

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6956107号

(P6956107)

(45) 発行日 令和3年10月27日 (2021. 10. 27)

(24) 登録日 令和3年10月6日 (2021.10. 6)

(51) Int. Cl. F I
G 0 6 F 16/25 (2019.01) G O 6 F 16/25
G 1 6 H 10/60 (2018.01) G 1 6 H 10/60

請求項の数 24 (全 20 頁)

(21) 出願番号	特願2018-553440 (P2018-553440)	(73) 特許権者	590000248
(86) (22) 出願日	平成29年4月19日 (2017. 4. 19)		コーニンクレッカ フィリップス エヌ ヴェ
(65) 公表番号	特表2019-514128 (P2019-514128A)		KONINKLIJKE PHILIPS N. V.
(43) 公表日	令和1年5月30日 (2019. 5. 30)		オランダ国 5656 アーヘー アイン ドーフエン ハイテック キャンパス 5 2
(86) 国際出願番号	PCT/EP2017/059266		
(87) 国際公開番号	W02017/182509	(74) 代理人	100122769
(87) 国際公開日	平成29年10月26日 (2017. 10. 26)		弁理士 笛田 秀仙
審査請求日	令和2年4月9日 (2020. 4. 9)	(74) 代理人	100163809
(31) 優先権主張番号	62/324, 363		弁理士 五十嵐 貴裕
(32) 優先日	平成28年4月19日 (2016. 4. 19)	(74) 代理人	100171701
(33) 優先権主張国・地域又は機関	米国 (US)		弁理士 浅村 敬一

最終頁に続く

(54) 【発明の名称】 明確な照合情報を持たない識別不能のヘルスケアデータベースの病院マッチング

(57) 【特許請求の範囲】

【請求項 1】

以下の動作を含むデータベース統合プロセスを実行することによって、 N (N は少なくとも3の値を有する正の整数) 個の匿名化されたヘルスケアデータベースを統合するようにプログラムされた少なくとも1つの電子プロセッサを有し、

前記動作は、

N 個の匿名化されたヘルスケアデータベースのうちの一組のデータベース i , j に関し、前記一組のデータベース i , j の両方のデータベース i , j に各々が含まれる特徴のセットを識別するとともに、前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベースの患者とマッチングする変換テーブルを生成する動作と、

$N(N-1)/2$ 個の変換テーブルを生成するために、前記 N 個の匿名化されたヘルスケアデータベースのデータベースの各一意の組について前記識別及び生成動作を繰り返す動作と、

を有し、

前記少なくとも1つの電子プロセッサは、前記 $N(N-1)/2$ 個の変換テーブルを用いて前記 N 個の匿名化されたヘルスケアデータベースに含まれる1又は複数の匿名化された患者について患者データを検索する動作を含む患者データ検索プロセスを実行するように更にプログラムされる、匿名化されたヘルスケアデータソース装置。

【請求項 2】

前記一組のデータベース i , j に関して前記特徴のセットを識別する動作が、特徴精度

10

20

メトリックが前記一組のデータベース i , j の各匿名化されたヘルスケアデータベースについての最小精度を満たす特徴を識別する動作を含む、請求項 1 記載の装置。

【請求項 3】

前記 N 個の匿名化されたデータベースに含まれる前記患者データを検索する動作が、クエリ特徴に関し、

前記クエリ特徴が、前記 N 個の匿名化されたヘルスケアデータベースの 1 つにのみ含まれる場合、前記クエリ特徴を含む匿名化されたヘルスケアデータベースから前記クエリ特徴を検索する動作、及び、

前記クエリ特徴が、前記 N 個の匿名化されたヘルスケアデータベースの 2 つ以上に含まれる場合、前記クエリ特徴を含む各々の匿名化されたヘルスケアデータベースにおける前記クエリ特徴に関する特徴精度メトリックに基づいて、前記クエリ特徴を含む前記 N の匿名化されたヘルスケアデータベースのうちの 2 つ以上において前記クエリ特徴の値から前記クエリ特徴に関する検索された値を生成する動作、

を含む、請求項 1 又は 2 に記載の装置。

【請求項 4】

前記変換テーブルを生成する動作が、 m (m は、前記一組のデータベース i , j においてマッチした患者の数) $\times 2$ の変換テーブルを生成する動作を含む、請求項 1 乃至 3 のいずれか 1 項に記載の装置。

【請求項 5】

前記データベース統合プロセスが、前記 $N(N - 1) / 2$ 個の変換テーブルの間でマッチングする患者の整合性に基づいて、前記 $N(N - 1) / 2$ 個の変換テーブルを改善する更なる動作を含む、請求項 1 乃至 4 のいずれか 1 項に記載の装置。

【請求項 6】

前記改善する動作が、前記識別された特徴のセットを使用しない、請求項 5 記載の装置。

【請求項 7】

前記データベース統合プロセスが、前記 N 個の匿名化されたヘルスケアデータベースのうちの少なくとも一組のデータベースに関し、各事象のタイムスタンプ間における時間間隔 t によって分離される一組のタイムスタンプされた事象によって規定される少なくとも 1 つの長さ方向の特徴を識別する動作、及び、

前記 2 つのデータベース i , j における患者に関する前記時間間隔 t の比較を含む前記長さ方向の特徴のマッチングに部分的に基づいて、前記一組のデータベースの患者にマッチングする前記変換テーブルを生成する動作、
を含む、請求項 1 乃至 6 のいずれか 1 項に記載の装置。

【請求項 8】

前記長さ方向の特徴のマッチングに部分的に基づいて前記一組のデータベースの患者にマッチングする前記変換テーブルを生成する動作が、前記 2 つのデータベース i , j における患者に関する事象のタイムスタンプの比較を含まない、請求項 7 記載の装置。

【請求項 9】

以下の動作を含むデータベース統合プロセスを実行することによって、ヘルスケアデータベース i とヘルスケアデータベース j とを統合するようにプログラムされた少なくとも 1 つの電子プロセッサを有し、

前記動作は、

一組のデータベース i , j に関し、事象のタイムスタンプ間の時間間隔 t によって分離される一組のタイムスタンプされた事象によって規定される少なくとも 1 つの長さ方向の特徴を含む前記一組のデータベース i , j の両方のデータベース i , j に各々含まれる特徴のセットを識別するとともに、前記 2 つのデータベース i , j における患者に関する時間間隔 t の比較を含む前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベース i , j の患者にマッチングする変換テーブルを生成する動作を含み、

10

20

30

40

50

前記少なくとも1つの電子プロセッサは、前記一組のデータベース*i*、*j*の患者とマッチングする前記変換テーブルを用いて、両方の匿名化されたヘルスケアデータベース*i*、*j*に含まれる1又は複数の匿名化された患者に関する患者データを検索する動作を含む患者データ検索プロセスを実行するように更にプログラムされる、匿名化されたヘルスケアデータソース装置。

【請求項10】

患者類似度に基づいて前記一組のデータベース*i*、*j*の患者とマッチングする前記変換テーブルを生成する動作が、前記2つのデータベース*i*、*j*における患者に関する事象のタイムスタンプの比較を含まない、請求項9記載の装置。

【請求項11】

前記特徴のセットを識別する動作が、前記一組のデータベース*i*、*j*の両方のデータベース*i*、*j*に含まれる非長さ方向の特徴のセットを識別する動作を含み、各データベース*i*、*j*における各患者に関し、前記患者に関する前記非長さ方向の特徴のセットの値の連結を有する患者に関する統一識別子を生成する動作を含み、

前記変換テーブルを生成する動作が、前記2つのデータベース*i*、*j*における患者に関する前記統一識別子の比較を更に含む前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベース*i*、*j*の患者とマッチングする前記変換テーブルを生成する動作を含む、請求項9又は10に記載の装置。

【請求項12】

前記特徴のセットを識別する動作が、前記特徴を抽出するための患者記録のテキストコンテンツについての自然言語処理を実行することによって、前記一組のデータベース*i*、*j*の少なくとも1つのデータベースにおける少なくとも1つの特徴を識別する動作を含む、請求項9乃至11のいずれか1項に記載の装置。

【請求項13】

前記一組のデータベース*i*、*j*の両方のデータベース*i*、*j*に各々含まれる前記特徴のセットを識別する動作が、特徴精度メトリックが、前記匿名化されたヘルスケアデータベース*i*と前記匿名化されたヘルスケアデータベース*j*との両方に関する最小精度を満たす特徴を識別する動作を含む、請求項9乃至12のいずれか1項に記載の装置。

【請求項14】

前記一組のデータベース*i*、*j*の患者にマッチングする前記変換テーブルを用いて両方の匿名化されたヘルスケアデータベース*i*、*j*に含まれる前記患者データを検索する動作が、クエリ特徴に関し、

前記クエリ特徴が、前記一組の匿名化されたヘルスケアデータベース*i*、*j*の1つのデータベースにのみ含まれる場合、前記クエリ特徴を含む前記匿名化されたヘルスケアデータベースから前記クエリ特徴を検索する動作、及び、

前記クエリ特徴が、前記一組の匿名化されたヘルスケアデータベース*i*、*j*の両方のデータベースに含まれる場合、前記クエリ特徴を含む各々匿名化されたヘルスケアデータベースにおける前記クエリ特徴に関する特徴精度メトリックに基づいて、前記一組の匿名化されたヘルスケアデータベース*i*、*j*における前記クエリ特徴の値から前記クエリ特徴に関する検索された値を生成する動作、

を含む、請求項9乃至13のいずれか1項に記載の装置。

【請求項15】

前記変換テーブルを生成する動作が、 m (m は、前記一組のデータベース*i*、*j*においてマッチした患者の数) $\times 2$ の変換テーブルを生成する動作を含む、請求項9乃至14のいずれか1項に記載の装置。

【請求項16】

前記少なくとも1つの電子プロセッサが、 $N(N-1)/2$ 個の変換テーブルを生成するために、前記 N 個の匿名化されたヘルスケアデータベースの各々一意の組のデータベースに関する識別及び生成動作を繰り返す更なる動作を含む前記データベース統合プロセスを実行することによって、前記匿名化されたヘルスケアデータベース*i*、前記匿名化され

10

20

30

40

50

たヘルスケアデータベース j 、及び、少なくとも 1 つの追加的な匿名化されたヘルスケアデータベースを含む N 個のデータベースを統合するようにプログラムされ、

前記少なくとも 1 つの電子プロセッサが、前記匿名化されたヘルスケアデータベース i 、 j の 1 つにおける患者の患者 ID を受信し、前記 $N(N - 1) / 2$ 個の変換テーブルを用いて前記 N 個の匿名化されたヘルスケアデータベースに含まれる前記患者に関する患者データを検索する動作を含む前記患者データ検索プロセスを実行するように更にプログラムされる、請求項 9 乃至 15 のいずれか 1 項に記載の装置。

【請求項 17】

N (N は、少なくとも 2 の値を持つ正の整数) 個の匿名化されたヘルスケアデータベースから匿名化された母集団イメージを再構成する匿名化された母集団イメージ再構成方法
10
を実行するためのコンピュータによって読み取り可能及び実行可能な命令を格納する非一時的な記録媒体であって、前記匿名化された母集団イメージ再構成方法は、

前記 N 個の匿名化されたヘルスケアデータベースの一組のデータベース i 、 j に関し、前記一組のデータベース i 、 j の両方のデータベース i 、 j に各々含まれる特徴のセットを識別するとともに、前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベースの患者にマッチングする変換テーブルを生成する動作と、

前記 $N(N - 1) / 2$ 個の変換テーブルによって統合される前記 N 個の匿名化されたヘルスケアデータベースのコンテンツを有する前記匿名化された母集団イメージを生成するために、前記 N 個の匿名化されたヘルスケアデータベースの各一意の組のデータベースに関して前記識別及び生成動作を繰り返す動作と、
20
を有する、非一時的な記録媒体。

【請求項 18】

前記格納された命令が、匿名化された母集団データクエリを受信するとともに、前記匿名化された母集団データクエリに応じて、前記 $N(N - 1) / 2$ 個の変換テーブルを用いて前記匿名化された母集団イメージから患者データを検索する動作を含む匿名化された母集団イメージデータ検索方法を更に実行するためのコンピュータによって読み取り可能及び実行可能である、請求項 17 記載の非一時的な記録媒体。

【請求項 19】

N が、少なくとも 3 の値を持つ正の整数である、請求項 17 又は 18 に記載の非一時的な記録媒体。
30

【請求項 20】

前記変換テーブルを生成する動作が、 m (m は、前記一組のデータベース i 、 j においてマッチした患者の数) $\times 2$ 個の変換テーブルを生成する動作を含み、これにより、前記 $N(N - 1) / 2$ 個の変換テーブルの各々は、 $m \times 2$ の変換テーブルである、請求項 19 記載の非一時的な記録媒体。

【請求項 21】

前記匿名化された母集団イメージ再構成方法が、前記 $N(N - 1) / 2$ 個の変換テーブルの間でマッチングする患者の整合性に基づいて、前記 $N(N - 1) / 2$ 個の変換テーブルを改善する更なる動作を含む、請求項 19 又は 20 に記載の非一時的な記録媒体。

【請求項 22】

前記改善する動作が、前記 $N(N - 1) / 2$ 個の変換テーブルについて作用し、前記識別された特徴のセットを使用しない、請求項 21 記載の非一時的な記録媒体。
40

【請求項 23】

前記匿名化された母集団イメージ再構成方法が、前記 N 個の匿名化されたヘルスケアデータベースの少なくとも一組のデータベースに関し、

事象のタイムスタンプ間の時間間隔 t によって分離される一組のタイムスタンプされる事象によって規定される少なくとも 1 つの長さ方向の特徴を識別する動作、及び、

前記 2 つのデータベース i 、 j における患者に関する前記時間間隔 t の比較を含む前記長さ方向の特徴のマッチングに部分的に基づいて、前記一組のデータベースの患者にマッチングする前記変換テーブルを生成する動作、
50

を含む、請求項 1 7 乃至 2 2 のいずれか 1 項に記載の非一時的な記録媒体。

【請求項 2 4】

前記長さ方向の特徴のマッチングに部分的に基づいて前記一組のデータベースの患者にマッチングする前記変換テーブルを生成する動作が、前記 2 つのデータベース i, j における患者に関する事象のタイムスタンプの比較を含まない、請求項 2 3 記載の非一時的な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

以下は、概して、医療研究及び開発技術、ヘルスケアデータベースのキュレーション技術、ヘルスケアデータマイニング技術、及び、関連技術に関する。

【背景技術】

【0002】

ヘルスケア研究開発の多くの領域は、医療患者のデータを含むヘルスケアデータベースを活用している。病院又は他の医療施設によって、及び／又は、心臓ケア機関（CCU）、集中治療室（CCU）、又は、緊急時アドミタンス部門などの個別の機関によって、病歴又は他の臨床データ、患者課金データ、病院ベッド占有率などの事項に関する管理記録が維持される。これらのデータベースには、（米国では）医療保険の相互運用性及び説明責任に関する法律（HIPAA）などの財政及び／又は医療プライバシー法の下で一般に秘密裏に維持されなければならない繊細な患者データが保存される。

【発明の概要】

【発明が解決しようとする課題】

【0003】

患者のプライバシーを維持しながら、患者データベースを臨床、病院管理、又は、他の目的のためのデータ分析に使用できるようにするために、患者識別情報（PII）を取り除くことによってデータベースを匿名化することが知られている。匿名化する必要のある情報には、患者名及び／又は医療識別番号（適切にランダムに割り当てられた番号などで置き換えられる）、住所などが含まれる。他の匿名化手段には、珍しい特徴の組み合わせによって識別可能な「まれな」患者の除去が含まれてもよい。例えば、特定の病気の102歳の患者は、その情報のみに基づいて識別され得るためである。

【0004】

まれな患者に加えて、患者は、患者記録に記録された事象のタイムスタンプ情報に基づいて識別可能であるかもしれない。例えば、患者が特定の状態で特定の日付に病院に入院した場合、その情報は、可能性のある患者識別の数を少数に減らすのに十分であり得る。しかしながら、長さ方向の情報、即ち、事象の時間的順序及び種々の事象間の時間間隔は、医療データ解析において有用であることがある。例えば、入院と退院との間の時間間隔は、病院の効率及び／又は特定の治療の有効性を分析するために有用であり、又は、さらに重要であり得る。タイムスタンプを使用して匿名化された患者を特定するための潜在的な可能性を減らし、医療データ分析の価値のある可能性のある長さ方向の情報を保持するために、一部の匿名化されたデータベースでは、タイムスタンプが、所与の患者の全てのタイムスタンプされた事象のための固定シフトを用いて、ランダムな量（一般に各患者ごとに異なる）だけシフトされる。タイムスタンプのランダムな固定時間シフトは、タイムスタンプによる患者の識別をより困難にするが、特に固定時間シフトの使用は、長さ方向の情報、即ち、事象間の時間間隔情報を保持する。

【課題を解決するための手段】

【0005】

本開示の一態様では、匿名化されたヘルスケアデータソース装置が、以下の動作を含むデータベース統合プロセスを実行することによって、 N （ N は少なくとも3の値を有する正の整数）個の匿名化されたヘルスケアデータベース10を統合するようにプログラムされた少なくとも1つの電子プロセッサを有し、前記動作は、 N 個の匿名化されたヘルスケ

データベースのうちの一組のデータベース i, j に関し、前記一組のデータベース i, j の両方のデータベース i, j に各々が含まれる特徴のセットを識別するとともに、前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベースの患者とマッチングする変換テーブルを生成する動作と、 $N(N-1)/2$ 個の変換テーブルを生成するために、前記 N 個の匿名化されたヘルスケアデータベースのデータベースの各一意の組について前記識別及び生成動作を繰り返す動作と、を有する。上記少なくとも 1 つ少なくとも 1 つの電子プロセッサは、 $N(N-1)/2$ 個の変換テーブルを用いて前記 N 個の匿名化されたヘルスケアデータベースに含まれる 1 又は複数の匿名化された患者について患者データを検索する動作を含む患者データ検索プロセスを実行するように更にプログラムされる。

10

【0006】

本開示の他の態様では、以下の動作を含むデータベース統合プロセスを実行することによって、ヘルスケアデータベース i とヘルスケアデータベース j とを統合するようにプログラムされた少なくとも 1 つの電子プロセッサを有し、前記動作は、前記一組のデータベース i, j に関し、事象のタイムスタンプ間の時間間隔 t によって分離される一組のタイムスタンプされた事象によって規定される少なくとも 1 つの長さ方向の特徴を含む前記一組のデータベース i, j の両方のデータベース i, j に各々含まれる特徴のセットを識別するとともに、前記 2 つのデータベース i, j における患者に関する時間間隔 t の比較を含む前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベース i, j の患者にマッチングする変換テーブルを生成する動作を含む。上記少なくとも 1 つの電子プロセッサは、一組のデータベース i, j の患者とマッチングする前記変換テーブルを用いて、両方の匿名化されたヘルスケアデータベース i, j に含まれる 1 又は複数の匿名化された患者に関する患者データを検索する動作を含む患者データ検索プロセスを実行するように更にプログラムされる。

20

【0007】

本開示の他の態様では、非一時的な記録媒体が、 N (N は、少なくとも 2 の値を持つ正の整数) 個の匿名化されたヘルスケアデータベースから匿名化された母集団イメージを再構成する匿名化された母集団イメージ再構成方法を実行するためのコンピュータによって読み取り可能及び実行可能な命令を格納する。匿名化された母集団イメージ再構成方法は、前記 N 個の匿名化されたヘルスケアデータベースの一組のデータベース i, j に関し、前記一組のデータベース i, j の両方のデータベース i, j に各々含まれる特徴のセットを識別するとともに、前記特徴のセットによって測定される患者類似度に基づいて、前記一組のデータベースの患者にマッチングする変換テーブルを生成する動作を有する。また、 $N(N-1)/2$ 個の変換テーブルによって統合される前記 N 個の匿名化されたヘルスケアデータベースのコンテンツを有する前記匿名化された母集団イメージを生成するために、前記 N 個の匿名化されたヘルスケアデータベースの各一意の組のデータベースに関する識別及び生成動作が繰り返される。

30

【0008】

1 つの利点は、2 つ、3 つ、4 つ、又は、それ以上の匿名化されたヘルスケアデータベースの統合を提供し、ヘルスケアデータ分析タスクのためにデータベースに含まれる結合データを活用することにある。

40

【0009】

別の利点は、1 又は複数の匿名化されたヘルスケアデータベースが非構造化ヘルスケアデータベースである前述のものを提供することにある。

【0010】

別の利点は、異なる匿名化されたヘルスケアデータベース内の匿名化された患者を照合する際に、長さ方向の情報、即ち事象間の時間間隔が活用される前述のものを提供することにある。

【0011】

所与の実施形態は、本開示を読んで理解すると当業者に明らかになるであろうように、

50

前述の利点のうちの1つ、2つ、それ以上、又は、全てを提供することができ、及び／又は、他の利点を提供し得る。

【図面の簡単な説明】

【0012】

本発明は、様々な構成要素及び構成要素の配置、並びに、様々なステップ及びステップの配置における形態をとることができる。図面は、好ましい実施形態を説明するためのものに過ぎず、本発明を限定するものとして解釈されるべきではない。

【図1】図1は、2以上の匿名化されたヘルスケアデータベースから統合された匿名化された患者データを活用する医療分析装置を概略的に図示している。

【図2】図2は、3以上の匿名化されたヘルスケアデータベースを統合するように構成された図1の装置によって実行されるデータベース統合プロセスの一実施形態を概略的に図示している。

【図3】図3は、異なる匿名化されたヘルスケアデータベースを統合する異なる特徴の選択の基準を模式的に示す表を示している。

【図4】図4は、図2のデータベース統合プロセスの実施形態の改善コンポーネントの動作を図式的に示している。

【図5】図5は、長さ方向の情報を利用する図1のデータベース統合プロセスの実施形態を図式的に示している。

【発明を実施するための形態】

【0013】

匿名化されたヘルスケアデータベースの統合には多くの課題がある。様々な匿名化されたヘルスケアデータベースは、任意の2つのデータベース間では一部のみが重複するデータであり、範囲が大きく異なる場合がある。実際、この部分的な重複は、複数の匿名化されたヘルスケアデータベースを統合したいという、1つのデータベースに欠落している情報を他のデータベースのコンテンツで「埋める」という願望の重要な動機となっている。例えば、本明細書で使用されるように、「匿名化されたヘルスケアデータベース」は、総合的な電子医療記録（EMR）あるいは心血管系情報システム（CVIS）又は集中治療機関（ICU）情報システムなどのドメイン特有の医療データベースから抽出された匿名化されたデータベースなどの医療記録データベース、病院課金部門データベースから抽出された匿名化データベース、医療保険会社データベースから抽出された匿名化されたデータベース、病院入院部門データベースから抽出された匿名化されたデータベースなどであってもよい。CVISから抽出された匿名化されたデータベースは、心臓血管疾患の診断及び治療に関する医療記録を含むことが期待されるが、それらの診断／治療のための保険適用範囲に関する情報は含まれない可能性がある。対照的に、病院課金部門から抽出された匿名化されたデータベースは、保険償還情報を含むが、医療診断／治療データは含まないことが期待される。これらのデータベースを組み合わせることで、患者集団の全体像を提供することができる。しかしながら、統合の動機づけを提供する2つのデータベース間の限られたコンテンツの重複もまた、そのような統合を困難にする。

【0014】

本明細書に開示される様々な実施形態において、複数（3以上）のヘルスケアデータベースの統合を活用することによって、これらの課題が克服される。これは、単一のプロセスでN個のデータベースの統合を行なうことへの動機付けとなる、全体的な重複の程度を大きくすることができる。しかしながら、逆説的に、統合を実行するためのより効率的で信頼性のある手法は、各組の変換テーブルを生成するために、最初に、一組の匿名化されたヘルスケアデータベースを統合し、その後、結果として生じる $N(N-1)/2$ 個の変換テーブルの間の患者マッチングの一貫性に基づいて、 $N(N-1)/2$ 個の変換テーブルを改善することが開示される。この手法は、N個のデータベース間の特徴の重なりが小さくなる可能性があることを認識し、さらに、重複が存在する場合でさえ、幾つかのデータベースでは特定の特徴が信頼できない可能性がある。一組のデータベースを最初に統合する開示された手法を採用することによって、その一組の匿名化されたヘルスケアデータ

ベースに対してうまく選択されたそのようなペアワイズ統合ごとに、一連の特徴を選択することができる。次いで、複数の ($N > 2$) データベースによって提供される追加情報は、幾つかの実施形態では、特徴に依存しない後続の精緻化ステップにおいて活用される。

【0015】

これに加えて、又は、これに代えて、本明細書に開示された実施形態では、これらの問題は、長さ方向の情報、即ち事象の時間的順序及び様々な事象間の時間間隔を利用することによって克服される。一般に、長さ方向の特徴は、事象のタイムスタンプ間の時間間隔 t によって分離された、匿名化されたヘルスケアデータベース内の単一の匿名化された患者に対する一組のタイムスタンプ付き事象によって定義される。このような長さ方向の特徴は、匿名化プロセスが各患者に関して全てのタイムスタンプのランダムではあるが固定のシフトを導入する匿名化されたヘルスケアデータベースにおいても明確に定義されている。これは、固定時間シフトは、事象間の時間間隔 t に影響を与えないためである。

【0016】

図1によれば、 N 個の匿名化されたヘルスケアデータベース 10 は、「データベース1」、「データベース2」、・・・、「データベース N 」とそれぞれ付されている。一般に、 N は、少なくとも2である、幾つかの実施形態では少なくとも3である、正の整数である。下限である $N = 2$ が、幾つかの実施形態で想定される。匿名化されたヘルスケアデータベース 10 は、幾つかの実施形態では100万人以上の患者エントリを有するような大規模データベースを匿名化するために、好ましくは自動化された（例えば、特定のクラス又はタイプのデータを除去するようにプログラムされたコンピュータによりコンピュータ実装された）適切な匿名化プロセス（図示せず）によって生成される。オプションで、匿名化は、例えば、特定の稀少な患者を除去するため、又は、他の異常な状況に対処するために、何らかの手動処理を含むこともできる。 N 個の匿名化されたデータベースを生成するために使用される匿名化プロセスは、一般に異なり、及び/又は、同じ情報を匿名化してもしなくてもよい。各匿名化プロセスは、好ましくは、患者名、患者住所、社会保障番号などの患者を直ちに識別することができる個人識別情報 (PII)、及び、病院の名前、郵便番号などの他の情報と組み合わせてPIIとなる可能性のある情報を匿名化する。情報が他の情報と組み合わせてPIIとなり得る場合、その組み合わせの一部のみを匿名化することで十分である。例えば、郵便番号、性別、及び、生年月日の組み合わせは、個人的に識別することができるが、郵便番号情報のみを匿名化することによって、受け入れ可能な患者の匿名性を達成することができる。匿名化プロセスは、任意に、ある特定の患者を識別することができるまれな情報、例えば、例えば、90歳などのある特定の最大年齢を超える任意の年齢、及び/又は、一般的な診断のリストの中にない診断を除去することもできる。

【0017】

一般に、特定のデータの匿名化は、データを削除する（修正する）か、又は、データをプレースホルダで置き換えることによって行なうことができる。後者は、匿名化ではあるが、その特定のタイプの情報との相関が望ましく保持される場合に好適である。例えば、医療機関（例えば、病院又はケア機関）のエントリは、データベースに対して内部的に一貫したプレースホルダによって置き換えられ得る。これらのプレースホルダは、所与のデータベース内で内部的に一貫しているが、基本的にデータベース間でランダムに変化する。例えば、データベース1では、病院「Blackacre総合病院」は、常にプレースホルダ（例えば、「8243」）で置き換えられ得る。一方、「Whiteacre地域医療センタ」は、常にプレースホルダ「1238」で置き換えられ得る。この例では、データベース1の医療機関「Blackacre総合病院」の全てのインスタンスが、（同じ）プレースホルダ医療機関「8243」に置き換えられ、データベース1の医療機関「Whiteacre地域医療センタ」の全てのインスタンスが、（同じ）プレースホルダ医療機関「1238」で置き換えられる。一方、データベース2の例を続けると、データベース2の医療機関「Blackacre総合病院」の各インスタンスは、（匿名化されたデータベース1におけるBlackacreに対して使用されたプレースホルダ「8243」とは異なる）同一のプレースホルダ医療機関「EADF」によって置き換

えられることができ、「Whiteacre地域医療センタ」の各インスタンスは、(匿名化されたデータベース1におけるWhiteacreに対して使用されたプレースホルダ「1238」とは異なる)同一のプレースホルダ医療機関「JSDF」によって置き換えられることができる。匿名化されたデータベース内で内部的に一致する医療機関プレースホルダによる医療機関の上記匿名化により、データベース上で動作するヘルスケアデータ分析プロセスが、患者の匿名性を維持しながら特定の医療機関との相関を識別することが可能になる。例えば、Blackacreが平均的な病院よりも心臓移植の成功率が統計的に有意に高い場合、これは、匿名化された病院「8243」で行なわれた心臓移植の統計的に有意に高い成功率として、データベース1(心臓移植結果データを格納していると仮定)に表示される。

【0018】

10

一方、一部の情報は、修正、つまり削除によって匿名化される場合がある。例えば、住宅アドレス情報は、完全に修正されることができ、これは識別性が高く、典型的なヘルスケアデータ分析プロセスでは、住宅アドレスとの有用な相関が期待できないためである。変形実施形態では、住宅アドレス相関が健康管理データ分析プロセスの有用な入力であると予想される場合、住所の匿名化は、各住宅アドレスをより地理的エリア(例えば、この都市が許容可能なレベルの匿名性を保証するのに十分な人口を有している場合には、住宅都市)で置き換えることによって実行されてもよい。人口が十分に少ない住宅都市又は郡は、個人的に識別できる「まれな」データを保持しないように完全に修正され、又は、居住国などの適切なより広い地理的单位で置き換えられてもよい。

【0019】

20

匿名化されたヘルスケアデータベース10は、一般に、例えばリレーショナルデータベースフォーマット又は他の構造化データベースフォーマットのような幾つかの構造化フォーマットで、スプレッドシート、検索可能な列区切りのリッチテキストファイルなどのようにそれぞれフォーマットされることが期待される。しかし、幾つかの実施形態では、データベース10のうちの1又は複数は、例えば患者に書かれたテキストレポートを格納するような非構造化データベースであってもよく、又は、限定された構造(例えば、患者名や人口統計情報などの情報に構造化されていないテキストレポートが続く構造化見出しなど)を有してもよい。そのような場合、自然言語処理(NLP)を使用して、データベースコンテンツの構造化された表現(例えば、テキスト文書の単語の表現)を抽出することができる。

30

【0020】

図1に示されるように、ヘルスケアデータ分析装置は、コンピュータ14(又は、より一般的には電子プロセッサ14)上で実装される匿名化されたヘルスケアデータソース装置12を含み、これは、例えば、ネットワークベースのサーバコンピュータ、クラウドコンピューティングリソース、サーバクラスタなどであってもよい。コンピュータ14は、データベース統合プロセス16及び患者データ検索プロセス18を実行するようにプログラムされ、後者は、 $N(N-1)/2$ 個の変換テーブル20のセットを使用する。本明細書の例示的な実施形態では、各変換テーブルは、 N 個のデータベース10のうちの一組のデータベースに関する $m \times 2$ 変換テーブルである。一般性を失うことなく、一組のデータベースをそれぞれデータベース i とデータベース j とし、集合的に一組のデータベース(i, j)を形成する。各変換テーブルは、データベース統合プロセス16によって一組のデータベース(i, j)にマッチングされた m 人の患者の行(又は、代替的に列)と、匿名化されたデータベース i における匿名化された患者IDと、匿名化されたデータベース j における匿名化された患者IDとをリスト化している2つの列(又は、代替的に行)と、を有する $m \times 2$ のテーブルである。 $N=2$ の場合、単一の組のデータベース(i, j)が存在する。 $N>2$ の場合、 $N(N-1)/2$ 個の一意なデータベースの組(i, j)が存在する。これは、 n のセットから k 個の要素をとる組み合わせ数の組み合わせ式を使用して取得され得る。

40

【数 1】

$$\binom{n}{k} = \frac{n(n-1) \cdots (n-k+1)}{k(k-1) \cdots 1}$$

現在のケースでは、組が図示されているため、 $k = 2$ であり、セットは、 N 個の匿名化されたヘルスケアデータベース 10 であるので、 $n = N$ である。従って、組み合わせは、 $N(N-1)/2$ まで減少する。一般に、 $N > 2$ である場合、データベースの組 (i, j) 毎に一致する患者数 m は異なるかもしれないが、データベースの組間の患者の重複は、3 以上の匿名化されたヘルスケアデータベースの有用な統合に対して期待される。

10

【0021】

$N(N-1)/2$ 個の変換テーブル 20 は単一のテーブルとして具体化されることが考えられる。例えば、単一の $m \times [N(N-1)]$ テーブルを形成するために、各々次元 $m \times 2$ の $N(N-1)/2$ 個のテーブルを連結したものである。この場合、全ての $N(N-1)/2$ 個の構成要素である $m \times 2$ 変換テーブルがマッチされる患者数 m を有すると仮定する。そうでない場合、パディングを使用して「欠落した」匿名化患者を説明することができる。例えば、データベース 1 の患者 49 がデータベース 3 に一致がない場合、組 $(i, j) = (1, 3)$ の構成要素である $m \times 2$ 変換テーブルは、<ヌル (null)>又はゼロ又は他のプレースホルダによって適切に埋められる。

【0022】

20

また、コンピュータ 14 は、匿名化された患者データを $N(N-1)/2$ 個の変換テーブル 20 を用いて N 個の匿名化されたヘルスケアデータベース 10 から検索するための患者データ検索プロセス 18 を実行するようにプログラムされる。例えば、クエリは、データベース 1 において使用される匿名化された患者 ID によって識別される所与の患者に関するクエリ特徴の値を取得するため、患者データ検索プロセス 18 の対象となり得る。この患者 ID は、データベース 1 からクエリ特徴の値を検索するために直接的に使用されるが、データベース $j = 2, \dots, N$ の各々に関し、データベースの組 $(1, j)$ のための適切な変換テーブルが、データベース j からクエリ特徴値を検索するために、データベース j において患者 ID をマッチングすべく使用される。

【0023】

30

しかしながら、一般的に、クエリ特徴は、 N 個全てのデータベースに含まれない可能性がある。クエリ特徴が、 N 個の匿名化されたヘルスケアデータベースのうちの 1 つのみに含まれる場合、クエリ特徴は、クエリ特徴を含む (単一の) 匿名化されたヘルスケアデータベースから検索される。一方、クエリ特徴が、 N 個の匿名化されたヘルスケアデータベースのうちの 2 以上に含まれる場合、検索された値は、クエリ特徴に関して、クエリ特徴を含む、 N 個の匿名化されたヘルスケアデータベースのうちの 2 以上のデータベースにおけるクエリ特徴の値から生成される。これは、例えば、クエリ特徴を含む各々の匿名化されたヘルスケアデータベースにおけるクエリ特徴に関する特徴精度メトリックを用いて実行され得る。例えば、クエリが、患者 49 に関する一次診断を要求するとともに、データベース 1, 2, 3 が一次診断フィールドをそれぞれ含む場合、(適切な $m \times 2$ 変換テーブルを用いたデータベース 2, 3 に関する匿名化された患者 ID 49 の変換後に) 患者 49 の一次診断に関する 3 つの値が提供される。データベース 1, 3 が、一次診断に関して 97% の精度を持つ一方、データベース 2 が、この特徴に関して極めて低い精度 (例えば 71%) を持つことが知られている場合、検索される値は、最も精度が高い可能性のあるデータベース 1, 3 から取得された一次診断として生成される。異なるデータベースが所与のクエリ特徴に関して異なる値を格納する場合、その特徴に関する最も高い精度メトリックを有する N 個のデータベース 10 のデータベースの値を取る、又は、最も一般的な値を取る (例えば、6 個のデータベースがその特徴に関する値を持ち、そのうち 5 個のデータベースが一致する場合、6 個のデータベースのうちの 5 個のデータベースにおいて現れる値が選択され得る)、又は、それらの値の平均値 (あるいは、その特徴の精度メトリック

40

50

が最も高いデータベースの一部のサブセットの値の平均値、又は、識別可能な異常値を削除した後の平均値)を有する数値をとるなど、様々な手法を使用して、検索値が生成される。

【0024】

患者データ検索プロセス18によって検索及び処理されるクエリは、クエリの目的に依存して変化し得る。例えば、30歳～50歳の全ての男性患者について一次診断を得ることが望ましい場合がある。この場合、クエリは、年齢及び性別による適切なフィルタリングの後で、一次診断のセットの要求として定式化(それぞれ異なる診断に関する列挙)され得る。この場合のクエリ結果は、各要素(診断、カウント)が診断を示す文字列と、その診断を有する患者数(年齢/性別フィルタリング後)を格納するデータの組{(診断、カウント)}であってもよい。N個のデータベース10がリレーショナルデータベースである場合、患者データ検索プロセス18は、SQLクエリを受け取るSQL(Structured Query Language)クエリエンジンとして実装され得る。

10

【0025】

引き続き図1を参照すると、ヘルスケアデータ分析装置は、コンピュータ24(又は、より一般的には電子プロセッサ24)上に実装されたヘルスケアデータ分析ツール22を更に含み、これは、例えば、ネットワークベースのサーバコンピュータ、クラウドコンピューティングリソース、サーバクラスタ、デスクトップコンピュータ(図示のように)などであってもよい。コンピュータ24は、例示的なキーボード28、マウス又は他のポインティングデバイス30、ディスプレイ26のタッチセンシティブオーバーレイなどの1又は複数のディスプレイコンポーネント/装置26並びに1又は複数のユーザ入力コンポーネント/装置を含むか、それらと動作可能に接続される。ヘルスケアデータ分析ツール22は、(一例として)特定の医療処置に関する保険適用範囲の評価、医療処置の生存率決定、患者に最も一般的に提供される医療ケアの種類との人口統計学的相関についての評価などの様々なヘルスケア分析を実行する。適切な実施形態では、ユーザは、ユーザ入力デバイス28, 30を操作して、実行される分析のタイプを構成する。例えば、ヘルスケアデータ分析ツール22は、匿名化されたヘルスケアデータソース装置12の患者データ検索プロセス18を介して匿名化されたデータベース10から適切なデータを取り出し、そのデータについて選択された分析分析を実行し、その結果は、ディスプレイコンポーネント26上に、図形表現などとして提示される。例えば、処置の保険適用範囲を日付間隔で格納されたヒストグラムとしてプロットするか、又は、異なる保険会社に対応するスライスを含む処置の保険適用範囲を示す円グラフとしてプロットするか、又は、地理的位置の関数として生存率をプロットするなどである。

20

30

【0026】

例示的な匿名化されたヘルスケアデータソース装置12が、コンピュータ14上に実装されるように図1に示されている。一方、ヘルスケアデータ分析ツール22が、異なるコンピュータ24上に実装されるように図1に示されている。しかしながら、他の実施形態では、匿名化されたヘルスケアデータソース装置とヘルスケアデータ分析ツールとが、単一のコンピュータ上に実装されてもよい。また、他のハードウェアセグメンテーショントポロジも考えられる。例えば、データベース統合プロセス16及び患者データ検索プロセス18は、異なるコンピュータ上で実装され得る。さらに、本明細書に記載のヘルスケアデータ解析装置の開示された機能は、開示された機能を実行するために電子プロセッサ14, 24によって読み取り可能且つ実行可能な命令を記憶する非一時的な記憶媒体として実施され得ることが理解されよう。非一時的な記憶媒体は、例えば、ハードディスクドライブ又は他の磁気記憶媒体、光ディスク又は他の光学記憶媒体、フラッシュメモリ、読み出し専用メモリ(ROM)、あるいは、他の電子記憶媒体、それらの様々な組み合わせなどを有していてもよい。

40

【0027】

図2を参照して、 $N > 2$ の場合のデータベース10に対するデータベース統合プロセス16の一実施形態を説明する。この実施形態では、Nは少なくとも3であり、より一般的

50

には、 N は3以上の任意の正の整数であり得る。動作40において、 N 個のデータベース10から、一組(第1の組)の匿名化されたヘルスケアデータベース(i, j)が選択される。一手法では、 i 及び j の値は、最初はそれぞれ1と2とに設定され、 i と j との全ての組の組み合わせが選択されるまで各次の反復して変更される。ここで、 $1 < i < N$ 、且つ、 $1 < j < N$ である(インデックス1, ..., N を用いて、 N 個のデータベース10の構成データベースを表している)。データベースの組(i, j)は2つの異なるデータベースを統合するので、その組は、 $i = j$ である全ての縮退ケースを除外する。

【0028】

以下、選択されたデータベース(i, j)における患者マッチングについての一例が説明される。動作42において、一致/除外基準を適用して、マッチさせるデータベース部分を選択する。データベース i 及びデータベース j からの患者記録を照合するために、おそらく関連性のある2つのデータベースのサブセットが抽出される。例えば、データベース i が医療外科及び熱傷ICU患者のデータのみをカバーする場合、データベース j から、入院時に医療外科及び熱傷ICU病棟に入院した患者のサブセットが抽出される(即ち、含まれる)一方、データベース i と重複しない他のエリアからのデータは除外される。除外された/包含されたデータは、特定のデータベースの組(i, j)についての重複によって決定され、異なる組に対して異なることがあることに留意されたい。

【0029】

動作44において、特徴のセットが、データベース組(i, j)を統合する際の使用のために識別される。ここで、データベース i 及びデータベース j を確実に統合するための一意でない識別の特徴のセットが選択される。選択された特徴は、それぞれ、データベース組(i, j)のデータベース i と j との両方に含まれる。さらに、選択された特徴は、信頼性に関する利用可能な情報に基づいて任意に選択される。例えば、データベースの1つが患者の性別記録に関して比較的正確であるが、データベース i とデータベース j との両方が体重記録の点で正確であることが分かっている場合、体重を特徴として適切に選択し、性別は、特徴として適切に選択されない。

【0030】

図3を簡単に参照すると、所与のデータベースの組(i, j)を統合するために選択される特徴のセットは、一般に、特定のデータベース i 及び j に依存することに留意されたい。例えば、図3は、3つの匿名化されたヘルスケアデータベース X 、 Y 、及び、 Z の特徴の表を示しており、各データベースの各特徴の確度をパーセンテージとして表にしている。図3に示す表の最後の3行は、各特徴を、指示されたデータベースの組み合わせ $i - j$ の特徴の集合として選択すべきかどうかを示している。例えば、図3は、データベース X 及び Y が、人種、死亡率、滞在の長さ、年齢、及び、体重の記録において正確であることを示しているので、これらの5つの特徴は、データベース X 及び Y の照合のために選択される。同様に、人種、滞在期間、年齢、一次診断、体重などの特徴セットが、データベース X とデータベース Z とを統合するために適切に選択される。また、性別、人種、滞在期間、年齢、体重などの特徴セットが、データベース Y とデータベース Z とを統合するために適切に選択される。図3の例では、精度パーセンテージは、特徴精度メトリックを形成し、サンプリング(患者の代表的なサンプリングを選択し、サンプルの特徴精度を検証すること)に基づいて、又は、明らかに誤った特徴値(例えば、年齢 = 0 又は年齢 = 200)に基づいて、又は、誤った特徴値(「エラー」として各々の欠落した特徴値をとる)に基づいて、生成され得る。

【0031】

再び図2によれば、動作46において、動作44で選択された特徴のセットが、データベース i 及び j 内の患者をマッチングするために使用される。様々な手法を用いることができる。単純な手法では、特徴のセットの特徴の使用可能な値の閾値分数(又は数)が一致する場合、データベース i とデータベース j とのそれぞれにおける2人の患者の間で一致が見つかる。オプションで、マッチングは、データベースに誤って記録された特徴値を有する可能性、特徴の選択性などの要因に基づいて、異なる特徴に異なる重みを適用する

10

20

30

40

50

ことができる。本質的に、データベース i の各患者は、動作 4 4 で選択された特徴のセットの値を記憶する要素を持つ特徴ベクトルによって表され、同様に、データベース j の各患者は、動作 4 4 で選択された特徴のセットの値を記憶する要素を持つ特徴ベクトルによって表される。これらの値の幾つかは空白であってもよい（例えば、ベクトルは $\langle \text{null} \rangle$ 又は他のプレースホルダを格納する）。2つのそのような特徴ベクトルの類似性を計算するための任意の手法を使用して、患者を比較し、2つのデータベース内の同様の患者を識別することができる。例えば、特徴の数が F である場合、適切な類似性計算は、次式によって与えられる2つの特徴ベクトル p_i と p_j との間の差であってもよい。

【数 2】

$$D(p_i, p_j) = \frac{1}{F} \sum_{f=1}^F w_f (p_i(f) - p_j(f))^2$$

10

ここで、 p_i 及び p_j は、データベース i において比較される患者とデータベース j において比較される患者とをそれぞれ表す特徴ベクトルであり、 $p_i(f)$ は、患者 p_i に関する f^{th} 番目の特徴の値を表し、同様に、 $p_j(f)$ は、患者 p_j に関する f^{th} 番目の特徴の値を表す。パラメータ w_f は、様々な特徴 $f = 1, \dots, F$ の相対的重要性を示すとともに、（必要に応じて）異なる特徴タイプを共通の単位に変換してその和を計算可能にするために選択される特徴重み及び/又は単位変換係数である。この式では、 $D(p_i, p_j)$ の値が小さいほど、より似ている患者を示し、結果、 $D(p_i, p_j)$ が幾らかの閾値よりも小さい場合に、2人の患者はマッチし得る。任意の欠落した特徴は、様々な方法で処理され得る。例えば、 $D(p_i, p_j)$ を形成する和からそれらを単純に省いたり（及び、然るべく $1/F$ スケーリングしたり）、又は、欠落している特徴 f の場合、 $p_i(f) - p_j(f)$ の幾らかの初期値を割り当てたりする。上記は単なる例示であり、実質的に他の比較形式を使用して、それぞれのデータベース i 及び j における一致する患者を識別することができることを理解されたい。

20

【0032】

動作 4 8 では、動作 4 6 で識別されたデータベース間の患者の一致が、データベース対 (i, j) の患者 ID 変換表に集計される。例えば、この表は、次のような $m \times 2$ の表であり得る。

30

表 1：データベース組 (i, j) に関する患者 ID 変換表

【表 1】

患者ID (データベースi)	患者ID(データベースj)
1	43
2	17
3	<ヌル>
4	98
5	2
5	3
6	6
...	
96	9
<ヌル>	6
<ヌル>	9
<ヌル>	23

ここで、この例では、データベース i の患者 ID = 3 はデータベース j に一致せず、同様に、患者 ID = 6、ID = 9、及び、ID = 23 はデータベース i に一致しないことに留意されたい。表 1 の例示的な例は、データベース i の患者 ID によってソートされているが、データベース j の患者 ID によるソートを実行することは自明である。そうすれば、表をより効率的に読み取ることができる（例えば、図 1 の患者データ検索プロセス 18 によって受信されるクエリがデータベース j 内の患者 ID によって索引付けされる）。

【0033】

幾つかの実施形態では、患者のマッチングは排他的ではないことに留意されたい。これは、データベース i の患者 ID = 5 がデータベース j の患者 2 とデータベース j の患者 3 との両方と一致する表 1 に示されている。このオプションの非排他性は、患者のマッチングにおける不確実性を捕捉することを可能にする。医療データ分析アプリケーションの場合、上記の非排他的マッチングは、そのような不確実なマッチングの数が比較的少ない場合には必ずしも問題ではなく、そのような場合、このように複数のマッチングを許容すると統計的な全体的な精度が向上する。表 1 に示されるデータベース (i, j) の例示的な変換テーブルでは、ストレージは、データベース i の患者 ID 又はデータベース j の患者 ID のうちのいずれかのテーブルのソートを容易にする利点を有する、データベース i の患者 ID 5 の重複エントリによるものである。

【0034】

決定動作 50 では、一意のデータベース組 (i, j) 毎に患者 ID 変換テーブルを生成するために、統合されている N 個のデータベース 10 のセット内の各一意の組のデータベース (i, j) について、処理を繰り返す。従って、このループは、N 個のデータベースから得られる $N(N-1)/2$ 個の一意のデータベース組に対する $N(N-1)/2$ 個の変換テーブルを生成するために $N(N-1)/2$ 回実行される。例えば、N = 3 の場合、組 (1, 2)、組 (1, 3)、組 (2, 3) の 3 つの繰り返しがある。他の例示的な例として、N = 5 の場合、(1, 2)、(1, 3)、(1, 4)、(1, 5)、(2, 3)、(2, 4)、(2, 5)、(3, 4)、(3, 5)、(4, 5) の 10 個の繰り返しがある。決定動作 50 により実施されるループは、例えば、ネストされたループ $i = 1$ to $(N-1)$; $j = i + 1$ to N (ここで、j は、内側のループである) によって実施され得る。

【 0 0 3 5 】

$N(N-1)/2$ 回のループ反復の出力は、 N 個のデータベース10のうちの $N(N-1)/2$ 個の一意的データベース組に対する $N(N-1)/2$ 個の変換テーブルとなる。幾つかの実施形態では、これは、患者データ検索プロセス18によって使用される $N(N-1)/2$ 個の変換テーブル20（それぞれ $m \times 2$ の次元）を提供する最終出力である。しかし、この時点でデータベース統合プロセス12が終了する場合、複数（3つ以上）のヘルスケアデータベース（即ち、 $N \geq 3$ ）からの情報は、個々の $m \times 2$ ペアワイズ変換テーブルを改善するために効果的に活用されない。

【 0 0 3 6 】

引き続き図2を参照するとともに図4を更に参照すると、例示的な実施形態では、 $N(N-1)/2$ 個の変換テーブルが構築された後に改善動作52が実行される。これは、 $N(N-1)/2$ 個の変換テーブル間の患者マッチングの整合性に基づいて、 $N(N-1)/2$ 個の変換テーブルを改善する。例示的な実施形態では、改善動作52は、動作44の反復で識別された特徴のセットを使用しない。むしろ、改善動作52は、 $N(N-1)/2$ 個の変換テーブル間で予想される整合性を考慮することによって、図4に図式的に示すように実行される。図4の例では、各円は、匿名化された患者ID（例えば、「患者1」は匿名化ID=1でラベル付けされている）及びデータベース（この例ではX、Y、又は、Z）によりラベル付けされた単一の匿名化された患者を表している。異なるデータベースの患者を結ぶ実線又は破線は、動作42、44、46、48のペアワイズマッチングによって見出される可能性のあるマッチングを示す。この例では、データベースXの患者1は、X-Y変換テーブルに基づいて、データベースYの患者22にリンクされる。一貫性を維持するために、データベースXの患者1とデータベースYの患者22との両方をデータベースZの同じ患者にリンクする必要がある。しかし、組 $i = X$ 、 $j = Z$ に対するペアワイズマッチングプロセスでは、データベースXの患者1は、データベースZの患者72及び患者31との両方にマッチングされる（このような非排他的マッチングは、表1の例におけるデータベースiの患者5の例について上述のように、許容され得る）。組 $i = Y$ 、 $j = Z$ に対するペアワイズマッチングプロセスでは、データベースYの患者22は、データベースZの患者72及び患者14の両方に一致した。自己整合性を維持するために、データベースXの患者1とデータベースYの患者22とがデータベースZの患者72と一致しなければならず、他の可能性のある一致は不一致である。従って、改善動作52では、データベースXの患者1とデータベースZの患者31との間の一致がXZ変換テーブルから除去され、同様にデータベースYの患者22とデータベースZの患者14との間のマッチングがYZ変換テーブルから除去される。

【 0 0 3 7 】

他の実施形態では、このような整合性解析は、反復ループ40、42、44、46、48、50の間に実行することができる。このアプローチは、既に作成された組毎の変換テーブルを活用することによって、後のループ反復を実行する処理時間を短縮する。例えば、X、Y、Zのインデックスを付けられたデータベースを有するとともに、X-Y、X-Z、Y-Z変換テーブルを作成するために反復ループ40、42、44、46、48、50がこの順序で実行されるケースを考える。X-Y及びX-Z変換テーブルの作成後、データベースXの患者10はデータベースYの患者11にリンクされ、データベースXの患者10はデータベースZの患者15にもリンクされることが分かる。次に、Y-Z変換テーブルを作成する最後の反復の間に、Y-Z変換テーブルと既に作成されたX-Y及びX-Z変換テーブルとの一貫性（整合性）を保証するために、データベースYの患者11をデータベースZの患者15にリンクする必要があることが既に知られている。

【 0 0 3 8 】

これに加えて、又は、これに代えて、本明細書に開示された幾つかの実施形態では、患者のマッチングを改善するために、長さ方向の情報が活用される。一般に、長さ方向の特徴は、事象のタイムスタンプ間の時間間隔 t によって分離された、匿名化されたヘルスケアデータベース内の単一の匿名化された患者に対するタイムスタンプ付き事象の組によ

10

20

30

40

50

って定義される。このような長さ方向の特徴は、匿名化された医療データベースにおいても明確に定義されており、厳密な時間シフトは事象間の時間間隔 t に影響を与えないので、匿名化プロセスは各患者の全てのタイムスタンプのランダムではあるが厳密なシフトを導入する。

【 0 0 3 9 】

図 5 を参照すると、長さ方向の特徴の一例が記載されている。長さ方向の特徴は、時間間隔 t で区切られた e 型の事象と、 f 型の事象とによって定義される。図 5 の例では、データベース X の患者 m は、時間間隔 t で分離された事象タイプ f の事象の発生に続く事象タイプ e の事象の発生を有する。同様に、データベース Y の患者 n は、同じ時間間隔 t で分離された事象タイプ f の事象の発生に続く事象タイプ e の事象の発生を有する。対照的に、データベース Z の患者 p は、事象タイプ e の事象とそれに続く事象タイプ f の事象を有するが、タイプ e 及び f の事象間の時間間隔は時間間隔 t よりもはるかに長い。従って、時間間隔 t で区切られた事象シーケンス $e \quad f$ の時間的特徴に基づいて、データベース X の患者 m は、データベース Y の患者 n と一致するが、データベース Z の患者 p とは一致しない。このような長さ方向の特徴を一致させる際に、例えば、タイムスタンプの入力の可能性のあるエラーを考慮して、異なるデータベースの患者のために t における或る程度の変動を許容することが考えられる。

【 0 0 4 0 】

より複雑な長さ方向の特徴、例えば、第 1 の時間間隔 t_1 で分離される事象 $g \quad e$ を有する事象タイプ $g \quad e \quad f$ と、第 2 の時間間隔 t_2 で分離される事象 $e \quad f$ とが想定される。他の考えられる長さ方向の特徴では、 t における許容可能な変動は、タイプ $e \quad f$ の事象がそれらの間の時間間隔に関係なく (t における許容可能な変動で規定されるある制限内で) 発生した場合に、実質的に長さ方向の特徴がマッチするように十分大きくてもよい。

【 0 0 4 1 】

例示的な長さ方向の特徴は、2つのデータベース (i, j) の患者の事象のタイムスタンプを比較するのではなく、事象間の時間間隔 t を使用する。前述したように、事象の絶対的なタイムスタンプに頼るのではなく、事象間の時間間隔に依存するこの手法は、患者のタイムラインが匿名化プロセスの一部としてランダムにシフトされた可能性に対して堅牢である。

【 0 0 4 2 】

幾つかの実施形態では、長さ方向の特徴は、動作 4 4 で識別され、動作 4 6 (図 2 参照) で使用される特徴のセットの他の特徴と同様に処理される。しかしながら、この手法は、長さ方向の特徴の識別及び処理が計算上複雑であり得るので、不要な計算上の複雑さを招く可能性がある。例えば、平均患者が $E = 30$ 事象を有する場合、 $e \quad f$ 形式の長さ方向の特徴を識別するために必要なペアワイズ事象比較の数は、 $E(E - 1) / 2 = 435$ 個の事象組である。一方、長さ方向の特徴のかなり高い特異性は、それらが合致する患者に対して非常に差別的であり得ることを意味する。従って、幾つかの実施形態では、患者のマッチング動作 4 6 は、長さ方向の特徴に依存することなく最初に行われ、長さ方向の特徴は、困難なマッチング (例えば、非長さ方向の特徴のみが使用される場合の、データベース Y の複数の患者と一致するデータベース X の患者) に対してのみ計算及び活用され得る。

【 0 0 4 3 】

幾つかの実施形態では、各患者の普遍的な患者 ID (又は UID) を使用して、非長さ方向の特徴マッチングが実行される (又は、部分的に実行される)。UID は、患者の性別、人種、年齢、及び、体重などの共通の特徴の集合の連結として構築される。例えば、患者の UID 1518170 は、以下の特徴を使用して生成される。男性又は性別 1 (1518170 の最初の桁)、ネイティブアメリカン又は人種 5 (1518170 の 2 桁目)、18 才 (1518170 の 3 桁目及び 4 桁目)、体重 170 ポンド (1518170 の 5 桁目 ~ 7 桁目)。従って、患者のために新しい記録 (医学報告又は請求記録) が生成されるたびに、UID が患者記録に割

10

20

30

40

50

り当てられる。U I Dは特徴ベースなので、異なる匿名化されたデータベースで同じにする必要がある。オプションで、幾つかの許容値が受け入れられる。例えば、年齢に対して ± 1 年の許容しきい値を使用した場合、データベースIIの80歳は、データベースIの79歳～81歳と同じとみなされる。特徴マッチングのためのそのようなU I D手法は、患者とのマッチングに使用される特徴のセットの全ての特徴に対して使用されてもよく、あるいは、より小さな特徴のサブセットが連結されてU I Dを形成してもよい。ここで、U I Dを形成している特徴のセットは、全てのN個のデータベース10に対して共通である。この後者の手法は、U I Dが一度計算され、図2の(i, j)ループの各反復に対して再使用されることを有利に可能にし、計算効率を高めることができる。この手法では、次の3段階マッチングプロセスが想定される。(1) U I Dに基づくマッチング、(2) 困難な場合

10

【0044】

所与の実施形態では、開示された態様の様々な組み合わせを採用することができる。が理解されよう。例えば、2つのデータベースの統合($N=2$)と複数のデータベースの統合($N=3$)との両方に、長さ方向の特徴マッチングが使用され得る。自然言語処理(NLP)を使用して、 $N=2$ 且つ $N=3$ の統合タスクのために、非構造化又は半構造化データベースから一連の特徴を生成できる。

【0045】

図1の開示されたヘルスケアデータ分析装置を見るための別の手法では、N個の匿名化されたヘルスケアデータベース10を統合するプロセスは、N個の匿名化されたヘルスケアデータベース10から匿名化された母集団イメージを再構築する匿名化母集団イメージ再構成方法と見ることができる。この代替図では、再構成された匿名化された母集団イメージは、 $N(N-1)/2$ 個の変換テーブル20によって統合されたN個の匿名化されたヘルスケアデータベース10の内容を含む。この代替的な視点では、匿名化された母集団イメージ再構成法は、N個の匿名化されたヘルスケアデータベース10の形式の母集団イメージデータを、 $N(N-1)/2$ 個の変換テーブル20によって統合されたN個の匿名化されたヘルスケアデータベース10の内容を有する匿名化された母集団イメージに再構築(又は変換)する。

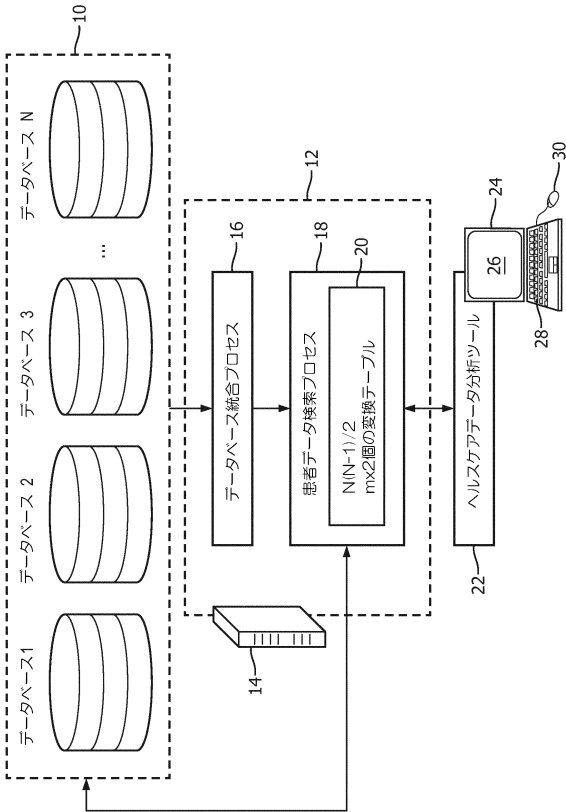
20

【0046】

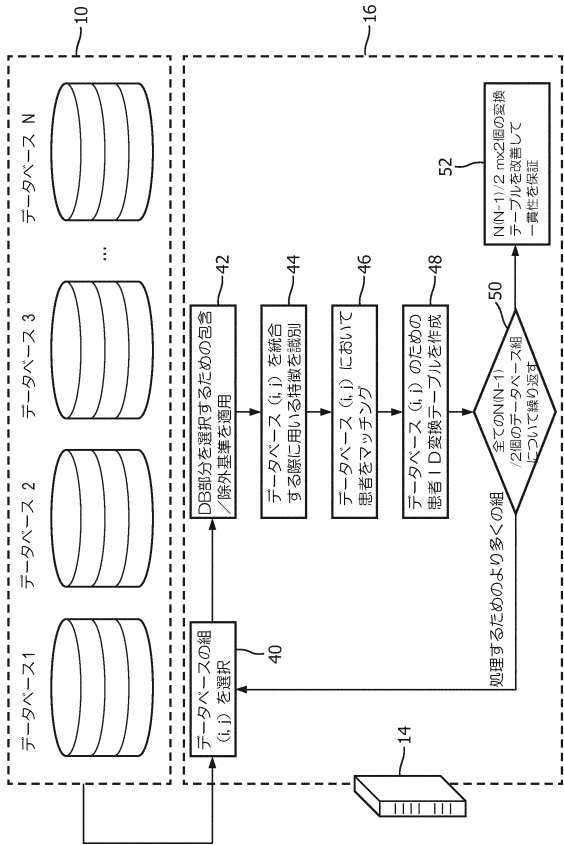
本発明を好ましい実施形態を参照して説明した。前述の詳細な説明を読んで理解すると、他への修正や変更がなされ得る。本発明は、添付の特許請求の範囲又はその均等の範囲内に入る限りにおいて、そのような改変及び変更の全てを含むと解釈されることが意図される。

30

【図 1】



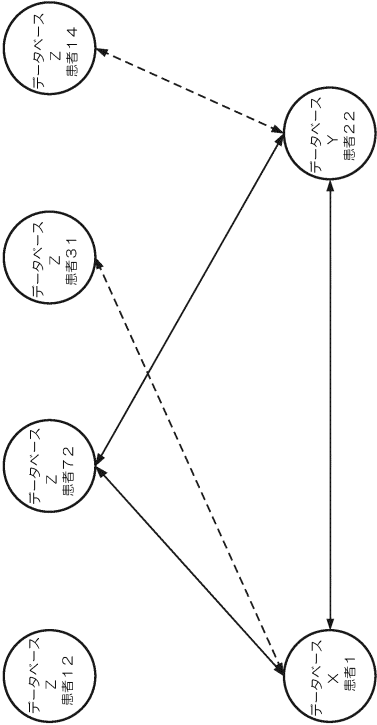
【図 2】



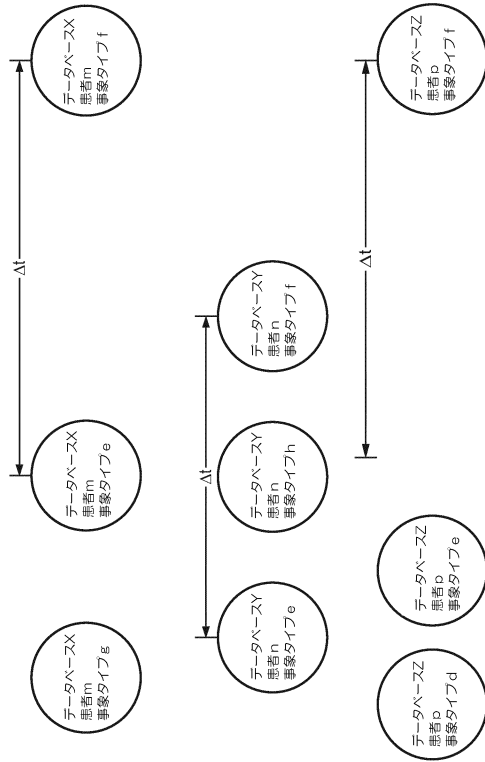
【図 3】

データベース	性別	人種	死亡率	在院日数	年齢	一次診断	体重
X	80%	99%	97%	95%	97%	97%	99%
Y	97%	97%	98%	94%	99%	71%	98%
Z	99%	96%	70%	98%	95%	97%	97%
X-Yの特徴?	いいえ	はい	はい	はい	はい	いいえ	はい
X-Zの特徴?	いいえ	はい	いいえ	はい	はい	はい	はい
Y-Zの特徴?	はい	はい	いいえ	はい	はい	いいえ	はい

【図 4】



【図 5】



フロントページの続き

- (72)発明者 シャリフィ セデュー レザ
オランダ国 5 6 5 6 アーエー アインドーフエン ハイ テック キャンパス 5
- (72)発明者 エルフォルト ダニエル ロベルト
オランダ国 5 6 5 6 アーエー アインドーフエン ハイ テック キャンパス 5
- (72)発明者 トゥルエン ロエル
オランダ国 5 6 5 6 アーエー アインドーフエン ハイ テック キャンパス 5

審査官 後藤 彰

- (56)参考文献 特開 2 0 1 6 - 0 3 8 7 8 0 (J P , A)
欧州特許出願公開第 0 2 8 7 9 0 6 9 (E P , A 2)
米国特許出願公開第 2 0 0 2 / 0 0 7 3 1 3 8 (U S , A 1)

- (58)調査した分野(Int.Cl. , D B 名)
- | | |
|---------|-------------------------|
| G 0 6 F | 1 6 / 0 0 - 1 6 / 9 5 8 |
| G 1 6 H | 1 0 / 0 0 - 8 0 / 0 0 |