タイプ5．その文は入院カルテの内容からだけでは書くことができない（情報が不足している）。分類作業は医療の知識がある4名で行い、不一致の場合は話し合いによって1つの分類に決定した。13の退院サマリを使用した。退院サマリは入院までの経過および入院中の経過を使用した。

結果、全体での各タイプの内訳は、タイプ1：43%、タイプ2：3%、タイプ3：9%、タイプ4：24%、タイプ5：21%、となった。入院までの経過における各タイプの内訳は、タイプ1：72%、タイプ2：1%、タイプ3：4%、タイプ4：10%、タイプ5：13%、となった。13の退院サマリのうち6の退院サマリでは、入院までの経過のすべての文がタイプ1であった。入院中の経過における各タイプの内訳は、タイプ1：24%、タイプ2：5%、タイプ3：12%、タイプ4：33%、タイプ5：26%、となった。入院中の経過ではすべての文がタイプ1の退院サマリはなかった。

入院までの経過は、前半部分に発症からの経過が、後半部分に入院を判断するに至った理由が書かれていることが多かった。入院までの経過は全体的に入院カルテから文をそのまま持って来ていること（タイプ1）が多かったが、すべてがタイプ1では

ない場合には、前半部分で特にその傾向が強くなり、一方で後半部分は医学的な知識がないと解釈ができない文（タイプ4）の割合が若干だが高かった。

入院中の経過は、入院中の症状と治療の経過が書かれ、その最後には退院をした旨と退院後の方針が書かれていることが多かった。退院後の方針は入院カルテの記載だけからでは書くことが難しいこと（タイプ5）が多い傾向にあった。

入院までの経過と入院中の経過を比較すると、入院中の経過はタイプ1の割合が低く、タイプ4と5の割合が高くなっていた。入院までの経過がタイプ1が72%だったのに対し、入院中の経過は逆にタイプ3～5が計71%となった。いずれの場合もタイプ2は非常に割合が低かった

5.退院サマリの自動生成に向けたアプローチの検討班

退院サマリの自動生成のため、主に自然言語処理の要素技術という観点から現状と全体像をつかむための予備的な調査研究および研究環境整備を、他グループと協調して行った。

研究環境整備については、電子カルテデータをセキュアな環境で扱えるようにするための環境設計と構築を行った。単にセキュアな環境を確保するだけでなく、効率的な研究開発を行えることが研究環境整備の目的のひとつである。そこで、電子カルテデータ自体は国立病院機構の厳重に管理された環境内にとどめつつ、プログラムの実行環境のほうを仮想マシンとして移動させることとし、環境の構築を行った。

他の分担研究により、本年度模擬カルテ

の提供と、その模擬カルテに基づいた、退院サマリ作成を考慮したアノテーション付与が行われた。

サマリを要約ととらえると、一般的な自動要約技術が適用できそうにも思われる。多くの自動要約は、トピックの統計的な解析を行ったうえで、文書中で重要なトピックが含まれるものを残す、という手法が骨格になっている。しかし退院サマリでは、統計的に重要でない、文書集合中で共通して頻出するトピックであっても、サマリとして残すべきことが多々ある。

また、入力にあたる電子カルテの文章中にない文章や表現が、サマリにどのくらい含まれているかという問題がある。入力の子セットでよいのであれば、切り貼りの範囲内におさまるが、現実には言い換えに始まり内容的にも新規な文章の挿入がある。

分担研究のデータによると、入院までの経過については7割以上のサマリ文がカルテの文章ほぼそのままであった。このことは、医師がサマリを作成する際に文の複製を使用しており、分量的な減少もあまりみられないことから、内容的にもあまり変更を必要としていないことを示唆している。ただし、入院前の記述は他の医師からのお願いの形式をとっており、そのままでは主語や言葉遣い、時制などが不適切なので、そうした部分の変換が必要かもしれない。入院中の経過については、2～3割程度がそのままの文であった。入院中のカルテの記載は文を完成させず断片的なスタイルのことが多く、一方サマリではきちんとした文にするため多かれ少なかれ文生成の要素

が必要と思われる。また、医師本人の記録なので、振り返ることで要素を追加したり、整理したりすることが想像される。

D. 考察

医療用情報システムの研究開発においては、医療現場に直接の恩恵が及ばないゴールが設定されることで、研究開発が現場のニーズから乖離するとともに、継続した開発投資に繋がらない悪循環が往々にして生じてきた。本研究提案は、医療現場における負担軽減策として期待が大きい退院サマリの自動要約技術の開発を目指す。これにより、紹介状の自動作成技術等、電子カルテの自動解析技術に関連する継続的な研究開発投資の実現が期待される。この研究開発サイクルを確立することにより、要素技術である電子カルテ上の記載から自動情報抽出における継続的な精度向上が期待される。

こうして確立する医療用自然言語処理技術は、大量の電子カルテからの効率的な情報抽出を実現し、健康医療政策に資する統計データの収集コストを劇的に低廉化することが期待される。とりわけ、様々な傷病や治療に関して、既存のDPCやレセプトには表れてこない深遠な実態を明らかとし、医療の質向上・均てん化・各種医療技術の臨床開発に必要なエビデンスを生み出すことが期待される。

また、国立病院機構の有する広域電子カルテ網は、各病院が独自に調達した電子カルテベンダー主要6社を網羅している。本研究によって、大口顧客としての交渉力を背景としたこれら主要ベンダーへの研究開発成果の技術移転が期待される。これらは、医療の情報化を進める厚生労働行政にとつ

て、新たな政策手段の実現をもたらす

加えて、本研究で作成した入院カルテ・退院サマリの高品質な個人情報を含まない模擬のデータセットは、研究成果として公表することにしており、電子カルテの現物を持たない情報系研究者が、容易に本領域の発展に貢献できる環境を提供出来る点も大きな成果になるものと考えている。

。

E. 結論

今年度の研究においては、NDCA データの利用に際した倫理審査の関係で、実カルテを対象とした解析は3月以降に限られる制約が生じた。しかしながら、高品質なテストデータの確保が可能であったことから、研究のコアである医療用自然言語処理部分の検討はほぼ計画通りに進めることが出来た。電子カルテの自由記載部分を自動解析する多施設構成での大規模データを対象とした研究としては、本邦初の試みとなろう。今後も本計画に沿った研究を遂行していきたい。

F . 研究発表

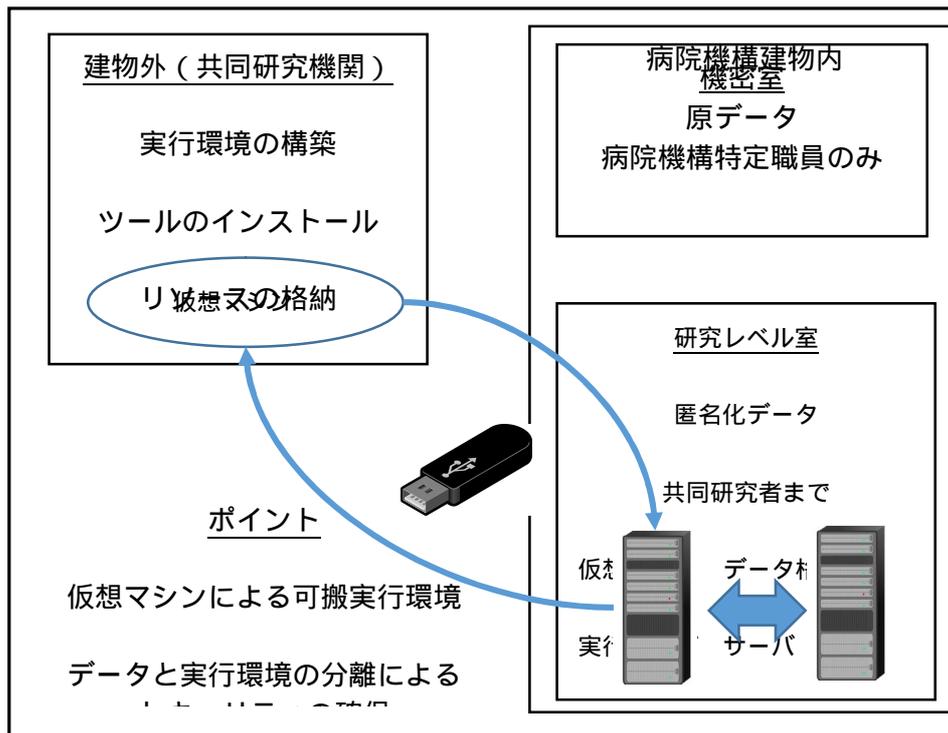
1 . 論文発表

なし

2 . 学会発表

なし

図 1



資料 1 NCDA における医師記録等の仕様書

趣旨

本事業では、各社の SS-MIX2 モジュールの拡張ストレージへの出力機能を利用し、以下の情報を出力することを求めている。その際、SS-MIX2 拡張ストレージ構成の説明と構築ガイドライン Ver.1.2d (以下、ガイドライン) に記載している仕様に対応していること。また、トランザクションストレージ、インデックスデータベースも同時に生成すること。

- 経過記録
- 退院時サマリー
- 診療情報提供書

以下に仕様を示す。

ドキュメントデータ 物理構造

```
|-- 拡張ストレージ ルートフォルダ
  |-- 患者 ID 先頭 3 文字
    |-- 患者 ID 4~6 文字
      |-- 患者 ID
        |-- 診療日
          |-- データ種別
            |-- コンテンツフォルダ
              |-- 主文書ファイル
```

診療日

特に指定しない。

データ種別

ガイドライン P4 (4)「データ種別フォルダ」について に則ること。

```
[ローカル文書コード]^ローカル文書名称^[ローカルコード体系コード]^標準文書コード^標準文書名称^標準コード体系コード
```

以下のように標準コードに対しローカルコードが複数あることは許容される。

```
L12345^ 入院診療録^99ZZZ^11506 -3^経過記録^LN
```

```
L12346^ 外来診療録^99ZZZ^11506 -3^経過記録^LN
```

コンテンツフォルダ

ガイドライン Ver.1.2d P5 (5)「コンテンツフォルダ」について に則ること。

```
患者 ID_診療日_データ種別コード_特定キー_発生日時_診療科コード_コンディションフラグ
```

いずれの文書も削除は想定していないが、電子カルテシステムによっては修正はあり得ると考える。その場合、ガイドライン P6 ④修正が発生する場合 に則り改版すること。

主文書ファイル

XML CDA R2 で出力すること。XML ファイル以外に画像ファイルや CSS ファイル等を出
力してもかまわない。

HEADER 部

いずれの文書も JAHIS 診療文書構造化記述規約 共通編 Ver.1.0 に則ること。

P27 6.3.11.検査・診療等行為 "documentationOf/ServiceEvent" によると、documentationOf
の制約・多重度は 0..1 となっているが、経過記録、退院時サマリについてはこれを 1..1 と
読み替えること。

経過記録は serviceEvent classCode(サービスイベントクラスコード)を ENC(診察)とし、
effectiveTime(実施日)は low value、high value とともに記録タイミングを出力すること。

退院時サマリは serviceEvent classCode(サービスイベントクラスコード)を ACCM(入院、
滞在)とし、effectiveTime(実施日)は low value に入院タイミング、high value に退院タイミ
ングを出力すること。

タイミングの粒度は日以上であれば良い。

BODY 部

診療情報提供書は、日本 HL7 協会 患者診療情報提供書 規格 Ver.1.00 に則ること。

診療情報提供書以外は、XML の文法に則ること。

参考資料

1. NCDA データベースの説明資料 (平成28年医療情報学会発表資料抜粋)

国立病院機構における 電子カルテデータ標準化について

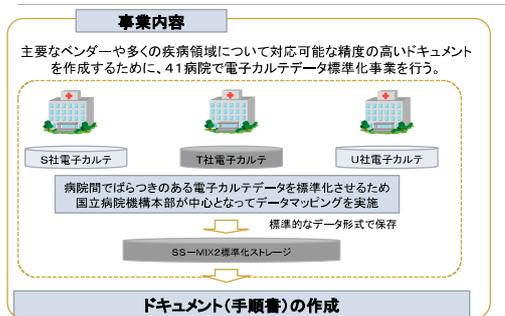
国立病院機構本部
IT推進部 医療情報データベース企画室長
堀口 裕正

※第36回医療情報学連合大会 COI開示
開示すべきCOIはありません。

補助金事業の事業背景

- 平成26年6月24日に閣議決定された「世界最先端IT国家創造宣言」では、地域を越えた国民への医療サービスの提供等を可能とする医療情報利活用基盤の構築を目指し、医療情報連携ネットワークについては、電子カルテを含めたデータやシステム仕様の標準化等を行い、平成30年度までに全国への普及・展開を図ることとされている。
- しかしながら、電子カルテについては、ベンダー毎で開発が行われ、各病院が使いやすいようにカスタマイズされるなど、電子カルテデータの形式が標準化されないまま普及したことから、電子カルテ上で使用されている病名や医薬品等のコードがベンダーや病院で異なり、標準化の課題となっている。
- 今回の『電子カルテデータ標準化等のためのIT基盤構築事業(13.0億円)』は、このような問題を解消するため、各病院の電子カルテデータを厚生労働省の定める標準コードに紐付けするデータマッピングを行い、SS-MIX2規格(標準化ストレージ機能)を用いて電子カルテデータの標準化を実施し、その工程を示したドキュメント(手順書)を作成・公開することを目的としている。

補助金事業の概要(課題・目的等)



事業の成果(標準化の普及促進関係)

- 最新のSS-MIX2Ver1.2cに完全準拠しているモジュールが41病院に導入
 - SS-MIX2 Ver1.2cモジュールの導入
 - SS-MIX2に完全準拠しているモジュール
- HOTコード・JLAC10・ICD10など標準コードを全面的に導入・活用
- 従前のモジュールで課題となっていたベンダー毎の表記ゆれ等の問題が解決され、データ形式の標準化が可能となります
- 本モジュールは6ベンダーから他の医療機関にも(有償にて)提供可能です。
- 他の医療機関が厚生労働省標準規格に準拠(SS-MIX2・標準コード等)したシステムを導入するに当たり、当該事業で作成したドキュメント(手順書)を活用することにより、専門的な知識を要することなく、簡単に導入することが可能となります。

事業の成果(標準コード及び標準化団体)

- 標準規格が持つ課題を標準化団体とともに解決
 - HOTコード・・・一般名処方用や持参薬用のコードの整備をMEDISIに依頼
 - JLAC10コード・・・体温等の検査コードの採番依頼
 - SS-MIX・・・各種規約の矛盾や、解釈について整理をJAMIに依頼

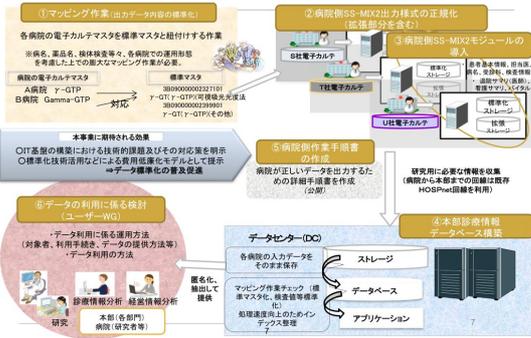
今回のプロジェクトのコンセプト

- 補助金事業として13ヶ月という短納期で仕上げる必要がある。
- 標準化の普及促進に資することを目標とする
- 以上の条件から以下のコンセプトで事業を実施した
- 検証環境での十分なテスト/検証を行い病院別の開発を極力行わない
- 病院における医療提供に係るユーザーインターフェイスは一切変更しない

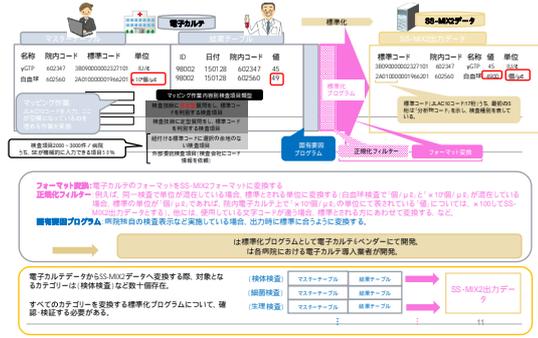
国立病院機構のDB事業概要(プロジェクト概要)

方針	
主な作業区分	内容
①マッピング作業(出力データ内容の標準化)	対象41病院を選定し、データマッピング作業を実施する
②病院側SS-MIX2出力様式の正規化(拡張部分を含む)	全てのSS-MIX機能(メッセージ)に対応できるよう、モジュールを各ベンダーで正規化(入力値の正規化・フルセット化等)する。併せて標準仕様以外の拡張データ(バイタル等)が出力できるようにする
③病院側SS-MIX2モジュールの導入	①で選定した対象病院に②で作成したSS-MIX2モジュールを導入する
④本部診療情報データベースシステム構築	データを収集する仕組みを検討し、外部データセンターにデータベースを構築する
⑤作業手順書の作成	本プロジェクト終了後、各病院がSS-MIX2を効率的に導入できるように、SS-MIX2モジュールを導入するベンダーが作業手順書を作成する(手順書は公開予定)
⑥データ利用に係る検討(ユーザーWG)	システム機能とユーザーの要望について調整する データベースの利用に係る規定(プロセスやルール)や具体的なデータ利用方法を検討する

SS-MIX2を用いた診療情報データベース構築プロジェクト 作業区分①～④

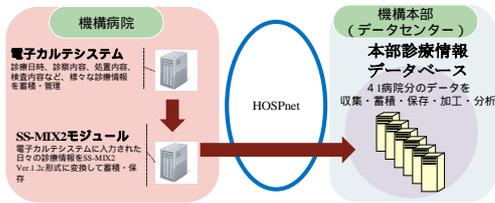


データ標準化のイメージ(SS-MIX2出力)

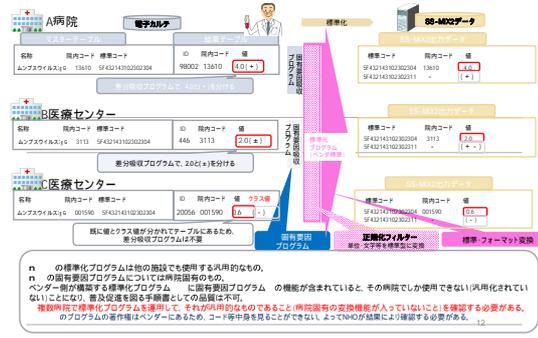


IT基盤構築の仕組み (標準化、IT基盤とは)

- 本事業では、各病院の電子カルテシステムのデータをSS-MIX2 Ver.1.2c形式に変換(標準化)し、診療情報データベースに収集する仕組み(IT基盤)を構築します。
- このIT基盤は、各病院に導入するSS-MIX2サーバと、機構本部(データセンター)に導入する本部診療情報データベースシステムから構成されます。



SS-MIX2変換プログラムの構成



検査結果の表記揺れについて

病院の部門システム等から発生する検査結果表記に関しては以下のような揺れがあります。今回、病院のモジュールでの揺れを標準的記法に統一しました。

- 検出限界超の値 $\{ > 100(\text{半角})、\text{スペース}、100\}$

100<	100<	100<	100<	>100	>100	>100	>100
------	------	------	------	------	------	------	------
- 検出限界以上の値 $\{ \geq 200(\text{半角})、\text{半角}=\text{スペース}、200\}$

200 ≤	200 ≤	200 ≤	200 ≤	200 ≤	200 ≤	200 ≤	200 ≤
200 <	200 <	200以上	200以上	≥ 200	≥ 200	≥ 200	≥ 200
⇒ 200	⇒ 200	⇒ 200	⇒ 200	⇒ 200	⇒ 200	200~	200
- 定性値 (-) $\{ \text{-(スペース}2\text{半角)-}\}$

-	(-)	フケンシュツ	インセイ	陰性	検出せず	検出せず
-	(-)	フケンシュツ	ミケンシュツ	ミケンシュツ	ケンシュツセズ	ケンシュツセズ
- 半定量値 $\{ 1+(\text{スペース}、1\text{スペース}、\text{半角}の+)\}$

1+	(2+)	(4+)	5+	(++)	+++	(++++)	+++++
----	------	------	----	------	-----	--------	-------

検査結果の変換について

各JLAC10においてどのような検査値をどのような単位で記載すべきか、どのように表記すべきかについて、マスターを提供・公表しました。

検査項目(参考)	JLAC10	単位
尿量	1A00500000192001	mL
末梢血液一般検査 - 赤血球数	2A990000001930951	$10^4/\mu\text{L}$
末梢血液一般検査 - 白血球数	2A990000001930952	$/\mu\text{L}$

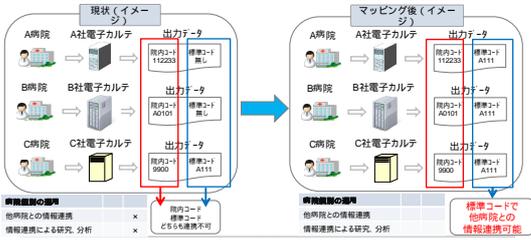
検査	分離前 JLAC10	検査値のバリエーション	分離後 JLAC10-1	分離後 JLAC10-2	検査結果	検査結果 2
麻疹ウイルス抗体IgG	SF431143102302304	数値(記号)文字列	2.0(-)未満	2.304	SF431143102302304	数値[+文字列]
鼻汁好酸球	2A300000006360411	(記号)数値	(2+)<60%	0.6360411	2A300000006360402	記号 数値

SS-MIX2をいれたIT基盤構築事業参加病院口覧

ベンダ・病院種別分布	350以上	350-499床	349床以下	複合(その他)	障害者病床中心	総計
富士通	5病院 滋賀医療、名古屋医療、大塚医療、九州医療、長崎医療、熊本医療	1病院 相模原	1病院 南和歌山医療	5病院 北海道医療、西野、馬、東京、長良医療、村山医療、福岡東医療	1病院 東三、張王、三里、広島西医療	22
日本電気	2病院 北海道がん、埼玉	2病院 旭川医療、帯広、釧路	1病院 仙台西多賀	1病院 島根	1病院 島根	7
ソフトウェア・サービス	5病院 高崎総合、四国がん、九州がん、湘野医療、鹿児島医療	1病院 氷子医療	1病院 崎北医療	1病院 崎北医療	1病院 崎北医療	7
亀田医療情報	1病院 静岡医療	1病院 西新潟中央、教養医療	1病院 大宅	1病院 大宅	1病院 大宅	2
SBS	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	2
日本IBM	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	1病院 仙台医療	1
地域毎	7	13	3	11	7	41
北海道	11病院	7病院	3病院	5病院	9病院	41病院

病院におけるマッピング作業

- 院内コードと標準コードを紐付ける対応表を作成します(マッピング作業)。
- 病院毎に異なる院内コードを、標準コードに変換することにより、他病院と連携した診療情報の分析等が可能になります。



検証環境におけるテスト

- 今回の事業を実施するため、本部に6ベンダーの電子カルテシステムをレンタルして各種検証を行った。
- すべての電子カルテに同じ処方オーダーを登録すると、そのデータがどう扱われて、SS-MIXデータに変換あされるのかについて、電子カルテの画面入力からスタートするテストを行って、各社の違いを確認した上で開発調整を行っている。
- これにより、NHO41病院に限らず、汎用的に使えるモジュールになるように開発を行った。

16

データベースについて

【国立病院機構 診療情報集積基盤】

(コクリツビョウインキコウ シンリョウジョウホクシユセキキバン)

英文表記 NHO Clinical Data Archives

省略形の記載法 「NCDA」

省略形の呼称 「クリニカルアーカイブス」

41病院で来院患者ベース 94万人/年 17,800床のデータベース

18

NCDAシステムフロー

各病院は標準化(拡張)ストレージ/標準化(拡張)トランザクションストレージ/インデックスDBの5つを作成

通常時、前日のTRストレージを本部システムが取りに行き、エラーチェックして本部DBに取り込む

不定期に本部のDB内のデータと病院の標準化(拡張)ストレージに齟齬が無いかバリデーションを行う(患者単位で実行可)

TRストレージでのデータ転送により大量の小さいファイルを転送せずに済み、通信コストが大幅に減る

1日~2日遅れでの取り込みとなり、リアルタイム性は望めない

SS-MIX2 標準化ストレージの構造

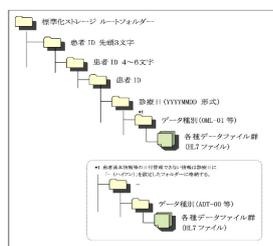


図 2.2-1 標準化ストレージ

(出典)SS-MIX2 標準化ストレージ構成の説明と構築ガイドライン

21

トランザクションストレージ

3.3 3.3.1

トランザクションストレージとは

- 「標準化ストレージ」は患者(患者ID)を特定してから、当該患者の診療情報を検索することに特化した物理構造を採用している。しかし、
- 1 何らかの理由で標準化ストレージ再作成しなければならない場合
 - 2 災害発生時への対策や地域医療連携の基盤として、外部接続回線を用いてデータセンター等の当該医療施設外に標準化ストレージの複製を作成する場合
 - 3 標準化ストレージ以外のシステムにおいて、本ガイドラインで定めた病院情報システムからの伝送データが再利用できると考えられる場合

上記のようなケースでは、診療情報がトランザクションとして標準化ストレージに記録された日時(以下「トランザクション発生日時」という)に着目して診療情報を参照することが必要であると考えられる。したがって、ここでは、病院情報システムから送られる標準化された診療情報そのものをデータソースとして再利用することによる便宜を考慮して、トランザクション発生日時により診療情報を参照することに特化したストレージとして、トランザクションストレージを規定する。

(出典)SS-MIX2 標準化ストレージ構成の説明と構築ガイドライン

22

トランザクションストレージ



図 3.3-1 標準化ストレージ

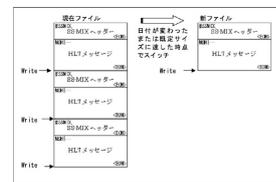


図 3.3-2 トランザクションデータファイルの切り替え

(出典)SS-MIX2 標準化ストレージ構成の説明と構築ガイドライン

23

NCDAにおけるエラーチェック

データ転送後、DB取り込み前に行うエラーチェックは大きく以下の2つ構成エラー

- SS-MIX2の仕様を満たしているかチェック
 - 型は正しいか、必須項目抜けは無い、行の順番は合っているか、項目毎のLengthは守られているか etc
- NHOの要求した仕様を満たしているかのチェック
 - 標準コードが入っているか
 - SN型にあった記載がされているか

毎日チェックした上でエラーを画面に報告、対処を行う

NHOのエラーについては、無視してDBに取り込むことが可能

NCDA本部DBでのデータ修正

基本SS-MIX2の使用に則って、DEL及びINSのSS-MIX(HL7)メッセージを作って登録することでエラー修正を行う。

これにより保存されているファイルからデータベースが再現可能となる

DB設計の際に考慮したこと

- HL7要素がすべて格納されること
 - オブジェクト型のDBを利用する
- 検査結果が抽出可能であること
 - 検査結果を単一の値では無く範囲で持つ
- いつ抽出作業を行っても同じ結果が出るようにすること
 - すべてのデータにデータ利用可能状態を示す情報に時間範囲を持たせることで実現
- 抽出処理を行う際にSQL言語が使えること

診療情報集積基盤における個人情報取り扱い

- 患者同意
 - 病院に提示されている「個人情報の利用目的」に「国立病院機構診療情報分析基盤での利用」を追加。(平成27年12月中に41病院で実施済み)
 - 併せて、ポスター・ちらしでの周知を開始
 - 患者の利用不可の申出には対応できるシステムとなっている
- 法令対応
 - 個人情報保護法・独立行政法人における個人情報保護法が来年施行見込みであり、今後出てくる政令・ガイドライン等に適切に対応していく
 - 研究の倫理指針の見直しがとりまとめられる方向なので、適切に対応していく
 - 医療等ID/代理機関等の法令改正が行われた場合にも適切に対応していく

31

参考) PostgreSQLにおける範囲型

2013年に組み込みの型としてリリース(23,34)と記載すると23より大きく34以下という意味となる中の値としては数値型とtimestamp型が利用可能
 '(3,7) '::int4rangeもしくはnumrange(1.0, 14.0, '[]')と書く

併せて専用の関数が定義されている。

診療情報集積基盤における利活用

- 患者に明示した個人情報の利用目的の範囲内で利活用を進める
- 利活用に際しては「利活用要項」を定め、それに従って利用を行う
- 利活用要項の骨子は以下の通り
 - データベース利用審査委員会を設置。データ利用について審議。
 - 利活用は匿名化後が原則
 - 研究における利用
 - 本要綱を遵守するとともに、倫理規定等の研究に関連する法令やルールを遵守する

32

参考) PostgreSQLにおける範囲関数

関数名	説明	戻り値	関数	戻り値	例	結果
range	range(下, 上) :: range	range				
range_lower	range_lower(下, 上) :: range	range	range_lower	range	range_lower(1, 10)	1
range_upper	range_upper(下, 上) :: range	range	range_upper	range	range_upper(1, 10)	10
range_min	range_min(下, 上) :: range	range	range_min	range	range_min(1, 10)	1
range_max	range_max(下, 上) :: range	range	range_max	range	range_max(1, 10)	10
range_contains	range_contains(範囲, 範囲) :: bool	bool	range_contains	bool	range_contains('1-10', '1-10')	true
range_overlap	range_overlap(範囲, 範囲) :: bool	bool	range_overlap	bool	range_overlap('1-10', '5-15')	true
range_disjoint	range_disjoint(範囲, 範囲) :: bool	bool	range_disjoint	bool	range_disjoint('1-10', '15-20')	true
range_adjacent	range_adjacent(範囲, 範囲) :: bool	bool	range_adjacent	bool	range_adjacent('1-10', '10-20')	true
range_contained_by	range_contained_by(範囲, 範囲) :: bool	bool	range_contained_by	bool	range_contained_by('1-10', '1-10')	true
range_containing	range_containing(範囲, 範囲) :: bool	bool	range_containing	bool	range_containing('1-10', '1-10')	true
range_same_as	range_same_as(範囲, 範囲) :: bool	bool	range_same_as	bool	range_same_as('1-10', '1-10')	true
range_equals	range_equals(範囲, 範囲) :: bool	bool	range_equals	bool	range_equals('1-10', '1-10')	true
range_is_empty	range_is_empty(範囲) :: bool	bool	range_is_empty	bool	range_is_empty('1-10')	false
range_is_nonempty	range_is_nonempty(範囲) :: bool	bool	range_is_nonempty	bool	range_is_nonempty('1-10')	true
range_is_properly_contained_by	range_is_properly_contained_by(範囲, 範囲) :: bool	bool	range_is_properly_contained_by	bool	range_is_properly_contained_by('1-10', '1-10')	false
range_is_properly_containing	range_is_properly_containing(範囲, 範囲) :: bool	bool	range_is_properly_containing	bool	range_is_properly_containing('1-10', '1-10')	false
range_is_subset	range_is_subset(範囲, 範囲) :: bool	bool	range_is_subset	bool	range_is_subset('1-10', '1-10')	true
range_is_superset	range_is_superset(範囲, 範囲) :: bool	bool	range_is_superset	bool	range_is_superset('1-10', '1-10')	true
range_is_disjoint	range_is_disjoint(範囲, 範囲) :: bool	bool	range_is_disjoint	bool	range_is_disjoint('1-10', '15-20')	true
range_is_adjacent	range_is_adjacent(範囲, 範囲) :: bool	bool	range_is_adjacent	bool	range_is_adjacent('1-10', '10-20')	true
range_is_contained	range_is_contained(範囲, 範囲) :: bool	bool	range_is_contained	bool	range_is_contained('1-10', '1-10')	true
range_is_containing	range_is_containing(範囲, 範囲) :: bool	bool	range_is_containing	bool	range_is_containing('1-10', '1-10')	true
range_is_equality	range_is_equality(範囲, 範囲) :: bool	bool	range_is_equality	bool	range_is_equality('1-10', '1-10')	true
range_is_subset_or_disjoint	range_is_subset_or_disjoint(範囲, 範囲) :: bool	bool	range_is_subset_or_disjoint	bool	range_is_subset_or_disjoint('1-10', '1-10')	true
range_is_superset_or_adjacent	range_is_superset_or_adjacent(範囲, 範囲) :: bool	bool	range_is_superset_or_adjacent	bool	range_is_superset_or_adjacent('1-10', '1-10')	true
range_is_subset_or_equality	range_is_subset_or_equality(範囲, 範囲) :: bool	bool	range_is_subset_or_equality	bool	range_is_subset_or_equality('1-10', '1-10')	true
range_is_superset_or_equality	range_is_superset_or_equality(範囲, 範囲) :: bool	bool	range_is_superset_or_equality	bool	range_is_superset_or_equality('1-10', '1-10')	true
range_is_subset_or_equality_or_disjoint	range_is_subset_or_equality_or_disjoint(範囲, 範囲) :: bool	bool	range_is_subset_or_equality_or_disjoint	bool	range_is_subset_or_equality_or_disjoint('1-10', '1-10')	true
range_is_superset_or_equality_or_adjacent	range_is_superset_or_equality_or_adjacent(範囲, 範囲) :: bool	bool	range_is_superset_or_equality_or_adjacent	bool	range_is_superset_or_equality_or_adjacent('1-10', '1-10')	true
range_is_subset_or_equality_or_adjacent_or_disjoint	range_is_subset_or_equality_or_adjacent_or_disjoint(範囲, 範囲) :: bool	bool	range_is_subset_or_equality_or_adjacent_or_disjoint	bool	range_is_subset_or_equality_or_adjacent_or_disjoint('1-10', '1-10')	true
range_is_superset_or_equality_or_adjacent_or_disjoint	range_is_superset_or_equality_or_adjacent_or_disjoint(範囲, 範囲) :: bool	bool	range_is_superset_or_equality_or_adjacent_or_disjoint	bool	range_is_superset_or_equality_or_adjacent_or_disjoint('1-10', '1-10')	true

病院におけるSS-MIX2のデータ精度について

- NHO内の研究チームにおいて本事業の開始「前」から導入されているSS-MIX2モジュールでのデータ精度を調査
- 電子カルテや検査部門システムに残っている検査結果のデータとSS-MIXのストレージ内のデータに齟齬が無い調査
- 4病院で、それぞれランダムに100人選んでカルテレビュー調査をおこなった。
- 結果、データの一致率は98%を超えた
- NHOとしてのデータ精度の結論
 - データは、間違い無く記載されている。表記の統一がきちんとされているかは別の話。
 - データを受け取るデータベースシステムがきちんと解釈できるかどうかの問題。
 - 上記2点をなるべく汎用的に解決することに取り組むべき

30

2 . NCDA システム仕様書

SS-MIX2 を用いた診療情報データベース構築の為の SS-MIX2 モジュール技術仕様書

1. システム要件

国立病院機構の各病院にて「国立病院機構診療情報分析基盤(NCDA)」に参加する為に調達する SS-MIX2 モジュールの機能は以下の通りである。 但し、本体の電子カルテシステム等の仕様上、作成が不可能であるものについては作成を要しない。その場合、何が不可能かを導入標準作業手順書に記載すること。

1.1 SS-MIX2 Ver.1.2d 機能

SS-MIX2 Ver.1.2d に準拠することとして、以下の機能を有すること。

- 日本医療情報学会発行の「SS-MIX2 標準化ストレージ構成の説明と構築ガイドライン Ver.1.2d」, 「SS-MIX2 拡張ストレージ構成の説明と構築ガイドライン Ver.1.2d」, 「SS-MIX2 標準化ストレージ仕様書 Ver.1.2d」, 「標準化ストレージ仕様書別紙：コード表 Ver.1.2d」, 「SS-MIX2 拡張ストレージ構成の説明と構築ガイドライン Ver.1.2d 別紙：標準文書コード表」に記載している仕様に対応していること。(尚、当初 Ver.1.2c 準拠としていたが、標準ストレージ部分では Ver.1.2c からの変更点について影響がないため Ver.1.2d 準拠ということとした。)
- 標準化ストレージ、拡張ストレージ、トランザクションストレージ、インデックスデータベースの4つのファイルを生成すること。
- 標準化ストレージにはデータ種別として36種のデータを出力すること。

(表 1-1 標準化ストレージ格納データ)

No	データ種別	種別名称	HL7メッセージ型
1	ADT-00	患者基本情報の更新	ADT^A08
2	ADT-00	患者基本情報の削除	ADT^A23

No	データ種別	種別名称	HL7メッセージ型
3	ADT-01	担当医の変更	ADT^A54
4	ADT-01	担当医の取消	ADT^A55
5	ADT-12	外来診察の受付	ADT^A04
6	ADT-21	入院予定	ADT^A14
7	ADT-21	入院予定の取消	ADT^A27
8	ADT-22	入院実施	ADT^A01
9	ADT-22	入院実施の取消	ADT^A11
10	ADT-31	外出泊実施	ADT^A21
11	ADT-31	外出泊実施の取消	ADT^A52
12	ADT-32	外出泊帰院実施	ADT^A22
13	ADT-32	外出泊帰院実施の取消	ADT^A53
14	ADT-41	転科・転棟(転室・転床)予定	ADT^A15
15	ADT-41	転科・転棟(転室・転床)予定の取消	ADT^A26
16	ADT-42	転科・転棟(転室・転床)実施	ADT^A02
17	ADT-42	転科・転棟(転室・転床)実施の取消	ADT^A12

No	データ種別	種別名称	HL7 メッセージ型
18	ADT-51	退院予定	ADT^A16
19	ADT-51	退院予定の取消	ADT^A25
20	ADT-52	退院実施	ADT^A03
21	ADT-52	退院実施の取消	ADT^A13
22	ADT-61	アレルギー情報の登録 / 更新	ADT^A60
23	PPR-01	病名 (歴) 情報の登録 / 更新	PPR^ZD1
24	OMD	食事オーダー	OMD^O03
25	OMP-01	処方オーダー	RDE^O 11
26	OMP-11	処方実施通知	RAS^O17
27	OMP-02	注射オーダー	RDE^O 11
28	OMP-12	注射実施通知	RAS^O17
29	OML-01	検体検査オーダー	OML^O33
30	OML-11	検体検査結果通知	OUL^R22
31	OMG-01	放射線検査オーダー	OMG^O19
32	OMG-11	放射線検査の実施通知	OMI^Z23

No	データ種別	種別名称	HL7 メッセージ型
33	OMG-02	内視鏡検査オーダ	OMG^O19
34	OMG-12	内視鏡検査の実施通知	OMI^Z23
35	OMG-03	生理検査オーダ	OMG^O19
36	OMG-13	生理検査結果通知	ORU^R01

「SS-MIX2 標準化ストレージ構成の説明と構築ガイドライン Ver.1.2d p11」

1.2 拡張ストレージへの出力機能

現在の SS-MIX2 モジュールでオプションとして既に導入している拡張ストレージへの出力機能は、そのまま提供すること。また、1.3.0 で規定する出力を行うこと。

1.3 NHO 対応としての設定

1.3.0 拡張ストレージへの出力機能

各社の SS-MIX2 モジュールの拡張ストレージへの出力機能を利用し、以下の情報を出力すること。その際、日本医療情報学会発行の「SS-MIX2 拡張ストレージ構成の説明と構築ガイドライン Ver.1.2d」に記載している仕様に対応していること。また、トランザクションストレージ、インデックスデータベースも同時に生成すること。

No	データ種別	種別名称	HL7 メッセージ型
1	L-OBSERVATIONS^OBSERVATION S^99ZL01	バイタル検査結果	HL7 V2.5 ORU^R30
2	^(ローカル名称)^11506-3^経過記	診療録(外来/入院含)	HL7 CDA R2

No	データ種別	種別名称	HL7メッセージ型
	録^LN	む)	
2.1	^(ローカル名称)^34108-1^外来診療録^LN	診療録(外来)(入院・外来が別の場合)	HL7 CDA R2
2.2	^(ローカル名称)^34112-3^入院診療録^LN	診療録(入院)(入院・外来が別の場合)	HL7 CDA R2
3	^(ローカル名称)^18842-5^退院時サマリー^LN	退院時サマリー	HL7 CDA R2
4	^(ローカル名称)^57133-1^紹介状^LN	診療情報提供書	HL7 CDA R2

1.3.1 バイタル検査結果通知の出力

(1) バイタル検査結果通知のデータを、別紙の形式で拡張ストレージに出力する。尚、「診療日」に出力する日付はOBX-14 トランザクション日時(測定した日)とする。

(2) ファイル作成の単位は、データの格納構造として日付の下にあるため、最大でも一日分が1ファイルにまとまっている形とする。一日の中で測定のたびに作成するのでも良い。一日1ファイルなら、特定キーは測定日を出力する。一日に複数回のデータを出力する場合は、特定キーに測定日の時間まで(YYYYMMDDHH)出力すること。

1.3.2 バイタルデータの項目及び形式等

(1) バイタルデータとして取得する項目は、「拡張期血圧、収縮期血圧、脈拍数、呼吸数、体温」の5項目とする。

(2) OBX-3 検査項目に出力するコードはJLAC10コードとする。バイタルデータを参考に適切なJLAC10を選択すること。

(3) 上記以外の項目を SS-MIX2 に出力することは問題ないが、今回の対応では扱わない。但し、今後の検討で仕様として扱うことになる場合は、JLAC10 コードを基準とした標準コードを必須とすることを想定している。この今後想定される検査項目は別表として提供する。

1.3.3 標準コード変換機能

SS-MIX2 データの出力に際しては、コードのマッピング表などに従って、院内のローカルコードを厚労省が定める標準コードに変換する機能を有すること。またマッピング表については、容易にその内容を変更できるマスターメンテナンスプログラム等の機能を有すること。

JLAC10 コード、JANIS コード、HOT コードについては、機構病院が NCDA 事業に参加する場合においては機構から提供する。

1.3.4 標準化ストレージにおける文字コードについて

メッセージの文字コードについては、「標準化ストレージガイドライン」で示されているとおり、1 バイト系文字は ISO IR-6 (ASCII)、2 バイト系文字は ISO IR87 (JIS X 0208 第一水準、第二水準) とする。ただし現実には上記以外の文字コードが電子カルテシステムに登録されている可能性があるため、以下のように対応することとする。

- 1 半角カナ文字 → 全角カナ文字に置き換えて SS-MIX2 に出力する。
- 2 外字 → ■で置き換えて SS-MIX2 に出力する。
- 3 環境依存文字については変換表を機構より提供するのでそれにより変換して SS-MIX2 に出力する。

1.3.5 単位の文字表記の統一

SS-MIX2 データの出力に際して、臨床検査データの OBX セグメントの 6 フィールド目の単位の文字表記を統一すること。

【単位の文字表記の統一ルール例】ASCII コードで表記すること

- ・かける → . (ドット)
- ・乗 → * (アスタリスク)
- ・μ → u (小文字ユー)
- ・語尾に名称 → () で

- → cel
- ‰ → permil
- 個 → pcs

【上記ルールの適用例】

- mL → mL (ASCIIコード)
- $\times 10^2 / \mu\text{l}$ → $\cdot 10^2 / \text{uL}$ (かける、乗、 μ)
- /HPF → /(hpf) (語尾に名称)

1.3.6 単位変換機能

SS-MIX2 データの出力に際して臨床検査データの単位に関しては、JLAC10 コードごとに、機構が定める単位に変換を行った上で SS-MIX2 データを生成すること。尚、JLAC10 コード別の単位表は別途機構から提供する。単位表は「SS-MIX2 標準化ストレージ仕様書 Ver.1.2」にも別表として添付する。

【単位変換例】

JLAC10 コード	数値	単位	→	JLAC10 コード	数値	単位
1A02500000127201	10.5	mg/l	→	1A02500000127201	1.05	mg/dL

1.3.7 計測値等の表記方法について

(1) 定性値・検出限界以下・検出限界以上の表記

- OBX (検体検査結果) セグメントの5フィールド目(検査値)に検査結果を記述する場合、現在そのデータ形式はOBX-2フィールドの説明にあるようにNM型、ST型、CWE型のうちいずれかの形式で記述することとなっている。
- 今回の仕様では、定性値・検出限界以下・検出限界以上のデータについては、SN型の表現方法を用いてSN型の”^”を” “(スペース)に置き換える。
- この件の説明は、「SS-MIX2 標準化ストレージ仕様書 Ver.1.2」 P104 表 3-77 検査結果セグメント(OBX)定義 のOBX-2の項目説明にも記述する。

(2) 複数の要素が一つの値で表現されている場合の表記

複数の要素が組み合わせられ一つの結果値として表記されている場合は、それぞれの要素に分離して表記すること。例えば定量値とクラス値が組み合わせられた結果値については、定量値とクラス値に分離する。

【定量値とクラス値の分離の例】

定量値とクラス値が組み合わせられた例

検査名称		院内コード	結果値
ムンプス Virus IgG		001591	2.3(±)
↓			
定量値とクラス値を分離した例			
SS-MIX2 標準コード	院内コード	結果値	備考
5F432143102302304	001591	2.3	
5F432143102302311	001591	+-	(半角スペース2つプラスマイナス)

1.3.8 トランザクションストレージのデータ保持期間

トランザクションストレージのデータ保持期間は、現在の標準化ストレージ及び拡張ストレージを作っているデータの再現に必要な分だけ保持しておくこと。

1.3.9 ST 型の長さ

- RXE-23(与薬速度)は ST 型で長さが 6 であるが、正負の記号と小数点を考慮し(例: +266.865)、本事業では 8 桁まで許容するものとする。

- CX型は先頭成分がST型で長さが15であるが、IN1-10(被保険者グループ雇用者ID)に長い名称の保険者が出力される場合などを考慮し、本事業ではCX型の先頭成分は30桁まで許容するものとする。
- XAD型は第8成分(その他地理表示)がST型で長さが50であるが、全角50文字(100バイト)と解釈しているシステムがあり半角文字で100文字登録出来るため、本事業ではXAD型の第8成分は100桁まで許容するものとする。

1.3.10 トランザクションストレージのファイル切り替え機能

SS-MIX2の仕様上、トランザクションストレージはカレントの日付が変わった時点、もしくは記録中のトランザクションデータファイルのファイルサイズが一定量を超えた時点で、新たなファイルを作成して記録先を切り替えるものとなっているが、同一日付内において一定時刻(例えば17:00)を経過した時点で記録先を切り替える機能を追加する。

3 . 倫理審査における計画書

**電子カルテ情報をセマンティクス（意味・内容）の標準化により分析
可能なデータに変換するための研究**

研究責任者：堀口 裕正

独立行政法人国立病院機構本部 総合研究センター
診療情報分析部 副部長

事務局/研究主催

独立行政法人国立病院機構本部 総合研究センター
診療情報分析部

堀口 水本

〒152-8621 目黒区東が丘 2 5 21

TEL: 03-5712-5133

FAX: 03-5712-5134

E-Mail : horiguchi-hiromasa@hosp.go.jp

第 1.0 版：2017 年 1 月 18 日

1. 背景

本研究では、電子カルテに実装可能な医療用語の標準化を行うシステムの開発に向け、そのステップとして、臨床現場からのニーズが極めて大きい退院サマリの自動生成技術の実現を目的とする。これは用語の標準化を目的とする研究として遠回りの課題設定である。しかし、電子カルテの自動解析は技術的な難易度が高く、実用的な精度を実現するためには多額の研究開発投資が求められる。そこで、本研究提案では、医療現場に直接的なメリットが生じる研究課題に取り組むことによって、現場の協力と今後の追加的な研究開発投資を呼び込み、その過程を通じて実用性の高い電子カルテの自動解析技術を実現する。

2. 目的

本研究は、電子カルテに実装可能な医療用語の標準化を行うシステムの開発に向け、そのステップとして、臨床現場からのニーズが極めて大きい退院サマリの自動生成技術の実現を目的とする

3. 研究方法

3 - 1. 研究実施場所

研究実施場所は、国立病院機構本部総合研究センター診療情報分析部(以下、診療情報分析部)研究室及び本部内分析室並びに静岡大学情報学部行動情報学科狩野研究室、岡山大学大学院医歯薬学総合研究科クリニカルバイオバンクネットワーク事業化研究講座研究室、国立保健医療科学院研究情報支援研究センター研究室とする。

3 - 2. 研究実施期間

研究実施期間は、倫理審査委員会承認後より2020年3月31日までとする。

3 - 3. 研究対象医療機関と対象患者

研究対象医療機関は、国立病院機構病院に所属するDPC病院のうち、診療情報集積基盤(以下、NCDA)を運用しデータ提供を行う医療機関とする。

対象患者は2016年1月1日から2019年12月31日までに入院し、退院時サマリを作成した全患者とする。

3 - 4 . 対象データ

研究に用いるデータは、研究対象医療機関より診療情報分析部に提供されたDPCデータおよびレセプトデータ、ならびにSS-MIX2ストレージに格納された情報から抽出した医師記録、退院サマリおよび入院中の検査結果、食事内容および処方内容である。

3 - 5 . 分析方法

(1) 対象

退院サマリを作成した全患者

(2) アウトカム

入院中に記載/記録された情報から退院サマリを自動生成する技術を開発すること

(3) 抽出する項目

入院中の医師記録・退院サマリ・入院中の検査結果、食事内容および処方内容

(4) 解析方法

入院中に記載/記録された情報を元データに、機械学習により自動的に情報収集を行い、退院サマリを自動で作成する。その作成結果と、実際の医師の書いた退院時サマリを比較/検討を行い、自動作成技術の能力評価を行い、またその能力の改善を行っていく。

4 . 倫理的配慮

本研究は、ヘルシンキ宣言、人を対象とする医学系研究に関する倫理指針(以下、倫理指針)に基づいて実施する。

4 - 1 . インフォームド・コンセント

本研究は既存試料・情報を用いて実施し、人体から取得された試料は用いない。研究対象者等からインフォームド・コンセントは受けないが、倫理指針「第12の1(2)イ」に則り、本計画書の4-3に記す通り、利用目的を含む本研究についての情報を研究対象者等に公開し、研究が実施されることについて研究

対象者が拒否できる機会を保障する。なお、NCDA 運用による診療情報の蓄積・利活用についての説明及び同意は、各施設での掲示で既に行われている。

4 - 2 . データ管理、個人情報等の取り扱いに関する配慮

研究の実施並びに種々のデータの収集及び取り扱いにおいては、国立病院機構診療情報データベース利活用規程に従うとともに、患者情報の機密保持に充分留意する。

本研究で用いるデータは、研究対象医療機関に 2016 年 1 月 1 日から 2019 年 12 月 31 日までに退院サマリを作成した全患者のデータであり、個人情報等を取り扱う。倫理指針「第 15 の 2 (1)」及び国立病院機構診療情報データベース利活用規程に則り、保有する個人情報等について、漏えい、滅失又はき損の防止その他の安全管理のため、下記の措置を講じる。

データは研究対象医療機関で収集され、本部 IT 推進部に提出される。データが保管されるサーバーを国立病院機構本部 2 階のセキュリティルームに設置し、セキュリティルーム内で IT 推進部システム開発専門職が匿名化処理を行う。研究者は匿名化後のデータを用いて本部内分析室において分析を実施する。

保有する個人情報に関する事項の公表等については、倫理指針「第 12 の 1 (2) イ」、「第 16 の 1 (1)」及び国立病院機構診療情報データベース利活用規程第 6 条第 3 項に則り、個人情報の取扱いを含む研究の実施についての情報を研究対象者等に公開する。

4 - 3 . 本研究における情報公開

本研究では、倫理審査委員会承認後、倫理指針「第 12 の 1 (2) イ」、「第 16 の 1 (1)」及び国立病院機構診療情報データベース利活用規程第 6 条第 3 項に則り、本部ホームページにおいて、本研究の意義、目的及び方法、研究機関、保有する個人情報に関して利用目的の通知、開示、訂正等又は利用停止の求めに応じる手続き並びに保有する個人情報に関する問い合わせや苦情等の窓口の連絡先に関する情報を公開する（公表する情報については別添資料を参照）。

4 - 4 . 研究成果の公表

本研究の成果は、報告書で公表するとともに、学会・論文で発表する。また、本研究結果を内包したソフトウェアの公表を実施する。データの集計・分析結果については、集団を記述した数値データもしくは機械学習の学習結果データとし、個人が同定されるデータの公表は行わない。

5 . 研究経費

本研究は、厚生労働科学研究費補助金（臨床研究等 I C T 基盤構築研究事業）「電子カルテ情報をセマンティクス（意味・内容）の標準化により分析可能なデータに変換するための研究」（代表 堀口裕正）を用いて研究を実施する

6 . 研究組織

総合研究センター診療情報分析部が主体となり、本部医療部、保険医療科学院、静岡大学、岡山大学等から協力を得て、研究を行う。

【研究代表者】

国立病院機構本部総合研究センター診療情報分析部 副部長 堀口 裕正

【共同研究者】

国立病院機構本部 企画役 岡田 千春
静岡大学情報学部行動情報学科 准教授 狩野 芳伸
岡山大学大学院医歯薬学総合研究科
クリニカルバイオバンクネットワーク
事業化研究講座研究室 准教授 森田 瑞樹
国立保健医療科学院研究情報支援研究センター
特命上席主任研究官 奥村 貴史

別添

「電子カルテ情報をセマンティクス（意味・内容）の標準化により分析可能なデータに変換するための研究」研究実施に関するお知らせ