

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6973636号  
(P6973636)

(45) 発行日 令和3年12月1日(2021.12.1)

(24) 登録日 令和3年11月8日(2021.11.8)

(51) Int.Cl. F 1  
G 0 6 F 21/62 (2013.01) G 0 6 F 21/62 3 4 5

請求項の数 6 (全 11 頁)

<p>(21) 出願番号 特願2020-518220 (P2020-518220)                  (86) (22) 出願日 平成31年4月17日 (2019.4.17)                  (86) 国際出願番号 PCT/JP2019/016447                  (87) 国際公開番号 W02019/216137                  (87) 国際公開日 令和1年11月14日 (2019.11.14)                  審査請求日 令和2年10月26日 (2020.10.26)                  (31) 優先権主張番号 特願2018-89641 (P2018-89641)                  (32) 優先日 平成30年5月8日 (2018.5.8)                  (33) 優先権主張国・地域又は機関                  日本国 (JP)</p>	<p>(73) 特許権者 000004226                  日本電信電話株式会社                  東京都千代田区大手町一丁目5番1号                  (74) 代理人 100121706                  弁理士 中尾 直樹                  (74) 代理人 100128705                  弁理士 中村 幸雄                  (74) 代理人 100147773                  弁理士 義村 宗洋                  (72) 発明者 長谷川 聡                  東京都千代田区大手町一丁目5番1号 日                  本電信電話株式会社内                   審査官 宮司 卓佳</p>
--	---

最終頁に続く

(54) 【発明の名称】 安全性評価装置、安全性評価方法、およびプログラム

(57) 【特許請求の範囲】

【請求項1】

複数のレコードからなる元データベースと上記元データベースを秘匿した秘匿データベースとを記憶するデータベース記憶部と、

上記秘匿データベースの各レコードについて、上記元データベースに対する近傍探索により所定の近傍数の近傍レコード集合を取得する近傍レコード探索部と、

上記秘匿データベースの各レコードについて上記近傍レコード集合の各レコードとの距離を計算し、当該レコードとの距離に基づいて最近傍レコードを取得する最近傍レコード計算部と、

上記秘匿データベースの各レコードについて、当該レコードに対応する上記元データベースのレコードが上記最近傍レコードと一致するか否かに基づいて当該レコードの再識別率を計算する再特定判定部と、

上記秘匿データベースの各レコードについて計算した再識別率に基づいて上記秘匿データベースの再識別率を計算する再識別率計算部と、  
 を含む安全性評価装置。

【請求項2】

請求項1に記載の安全性評価装置であって、

上記近傍レコード集合のレコード数が所定の閾値よりも多い場合に上記近傍レコード集合中の重複レコードを排除する重複排除部をさらに含む、  
 安全性評価装置。

**【請求項 3】**

請求項 2 に記載の安全性評価装置であって、

上記近傍レコード探索部は、上記近傍数を上記元データベースのレコード件数の対数として上記近傍レコード集合を取得するものであり、

上記重複排除部は、上記閾値を上記近傍数の 2 倍として上記近傍レコード集合中の重複レコードを排除するものである、

安全性評価装置。

**【請求項 4】**

請求項 1 から 3 のいずれかに記載の安全性評価装置であって、

上記近傍レコード探索部は、kd木を用いる近傍探索により上記近傍レコード集合を取得するものである、

安全性評価装置。

10

**【請求項 5】**

データベース記憶部に、複数のレコードからなる元データベースと上記元データベースを秘匿した秘匿データベースとが記憶されており、

近傍レコード探索部が、上記秘匿データベースの各レコードについて、上記元データベースに対する近傍探索により所定の近傍数の近傍レコード集合を取得し、

最近傍レコード計算部が、上記秘匿データベースの各レコードについて上記近傍レコード集合の各レコードとの距離を計算し、当該レコードとの距離に基づいて最近傍レコードを取得し、

20

再特定判定部が、上記秘匿データベースの各レコードについて、当該レコードに対応する上記元データベースのレコードが上記最近傍レコードと一致するか否かに基づいて当該レコードの再識別率を計算し、

再識別率計算部が、上記秘匿データベースの各レコードについて計算した再識別率に基づいて上記秘匿データベースの再識別率を計算する、

安全性評価方法。

**【請求項 6】**

請求項 1 から 4 のいずれかに記載の安全性評価装置としてコンピュータを機能させるためのプログラム。

**【発明の詳細な説明】**

30

**【技術分野】****【0001】**

この発明は、データベースに対して決定的手法もしくは確率的手法により個別データを秘匿したデータベースの安全性を評価する技術に関する。

**【背景技術】****【0002】**

データベース（以下、「元データベース」と呼ぶ）に対して決定的手法により個別データを秘匿する技術として、k-匿名法（非特許文献 1 および 2 参照）がある。また、確率的手法により秘匿する技術として、Pk-匿名法（非特許文献 3 および 4 参照）がある。これらの秘匿処理を施したデータベース（以下、「秘匿データベース」と呼ぶ）の安全性を評価するために、レコードリンケージと呼ばれる手法（非特許文 5 および 6 参照）が用いられる。レコードリンケージとは、あるレコードを再特定しようとすることで、どれだけそのレコードが秘匿できているかを測定する方法である。従来技術では、再特定を試みる秘匿データベースの対象レコードと元データベースの全レコードとの距離を計算し、最近傍レコードと対象レコードとが一致したら再特定できたとして、最近傍レコード数の逆数を対象レコードの再識別率とする。これを秘匿データベースの全レコードについて実施し、各レコードの再識別率を合計した値をデータベースの再識別率として評価する。

40

**【先行技術文献】****【非特許文献】****【0003】**

50

【非特許文献1】Kristen LeFevre, David J DeWitt, and Raghu Ramakrishnan, "Incognito: Efficient full-domain k-anonymity", Proceedings of the 2005 ACM SIGMOD international conference on Management of data, pp. 49-60, 2005.

【非特許文献2】Florian Kohlmayer, Fabian Prasser, Claudia Eckert, Alfons Kemper, and Klaus A Kuhn, "Flash: efficient, stable and optimal k-anonymity", Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Conference on Social Computing (SocialCom), pp. 708-717, 2012.

【非特許文献3】五十嵐大, 千田浩司, 高橋克巳, "数値属性における, k-匿名性を満たすランダム化手法", コンピュータセキュリティシンポジウム2011, pp. 450-455, 2011年

【非特許文献4】五十嵐大, 千田浩司, 高橋克巳, "k-匿名性の確率的指標への拡張とその適用例", コンピュータセキュリティシンポジウム2009, pp. 1-6, 2009年

【非特許文献5】Vicenc Torra, John M Abowd, and Josep Domingo-Ferrer, "Using mahalanobis distance-based record linkage for disclosure risk assessment", International Conference on Privacy in Statistical Databases, pp. 233-242, 2006.

【非特許文献6】Josep Domingo-Ferrer and Vicenc Torra, "Distance-based and probabilistic record linkage for re-identification of records with categorical variables", Butlletí de IACIA, Associació Catalana d'Intel·ligència Artificial, pp. 243-250, 2002.

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0004】

近年ビッグデータの利活用が注目されており、匿名化の対象となるデータも大規模データとなることが想定される。従来技術では、レコードリンケージの際に、レコード数が増えるに連れて処理時間が増えることが問題であった。より具体的には、レコード数の線形的な増加に伴い、処理時間が2乗で増えてしまう。したがって、大規模なデータに対し、実用的な処理時間でレコードリンケージを行うことが課題であった。

#### 【0005】

この発明は、上記のような技術的課題に鑑みて、大規模なデータを秘匿したデータベースの安全性を効率的に評価することを目的とする。

#### 【課題を解決するための手段】

#### 【0006】

上記の課題を解決するために、この発明の一態様の安全性評価装置は、複数のレコードからなる元データベースと元データベースを秘匿した秘匿データベースとを記憶するデータベース記憶部と、秘匿データベースの各レコードについて、元データベースに対する近傍探索により所定の近傍数の近傍レコード集合を取得する近傍レコード探索部と、秘匿データベースの各レコードについて近傍レコード集合の各レコードとの距離を計算し、当該レコードとの距離に基づいて最近傍レコードを取得する最近傍レコード計算部と、秘匿データベースの各レコードについて、当該レコードに対応する元データベースのレコードが最近傍レコードと一致するか否かに基づいて当該レコードの再識別率を計算する再特定判定部と、秘匿データベースの各レコードについて計算した再識別率に基づいて秘匿データベースの再識別率を計算する再識別率計算部と、を含む。

#### 【発明の効果】

#### 【0007】

この発明によれば、レコードリンケージを行う際に、従来技術では $O(N^2)$ の計算量を要する処理が、近傍数を $\log N$ とした場合には $O(N \log N)$ の計算量となる。そのため、大規模なデータに対し、実用的な処理時間でレコードリンケージを行うことができる。したがって、大規模なデータを秘匿したデータベースの安全性を効率的に評価することができる。

#### 【図面の簡単な説明】

【0008】

【図1】図1は、本発明で対象とするデータベースの定義を説明するための図である。

【図2】図2は、従来のレコードリンケージを説明するための概念図である。

【図3】図3は、本発明のレコードリンケージを説明するための概念図である。

【図4】図4は、実施形態の安全性評価装置の機能構成を例示する図である。

【図5】図5は、実施形態の安全性評価方法の処理手続きを例示する図である。

【発明を実施するための形態】

【0009】

以下、この発明の実施の形態について詳細に説明する。なお、図面中において同じ機能を有する構成部には同じ番号を付し、重複説明を省略する。

10

【0010】

[記号]

ある属性の集合を大文字 $X$ と表現し、属性 $X$ の値を小文字 $x$   $X$ と表現する。

【0011】

データベースの1レコードを横ベクトルとして表現する。例えば、 $M$ 属性あるデータベースの $i$ 番目のレコードは、 $x_i = \{x_{i1}, \dots, x_{ij}, \dots, x_{iM}\}$ とする。

【0012】

複数のレコードからなる集合をデータベース $X$ とする。例えば、レコード数 $N$ のデータベースは、 $X = \{x_1, \dots, x_N\}$ とする。

【0013】

データベース $X$ の各レコードを決定的手法もしくは確率的手法により秘匿したレコードからなる集合を秘匿データベース $Y$ とする。例えば、レコード数 $N$ の秘匿データベースは、 $Y = \{y_1, \dots, y_N\}$ とする。

20

【0014】

秘匿データベース $Y$ はレコードの順番がシャッフルされている場合もある。そこで、秘匿データベース $Y$ の行番号と元データベース $X$ の行番号(以下、「真の行番号」と呼ぶこともある)とを対応付ける行番号対応関数 $f_y: R \rightarrow R$ を定義する。

【0015】

図1に、元データベース $X$ 、秘匿データベース $Y$ 、および行番号対応関数 $f_y$ の例を示す。元データベース $X$ は、 $M$ 属性からなる平文のレコードを $N$ レコード含むデータベースである。秘匿データベース $Y$ は、元データベースの各レコードが秘匿され、かつ、順番がシャッフルされたデータベースである。行番号対応関数 $f_y$ は、元データベースの行番号と秘匿データベースの行番号との対応が表された参照表である。

30

【0016】

[処理の概要]

本発明の安全性評価技術では、元データベース $X$ と秘匿データベース $Y$ と行番号対応関数 $f_y$ とを用いて、データベース全体の再識別率を計算し、安全性の評価を行う。本発明では、あるレコードの再識別率を計算するにあたり、大まかに以下の2つの処理を行う。

【0017】

処理1. 近傍探索の対象となるレコードの近傍レコードを、指定した近傍数分取得する。近傍レコードの探索は、木構造を用いたもの(参考文献1参照)や、ハッシングを用いたもの(参考文献2参照)があり、それらを用いて近傍レコードを取得する。木構造としては、例えばkd木等が挙げられる。

40

【0018】

[参考文献1] Jon Louis Bentley, "Multidimensional binary search trees used for associative searching", Communications of the ACM, Vol. 18, No. 9, pp. 509-517, 1975.

[参考文献2] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni, "Locality-sensitive hashing scheme based on p-stable distributions", In Proceedings of the twentieth annual symposium on Computational geometry, pp. 253-262,

50

2004.

【 0 0 1 9 】

処理 2 . 近傍レコードに基づくレコードリンケージの対象レコードと近傍レコードとの距離を計算し、最も距離が近い近傍レコードの行番号とレコードリンケージの対象レコードの真の行番号とが一致したら再特定できたとする。計算する距離としては、例えば、ユークリッド距離、ハミング距離、マンハッタン距離等、適切な距離を用いることができる。レコードの属性値が重複している場合、近傍探索で指定した近傍数以上の近傍レコードが取得される。その場合には、近傍レコード中の重複レコードを排除した上で、近傍レコードとの距離を計算する。

【 0 0 2 0 】

上記処理 1 , 2 を秘匿データベースの各レコードについて行い、各レコードの再識別率の合計値をデータベースの再識別率とし、秘匿データベース全体の安全性を評価する。

【 0 0 2 1 】

図 2 は、従来技術によるレコードリンケージを表す概念図であり、図 3 は、本発明によるレコードリンケージを表す概念図である。従来技術は秘匿データベースのあるレコードについて元データベースの全レコードとの距離を計算し、最も近いレコードの行番号がそのレコードの真の行番号と一致した場合に、再特定できたものと判定する。一方、本発明は秘匿データベースのあるレコードについて元データベースから近傍探索により取得した所定の近傍数の近傍レコードとの距離を計算し、最も近いレコードの行番号がそのレコードの真の行番号と一致した場合に、再特定できたものと判定する。

【 0 0 2 2 】

本発明では、木構造を用いた近傍探索もしくはハッシングを用いた近傍探索のどちらかを用いることとする。本発明の具体的な処理を<Algorithm 1>に示す。上記の処理 1 ( 近傍探索 ) は 2 ~ 7 行目に対応し、処理 2 ( 近傍レコードに基づくレコードリンケージ ) は 8 ~ 22 行目に対応する。なお、 $|\cdot|$  は集合  $\cdot$  の要素数を表す。

【 0 0 2 3 】

<Algorithm 1>近傍探索を用いたレコードリンケージ

Input: レコード数Nの元データベース  $X = \{x_1, \dots, x_N\}$ , レコード数Nの秘匿データベース  $Y = \{y_1, \dots, y_N\}$ , 行番号対応関数  $f_y: R \rightarrow R$ , 近傍数  $K (1 < K < N)$ , 許容範囲  $(\epsilon > 1)$

Output: 再識別率  $r$

```

1:  $r = 0$ 
2: for  $i=1$  to  $N$  do
3:    $y_i$  に対する元データベース  $X$  の近傍レコード集合  $X_i^{near} = \{x_j\}$  (ただし、 $|X_i^{near}| \geq K$ ) を近傍探索により取得する
4: end for
5: if  $|X_i^{near}| \geq K$  となる  $i$  が存在する場合 then
6:    $|X_i^{near}| \geq K$  となる  $i$  に対して、 $X_i^{near}$  のうち重複を除いたレコード集合  $X_i^{uniq} = \{x_j\}$  とし、各  $x_j$  に対応する重複レコードの行番号集合を返す関数  $f_{dup}^i$  を保持する
7: end if
8: for  $i=1$  to  $N$  do
9:   if  $|X_i^{near}| \geq K$  となる場合 then
10:     $y_i$  と  $X_i^{uniq}$  の各レコードとの距離を求め、 $y_i$  に最も距離の近いレコードを最近傍レコード  $Z$  とする
11:    if  $Z$  に  $f_y(k)=i$  (ただし  $k: x_k = Z$ ) となるレコードが存在する場合 then
12:      for  $x_j \in Z$  do
13:         $r = r + 1 / (|f_{dup}^i(j)| |Z|)$ 
14:      end for

```

10

20

30

40

50

```

15:         end if
16:     else
17:          $y_i$ と $X_i^{near}$ の各レコードとの距離を求め、 $y_i$ に最も距離の近いレコー
ドを最近傍レコード $Z$ とする
18:         if  $Z$ に $f_y(k)=i$ (ただし $k:x_k Z$ )となるレコードが存在する場合 the
n
19:              $r = r+1/|Z|$ 
20:         end if
21:     end if
22: end for

```

10

## 【0024】

まず、秘匿データベース $Y$ のレコード $y_i$ ごとに近傍探索を用いて近傍レコード集合 $X_i^{near}$ を取得する(2~4行目に対応)。元データベース $X$ 中に重複したレコードが少なければ、取得した近傍レコード集合 $X_i^{near}$ が指定した近傍数 $K$ 以下となる。しかしながら、元データベース $X$ 中に重複するレコードが多い場合、取得した近傍レコード集合 $X_i^{near}$ が指定した近傍数 $K$ を超えることがあり、近傍探索した効果がなくなってしまう。そこで、近傍レコード集合 $X_i^{near}$ が $K$ 件を超えた場合は、近傍レコード集合 $X_i^{near}$ 中の重複するレコードを排除したレコード集合 $X_i^{uniq}$ を生成し、重複するレコードの行番号を返す関数 $f_{dup}^i$ を保持する(5~7行目に対応)。なお、許容範囲と近傍数 $K$ は、例えば、 $=2.0$ 、 $K=\log N$ などに設定するとよい。

20

## 【0025】

次に、秘匿データベース $Y$ のレコード $y_i$ ごとに近傍レコード集合の各レコードとの距離を計算し、最も距離が近いレコードの真の行番号(すなわち、元データベース上の行番号)が現在のレコードの真の行番号(すなわち、現在のレコードの秘匿データベース上の行番号に対応付けられた元データベース上の行番号)と一致した場合、再識別成功として再識別率 $r$ を加算する。

## 【0026】

より具体的には、近傍レコード集合 $X_i^{near}$ のレコード数が $K$ 件を超えていた場合は、まず重複を排除した近傍レコード集合 $X_i^{uniq}$ の各レコードとの距離計算を行う。そして、最も距離が近いレコードの重複するレコードの真の行番号を探索し、現在のレコードの行番号と一致した場合、重複レコード間で平均して再識別できたとして、 $1/\text{重複レコード数}$ ( $1/|f_{dup}^i|$ )を再識別率として加算する。その際、最近傍レコードが複数ある場合は、それらも平均して( $1/|Z|$ )再識別できたとして加算する(9~15行目に対応)。

30

## 【0027】

もし近傍レコード集合 $X_i^{near}$ のレコード数が $K$ 件以下であった場合は、まず近傍レコード集合 $X_i^{near}$ の各レコードとの距離計算を行う。そして、最も距離が近いレコードの重複するレコードの真の行番号を探索し、現在のレコードの行番号と一致した場合、 $1$ を再識別率として加算する。その際、最近傍レコードが複数ある場合は、それらを平均して( $1/|Z|$ )再識別できたとして加算する(16~21行目に対応)。

## 【0028】

## [実施形態]

実施形態の安全性評価装置および方法は、上記<Algorithm 1>を実行して秘匿データベースの安全性を評価する。実施形態の安全性評価装置1は、図4に例示するように、データベース記憶部10、近傍レコード探索部11、重複排除部12、最近傍レコード計算部13、再特定判定部14、および再識別率計算部15を備える。この安全性評価装置1が、図5に例示する各ステップの処理を行うことにより実施形態の安全性評価方法が実現される。

40

## 【0029】

安全性評価装置1は、例えば、中央演算処理装置(CPU: Central Processing Unit)、主記憶装置(RAM: Random Access Memory)などを有する公知又は専用のコンピュータに

50

特別なプログラムが読み込まれて構成された特別な装置である。安全性評価装置 1 は、例えば、中央演算処理装置の制御のもとで各処理を実行する。安全性評価装置 1 に入力されたデータや各処理で得られたデータは、例えば、主記憶装置に格納され、主記憶装置に格納されたデータは必要に応じて中央演算処理装置へ読み出されて他の処理に利用される。安全性評価装置 1 の各処理部は、少なくとも一部が集積回路等のハードウェアによって構成されていてもよい。安全性評価装置 1 が備える各記憶部は、例えば、RAM (Random Access Memory) などの主記憶装置、ハードディスクや光ディスクもしくはフラッシュメモリ (Flash Memory) のような半導体メモリ素子により構成される補助記憶装置、またはリレーショナルデータベースやキーバリューストアなどのミドルウェアにより構成することができる。

10

## 【0030】

以下、図 5 を参照して、実施形態の安全性評価装置 1 が実行する安全性評価方法について説明する。

## 【0031】

データベース記憶部 10 には、平文のレコード  $x_i$  ( $i=1, \dots, N, N+2$ ) からなる元データベース  $X = \{x_1, \dots, x_N\}$  と、元データベース  $X$  を秘匿した秘匿データベース  $Y = \{y_1, \dots, y_N\}$  と、元データベース  $X$  の行番号と秘匿データベース  $Y$  の行番号とを対応付ける行番号対応関数  $f_y$  が記憶されている。

## 【0032】

ステップ S 11 において、近傍レコード探索部 11 は、秘匿データベース  $Y$  の各レコード  $y_i$  について、元データベース  $X$  に対する近傍探索により所定の近傍数  $K$  の近傍レコード集合  $X_i^{near} = \{x_j\}$  ( $j \in \{1, \dots, N\}$ ) を取得する。このとき、近傍数  $K$  は、例えば、元データベース  $X$  のレコード件数  $N$  の対数  $\log N$  とする。近傍探索は、木構造もしくはハッシングを用いた近傍探索のどちらかを用い、例えば、kd木を用いる手法を用いる。近傍レコード探索部 11 は、取得した近傍レコード集合  $X_i^{near}$  を重複排除部 12 へ出力する。

20

## 【0033】

ステップ S 12 において、重複排除部 12 は、近傍レコード集合  $X_i^{near}$  のレコード数が所定の閾値  $K$  よりも多い場合に、近傍レコード集合  $X_i^{near}$  中の重複レコードを排除して、重複排除済み近傍レコード集合  $X_i^{uniq}$  を生成する。このとき、閾値  $K$  は、例えば、近傍数の 2 倍 (すなわち、許容範囲  $=2.0$ ) とする。重複排除部 12 は、重複排除済み近傍レコード集合  $X_i^{uniq}$  を最近傍レコード計算部 13 へ出力する。近傍レコード集合  $X_i^{near}$  のレコード数が閾値  $K$  以下だった場合は、近傍レコード集合  $X_i^{near}$  を最近傍レコード計算部 13 へ出力する。

30

## 【0034】

ステップ S 13 において、最近傍レコード計算部 13 は、秘匿データベース  $Y$  の各レコード  $y_i$  について、近傍レコード集合  $X_i^{near}$  のレコード数が閾値  $K$  よりも多かった場合には、重複排除済み近傍レコード集合  $X_i^{uniq}$  の各レコードとの距離を計算し、近傍レコード集合  $X_i^{near}$  のレコード数が閾値  $K$  以下だった場合には、近傍レコード集合  $X_i^{near}$  の各レコードとの距離を計算し、当該レコード  $y_i$  との距離が最も近い最近傍レコード  $Z = \{x_k\}$  ( $k \in \{1, \dots, N\}$ ) を取得する。最近傍レコード計算部 13 は、取得した最近傍レコード  $Z$  を再特定判定部 14 へ出力する。

40

## 【0035】

ステップ S 14 において、再特定判定部 14 は、秘匿データベース  $Y$  の各レコード  $y_i$  について、当該レコード  $y_i$  に対応付けられた元データベース  $X$  のレコード  $x_j$  が最近傍レコード  $Z$  中に存在するか否かに基づいて当該レコード  $y_i$  の再識別率  $r_i$  を計算する。レコード  $y_i$  に対応付けられたレコード  $x_j$  は行番号対応関数  $f_y$  を用いて求めることができる。再特定判定部 14 は、計算したレコード  $y_i$  の再識別率  $r_i$  を再識別率計算部 15 へ出力する。

## 【0036】

50

ステップS15において、再識別率計算部15は、秘匿データベースYの各レコード $y_i$ について計算した再識別率 $r_i$ に基づいて秘匿データベースYの再識別率 $r$ を計算する。例えば、秘匿データベースYの各レコード $y_i$ の再識別率 $r_i$ の総和 $\sum_{i=1}^N r_i$ を秘匿データベースYの再識別率 $r$ とする。再識別率計算部15は、秘匿データベースYの再識別率 $r$ を安全性評価装置1の出力とする。

【0037】

本形態のポイントは、レコードリンケージに対して近傍探索を単に組み合わせただけでは解決できない課題、すなわち、近傍レコードが大量に出現した場合の問題を解決したことである。具体的には、大量の近傍レコードに対して重複排除処理を加えることで、処理時間を抑えたことである。近傍レコードを取得する際に属性値が重複等している場合、指定した近傍数以上の近傍レコードを取得するため、そのまま処理を行うと実行時間が長くなる。最悪の場合、近傍レコードがデータベース中のレコード数分出力されてしまい、結果として近傍レコードを探索した効果がなくなってしまう。本形態では、近傍レコードを取得する際に重複排除の処理を加えていることから、上記問題を回避でき、高速な実行が可能となっている。

10

【0038】

以上、この発明の実施の形態について説明したが、具体的な構成は、これらの実施の形態に限られるものではなく、この発明の趣旨を逸脱しない範囲で適宜設計の変更等があっても、この発明に含まれることはいうまでもない。実施の形態において説明した各種の処理は、記載の順に従って時系列に実行されるのみならず、処理を実行する装置の処理能力あるいは必要に応じて並列的にあるいは個別に実行されてもよい。

20

【0039】

[プログラム、記録媒体]

上記実施形態で説明した各装置における各種の処理機能をコンピュータによって実現する場合、各装置が有すべき機能の処理内容はプログラムによって記述される。そして、このプログラムをコンピュータで実行することにより、上記各装置における各種の処理機能がコンピュータ上で実現される。

【0040】

この処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、例えば、磁気記録装置、光ディスク、光磁気記録媒体、半導体メモリ等のようなものでもよい。

30

【0041】

また、このプログラムの流通は、例えば、そのプログラムを記録したDVD、CD-ROM等の可搬型記録媒体を販売、譲渡、貸与等することによって行う。さらに、このプログラムをサーバコンピュータの記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピュータにそのプログラムを転送することにより、このプログラムを流通させる構成としてもよい。

【0042】

このようなプログラムを実行するコンピュータは、例えば、まず、可搬型記録媒体に記録されたプログラムもしくはサーバコンピュータから転送されたプログラムを、一旦、自己の記憶装置に格納する。そして、処理の実行時、このコンピュータは、自己の記憶装置に格納されたプログラムを読み取り、読み取ったプログラムに従った処理を実行する。また、このプログラムの別の実行形態として、コンピュータが可搬型記録媒体から直接プログラムを読み取り、そのプログラムに従った処理を実行することとしてもよく、さらに、このコンピュータにサーバコンピュータからプログラムが転送されるたびに、逐次、受け取ったプログラムに従った処理を実行することとしてもよい。また、サーバコンピュータから、このコンピュータへのプログラムの転送は行わず、その実行指示と結果取得のみによって処理機能を実現する、いわゆるASP(Application Service Provider)型のサービスによって、上述の処理を実行する構成としてもよい。なお、本形態におけるプログラムには、電子計算機による処理の用に供する情報であってプログラムに準ずるもの(コンピ

40

50

ュータに対する直接の指令ではないがコンピュータの処理を規定する性質を有するデータ等)を含むものとする。

【0043】

また、この形態では、コンピュータ上で所定のプログラムを実行させることにより、本装置を構成することとしたが、これらの処理内容の少なくとも一部をハードウェア的に実現することとしてもよい。

【図1】

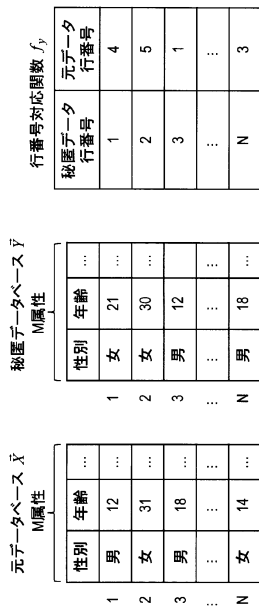


図1

【図2】

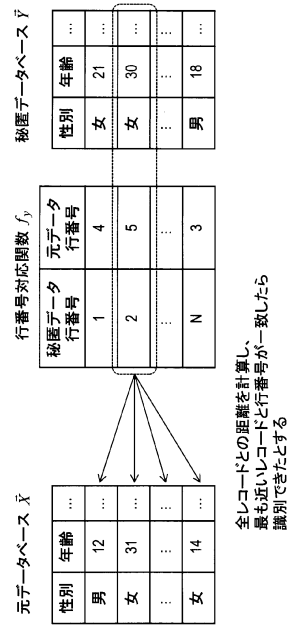


図2

【 図 3 】

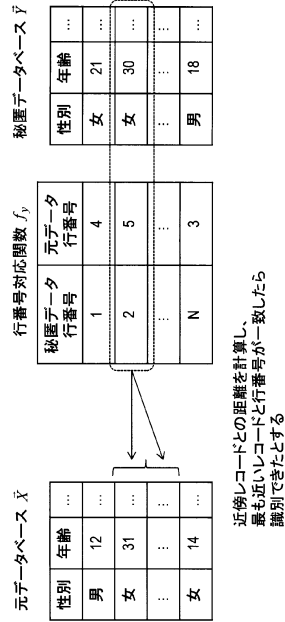


図3

【 図 4 】

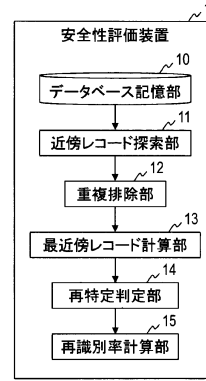


図4

【 図 5 】

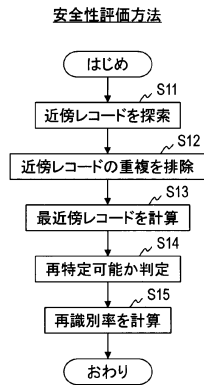


図5

---

フロントページの続き

(56)参考文献 特開2018-49437(JP,A)

国際公開第2007/132564(WO,A1)

伊藤聡志, 菊池浩明, ユークリッド距離を用いた再識別手法とPWSCup2015の匿名加工データを用いた評価, 電子情報通信学会技術研究報告 Vol.116 No.65 IEICE Technical Report, 日本, 一般社団法人電子情報通信学会, 2016年05月19日, 第116巻, 第65号, p.145-p.152

(58)調査した分野(Int.Cl., DB名)

G06F 21/62