

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7363662号
(P7363662)

(45)発行日 令和5年10月18日(2023.10.18)

(24)登録日 令和5年10月10日(2023.10.10)

(51)国際特許分類 F I
G 0 6 F 21/62 (2013.01) G 0 6 F 21/62 3 5 4

請求項の数 6 (全20頁)

(21)出願番号	特願2020-79550(P2020-79550)	(73)特許権者	000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号
(22)出願日	令和2年4月28日(2020.4.28)	(74)代理人	100092978 弁理士 真田 有
(65)公開番号	特開2021-174390(P2021-174390 A)	(74)代理人	100189201 弁理士 横田 功
(43)公開日	令和3年11月1日(2021.11.1)	(72)発明者	福岡 尊 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
審査請求日	令和5年1月12日(2023.1.12)	(72)発明者	山岡 裕司 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	宮司 卓佳

最終頁に続く

(54)【発明の名称】 生成方法、情報処理装置及び生成プログラム

(57)【特許請求の範囲】

【請求項1】

複数の項目値を含む複数の個人情報を受け付け、

前記複数の項目値それぞれに対応付けられた第1のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第1の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第1の匿名情報を生成し、

前記複数の項目値それぞれに対応付けられた前記第1のパラメータとは異なる第2のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第2の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第2の匿名情報を生成し、

前記第1の匿名情報及び前記第2の匿名情報を含む匿名情報を生成する、
処理をコンピュータが実行することを特徴とする生成方法。

【請求項2】

前記第1及び第2のパラメータについての多様性に関する指標が最大化されるように、当該第1及び第2のパラメータを決定する、
処理を前記コンピュータに実行させることを特徴とする、請求項1に記載の生成方法。

【請求項3】

前記第1及び第2のパラメータは、前記複数の項目値間の優先順序である、
ことを特徴とする、請求項1又は2に記載の生成方法。

【請求項 4】

前記第 1 及び第 2 のパラメータは、前記複数の項目値それぞれの重み付け値である、ことを特徴とする、請求項 1 又は 2 に記載の生成方法。

【請求項 5】

複数の項目値を含む複数の個人情報を受け付ける受付処理部と、

前記複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 1 の匿名情報を生成する第 1 匿名情報生成部と、

前記複数の項目値それぞれに対応付けられた前記第 1 のパラメータとは異なる第 2 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 2 の匿名情報を生成する第 2 匿名情報生成部と、

前記第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する結合情報生成部と、

を備えることを特徴とする情報処理装置。

【請求項 6】

複数の項目値を含む複数の個人情報を受け付け、

前記複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 1 の匿名情報を生成し、

前記複数の項目値それぞれに対応付けられた前記第 1 のパラメータとは異なる第 2 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 2 の匿名情報を生成し、

前記第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する、処理をコンピュータに実行させることを特徴とする生成プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、生成方法、情報処理装置及び生成プログラムに関する。

【背景技術】

【0002】

顧客から預かったデータをもとに、人工知能（AI）製品を展開するビジネスが存在する。このような AI 製品としては、例えば、顧客から預かったデータを利用して機械学習し、事象を予測するモデルがある。

【0003】

図 1 は、AI 製品の機械学習を例示する図である。

【0004】

図 1 の A 1 に示す Id と職業と性別と年収とが対応付けられたデータに対して、年収を目的変数として機械学習を実行することにより、符号 A 2 に示すように、職業及び性別から、年収を予測するモデルが生成される。

【0005】

利用するデータが個人情報である場合には、匿名化処理が実行されることがある。

【先行技術文献】

【特許文献】

【0006】

【文献】特開 2017 - 182508 号公報

再表 2013 - 114445 号公報

10

20

30

40

50

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、匿名化処理された匿名化データを学習で利用すると、匿名化データはオリジナルのデータよりも情報量が低下するため、生成されるモデルの精度が低下するおそれがある。

【0008】

1つの側面では、機械学習によって生成する学習モデルの精度を向上させることを目的とする。

【課題を解決するための手段】

【0009】

1つの側面では、生成方法は、複数の項目値を含む複数の個人情報を受け付け、前記複数の項目値それぞれに対応付けられた第1のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第1の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第1の匿名情報を生成し、前記複数の項目値それぞれに対応付けられた前記第1のパラメータとは異なる第2のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第2の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第2の匿名情報を生成し、前記第1の匿名情報及び前記第2の匿名情報を含む匿名情報を生成する、処理をコンピュータが実行する。

【発明の効果】

【0010】

1つの側面では、機械学習によって生成する学習モデルの精度を向上できる。

【図面の簡単な説明】

【0011】

【図1】AI製品の機械学習を例示する図である。

【図2】実施形態の一例におけるk-匿名化を例示するテーブルである。

【図3】実施形態の一例における情報処理装置のハードウェア構成例を模式的に示すブロック図である。

【図4】図3に示した情報処理装置のソフトウェア構成例を模式的に示すブロック図である。

【図5】関連例における抑制k-匿名化を例示するテーブルである。

【図6】関連例における一般化匿名化を例示するテーブルである。

【図7】図6に示したテーブルに対応する一般化木を例示する図である。

【図8】実施形態の一例におけるk-匿名化データの結合例を示す図である。

【図9】実施形態の一例における匿名化データの生成処理を説明するフローチャートである。

【図10】実施形態の一例における抑制k-匿名化処理を利用した匿名化処理の生成処理を説明するフローチャートである。

【図11】図10に示したオリジナルデータからの匿名化データの結合例を示す図である。

【図12】実施形態の一例における一般化匿名化処理を利用した匿名化処理の生成処理を説明するフローチャートである。

【図13】図12に示したオリジナルデータからの匿名化データの結合例を示す図である。

【図14】関連例における学習モデルの生成処理を模式的に示す図である。

【図15】実施形態の一例における学習モデルの生成処理を模式的に示すブロック図である。

【図16】関連例における学習モデルと実施形態の一例における学習モデルとの精度を比較するテーブルである。

【図17】実施形態の一例におけるデータの削除箇所を例示する図である。

10

20

30

40

50

【図18】実施形態の一例におけるデータの削除箇所の決定アルゴリズムの第1の例を説明するテーブルである。

【図19】学習モデルの生成処理の違いによる学習モデルの精度を比較するテーブルである。

【発明を実施するための形態】

【0012】

〔A〕実施形態

以下、図面を参照して一実施の形態を説明する。ただし、以下に示す実施形態はあくまでも例示に過ぎず、実施形態で明示しない種々の変形例や技術の適用を排除する意図はない。すなわち、本実施形態を、その趣旨を逸脱しない範囲で種々変形して実施することができる。また、各図は、図中に示す構成要素のみを備えるという趣旨ではなく、他の機能等を含むことができる。

10

【0013】

以下、図中において、同一の各符号は同様の部分を示しているため、その説明は省略する。

【0014】

〔A-1〕概要

図2は、実施形態の一例におけるk-匿名化を例示するテーブルである。

【0015】

実施形態の一例においては、分類モデルの精度を下げないk-匿名化データが構築される。k-匿名化データとは、同一レコードが少なくともk個ある個人を特定できないデータである。ただし、全て欠損しているレコードは無視される。k-匿名化が施されているか否かは容易にチェックできる。

20

【0016】

図2においては、符号B1に示すオリジナルデータを入力として、2-匿名化により、符号B2に示す匿名化データ（別言すれば、匿名情報）が構築される。

【0017】

k-匿名化においては、「どの属性に関する情報を残したいか」といった属性に関する設定が、入力として求められる場合がある。図示する例では、年収及び性別ができるだけ残されるように、匿名化が実行されている。

30

【0018】

図3は、実施形態の一例における情報処理装置1のハードウェア構成例を模式的に示すブロック図である。

【0019】

図3に示すように、情報処理装置1は、Central Processing Unit (CPU) 11, メモリ部12, 表示制御部13, 記憶装置14, 入力Interface (IF) 15, 外部記録媒体処理部16及び通信IF17を備える。

【0020】

メモリ部12は、記憶部の一例であり、例示的に、Read Only Memory (ROM) 及びRandom Access Memory (RAM) などである。メモリ部12のROMには、Basic Input/Output System (BIOS) 等のプログラムが書き込まれてよい。メモリ部12のソフトウェアプログラムは、CPU11に適宜に読み込まれて実行されてよい。また、メモリ部12のRAMは、一時記録メモリあるいはワーキングメモリとして利用されてよい。

40

【0021】

表示制御部13は、表示装置130と接続され、表示装置130を制御する。表示装置130は、液晶ディスプレイやOrganic Light-Emitting Diode (OLED) ディスプレイ, Cathode Ray Tube (CRT), 電子ペーパーディスプレイ等であり、オペレータ等に対する各種情報を表示する。表示装置130は、入力装置と組み合わせられたものでもよく、例えば、タッチパネルでもよい。

50

【 0 0 2 2 】

記憶装置 1 4 は、高 I O 性能の記憶装置であり、例えば、Dynamic Random Access Memory (D R A M) や Solid State Drive (S S D) , Storage Class Memory (S C M) , Hard Disk Drive (H D D) が用いられてよい。

【 0 0 2 3 】

入力 I F 1 5 は、マウス 1 5 1 やキーボード 1 5 2 等の入力装置と接続され、マウス 1 5 1 やキーボード 1 5 2 等の入力装置を制御してよい。マウス 1 5 1 やキーボード 1 5 2 は、入力装置の一例であり、これらの入力装置を介して、オペレータが各種の入力操作を行なう。

【 0 0 2 4 】

外部記録媒体処理部 1 6 は、記録媒体 1 6 0 が装着可能に構成される。外部記録媒体処理部 1 6 は、記録媒体 1 6 0 が装着された状態において、記録媒体 1 6 0 に記録されている情報を読み取り可能に構成される。本例では、記録媒体 1 6 0 は、可搬性を有する。例えば、記録媒体 1 6 0 は、フレキシブルディスク、光ディスク、磁気ディスク、光磁気ディスク、又は、半導体メモリ等である。

【 0 0 2 5 】

通信 I F 1 7 は、外部装置との通信を可能にするためのインタフェースである。

【 0 0 2 6 】

C P U 1 1 は、種々の制御や演算を行なう処理装置であり、メモリ部 1 2 に格納された Operating System (O S) やプログラムを実行することにより、種々の機能を実現する。

【 0 0 2 7 】

情報処理装置 1 全体の動作を制御するための装置は、C P U 1 1 に限定されず、例えば、M P U や D S P , A S I C , P L D , F P G A のいずれか 1 つであってもよい。また、情報処理装置 1 全体の動作を制御するための装置は、C P U , M P U , D S P , A S I C , P L D 及び F P G A のうちの 2 種類以上の組み合わせであってもよい。なお、M P U は Micro Processing Unit の略称であり、D S P は Digital Signal Processor の略称であり、A S I C は Application Specific Integrated Circuit の略称である。また、P L D は Programmable Logic Device の略称であり、F P G A は Field Programmable Gate Array の略称である。

【 0 0 2 8 】

図 4 は、図 3 に示した情報処理装置 1 のソフトウェア構成例を模式的に示す図である。

【 0 0 2 9 】

情報処理装置 1 は、受付処理部 1 1 1 , 匿名情報生成部 1 1 2 及び結合情報生成部 1 1 3 として機能する。

【 0 0 3 0 】

受付処理部 1 1 1 は、複数の項目値を含む複数の個人情報を受け付ける。

【 0 0 3 1 】

匿名情報生成部 1 1 2 は、複数の項目値それぞれに対応付けられたパラメータに応じて複数の項目値それぞれの秘匿化されやすさが決まる秘匿化アルゴリズムを、受け付けた複数の個人情報に適用して、複数の項目値の少なくとも何れかの項目値が匿名化された匿名情報を生成する。具体的には、匿名情報生成部 1 1 2 は、複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて複数の項目値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた複数の個人情報に適用して、複数の項目値の少なくとも何れかの項目値が匿名化された第 1 の匿名情報を生成する第 1 匿名情報生成部の一例として機能する。また、匿名情報生成部 1 1 2 は、複数の項目値それぞれに対応付けられた第 1 のパラメータとは異なる第 2 のパラメータに応じて複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化アルゴリズムを、受け付けた複数の個人情報に適用して、複数の項目値の少なくとも何れかの項目値が匿名化された第 2 の匿名情報を生成する第 2 匿名情報生成部の一例として機能する。なお、詳細は後述されるが、開示の技術におけるパラメータまたは秘匿化アルゴリズムは 2 つに限定されるわけではなく、3 つ以上で

10

20

30

40

50

あってもよい。

【 0 0 3 2 】

結合情報生成部 1 1 3 は、第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する。

【 0 0 3 3 】

〔 A - 2 〕匿名化処理

k - 匿名化においては、データ品質を表す量を定義し、その量が最大化される。データの品質を表す量としては、例えば、加工（削除等）されなかったセルの数やエントロピーがある。

【 0 0 3 4 】

しかしながら、データ品質を表す量が大きくても、うまく学習モデルを構成できる保証はないため、出力を学習データとして用いると精度劣化が大きくなるおそれがある。モデル学習に最適な匿名化データを見つけることは計算量的に容易でなく、現実的ではない。

【 0 0 3 5 】

図 5 は、関連例における抑制 k - 匿名化を例示するテーブルである。

【 0 0 3 6 】

抑制 k - 匿名化においては、セルを削除することにより、k - 匿名化を実現することができる。符号 C 1 に示すような表形式のオリジナルデータと、k（正の整数）と、属性の順番付け（属性優先順序と称されてもよい。）とが入力されると、符号 C 2 及び C 3 に示すような表形式の k - 匿名化データが出力される。

【 0 0 3 7 】

符号 C 2 に示す k - 匿名化データでは、属性優先順序が「年収 > 性別 > 職業」に設定されている。また、符号 C 3 に示す k - 匿名化データでは、属性優先順序が「年収 > 職業 > 性別」に設定されている。

【 0 0 3 8 】

抑制 k - 匿名化においては、匿名化の際に、属性優先順序が「属性に関する設定」として使用される。これにより、属性優先順序が高い属性ほど、データが残りやすいように匿名化が実行される。

【 0 0 3 9 】

図 6 は、関連例における一般化匿名化を例示するテーブルである。図 7 は、図 6 に示したテーブルに対応する一般化木を例示する図である。

【 0 0 4 0 】

一般化 k - 匿名化においては、抑制 k - 匿名化におけるセルの削除に加えて、セルの置き換えも実施される。

【 0 0 4 1 】

図 6 の符号 D 1 におけるオリジナルデータは、符号 D 2 及び D 3 に示すように、一般化された k - 匿名化データとして出力される。符号 D 2 に示す k - 匿名化データでは、職業の一般化が行なわれている。また、符号 D 3 に示す k - 匿名化データでは、住所の一般化が行なわれている。

【 0 0 4 2 】

図 7 の符号 E 1 に示す一般化木においては、「杉並区」、 「世田谷区」及び「目黒区」の一般化概念として「東京都」が定義されており、「横浜市」及び「川崎市」の一般化概念として「神奈川県」が定義されている。

【 0 0 4 3 】

また、図 7 の符号 E 2 に示す一般化木においては、「飲食店」及び「美容師」の一般化概念として「自営業」が定義されており、「開発」及び「営業」の一般化概念として「会社員」が定義されており、「教授」及び「教諭」の一般化概念として「教育職」が定義されている。

【 0 0 4 4 】

なお、図 7 の符号 E 3 に示す一般化木において、年収の「500 万円以上」及び「50

10

20

30

40

50

0万円未満」については、一般化概念が定義されていない。

【0045】

匿名化は、図7に示した一般化木に加えて、各属性の重みも入力とし、次式のNCPを最小化するように実行されてよい。

【0046】

【数1】

$$NCP = \sum_A (Aの重み) \sum_T \frac{TのAにおける値を先祖とするリーフの数}{Aに対応する一般化木のリーフの数}$$

10

なお、Aは属性であり、Tはレコードである。重みが大きい（別言すれば、重要な）属性は、匿名化されづらくなる。

【0047】

ここで、重みを(住所,職業,年収)=(a,b,c)とすると、図6の符号D2に示した匿名化データのNCPは $5a+(5/3)b+4c$ となり、図6の符号D3に示した匿名化データのNCPは $2a+5b+3c$ となる。

【0048】

(a,b,c)=(0.1,0.6,0.3)とすると、図6の符号D2に示した匿名化データのNCPは $5*0.1+(5/3)*0.6+4*0.3=2.7$ となり、図6の符号D3に示した匿名化データのNCPは $2*0.1+5*0.6+3*0.3=4.1$ となる。すなわち、図6の符号D2に示した匿名化データの方がNCPが低くなる。

20

【0049】

一方、(a,b,c)=(0.3,0.1,0.6)とすると、図6の符号D2に示した匿名化データのNCPは $5*0.3+(5/3)*0.1+4*0.6=4.0666\dots$ となり、図6の符号D3に示した匿名化データのNCPは $2*0.3+5*0.1+3*0.6=2.9$ となる。すなわち、図6の符号D3に示した匿名化データの方がNCPが低くなる。

【0050】

〔A-3〕結合処理

図8は、実施形態の一例におけるk-匿名化データの結合例を示す図である。

30

【0051】

実施形態の一例において、オリジナルデータを機械学習するにあたって、レコード数は保たれなくてもよい。そこで、複数の異なるk-匿名化データが結合されることで、機械学習に適した匿名化データが生成される。

【0052】

符号F1に示すオリジナルデータを入力として、符号F2に示すように2つのk-匿名化データが出力される。そして、符号F3に示すように、2つのk-匿名化データが結合されることにより、学習モデルが生成される。

【0053】

図9は、実施形態の一例における匿名化データの生成処理を説明するフローチャートである。図9に示すフローチャート(ステップS1~S3)に従って、匿名化データの生成処理を説明する。

40

【0054】

必要に応じて優先順序等の補助入力を受け付け、属性に関するパラメータとして、匿名化アルゴリズム設定#1~#nが生成される(ステップS1)。

【0055】

表形式のオリジナルデータの入力を受け付け、各匿名化アルゴリズム設定#1~#nに応じた匿名化アルゴリズムで匿名化処理が実行され、匿名化データ#1~#nが出力される(ステップS2)。

【0056】

50

匿名化データ # 1 ~ # n について互いに結合処理が実行され、結合された匿名化データが出力される (ステップ S 3)。そして、匿名化データの生成処理は終了する。

【 0 0 5 7 】

図 1 0 は、実施形態の一例における抑制 k - 匿名化処理を利用した匿名化処理の生成処理を説明するフローチャートである。図 1 1 は、図 1 0 に示したオリジナルデータからの匿名化データの結合例を示す図である。図 1 0 に示すフローチャート (ステップ S 1 1 ~ S 1 3) に従って、抑制 k - 匿名化処理を利用した匿名化処理の生成処理を説明する。

【 0 0 5 8 】

表形式のオリジナルデータが入力として受け付けられる (ステップ S 1 1)。

【 0 0 5 9 】

属性に関するパラメータ (別言すれば、属性優先順序) が、匿名化データの生成個数と順序を固定したい属性とを追加入力とした上で、ランダムに生成される (ステップ S 1 2)。

【 0 0 6 0 】

生成した各属性優先順序の設定の元でそれぞれの匿名化データが出力され、得られた匿名化データが結合される (ステップ S 1 3)。そして、抑制 k - 匿名化処理を利用した匿名化処理の生成処理は終了する。

【 0 0 6 1 】

図 1 1 に示す例では、生成個数が「2」に設定され、「年収」の属性優先順序が一番として固定され、残った2つの属性について順序をランダムにして匿名化データが生成される。符号 G 1 に示す例では、「年収 職業 性別」及び「年収 性別 職業」を属性優先順序とする匿名化データが生成される。そして、符号 G 2 に示すように、2つの匿名化データが結合されて学習モデルが出力される。

【 0 0 6 2 】

図 1 2 は、実施形態の一例における一般化匿名化処理を利用した匿名化処理の生成処理を説明するフローチャートである。図 1 3 は、図 1 2 に示したオリジナルデータからの匿名化データの結合例を示す図である。図 1 2 に示すフローチャート (ステップ S 2 1 ~ S 2 3) に従って、一般化匿名化処理を利用した匿名化処理の生成処理を説明する。

【 0 0 6 3 】

表形式のオリジナルデータが入力として受け付けられる (ステップ S 2 1)。

【 0 0 6 4 】

属性に関するパラメータ (別言すれば、属性に対する重み) がランダムに生成される (ステップ S 2 2)。

【 0 0 6 5 】

生成した各属性優先順序の設定の元でそれぞれの匿名化データが出力され、得られた匿名化データが結合される (ステップ S 2 3)。そして、一般化匿名化処理を利用した匿名化処理の生成処理は終了する。

【 0 0 6 6 】

図 1 3 に示す例では、属性に対する重みとして、(住所, 職業, 年収) = (0.1, 0.6, 0.3), (0.3, 0.1, 0.6) が生成されると、符号 H 1 に示す匿名化データが生成される。そして、符号 H 2 に示すように、2つの匿名化データが結合されて学習モデルが出力される。

【 0 0 6 7 】

〔 A - 4 〕 関連例との比較

図 1 4 は、関連例における学習モデルの生成処理を模式的に示す図である。

【 0 0 6 8 】

関連例においては、図 1 4 の符号 I 1 に示すオリジナルデータを入力として、符号 I 2 に示す1つの匿名化データが生成される。そして、1つの匿名化データに対して機械学習が実施されることにより、符号 I 3 に示すように、学習モデルが生成される。

【 0 0 6 9 】

図 1 5 は、実施形態の一例における学習モデルの生成処理を模式的に示すブロック図で

10

20

30

40

50

ある。

【 0 0 7 0 】

一方、実施形態の一例においては、図 1 5 の符号 J 1 に示すオリジナルデータを入力として、符号 J 2 に示す複数の匿名化データ # 1 ~ # n が生成される。複数の匿名化データ # 1 ~ # n が結合されて、符号 J 3 に示すように、結合匿名化データが生成される。そして、結合匿名化データに対して機械学習が実施されることにより、符号 J 4 に示すように、学習モデルが生成される。

【 0 0 7 1 】

図 1 6 は、関連例における学習モデルと実施形態の一例における学習モデルとの精度を比較するテーブルである。

【 0 0 7 2 】

図 1 6 に示す例では、単一の匿名化データで学習したモデルと、複数の匿名化データで学習したモデルとの精度が比較されている。

【 0 0 7 3 】

実験方法として、2 から 5 0 までの k に対して、8 個の k - 匿名化データを生成する。次に、2 から 5 0 までの k に対し、8 個の k - 匿名化データそれぞれで学習した 8 個の学習モデルによる精度のうち最高のものと、8 個の k - 匿名化データを結合したデータで学習した学習モデルの精度とを記録する。そして、それぞれの精度について、k に関する平均値、最小値をとった。

【 0 0 7 4 】

このような実験の結果、図 1 6 に示すように、複数の匿名化データから複数のモデルを作るよりも、複数の匿名化データを結合して一つのモデルを作った方が、平均値及び最小値の両方の精度が高くなる。

【 0 0 7 5 】

〔 A - 5 〕匿名化データの生成処理の詳細

匿名化データの生成処理の具体例としては、データオーギュメンテーションとランダム生成とが想定される。

【 0 0 7 6 】

データオーギュメンテーションにおいては、安定性は高くなるものの、精度が低くなるおそれがある。例えば、よく似た匿名化データが結合されることにより、多様性が失われてしまい、学習モデルに汎用性がなくなる。特に、表形式のオリジナルデータの場合には、画像データとは異なり、似たようなデータが入力されても、機械学習の効果が薄くなる。また、匿名化データが少しずつ変更されるため、結合匿名化データが、匿名化アルゴリズムのハイパーパラメータの初期値に強く依存するおそれがある。

【 0 0 7 7 】

一方、ランダム生成においては、安定性が低くなるおそれがある。匿名化データがランダムに生成されてしまうので、学習に適した匿名化データが生成されないことがある。また、大量に生成して多様性を担保することは、学習コストの面で非効率になる。

【 0 0 7 8 】

図 1 7 は、実施形態の一例におけるデータの削除箇所を例示する図である。

【 0 0 7 9 】

符号 K 1 に示す様に、オリジナルデータから複数の匿名化データ # 1 ~ # m 間で似たような部分を削除するのは学習に適さないと想定される一方、符号 K 2 に示す様に、オリジナルデータから複数の匿名化データ # 1 ~ # m 間で大きく異なる部分を削除するのが学習に適すると想定される。

【 0 0 8 0 】

すなわち、“大いに異なる”匿名化データを構成することができれば、データを補完し合えるので、学習に適したデータを生成できると想定される。

【 0 0 8 1 】

匿名化データの生成処理の第 1 の具体例として、匿名化データの水増し件数 m 及び順序

10

20

30

40

50

同士の距離関数を入力とする。順序同士の距離関数は、Kendallの距離やCayley距離等の任意の関数でよい。また、属性数を n としたとき、 $m \leq n!$ を満たすものとする。

【0082】

与えられた距離関数によって、順序の間の距離の総和といった多様性を表す指標が最大となる m 個の異なる属性順序が、全ての組み合わせを調べることによって決定される。

【0083】

決定された m 個の異なる属性順序を使って、属性に関する設定が m 個作成され、 k - 匿名化データが生成・結合される。

【0084】

図18は、実施形態の一例におけるデータの削除箇所の決定アルゴリズムの第1の例を説明するテーブルである。

10

【0085】

ここで、匿名化データの生成処理の第1の具体例において、水増し件数を $m=2$ 、属性数は $n=3$ 、距離関数を Kendall の距離関数とする。なお、Kendall の距離は、二つの1から n の整数からなる配列 a, b が与えられたとき、 $i < j$ を満たす1から n の整数の組 (i, j) であって、 a の i 番目の数と a の j 番目の数の間の大小関係が、 b のそれと食い違っているものを数え上げた数である。

【0086】

順序間の距離は、図18に示すようになる。

【0087】

多様性を表す指標として、距離が最大になる異なる2つの属性優先順序を、全ての組み合わせを調べ決定する。この場合の異なる2つの属性優先順序は、 $\{(123), (132)\}$, $\{(123), (213)\}$, $\{(123), (231)\}$, $\{(123), (312)\}$, $\{(123), (321)\}$, $\{(132), (213)\}$, $\{(132), (231)\}$, $\{(132), (312)\}$, $\{(132), (321)\}$, $\{(213), (231)\}$, $\{(213), (312)\}$, $\{(213), (321)\}$, $\{(231), (312)\}$, $\{(231), (321)\}$, $\{(312), (321)\}$ となる。

20

【0088】

2つの属性優先順序の間の距離を、図18を用いて計算すると、それぞれ1,1,2,2,3,2,1,3,2,3,1,2,2,1,1となる。

【0089】

そして、最大となるものが選択される。複数ある場合はランダムにとることで、一つ選択される。本例では、 $\{(123), (321)\}$, $\{(132), (312)\}$, $\{(213), (231)\}$ から一つが選択される。

30

【0090】

匿名化データの生成処理の第2の具体例として、匿名化データの水増し件数 m を入力とする。また、属性数を n としたとき、 $m \leq 2n$ を満たすものとする。

【0091】

以下、匿名化データの生成処理の第2の具体例におけるアルゴリズムを説明する。

【0092】

まず、初期値である属性優先順序に対して、順序を固定する属性が選ばれる。その後、動かす順序が決定され、その並びを $[1, 2, \dots, n]$ とおく。

40

【0093】

次に、数列 $(0, 1/2, 1/3, 2/3, 1/4, 2/4, 3/4, \dots)$ に n を掛け、整数に切り下げ、前から見て重複して現れた番号は除くことでできる数列を (a_1, \dots, a_n) とする。

【0094】

次に、1から n の各 i に対して、属性優先順序 b_i を $[a_{i+1}, a_{i+2}, \dots, n, 1, \dots, a_i]$ と置く。

【0095】

次に、属性優先順序の列 $B = (b_1, r(b_1), b_2, r(b_2), \dots, b_n, r(b_n))$ が生成される。ここで、優先順序 b に対して、 $r(b)$ はそれをひっくり返したものを表す。例えば、 $b = [3, 4, 1, 2]$ ならば $r(b) = [2, 1, 4, 3]$ である。

【0096】

50

そして、Bの先頭からm項をとり、それらの属性優先順序を用いて、k-匿名化データが生成・結合される。

【0097】

匿名化データの生成処理の第2の具体例におけるアルゴリズムでは、計算量は非常に少なくなる。またアルゴリズムの2~3が作用して、データ間のKendallの距離の総和は、mが偶数なら最大で、奇数の場合もランダムに比べ大きくなると期待できる。

【0098】

ここで、匿名化データの生成処理の第2の具体例において、水増し件数をm=3とする。属性は{年収、職業、住所、性別、学歴}の5種類で、さらに年収は必ず属性優先順序の最初に置くと決める。すなわち、動かす属性は{職業、住所、性別、学歴}の4つであるため、n=4の場合に対応する。

【0099】

初期値である属性優先順序を固定し、その並びを[1,2,3,4]とおく。なお、実際には[職業、学歴、性別、住所]などと並びが、便宜上数値とする。

【0100】

数列(0,1/2,1/3,2/3,1/4,2/4,3/4,...)に4を掛け、整数に切り下げ、前から見て重複して現れた番号は除くことのできる数列は(0,2,1,3)となる。

【0101】

1から4の各iに対して、属性優先順序biは、b1=[1,2,3,4], b2=[3,4,1,2], b3=[2,3,4,1], b4=[4,1,2,3]となる。

【0102】

属性優先順序の列B=(b1,r(b1),b2,r(b2),...,bn,r(bn))は、([1,2,3,4],[4,3,2,1],[3,4,1,2],[2,1,4,3],[2,3,4,1],[1,4,3,2],[4,1,2,3],[3,2,1,4])となる。

【0103】

Bの先頭からm=3項をとると、[1,2,3,4], [4,3,2,1], [3,4,1,2]が生成され、それらの属性優先順序を用いて、k-匿名化データが生成・結合される。

【0104】

本例では、結果として、[職業、学歴、性別、住所]、[住所、性別、学歴、職業]、[性別、住所、職業、学歴]の3つが生成される。

【0105】

図19は、学習モデルの生成処理の違いによる学習モデルの精度を比較するテーブルである。

【0106】

図19においては、データオーギュメンテーション及びランダム生成による匿名化データで学習した学習モデルと、匿名化データの生成処理の第2の具体例による匿名化データで学習した学習モデルとが比較されている。

【0107】

実験方法として、データオーギュメンテーション、ランダム生成A、ランダム生成B、匿名化データの生成処理の第2の具体例により、2から15までのkに対して、k-匿名化データを8個作成する。なお、ランダム生成A,Bは、異なるシードによるランダム生成を意味する。また、目的変数は必ず優先順序を1位とした。初期値となる属性優先順序は、学習器の特徴量重要度を用いた。

【0108】

それぞれの場合で、データを結合し、学習したモデルの精度を各kについて比較し、kに関する精度の平均、最小値、最大値を記録すると、図19に示すテーブルが得られた。

【0109】

図19に示す実験結果において、匿名化データの生成処理の第2の具体例は、データオーギュメンテーションよりも精度が高くなる。また、ランダム生成は精度が良い場合もあるが、シードによって値がばらつくため安定しないことが確認される。匿名化データの生成処理の第2の具体例はシードに寄らず、平均値は最大となる。

10

20

30

40

50

【 0 1 1 0 】

〔 A - 6 〕 効果

上述した実施形態の一例における生成方法，情報処理装置 1 及び生成プログラムによれば、例えば、以下の作用効果を奏することができる。

【 0 1 1 1 】

受付処理部 1 1 1 は、複数の項目値を含む複数の個人情報を受け付ける。匿名情報生成部 1 1 2 は、複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて複数の項目値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた複数の個人情報に適用して、複数の項目値の少なくとも何れかの項目値が匿名化された第 1 の匿名情報を生成する。また、匿名情報生成部 1 1 2 は、複数の項目値それぞれに対応付けられた第 1 のパラメータとは異なる第 2 のパラメータに応じて複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化アルゴリズムを、受け付けた複数の個人情報に適用して、複数の項目値の少なくとも何れかの項目値が匿名化された第 2 の匿名情報を生成する。そして、結合情報生成部 1 1 3 は、第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する。

10

【 0 1 1 2 】

これにより、個人情報の匿名性を確保しつつ、機械学習によって生成する学習モデルの精度を向上できる。

【 0 1 1 3 】

第 1 及び第 2 のパラメータについての多様性に関する指標が最大化されるように、当該第 1 及び第 2 のパラメータを決定する。これにより、匿名情報の多様性を向上させることができる。

20

【 0 1 1 4 】

第 1 及び第 2 のパラメータは、例えば、前記複数の項目値間の優先順序である。第 1 及び第 2 のパラメータは、例えば、前記複数の項目値それぞれの重み付け値であってもよい。これにより、多様性を有する匿名情報を容易に生成することができる。

【 0 1 1 5 】

〔 B 〕 その他

開示の技術は上述した実施形態に限定されるものではなく、本実施形態の趣旨を逸脱しない範囲で種々変形して実施することができる。本実施形態の各構成及び各処理は、必要に応じて取捨選択することができ、あるいは適宜組み合わせてもよい。

30

【 0 1 1 6 】

〔 C 〕 付記

以上の実施形態に関し、更に以下の付記を開示する。

【 0 1 1 7 】

(付記 1)

複数の項目値を含む複数の個人情報を受け付け、

前記複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 1 の匿名情報を生成し、

40

前記複数の項目値それぞれに対応付けられた前記第 1 のパラメータとは異なる第 2 のパラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化アルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された第 2 の匿名情報を生成し、

前記第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する、
処理をコンピュータが実行することを特徴とする生成方法。

【 0 1 1 8 】

(付記 2)

前記第 1 及び第 2 のパラメータについての多様性に関する指標が最大化されるように、

50

当該第 1 及び第 2 のパラメータを決定する、
処理を前記コンピュータに実行させることを特徴とする、付記 1 に記載の生成方法。

【 0 1 1 9 】

(付記 3)

前記第 1 及び第 2 のパラメータは、前記複数の項目値間の優先順序である、
ことを特徴とする、付記 1 又は 2 に記載の生成方法。

【 0 1 2 0 】

(付記 4)

前記第 1 及び第 2 のパラメータは、前記複数の項目値それぞれの重み付け値である、
ことを特徴とする、付記 1 又は 2 に記載の生成方法。

10

【 0 1 2 1 】

(付記 5)

複数の項目値を含む複数の個人情報を受け付ける受付処理部と、
前記複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて前記複数の項目
値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた前記複
数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された
第 1 の匿名情報を生成する第 1 匿名情報生成部と、

前記複数の項目値それぞれに対応付けられた前記第 1 のパラメータとは異なる第 2 のパ
ラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化ア
ルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくと
も何れかの項目値が匿名化された第 2 の匿名情報を生成する第 2 匿名情報生成部と、

20

前記第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する結合情報生成部
と、
を備えることを特徴とする情報処理装置。

【 0 1 2 2 】

(付記 6)

前記第 1 及び第 2 のパラメータについての多様性に関する指標が最大化されるように、
当該第 1 及び第 2 のパラメータを決定する、
ことを特徴とする、付記 5 に記載の情報処理装置。

【 0 1 2 3 】

(付記 7)

前記第 1 及び第 2 のパラメータは、前記複数の項目値間の優先順序である、
ことを特徴とする、付記 5 又は 6 に記載の情報処理装置。

30

【 0 1 2 4 】

(付記 8)

前記第 1 及び第 2 のパラメータは、前記複数の項目値それぞれの重み付け値である、
ことを特徴とする、付記 5 又は 6 に記載の情報処理装置。

【 0 1 2 5 】

(付記 9)

複数の項目値を含む複数の個人情報を受け付け、
前記複数の項目値それぞれに対応付けられた第 1 のパラメータに応じて前記複数の項目
値それぞれの秘匿化されやすさが決まる第 1 の秘匿化アルゴリズムを、受け付けた前記複
数の個人情報に適用して、前記複数の項目値の少なくとも何れかの項目値が匿名化された
第 1 の匿名情報を生成し、

40

前記複数の項目値それぞれに対応付けられた前記第 1 のパラメータとは異なる第 2 のパ
ラメータに応じて前記複数の項目値それぞれの秘匿化されやすさが決まる第 2 の秘匿化ア
ルゴリズムを、受け付けた前記複数の個人情報に適用して、前記複数の項目値の少なくと
も何れかの項目値が匿名化された第 2 の匿名情報を生成し、

前記第 1 の匿名情報及び前記第 2 の匿名情報を含む匿名情報を生成する、
処理をコンピュータに実行させることを特徴とする生成プログラム。

50

【符号の説明】

【 0 1 2 6 】

1	: 情報処理装置	
1 1	: C P U	
1 1 1	: 受付処理部	
1 1 2	: 匿名情報生成部	
1 1 3	: 結合情報生成部	
1 2	: メモリ部	
1 3	: 表示制御部	
1 4	: 記憶装置	10
1 5	: 入力 I F	
1 5 1	: マウス	
1 5 2	: キーボード	
1 6	: 外部記録媒体処理部	
1 6 0	: 記録媒体	
1 3 0	: 表示装置	
1 7	: 通信 I F	

20

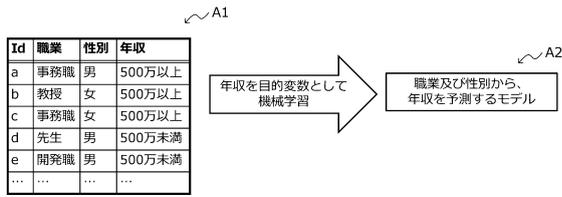
30

40

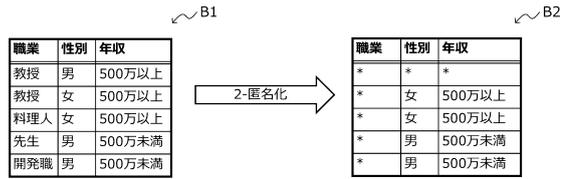
50

【図面】

【図 1】

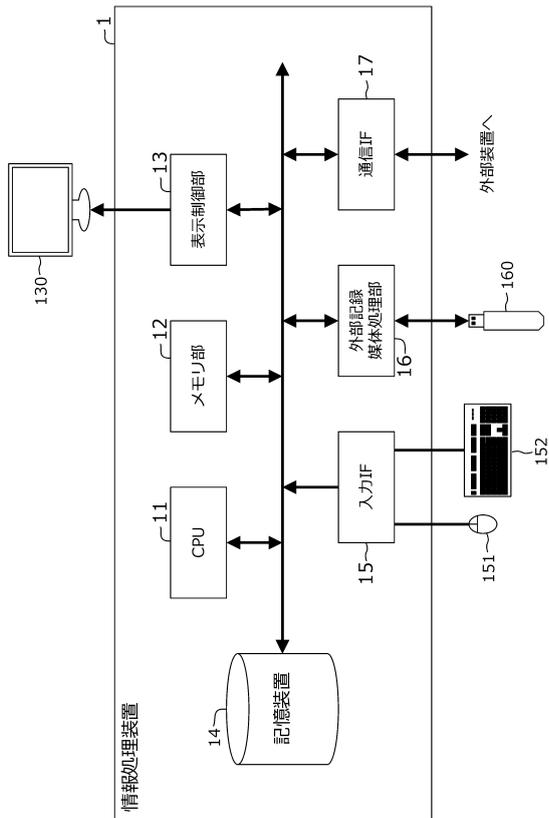


【図 2】

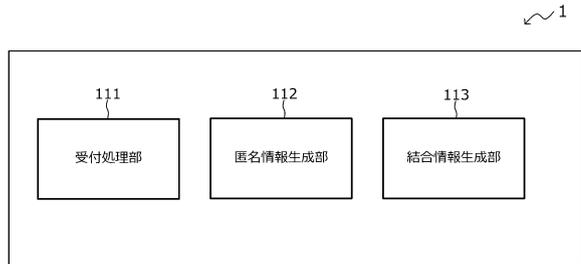


10

【図 3】



【図 4】



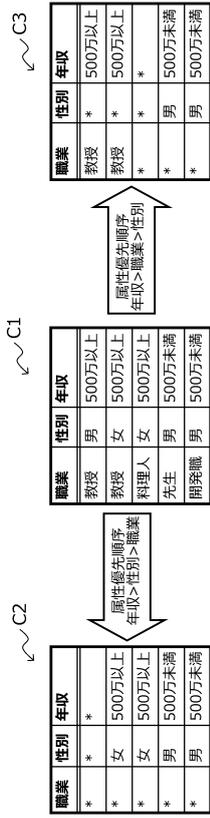
20

30

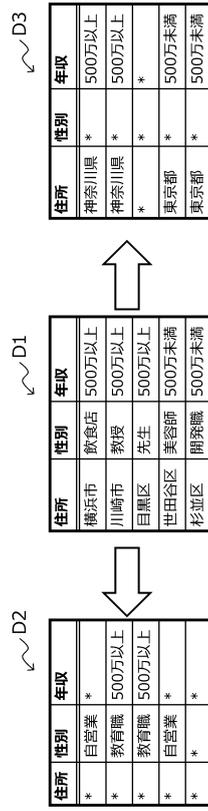
40

50

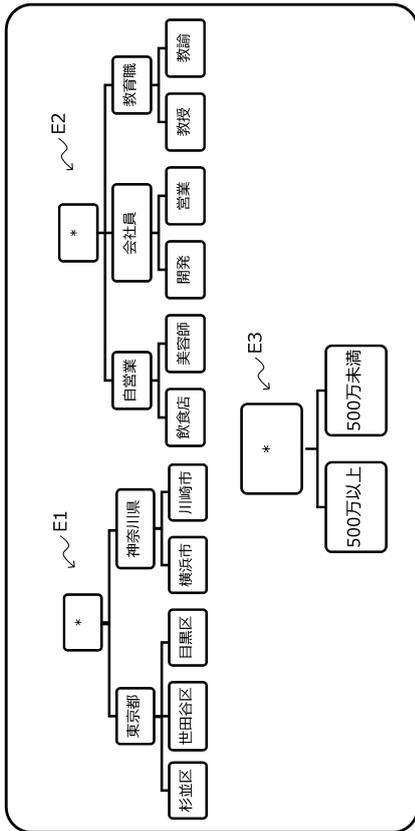
【 図 5 】



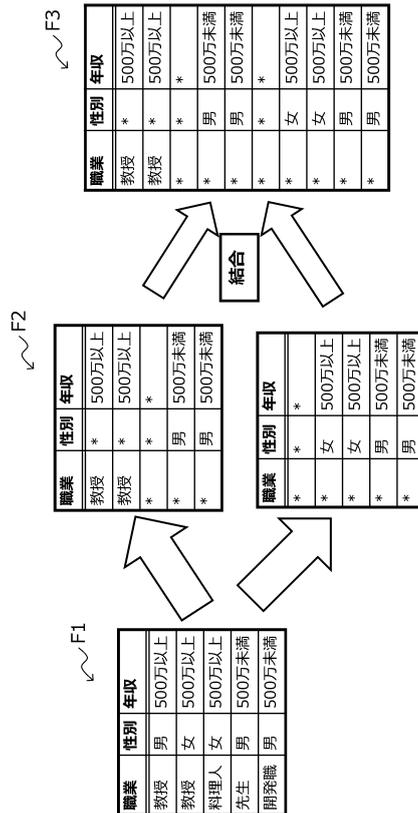
【 図 6 】



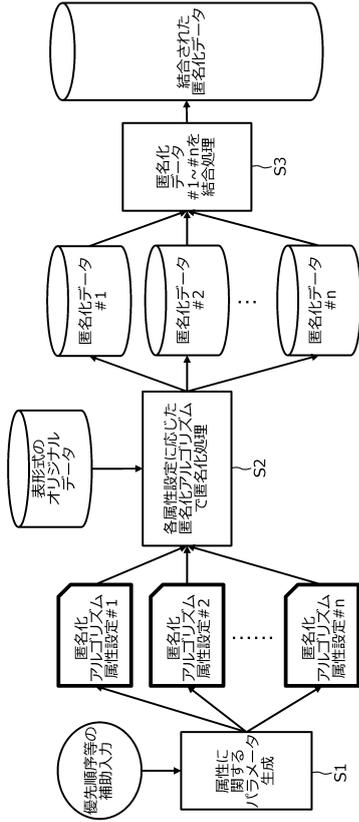
【 図 7 】



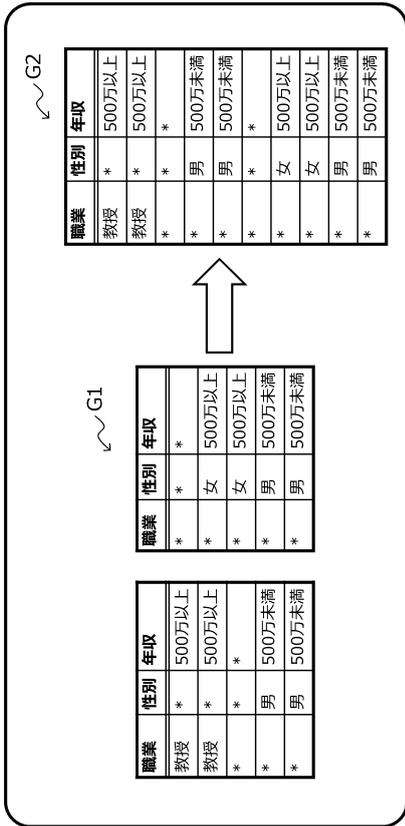
【 図 8 】



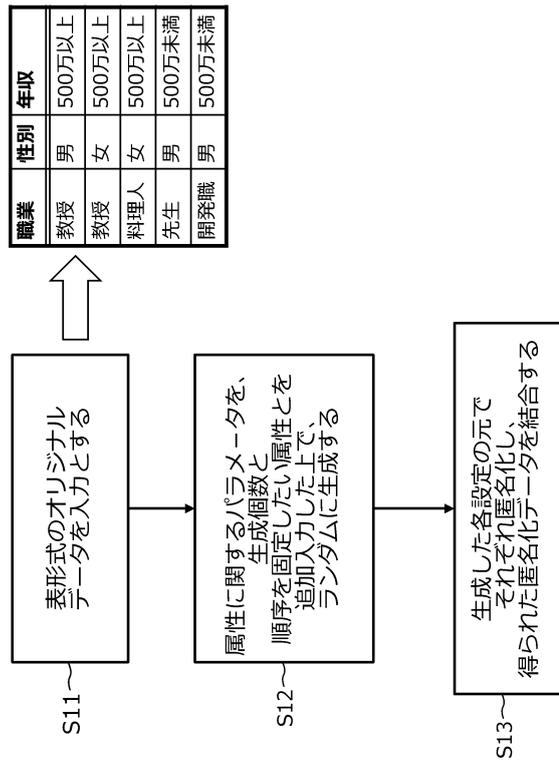
【 図 9 】



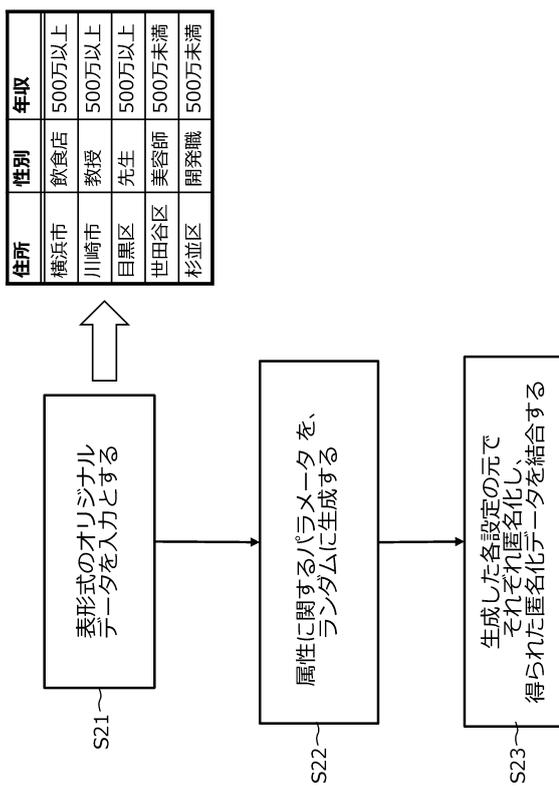
【 図 1 1 】



【 図 1 0 】



【 図 1 2 】



10

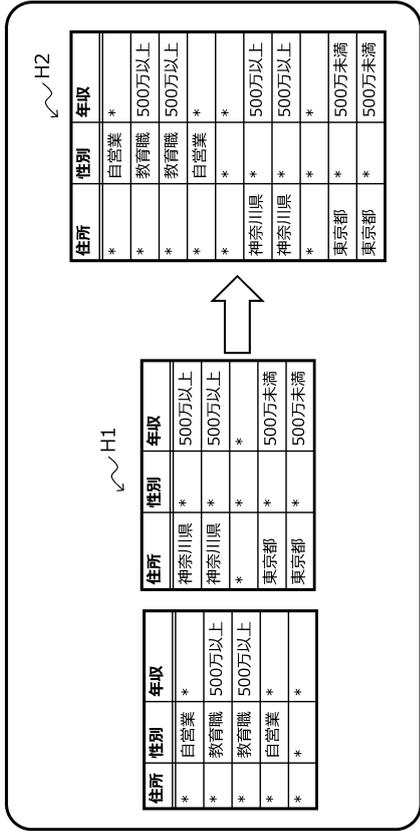
20

30

40

50

【 図 1 3 】



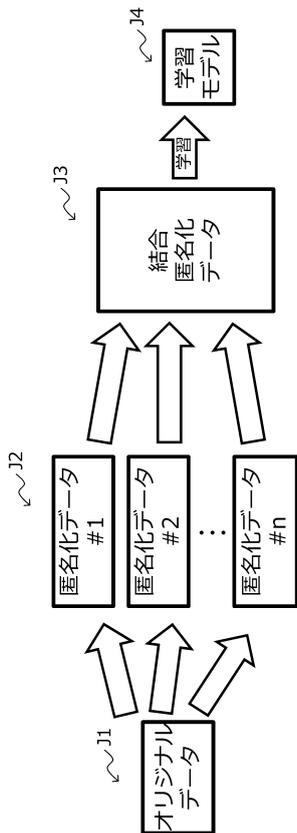
【 図 1 4 】



10

20

【 図 1 5 】



【 図 1 6 】

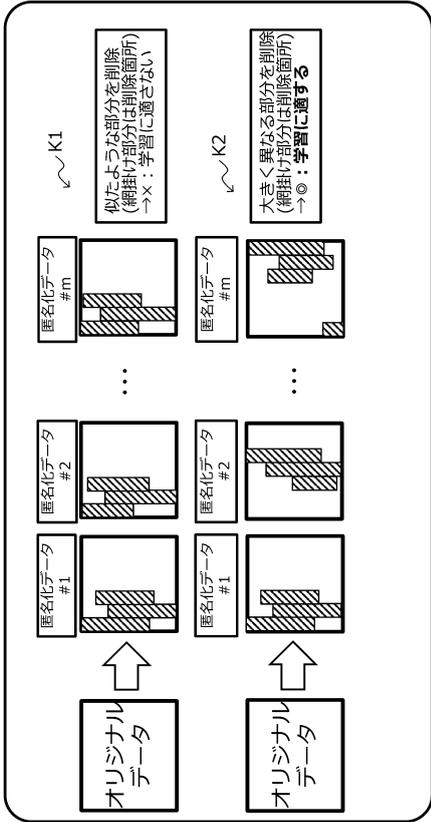
精度のkに関する要約統計量	平均	最小値
それぞれ学習した場合の最大精度	0.8270	0.8236
結合したデータで学習した場合の精度	0.8275	0.8248

30

40

50

【図 17】



【図 18】

順序間の距離	(123)	(132)	(213)	(231)	(312)	(321)
(123)	0	1	1	2	2	3
(132)	1	0	2	1	3	2
(213)	1	2	0	3	1	2
(231)	2	1	3	0	2	1
(312)	2	3	1	2	0	1
(321)	3	2	2	1	1	0

10

20

【図 19】

精度のkに関する要約統計量	平均	最小値	最大値
データオーギュメンテーション	0.8271	0.8256	0.8278
ランダム生成A	0.8274	0.8262	0.8288
ランダム生成B	0.8277	0.8262	0.8293
赤: 補助手段その2の実施例2	0.8281	0.8265	0.8293

30

40

50

フロントページの続き

- (56)参考文献 特表 2019 - 526851 (JP, A)
特開 2014 - 229039 (JP, A)
国際公開第 2012 / 067213 (WO, A1)
特開 2011 - 209800 (JP, A)
米国特許出願公開第 2018 / 0004978 (US, A1)
米国特許出願公開第 2016 / 0132697 (US, A1)
- (58)調査した分野 (Int.Cl., DB名)
G06F 21 / 62