

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6248774号  
(P6248774)

(45) 発行日 平成29年12月20日(2017.12.20)

(24) 登録日 平成29年12月1日(2017.12.1)

(51) Int.Cl. F I  
**G 0 6 F 17/30 (2006.01)**  
 G 0 6 F 17/30 3 5 0 C  
 G 0 6 F 17/30 2 2 0 Z

請求項の数 5 (全 16 頁)

(21) 出願番号	特願2014-85624 (P2014-85624)	(73) 特許権者	000005223
(22) 出願日	平成26年4月17日(2014.4.17)		富士通株式会社
(65) 公開番号	特開2015-207055 (P2015-207055A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成27年11月19日(2015.11.19)	(74) 代理人	100103528
審査請求日	平成29年1月10日(2017.1.10)		弁理士 原田 一男
		(72) 発明者	武部 浩明
			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
		(72) 発明者	上原 祐介
			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	齊藤 貴孝

最終頁に続く

(54) 【発明の名称】 最近傍探索のための情報処理装置及び方法

(57) 【特許請求の範囲】

【請求項 1】

複数のデータサンプルを格納するデータ格納部と、

前記データ格納部に格納された前記複数のデータサンプル間についての類似度行列に、  
 前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記  
 複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽  
 出する抽出部と、

を有する情報処理装置。

【請求項 2】

前記抽出部が、

前記複数のデータサンプル間についての類似度行列を算出する類似度行列算出部と、  
 算出された前記類似度行列に対して固有値分解を行って固有ベクトルを生成する固有値  
 分解処理部と、

生成された前記固有ベクトルの各成分の絶対値に基づき、前記複数のデータサンプルか  
 ら前記特定のデータサンプルを特定する特定部と、

を有する請求項 1 記載の情報処理装置。

【請求項 3】

前記特定のデータサンプルを用いて前記複数のデータサンプルの各々についてハッシュ  
 ベクトルを算出する算出部と、

算出された前記ハッシュベクトルと、前記データサンプルについてのカテゴリ名とを対

応付けてデータベースに登録する登録部と、

あるデータについて前記特定のデータサンプルを用いて算出されたハッシュベクトルと前記データベースに登録されたハッシュベクトルとの距離を算出して、最短距離のハッシュベクトルに対応付けられたカテゴリ名を特定する処理部と、

をさらに有する請求項 1 又は 2 記載の情報処理装置。

【請求項 4】

複数のデータサンプルを格納するデータ格納部に格納された前記複数のデータサンプル間についての類似度行列に、前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽出する

10

処理を含み、コンピュータが実行する情報処理方法。

【請求項 5】

複数のデータサンプルを格納するデータ格納部に格納された前記複数のデータサンプル間についての類似度行列に、前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽出する

処理を、コンピュータに実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

20

本発明は、最近傍探索技術に関する。

【背景技術】

【0002】

最近傍探索は、データが分布する  $n$  次元ベクトル空間において、与えられた点に最も近いデータを探す技術である。

【0003】

高次元データ空間における最近傍探索技術として、高次元データ空間を低次元データ空間に写像することにより、大規模データを小さく圧縮してから最近傍探索を行う方法が有用であることが近年示されている。

【0004】

30

このような方法において、元のデータ集合における距離関係を保存したまま低次元データ空間に写像するラプラシアン固有マップ法（例えば非特許文献 1）を用いることが好ましい。

【0005】

しかし、ラプラシアン固有マップ法を用いる場合、入力される未知のデータを低次元データ空間に写像するときに、未知のデータが低次元データ空間においてどこに写像されるかを知るために、未知のデータと既知のすべてのデータサンプルとの類似度計算を行うことになる。そのため、ラプラシアン固有マップ法を最近傍探索にそのまま適用すると、処理時間が膨大になってしまうという問題があった。

【0006】

40

これに対して、Anchor Graph Hashing（例えば非特許文献 2）という方法は、データ集合からアンカーと呼ばれる  $k$  個のデータサンプルを選んだ後に、各データサンプルを、それから最も近いアンカーからの距離（又は類似度）で表す。これにより、入力された未知のデータに対してもアンカーとの距離を計測することで低次元データ空間に写像することができる。すなわち、未知のデータを低次元データ空間に写像する場合、未知のデータとアンカーとの類似度計算を行えばよいので、処理時間の問題を解決できる。

【0007】

しかしながら、アンカーを  $k$ -means 法（ $K$ -平均法）によって選択しており、データの分布の形状によっては最近傍探索の精度が低下する問題があった。

【0008】

50

k - m e a n s 法は、データサンプル群をクラスタに分類する代表的な方法として知られている。具体的なアルゴリズムは以下のとおりである。

1. 各データサンプルに対してランダムにクラスタを割り振る。
2. 割り振ったデータサンプルを基に各クラスタの中心を計算する。
3. 各データサンプルについて各クラスタ中心との距離を求め、最も近いクラスタ中心のクラスタに割り当て直す。
4. 上記の処理で全てのデータサンプルについてクラスタの割り当てが変化しなかったら終了する。そうでない場合は2及び3の処理を繰り返す。

【0009】

この方法によってアンカーが選ばれる場合、データサンプルの分布の形状によっては最近傍探索の精度が低下する問題がある。

【0010】

例えば、図1に模式的に示すような特徴空間におけるデータサンプル群（×印）があったとき、k - m e a n s 法でアンカーを算出すると、アンカー a は分布の中心付近となる。そうすると、各データサンプルはアンカー a からの距離がほぼ等しくなるので、アンカーを経由した類似度は、どれもほぼ同じになる。従って、図2に模式的に示すように、各データサンプルを低次元特徴空間に写像すると、これらのデータサンプルはおおよそ一カ所に密集した状態となる。

【0011】

一方、図3に示すように、特徴空間において未知のデータ X が入力されて、未知のデータ X に対する最近傍点を求めるものとする。このとき、図4に模式的に示すように、未知のデータ X を低次元特徴空間に写像すると、密集したデータ集合の中に写像されるため、誤った点を最近傍点として抽出する可能性が高くなる。図5に示すように、本来特徴空間では、未知のデータ X はデータサンプル Z に最も近いが、低次元特徴空間では、アンカー a を経由した距離が同じであるデータサンプル Y が最も近いデータサンプルとして抽出されてしまう。

【先行技術文献】

【非特許文献】

【0012】

【非特許文献1】“Laplacian Eigenmaps and Spectral Techniques for Embedding and Clustering” M. Belkin and P. Niyogi. NIPS 14, 2002

【非特許文献2】“Hashing with Graphs” W.Liu, J.Wang, S.Kumar and S.F.Chang. I CML ' 11, 2011

【発明の概要】

【発明が解決しようとする課題】

【0013】

従って、本発明の目的は、一側面によれば、最近傍探索における未知のデータに対する識別精度を向上させるための技術を提供することである。

【課題を解決するための手段】

【0014】

本発明に係る情報処理装置は、(A)複数のデータサンプルを格納するデータ格納部と、(B)データ格納部に格納された複数のデータサンプル間についての類似度行列に、複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく複数のデータサンプル間の類似度行列が近似するように、特定のデータサンプルを抽出する抽出部とを有する。

【発明の効果】

【0015】

一側面としては、最近傍探索における未知のデータに対する識別精度を向上させることができるようになる。

【図面の簡単な説明】

10

20

30

40

50

## 【 0 0 1 6 】

【図 1】図 1 は、従来技術の問題を説明するための図である。

【図 2】図 2 は、従来技術の問題を説明するための図である。

【図 3】図 3 は、従来技術の問題を説明するための図である。

【図 4】図 4 は、従来技術の問題を説明するための図である。

【図 5】図 5 は、従来技術の問題を説明するための図である。

【図 6】図 6 は、本実施の形態の概要を説明するための図である。

【図 7】図 7 は、本実施の形態の概要を説明するための図である。

【図 8】図 8 は、本実施の形態の概要を説明するための図である。

【図 9】図 9 は、本実施の形態の概要を説明するための図である。

10

【図 10】図 10 は、本実施の形態の概要を説明するための図である。

【図 11】図 11 は、本実施の形態の概要を説明するための図である。

【図 12】図 12 は、本実施の形態の概要を説明するための図である。

【図 13】図 13 は、本実施の形態の概要を説明するための図である。

【図 14】図 14 は、本実施の形態に係る情報処理装置の構成例を示す図である。

【図 15】図 15 は、学習処理の処理フローを示す図である。

【図 16】図 16 は、アンカー抽出処理の処理フローを示す図である。

【図 17】図 17 は、初期類似度行列を説明するための図である。

【図 18】図 18 は、識別処理の処理フローを示す図である。

【図 19】図 19 は、コンピュータの機能ブロック図である。

20

【発明を実施するための形態】

## 【 0 0 1 7 】

まず、図 6 乃至図 13 を用いて、本実施の形態の概要を説明しておく。

## 【 0 0 1 8 】

本実施の形態では、A G H 法におけるアンカーを *k - m e a n s* 法で抽出するのではなく、特徴空間のデータ分布を近似するアンカーを算出する。

## 【 0 0 1 9 】

まず、アンカーを設定することによってデータ分布がどのように変化するのかを図 6 を用いて模式的に説明する。

## 【 0 0 2 0 】

30

A G H 法では、アンカーを設定した後、各データサンプル間の距離又は類似度を、アンカーを経由して計測する。アンカー *a* を経由して計測したときの 2 点 A 及び B 間の距離を  $|AB|_{\text{anchor}}$  と表すと、 $|AB|_{\text{anchor}}$  及び  $|AC|_{\text{anchor}}$  は、元の特徴空間における 2 点間距離  $|Aa|$ 、 $|Ba|$  及び  $|Ca|$  によって以下のように表される。

## 【 0 0 2 1 】

【数 1】

$$\begin{aligned} |AB|_{\text{anchor}} &\approx ||Aa| - |Ba|| \\ |AC|_{\text{anchor}} &\approx ||Aa| - |Ca|| \end{aligned}$$

40

## 【 0 0 2 2 】

特徴空間のデータ分布を近似するアンカーを算出するということは、データサンプル間の距離をできるだけ保存するアンカーを求めるということである。具体的には、アンカーを経由した場合の距離の比  $||Aa| - |Ba|| : ||Aa| - |Ca||$  を、元の特徴空間における距離の比  $|AB| : |AC|$  と可能な限り同じにするようなアンカーを抽出する。

## 【 0 0 2 3 】

すなわち、データサンプル間の距離を可能な限り保存するアンカーを求めるという問題は、類似度行列によって定式化できる。すなわち、図 7 に模式的に示すように、アンカー

50

を設定したときのデータ分布を表す類似度行列  $Z^{\wedge}_{\text{anchor}}$  ( " ^ " は、以下、 $Z$  の上に置かれたものとする ) が、図 8 に模式的に示すような元の特徴空間におけるデータ分布を表す類似度行列  $Z^{\wedge}$  に可能な限り一致するようにアンカー  $a$  を設定する。

【 0 0 2 4 】

類似度行列とは、データサンプルが  $m$  個与えられた場合は  $m \times m$  の行列であり、行列の  $(i, j)$  成分はデータサンプル  $i$  の特徴ベクトル  $s_i$  とデータサンプル  $j$  の特徴ベクトル  $s_j$  の類似度を表すものである。

【 0 0 2 5 】

アンカーを設定したときのデータ分布を表す類似度行列  $Z^{\wedge}_{\text{anchor}}$  は、アンカーを 1 個として、それらの特徴ベクトルを  $a_1, a_2, \dots, a_l$  と表すと、 $a_1, a_2, \dots, a_l$  の関数となるが、類似度行列の定義からして非線形の極めて複雑な関数となる。具体的には、A G H 法における類似度の算出方法は足切りなどの演算が存在するので、 $a_1, a_2, \dots, a_l$  を陽に変数とした関数として表すことができず、以下のような方程式を立てて解くことは困難である。

【 0 0 2 6 】

【 数 2 】

$$\hat{Z} = \hat{Z}_{\text{anchor}}(a_1, a_2, \dots, a_l)$$

【 0 0 2 7 】

そこで、類似度行列  $Z^{\wedge}$  を固有値分解して、 $Z^{\wedge}$  を近似するアンカーを算出する。 $Z^{\wedge}$  を固有値分解すると、 $Z^{\wedge}$  は射影行列の和に近似でき、「どのデータサンプルをアンカーとして用いれば類似度行列に近い行列を得られるか」という問題が、「どのデータサンプルをアンカーとして用いれば射影行列に近い行列を得られるか」という問題に分解できる。そして、「どのデータサンプルをアンカーとして用いれば射影行列に近い行列を得られるか」という問題については、近似的な解を得ることができる。

【 0 0 2 8 】

すなわち、データサンプル  $s_j$  の射影行列  $v_j v_j^T$  に対する寄与度が固有ベクトル  $v_i$  の成分の絶対値で測れるので、固有ベクトルの成分の絶対値が大きいデータサンプルをアンカーとして選択する。

【 0 0 2 9 】

具体的には、本実施の形態におけるアンカーには、図 9 に模式的に示すような特徴空間においては、データサンプル  $b$  のようなデータサンプルが選択される。そうすると、低次元特徴空間においては、図 10 に模式的に示すように、1 箇所集中することなく、分散するようになる。このような場合、図 11 に模式的に示すように特徴空間において未知のデータ  $X$  が入力された場合、図 12 に模式的に示すように低次元特徴空間に写像されるわけであるが、分布に応じた低次元特徴空間となっているので、アンカー  $b$  からの距離に近い他のデータサンプルが多数となることはない。従って、図 13 に模式的に示すように、未知のデータ  $X$  に近いデータサンプル  $Z$  が正しく見つかるようになる。

【 0 0 3 0 】

次に、このような前提の下処理を行う、本実施の形態における情報処理装置の構成例を図 14 に示す。

【 0 0 3 1 】

本実施の形態における情報処理装置 100 は、第 1 データ格納部 101 と、アンカー抽出部 102 と、アンカー格納部 103 と、ハッシュベクトル算出部 104 と、登録部 105 と、データベース 106 と、データ入力部 107 と、第 2 データ格納部 108 と、照合処理部 109 と、データ出力部 110 とを有する。

【 0 0 3 2 】

第 1 データ格納部 101 は、例えば  $m$  個のデータサンプルの特徴ベクトル及びカテゴリ名が格納されている。

## 【 0 0 3 3 】

アンカー抽出部 1 0 2 は、類似度行列算出部 1 0 2 1 と、固有値分解処理部 1 0 2 2 と、アンカー特定部 1 0 2 3 とを有する。

## 【 0 0 3 4 】

類似度行列算出部 1 0 2 1 は、第 1 データ格納部 1 0 1 に格納されているデータサンプルの特徴ベクトルから類似度行列を算出する。

## 【 0 0 3 5 】

固有値分解処理部 1 0 2 2 は、類似度行列に対して固有値分解処理を実行し、固有ベクトルを生成する。

## 【 0 0 3 6 】

アンカー特定部 1 0 2 3 は、固有ベクトルを用いて、第 1 データ格納部 1 0 1 に格納されているデータサンプルのうちアンカーとして採用するデータサンプルを選定し、アンカー格納部 1 0 3 に格納する。

## 【 0 0 3 7 】

ハッシュベクトル算出部 1 0 4 は、第 1 データ格納部 1 0 1 に格納されているデータサンプルの特徴ベクトルから、アンカー格納部 1 0 3 に格納されているアンカーに基づきハッシュベクトルを算出し、登録部 1 0 5 に出力する。ハッシュベクトル算出部 1 0 4 によるハッシュベクトル算出処理については、従来と同じなので詳細には述べない。

## 【 0 0 3 8 】

登録部 1 0 5 は、ハッシュベクトルと、当該ハッシュベクトルの元となったデータサンプルのカテゴリ名とを、データベース 1 0 6 に登録する。

## 【 0 0 3 9 】

データ入力部 1 0 7 は、ユーザ又は他のコンピュータから、未知のデータの特徴ベクトルの入力を受け付け、第 2 データ格納部 1 0 8 に格納する。

## 【 0 0 4 0 】

照合処理部 1 0 9 は、ハッシュベクトル算出部 1 0 4 に、第 2 データ格納部 1 0 8 に格納されているデータの特徴ベクトルについてのハッシュベクトルを算出させる。そして、照合処理部 1 0 9 は、当該ハッシュベクトルと、データベース 1 0 6 に格納されている各ハッシュベクトルとのハミング距離を算出して、最短のハミング距離が算出されたハッシュベクトルに対応付けて格納されているカテゴリ名を特定し、データ出力部 1 1 0 に出力する。

## 【 0 0 4 1 】

データ出力部 1 1 0 は、出力装置（表示装置、印刷装置や他のコンピュータなど）に、特定されたカテゴリ名を出力する。

## 【 0 0 4 2 】

次に、図 1 5 乃至図 1 8 を用いて、情報処理装置 1 0 0 の処理内容について説明する。

## 【 0 0 4 3 】

本実施の形態では、A G H 法と同様に、予め実施しておく学習処理と、学習処理の結果に基づき未知のデータに対してカテゴリ名を識別する識別処理とが行われる。

## 【 0 0 4 4 】

まず、学習処理について、図 1 5 乃至図 1 7 を用いて説明する。

## 【 0 0 4 5 】

アンカー抽出部 1 0 2 は、本実施の形態における主要部であるアンカー抽出処理を実行する（図 1 5：ステップ S 1）。アンカー抽出処理については、図 1 6 を用いて説明する。

## 【 0 0 4 6 】

まず、類似度行列算出部 1 0 2 1 は、第 1 データ格納部 1 0 1 に格納されているデータサンプルの特徴ベクトルから、類似度行列を算出する（図 1 6：ステップ S 1 1）。

## 【 0 0 4 7 】

（ A ）初期類似度行列 Z の算出

10

20

30

40

50

類似度行列算出部 1021 は、各データサンプルの特徴ベクトル  $s_i$  に対して、自データサンプル以外の他のデータサンプルの特徴ベクトル  $s_j$  とのユークリッド距離に基づく類似度  $h(s_i, s_j)$  を算出する。具体的には、以下のような式に従って類似度を算出する。

【0048】

【数3】

$$h(s_i, s_j) = \exp\left(-\frac{D^2(s_i, s_j)}{T}\right)$$

10

【0049】

なお、 $D(A, B)$  は、 $A$  と  $B$  とのユークリッド距離を表しており、 $D^2(A, B)$  は、ユークリッド距離の二乗を表す。また、 $T$  は、比例定数で、予め設定された値である。

【0050】

このような計算をデータサンプル  $i$  及び  $j$  のペアの各々について算出すれば、図 17 に示すように、初期類似度行列  $Z$  の  $i, j$  成分の値が算出される。

【0051】

(B) 足切り処理

類似度  $h(s_i, s_j)$  については、以下のような足切りを行う。

20

・  $\{h(s_i, s_j) \mid i: \text{fixed}, 1 \leq j \leq m, j \neq i\}$  において、上位  $r$  個以外は  $h(s_i, s_j) = 0$

【0052】

すなわち、初期類似度行列  $Z$  の行毎に、類似度が大きい順に  $r + 1$  番目以降は類似度を 0 に設定する。図 17 の例では、点線矩形で囲まれた  $i$  行目において、類似度を大きい順にソートして、 $r + 1$  番目以降の類似度を 0 に設定する。

【0053】

(C) 正規化処理

足切り処理後の初期類似度行列  $Z$  の行毎に、成分の総和が 1 となるように正規化する。図 17 の例では、点線矩形で囲まれた  $i$  行目において、類似度の総和が 1 となるように正規化する。

30

【0054】

(D) 対称行列生成

正規化処理後の初期類似度行列  $Z$  から、類似度行列  $Z^{\wedge}$  を算出する。これによって、 $m \times m$  の対称行列となった類似度行列  $Z^{\wedge}$  が得られる。

【0055】

【数4】

$$\hat{Z} = ZZ^T$$

【0056】

40

次に、固有値分解処理部 1022 は、類似度行列  $Z^{\wedge}$  に対して固有値分解処理を実行する (ステップ S13)。

【0057】

本実施の形態では、固有値のうち大きい順に上位  $k$  個の固有値  $\sigma_i$  を採用して、類似度行列  $Z^{\wedge}$  を以下のように近似する。

【0058】

【数5】

$$\hat{Z} \approx \mathbf{U}_k \Sigma_k \mathbf{U}_k^T = \sigma_1 \mathbf{v}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{v}_2 \mathbf{v}_2^T \cdots + \sigma_k \mathbf{v}_k \mathbf{v}_k^T$$

50

【 0 0 5 9 】

【 数 6 】

$$\Sigma_k = \text{diag}(\sigma_t) = \begin{pmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_k \end{pmatrix}$$

【 0 0 6 0 】

10

【 数 7 】

$$\mathbf{U}_k = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k)$$

なお、 $\mathbf{v}_i$ は、固有値  $\lambda_i$ に対応する固有ベクトルである。また、 $\mathbf{v}_i \mathbf{v}_i^T$ は、射影行列と呼ばれる。

【 0 0 6 1 】

その後、アンカー特定部 1 0 2 3 は、各固有ベクトル  $\mathbf{v}_t$  から、アンカーとして用いるべきデータサンプルを特定し、当該データサンプルの特徴ベクトルをアンカー格納部 1 0 3 に格納する（ステップ S 1 5）。そして呼び出し元の処理に戻る。

20

【 0 0 6 2 】

具体的には、固有ベクトル  $\mathbf{v}_t$  ( $1 \leq t \leq k$ ) の各々について、当該固有ベクトル  $\mathbf{v}_t$  において絶対値が最も大きい成分を特定する。そして、 $k > 1$  であれば、 $k$  個の固有ベクトルのうち、特定された成分の絶対値が大きい順に 1 個の固有ベクトルを選択して、選択された固有ベクトルについて特定された成分の順番  $i_u$  のデータサンプル  $s_{i_u}$  を、アンカー  $a_{i_u}$  として採用する。

【 0 0 6 3 】

すなわち、固有ベクトル  $\mathbf{v}_t$  を以下のように表す。

【 0 0 6 4 】

【 数 8 】

30

$$\mathbf{v}_t = (v_1, v_2, \dots, v_m)^T$$

【 0 0 6 5 】

そして、固有ベクトル  $\mathbf{v}_t$  において絶対値が最も大きい成分を特定するということは、以下のように表される。

【 0 0 6 6 】

【 数 9 】

$$i_u = \arg \max_{i \in \{1, 2, \dots, m\} \setminus \{i_1, \dots, i_{u-1}\}} (|v_i|)$$

40

なお、既来选择された  $i_1$  乃至  $i_{u-1}$  と同じ  $i_u$  については選択できないので、次に絶対値が大きい成分の番号を特定する。

【 0 0 6 7 】

このような処理を実行すれば、 $m$  個のデータサンプル間についての類似度行列  $\hat{Z}$  に、 $m$  個のデータサンプルから選択されたアンカーとの類似度に基づく  $m$  個のデータサンプル間の類似度行列  $\hat{Z}_{\text{anchor}}$  が近似するように、アンカーが選択される。

【 0 0 6 8 】

すなわち、A G H 法における識別処理の高速性を保持しつつ、識別精度を向上させるこ

50



とができるようになる。

【0069】

すなわち、本実施の形態では、高速且つ高精度な最近傍探索を行うことができる。特に、高次元空間内においてデータ分布がいかなる状態であっても高精度な最近傍探索を行うことができる。

【0070】

図15の処理の説明に戻って、ハッシュベクトル算出部104は、第1データ格納部101に格納されている各データサンプルの特徴ベクトルから、アンカー格納部103に格納されているアンカーのデータに基づき、ハッシュベクトルを算出し、登録部105に出力する(ステップS3)。ハッシュベクトルは、 $n$ 次元のベクトルであって、個々の成分の値は0又は1である。この処理は、従来と同じであるから、詳細な説明については省略する。

10

【0071】

そして、登録部105は、ハッシュベクトル算出部104によって算出されたハッシュベクトルと、当該ハッシュベクトルの元となったデータサンプルのカテゴリ名とを、対応付けてデータベース106に登録する(ステップS5)。

【0072】

このようにすれば、学習処理が完了して、識別処理を行う準備ができたことになる。

【0073】

次に、図18を用いて、識別処理の処理フローを説明する。

20

【0074】

まず、データ入力部107は、ユーザ又は他のコンピュータから、未知のデータの入力を受け付け、第2データ格納部108に格納する(図18:ステップS21)

【0075】

そうすると、照合処理部109は、第2データ格納部108に格納されている未知のデータの特徴ベクトルに対するハッシュベクトルを、ハッシュベクトル算出部104に算出させる(ステップS23)。この処理も、従来と同じであり、演算内容としてはステップS3と同様であり、詳細な説明は省略する。

【0076】

そして、照合処理部109は、未知のデータに対するハッシュベクトルと、データベース106に格納されている各ハッシュベクトルとのハミング距離を算出し、最もハミング距離が短いハッシュベクトルに対応付けられているカテゴリ名を特定する照合処理を実行し、処理結果をデータ出力部110に出力する(ステップS25)。

30

【0077】

データ出力部110は、照合処理部109から処理結果であるカテゴリ名を受け取ると、ユーザ又は他のコンピュータに対してカテゴリ名を出力する(ステップS27)。

【0078】

以上のような処理を行うことで、アンカーが適切に設定されているので、高精度なデータ識別を行うことができるようになる。

【0079】

以上本発明の実施の形態を説明したが、本発明はこれに限定されるものではない。

40

【0080】

例えば、図14の情報処理装置100の機能ブロック構成は、プログラムモジュール構成とは一致しない場合がある。また、処理フローについても、処理結果が変わらない限り、順番を入れ替えたり、並列に実行する場合もある。

【0081】

なお、本実施の形態におけるデータは、画像データであったり、文字データであったり、個人情報であったりする。

【0082】

さらに、情報処理装置100は1台のコンピュータではなく、複数台のコンピュータで

50

機能分担する場合もある。

【 0 0 8 3 】

なお、上で述べた情報処理装置 1 0 0 は、コンピュータ装置であって、図 1 9 に示すように、メモリ 2 5 0 1 と C P U (Central Processing Unit) 2 5 0 3 とハードディスク・ドライブ (H D D : Hard Disk Drive) 2 5 0 5 と表示装置 2 5 0 9 に接続される表示制御部 2 5 0 7 とリムーバブル・ディスク 2 5 1 1 用のドライブ装置 2 5 1 3 と入力装置 2 5 1 5 とネットワークに接続するための通信制御部 2 5 1 7 とがバス 2 5 1 9 で接続されている。オペレーティング・システム (O S : Operating System) 及び本実施例における処理を実施するためのアプリケーション・プログラムは、H D D 2 5 0 5 に格納されており、C P U 2 5 0 3 により実行される際には H D D 2 5 0 5 からメモリ 2 5 0 1 に読み出される。C P U 2 5 0 3 は、アプリケーション・プログラムの処理内容に応じて表示制御部 2 5 0 7 、通信制御部 2 5 1 7 、ドライブ装置 2 5 1 3 を制御して、所定の動作を行わせる。また、処理途中のデータについては、主としてメモリ 2 5 0 1 に格納されるが、H D D 2 5 0 5 に格納されるようにしてもよい。本技術の実施例では、上で述べた処理を実施するためのアプリケーション・プログラムはコンピュータ読み取り可能なリムーバブル・ディスク 2 5 1 1 に格納されて頒布され、ドライブ装置 2 5 1 3 から H D D 2 5 0 5 にインストールされる。インターネットなどのネットワーク及び通信制御部 2 5 1 7 を経由して、H D D 2 5 0 5 にインストールされる場合もある。このようなコンピュータ装置は、上で述べた C P U 2 5 0 3 、メモリ 2 5 0 1 などのハードウェアと O S 及びアプリケーション・プログラムなどのプログラムとが有機的に協働することにより、上で述べたような各種機能を実現する。

【 0 0 8 4 】

以上述べた本実施の形態をまとめると、以下のようになる。

【 0 0 8 5 】

本実施の形態に係る情報処理装置は、( A ) 複数のデータサンプルを格納するデータ格納部と、( B ) データ格納部に格納された複数のデータサンプル間についての類似度行列に、複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく複数のデータサンプル間の類似度行列が近似するように、特定のデータサンプルを抽出する抽出部とを有する。

【 0 0 8 6 】

このような特定のデータサンプルを抽出して A G H 法におけるアンカーとして用いれば、識別精度を向上させることができるようになる。

【 0 0 8 7 】

なお、上で述べた抽出部が、( b 1 ) 複数のデータサンプル間についての類似度行列を算出する類似度行列算出部と、( b 2 ) 算出された類似度行列に対して固有値分解を行って固有ベクトルを生成する固有値分解処理部と、( b 3 ) 生成された固有ベクトルの各成分の絶対値に基づき、複数のデータサンプルから上記特定のデータサンプルを特定する特定部とを有する。

【 0 0 8 8 】

例えば、固有ベクトルの各成分の絶対値のうち最も大きい絶対値が得られた成分の順番を用いて、複数のデータサンプルから上記順番における特定のデータサンプルを特定するものである。

【 0 0 8 9 】

さらに、上で述べた情報処理装置は、( C ) 特定のデータサンプルを用いて複数のデータサンプルの各々についてハッシュベクトルを算出する算出部と、( D ) 算出されたハッシュベクトルと、データサンプルについてのカテゴリ名とを対応付けてデータベースに登録する登録部と、( E ) あるデータについて特定のデータサンプルを用いて算出されたハッシュベクトルとデータベースに登録されたハッシュベクトルとの距離を算出して、最短距離のハッシュベクトルに対応付けられたカテゴリ名を特定する処理部とをさらに有するようにしても良い。このようにすれば、A G H 法で最近傍探索が行われる。

## 【 0 0 9 0 】

なお、上で述べたような処理をプロセッサ又はコンピュータに実行させるためのプログラムを作成することができ、当該プログラムは、例えばフレキシブル・ディスク、CD-ROMなどの光ディスク、光磁気ディスク、半導体メモリ（例えばROM）、ハードディスク等のコンピュータ読み取り可能な記憶媒体又は記憶装置に格納される。なお、処理途中のデータについては、RAM等の記憶装置に一時保管される。

## 【 0 0 9 1 】

以上の実施例を含む実施形態に関し、さらに以下の付記を開示する。

## 【 0 0 9 2 】

（付記 1）

複数のデータサンプルを格納するデータ格納部と、

前記データ格納部に格納された前記複数のデータサンプル間についての類似度行列に、前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽出する抽出部と、

を有する情報処理装置。

## 【 0 0 9 3 】

（付記 2）

前記抽出部が、

前記複数のデータサンプル間についての類似度行列を算出する類似度行列算出部と、

算出された前記類似度行列に対して固有値分解を行って固有ベクトルを生成する固有値分解処理部と、

生成された前記固有ベクトルの各成分の絶対値に基づき、前記複数のデータサンプルから前記特定のデータサンプルを特定する特定部と、

を有する付記 1 記載の情報処理装置。

## 【 0 0 9 4 】

（付記 3）

前記特定のデータサンプルを用いて前記複数のデータサンプルの各々についてハッシュベクトルを算出する算出部と、

算出された前記ハッシュベクトルと、前記データサンプルについてのカテゴリ名とを対応付けてデータベースに登録する登録部と、

あるデータについて前記特定のデータサンプルを用いて算出されたハッシュベクトルと前記データベースに登録されたハッシュベクトルとの距離を算出して、最短距離のハッシュベクトルに対応付けられたカテゴリ名を特定する処理部と、

をさらに有する付記 1 又は 2 記載の情報処理装置。

## 【 0 0 9 5 】

（付記 4）

複数のデータサンプルを格納するデータ格納部に格納された前記複数のデータサンプル間についての類似度行列に、前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽出する

処理を含み、コンピュータが実行する情報処理方法。

## 【 0 0 9 6 】

（付記 5）

前記抽出する処理が、

前記複数のデータサンプル間についての類似度行列を算出し、

算出された前記類似度行列に対して固有値分解を行って固有ベクトルを生成し、

生成された前記固有ベクトルの各成分の絶対値に基づき、前記複数のデータサンプルから前記特定のデータサンプルを特定する

処理を含む付記 4 記載の情報処理方法。

10

20

30

40

50

## 【 0 0 9 7 】

( 付 記 6 )

複数のデータサンプルを格納するデータ格納部に格納された前記複数のデータサンプル間についての類似度行列に、前記複数のデータサンプルから選択された特定のデータサンプルとの類似度に基づく前記複数のデータサンプル間の類似度行列が近似するように、前記特定のデータサンプルを抽出する

処理を、コンピュータに実行させるためのプログラム。

## 【 0 0 9 8 】

( 付 記 7 )

前記抽出する処理が、

10

前記複数のデータサンプル間についての類似度行列を算出し、

算出された前記類似度行列に対して固有値分解を行って固有ベクトルを生成し、

生成された前記固有ベクトルの各成分の絶対値に基づき、前記複数のデータサンプルから前記特定のデータサンプルを特定する

処理を含む付記 6 記載のプログラム。

## 【 符号の説明 】

## 【 0 0 9 9 】

1 0 1 第 1 データ格納部

1 0 2 アンカー抽出部

1 0 3 アンカー格納部

20

1 0 4 ハッシュベクトル算出部

1 0 5 登録部

1 0 6 データベース

1 0 7 データ入力部

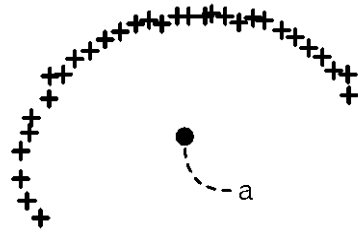
1 0 8 第 2 データ格納部

1 0 9 照合処理部

1 1 0 データ出力部

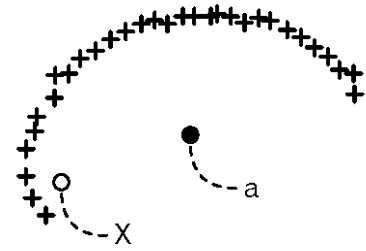
【図 1】

特徴空間



【図 3】

特徴空間



【図 2】

低次元特徴空間



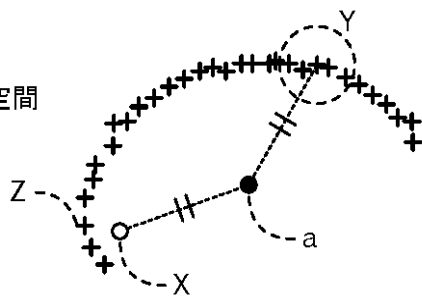
【図 4】

低次元特徴空間



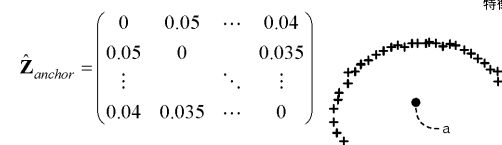
【図 5】

特徴空間



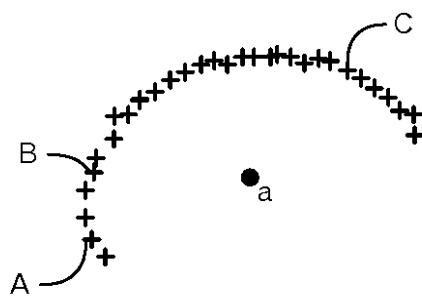
【図 7】

特徴空間



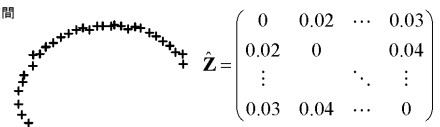
【図 6】

特徴空間



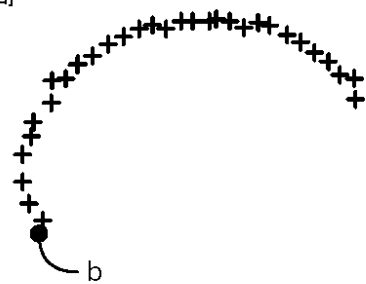
【図 8】

特徴空間



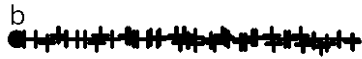
【図 9】

特徴空間



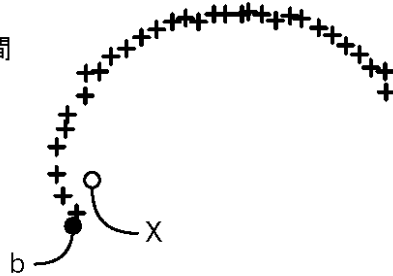
【図 10】

低次元特徴空間



【図 11】

特徴空間

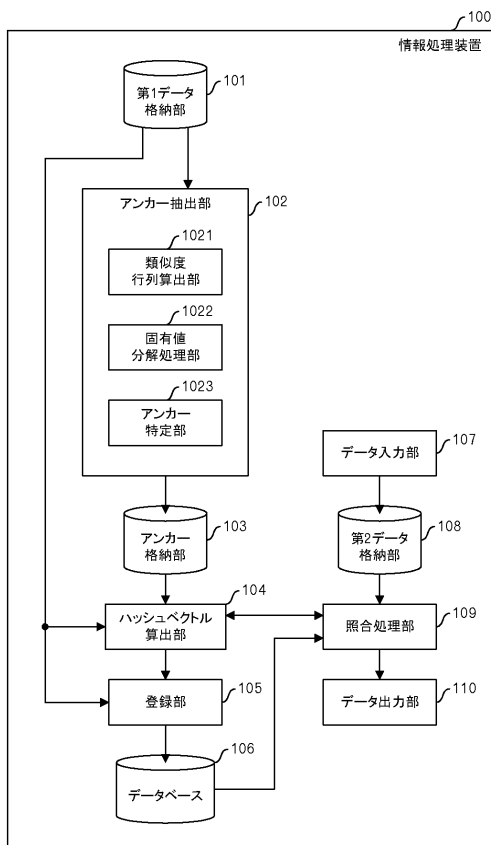


【図 12】

低次元特徴空間

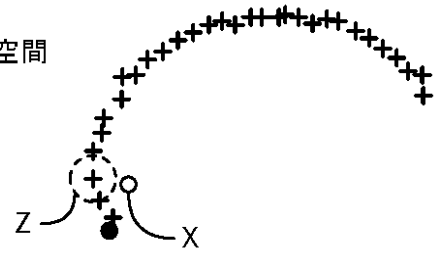


【図 14】

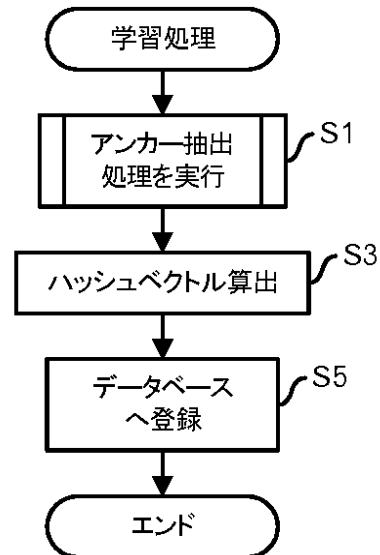


【図 13】

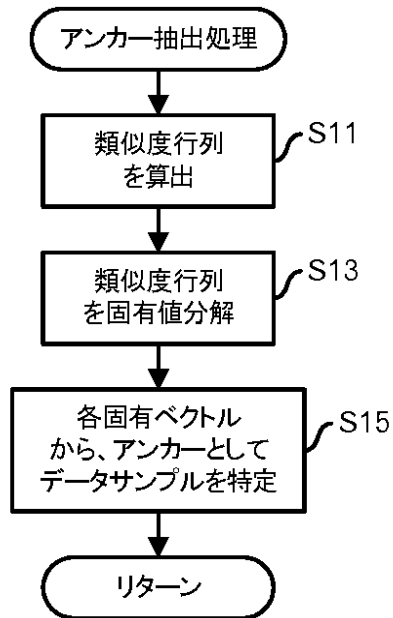
特徴空間



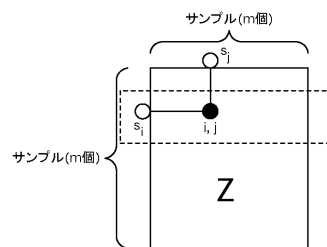
【図 15】



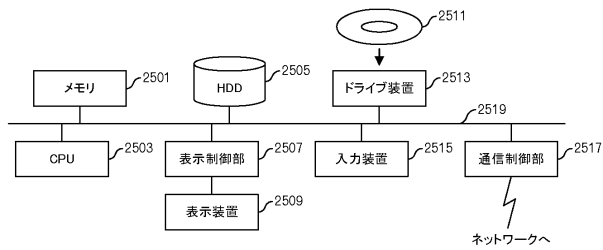
【図 16】



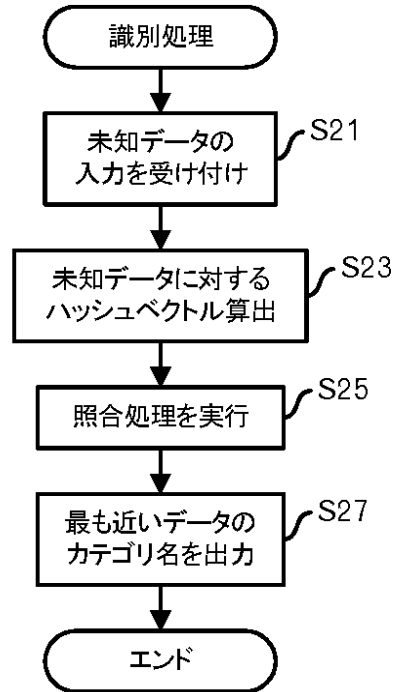
【図 17】



【図 19】



【図 18】



---

フロントページの続き

- (56)参考文献 特開2012-038272(JP,A)  
国際公開第2012/050952(WO,A1)  
米国特許出願公開第2013/0086553(US,A1)  
特開2010-250377(JP,A)  
特開2012-088972(JP,A)  
特開2013-020290(JP,A)  
田中 竜仁、外4名、スケッチを入力とした3次元モデル検索の高速化, Media Computing Conference 2012 2012年度 画像電子学会第40回年次大会 予稿集 Visual Computing / グラフィクスとCAD合同シンポ [DVD-ROM] Media Computing Conference, 日本, 一般社団法人画像電子学会, 2012年 7月 2日, p. 1 - 4  
鈴木 郁美、外3名、「ハブの出現しやすさ」から見たラプリアンベースカーネル, 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2011年11月 2日, 第111巻, 第275号, p. 257 - 262  
三川 健太、外3名、テキスト分類問題におけるカテゴリ情報を用いた適応的距離学習に関する一考察, 電子情報通信学会技術研究報告, 日本, 一般社団法人電子情報通信学会, 2012年10月31日, 第112巻, 第279号, p. 83 - 88

- (58)調査した分野(Int.Cl., DB名)  
G06F 17/30