

(19) 日本国特許庁(JP)

再公表特許(A1)

(11) 国際公開番号

W02012/077818

発行日 平成26年5月22日 (2014. 5. 22)

(43) 国際公開日 平成24年6月14日 (2012. 6. 14)

(51) Int.Cl. F 1 テーマコード (参考)  
**G 0 6 F 1 7 / 3 0 ( 2 0 0 6 . 0 1 )**  
 G 0 6 F 1 7 / 3 0 3 5 0 C  
 G 0 6 F 1 7 / 3 0 4 1 2

審査請求 未請求 予備審査請求 未請求 (全 22 頁)

<p>出願番号 特願2012-547940 (P2012-547940)</p> <p>(21) 国際出願番号 PCT/JP2011/078702</p> <p>(22) 国際出願日 平成23年12月12日 (2011. 12. 12)</p> <p>(31) 優先権主張番号 特願2010-276013 (P2010-276013)</p> <p>(32) 優先日 平成22年12月10日 (2010. 12. 10)</p> <p>(33) 優先権主張国 日本国 (JP)</p>	<p>(71) 出願人 304027349                  国立大学法人豊橋技術科学大学                  愛知県豊橋市天伯町雲雀ヶ丘 1-1</p> <p>(74) 代理人 100095577                  弁理士 小西 富雅</p> <p>(72) 発明者 青野 雅樹                  愛知県豊橋市天伯町雲雀ヶ丘 1-1 国立                  大学法人豊橋技術科学大学内</p> <p>(72) 発明者 立間 淳司                  愛知県豊橋市天伯町雲雀ヶ丘 1-1 国立                  大学法人豊橋技術科学大学内</p>
--	---

最終頁に続く

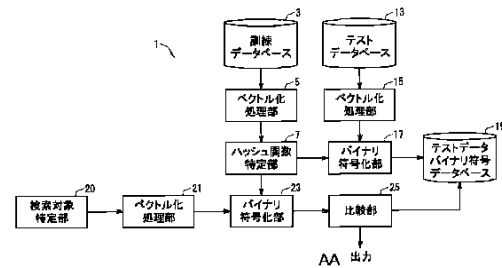
(54) 【発明の名称】 ハッシュ関数の変換行列を定める方法、該ハッシュ関数を利用するハッシュ型近似最近傍探索方法、その装置及びそのコンピュータプログラム

(57) 【要約】

インターネット上に存在するデータの量は拡大の一途をたどっているため、検索サイトに用いられる検索アルゴリズムには更なる精度の向上が期待されている。

データによる多様体上の局所的な近傍構造で表される非線形な関係を保持し、短いバイナリ符号に変換する新しいハッシュ型近似最近傍探索手法を提案する。

【図1】



- 3 Training database
- 5, 15, 21 Vectorization processing unit
- 7 Hash function designation unit
- 13 Test database
- 17, 23 Binary encoding unit
- 19 Test data binary code database
- 20 Search target designation unit
- 25 Comparison unit
- AA Output

【特許請求の範囲】

【請求項 1】

データベースに含まれる第 1 のベクトルデータ  $\mathbf{x}$  (  $n$  次元 ) をバイナリ符号である  $y = [y_1, y_2, \dots, y_d]$ 、ただし、 $n \gg d$  に変換するハッシュ型近似最近傍探索方法において、下記式 ( A ) ~ 式 ( C ) より前記バイナリ符号  $y$  を得る下記ハッシュ関数  $h(\mathbf{x})$  に適用する変換行列  $\mathbf{r}$  を定める方法であって、

【数 6】

$$h_i(\mathbf{x}) = \text{sign}(\mathbf{r}_i^T \mathbf{x}) \quad \text{式 (A)}$$

10

ここに、前記バイナリ符号  $y$  は次のように表わされる、

【数 7】

$$\mathbf{y} = [y_1, y_2, \dots, y_d]^T \quad \text{式 (B)}$$

【数 8】

$$y_i = (1 + h_i(\mathbf{x}))/2 \quad \text{式 (C)}$$

前記第 1 のベクトルデータ  $\mathbf{x}$  を前記バイナリ符号  $y$  のビット数である前記  $d$  の次元に射影したときの第 2 のベクトルデータ  $\mathbf{h}$  を下記式 ( D ) で規定したとき、

【数 9】

$$\mathbf{h} = F^T(\mathbf{x} - \bar{\mathbf{x}}) \quad \text{式 (D)}$$

20

ただし、

【数 10】

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}_i$$

下記式 ( E ) が最小になる変換行列  $F$  を求め、この変換行列  $F$  を式 ( A ) の変換行列  $\mathbf{r}$  とする、

【数 11】

$$\Phi(W) = \sum_i \left\| \mathbf{h}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{h}_j \right\|^2 \quad \text{式 (E)}$$

30

ただし、 $w_{ij}$  は前記データベースの第 1 のベクトルデータ  $\mathbf{x}$  につき、

【数 12】

$$\mathcal{E}(W) = \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j \right\|^2 \quad \text{式 (F)}$$

式 ( F ) を最小とする重みである、ハッシュ関数の変換行列を定める方法。

40

【請求項 2】

前記変換行列  $F$  を求める際に前記重み  $w_{ij}$  を正規化する、請求項 1 に規定の方法。

【請求項 3】

前記正規化は、前記データベースにおける第 1 のベクトルデータ  $\mathbf{x}$  での近傍の分布を下記式 ( G ) のスコアで定義し、

【数 35】

$$d_i = \sum_{j \in \mathcal{N}_i} \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|/\sigma^2)$$

この分布スコアの対角行列  $D = \text{diag}[d_1, \dots, d_m]$  をもとめて、グラフラプラシアン理論

50

に基づき行なう、請求項 2 に記載の方法。

【請求項 4】

原データベースの原データから訓練データを抽出する訓練データ生成部と、  
前記訓練データに基づき第 1 のハッシュ関数を特定するハッシュ関数特定部と、  
前記第 1 のハッシュ関数を用いて前記原データベースの原データをバイナリ符号化する  
バイナリ符号化部と、

バイナリ符号化された前記原データを保存するバイナリ符号化データ保存部と、を備える  
検索サーバと、

入力された検索対象データを前記第 1 のハッシュ関数を用いてバイナリ符号化する第 2  
のバイナリ符号化部を備えるクライアント端末と、

前記クライアント端末の第 2 のバイナリ符号化部でバイナリ符号化された検索対象デー  
タと前記バイナリ符号化された原データとを比較する比較部と、

を備える、ハッシュ型近似最近傍探索装置において、

前記第 1 のハッシュ関数は請求項 1 ~ 請求項 3 にいずれかの方法で定められる、ハッシ  
ュ型近似最近傍探索装置。

【請求項 5】

ハッシュ型近似最近傍探索装置に用いられる検索サーバであって、

原データベースの原データから訓練データを抽出する訓練データ生成部と、

前記訓練データに基づき第 1 のハッシュ関数を特定するハッシュ関数特定部と、

前記第 1 のハッシュ関数を用いて前記原データベースの原データをバイナリ符号化する  
バイナリ符号化部と、

バイナリ符号化された前記原データを保存するバイナリ符号化データ保存部と、を備え

、  
前記第 1 のハッシュ関数は請求項 1 ~ 請求項 3 にいずれかの方法で定められる、検索サ  
ーバ。

【請求項 6】

原データベースの原データから訓練データを抽出する訓練データ生成部と、

前記訓練データを第 1 のベクトル化方法に基づきベクトル化処理する第 1 のベクトル化  
処理部と、

ベクトル化処理された前記訓練データに基づき第 1 のハッシュ関数を特定するハッシュ  
関数特定部と、

前記原データを前記第 1 のベクトル化方法に基づきベクトル化処理する第 2 のベクトル  
化処理部と、

前記第 1 のハッシュ関数を用いて前記ベクトル化処理された原データをバイナリ符号化  
するバイナリ符号化部と、

バイナリ符号化された前記原データを保存するバイナリ符号化データ保存部と、を備え  
る検索サーバと、

入力された検索対象データを前記第 1 のベクトル化方法に基づきベクトル化処理する第  
3 のベクトル化処理部と、

前記ベクトル化処理された検索対象データを前記第 1 のハッシュ関数を用いてバイナリ  
符号化する第 2 のバイナリ符号化部を備えるクライアント端末と、

前記クライアント端末の第 2 のバイナリ符号化部でバイナリ符号化された検索対象デー  
タと前記バイナリ符号化された原データとを比較する比較部と、

を備える、ハッシュ型近似最近傍探索装置において、

前記第 1 のハッシュ関数は請求項 1 ~ 請求項 3 にいずれかの方法で定められるハッシ  
ュ型近似最近傍探索装置。

【請求項 7】

ハッシュ型近似最近傍探索装置に用いられる検索サーバであって、

原データベースの原データから訓練データを抽出する訓練データ生成部と、

前記訓練データを第 1 のベクトル化方法に基づきベクトル化処理する第 1 のベクトル化

10

20

30

40

50

処理部と、

ベクトル化処理された前記訓練データに基づき第1のハッシュ関数を特定するハッシュ関数特定部と、

前記原データを前記第1のベクトル化方法に基づきベクトル化処理する第2のベクトル化処理部と、

前記第1のハッシュ関数を用いて前記ベクトル化処理された原データをバイナリ符号化するバイナリ符号化部と、

バイナリ符号化された前記原データを保存するバイナリ符号化データ保存部と、を備え、

前記第1のハッシュ関数は請求項1～請求項3にいずれかの方法で定められる、検索サーバ。 10

【請求項8】

請求項1～請求項3のいずれかの方法により定められた変換行列を備えるハッシュ関数を利用するハッシュ型近似最近傍探索方法。

【請求項9】

前記データベースとして所定の訓練データベースを用いて、請求項1～請求項3のいずれかの方法により定められた変換行列を備えるハッシュ関数を特定するステップと、

特定されたハッシュ関数をテストデータベースに適用し、該テストデータベースのベクトルデータをテストデータバイナリ符号に変換するステップと、

検索対象のベクトルデータへ前記ハッシュ関数を適用して検索対象バイナリ符号を作成するステップと、 20

該検索対象バイナリ符号を前記テストデータバイナリ符号と比較するステップと、  
備えるハッシュ型近似最近傍探索方法。

【請求項10】

前記データベースとして所定の訓練データベースを用いて、請求項1～請求項3のいずれかの方法により定められた変換行列を備えるハッシュ関数を特定するハッシュ関数特定部と、

特定されたハッシュ関数をテストデータベースに適用し、該テストデータベースのベクトルデータをテストデータバイナリ符号に変換する変換部と、

検索対象のベクトルデータへ前記ハッシュ関数を適用して検索対象バイナリ符号を作成するバイナリ符号作成部と、 30

該検索対象バイナリ符号を前記テストデータバイナリ符号と比較する比較部と、  
備えるハッシュ型近似最近傍探索装置。

【請求項11】

コンピュータを、

前記データベースとして所定の訓練データベースを用いて、請求項1～請求項3のいずれかの方法により定められた変換行列を備えるハッシュ関数を特定するハッシュ関数特定部と、

特定されたハッシュ関数をテストデータベースに適用し、該テストデータベースのベクトルデータをテストデータバイナリ符号に変換する変換部と、 40

検索対象のベクトルデータへ前記ハッシュ関数を適用して検索対象バイナリ符号を作成するバイナリ符号作成部と、

該検索対象バイナリ符号を前記テストデータバイナリ符号と比較する比較部と、  
備えるハッシュ型近似最近傍探索装置として機能させるコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はハッシュ関数の変換行列を定める方法、該ハッシュ関数を利用するハッシュ型近似最近傍探索方法、その装置及びそのコンピュータプログラムに関する。 50

## 【背景技術】

## 【0002】

現在、インターネット上には、文書、画像、音楽、動画など様々なデータが大量に存在している。これら大量のデータを有効利用するため、大規模データベースを対象とし、高速に検索質問と類似するものを見つけ出す技術が、コンピュータビジョンやテキストマイニングの分野で注目されている。例えば、画像であれば、カメラ付き携帯電話で商品を撮影し、それと見た目が類似した商品を、大量の商品データから瞬時に検索することができる。また、特定の風景画像を検索したい場合に、大規模画像データベースから類似する画像を高速に検索できれば、どこで撮影された画像かを即座に判定することができる。

## 【0003】

一般的に、文書ベクトルや画像の特徴ベクトルは、数百から数千の高次元なものとなる。高次元特徴ベクトルで、大規模データベースを検索対象とした場合、線形探索では実用的な検索速度を得ることは難しい。この問題に対して、近似最近傍探索という、大規模データベースで高速な検索を実現する技術が注目されている。

近似最近傍探索は、木構造型とハッシュ型（非特許文献1）の二つに大きく分類される。木構造型近似最近傍探索は、特徴空間上に張られる軸の分割を繰り返して木構造を生成し、探索の際に探索範囲を狭めることで高速に探索を行う。探索範囲は、検索クエリとの暫定的な距離と許容誤差で定義される半径による超球で定義されるが、高次元ベクトルデータを対象とした場合に、次元の呪いの影響をうける。ハッシュ型近似最近傍探索は、高次元ベクトルデータを、短いバイナリ符号に変換し、これをハッシュテーブルのキーとすることで、高速に探索を行う。特徴空間上での距離関係をとらえ、類似するベクトルデータ同士では、バイナリ符号間のハミング距離が小さくなるように変換することで、次元の呪いの影響を小さくできる。また、短いバイナリ符号に変換することで、検索インデックスに必要な容量を抑えることもできる。

## 【0004】

ハッシュ型近似最近傍探索で焦点となるのは、高次元ベクトルデータをバイナリ符号に変換するアルゴリズムである。その目的は、 $n$ 次元の  $m$ 個のベクトルデータ集合  $X = [x_1, x_2, \dots, x_m] \in \mathbb{R}^n \times \mathbb{R}^m$  が与えられた場合に、ハッシュ関数  $h$  を用いて、ベクトルデータ間の類似関係を保持したまま、 $d$ ビットのバイナリ符号集合  $Y = [y_1, y_2, \dots, y_m] \in \mathbb{B}^d \times \mathbb{B}^m$

へと変換することである。

## 【0005】

Locality Sensitive Hashing (LSH) は、最も知られているハッシュ型近似最近傍探索アルゴリズムである。LSHのハッシュ関数は以下の特性を満たすことを条件としている。

$$\Pr[h(x_i) = h(x_j)] = \text{sim}(x_i, x_j)$$

ここで、 $\text{sim}(x_i, x_j) \in [0, 1]$  は、類似度を表す関数である。これは、類似するベクトルデータ同士は、同じハッシュ値になることを示している。Charikarは、内積による類似度  $\text{sim}(x_i, x_j) = x_i^T x_j$  を考え、データ  $x$  と同じ次元の標準正規分布  $N(0, 1)$  によるランダムな超平面（変換ベクトル） $r$  との積によるハッシュ関数を提案した。

$h_i(x) = \text{sign}(r_i^T x)$  ここで、 $\text{sign}$  は、与えられた数値の符号を返す関数である。バイナリ符号  $y$  は、以下のようにして得る。

$$y = [y_1, y_2, \dots, y_d]^T$$

$$y_i = (1 + h_i(x)) / 2$$

このハッシュ関数が、LSHの特性を満たすことは、最大カット問題の近似解法で示される。

## 【数1】

$$\Pr[\text{sign}(r^T x_i) = \text{sign}(r^T x_j)] = 1 - \frac{1}{\pi} \cos^{-1} \left( \frac{x_i^T x_j}{\|x_i\| \|x_j\|} \right)$$

10

20

30

40

50

## 【 0 0 0 6 】

また、Kulisらは、非線形な写像  $(x)$  の内積による類似度  $\text{sim}(x_i, x_j) = k(x_i, x_j) = (x_i)^T (x_j)$  を用いた Kernelized Locality Sensitive Hashing (KLSH) を提案した。

Salakhutdinovらは、ユニット数を徐々に減少させた複数の Restricted Boltzmann Machines (RBM) によるネットワーク構造を用いてバイナリ符号を得る Semantic Hashing を提案した。Semantic Hashing のアルゴリズムでは、教師なしの事前訓練フェーズと、教師ありの微調整フェーズの二段階の学習からなる。事前訓練フェーズでは、ある層での出力は次の層の入力となるように、段階を追って各層ごとに訓練が実施される。微調整フェーズでは、ラベル付きデータを用いて、誤差逆伝搬法により、事前訓練フェーズで得られた重みを調整する。Torralbaらは、この Semantic Hashing を類似画像検索に応用し、LSH よりも高い検索精度を得た。

10

## 【 0 0 0 7 】

Weissらが提案した Spectral Hashing (SH) は、グラフの分割問題を応用してバイナリ符号を求める (非特許文献 2)。Weissらは、ハッシュ型近似最近傍探索における有効なバイナリ符号を求めるため、(1) 新規データのバイナリ符号計算が容易であること (2) わずかなビット数で全データセットを表現すること (3) 類似するデータは類似するバイナリ符号となること、の三つの条件を設定した。これら条件を満たすバイナリ符号を求めるため、Weissらは、以下の最小化問題を考えた。

目的関数：

【数 2】

20

$$\sum_{i,j} W_{i,j} \| \mathbf{h}_i - \mathbf{h}_j \|^2 \quad \text{式 (1)}$$

制約条件

【数 3】

$$\mathbf{h}_i = [h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_d(\mathbf{x}_i)]^T, \quad h_k(\mathbf{x}_i) \in \{-1, 1\}, \quad 1 \leq k \leq d \quad \text{式 (2)}$$

【数 4】

$$\sum_i \mathbf{h}_i = 0 \quad \text{式 (3)}$$

30

【数 5】

$$\frac{1}{n} \sum_i \mathbf{h}_i \mathbf{h}_i^T = I \quad \text{式 (4)}$$

ここで、 $W_{i,j} = \exp(-\|x_i - x_j\|^2 / \sigma^2)$  であり、式 (1) は、特徴空間の局所的な類似関係をバイナリ符号に反映することを示す。制約条件において、式 (2) は、バイナリ符号が -1 と 1 からなることを、式 (3) は、各ビットは偏りなく -1 もしくは 1 を取り得ることを、式 (4) は、異なるビット間は互いに独立であることを表す。Weissらは、式 (2) の制約条件を緩和することで、グラフラプラシアン固有ベクトルを求める問題とした。SH は、LSH、KLSH、Semantic Hashing などよりも検索精度が高いことが知られている。

40

## 【 0 0 0 8 】

この他に、逐次学習を応用して頑健なバイナリ符号を得る手法、ラベルが付与されたデータを用いてデータの類似・相違を教師する手法、平行移動不変カーネルで表されるデータ間の関係を保持したバイナリ符号を得る手法、ベクトルデータの分布とバイナリ符号の分布によるカルバック・ライブラー情報量が最小となるようにバイナリ符号を求める手法などがある。

【先行技術文献】

【特許文献】

50

【 0 0 0 9 】

【 特許文献 1 】 特開 2 0 1 0 - 3 9 7 7 8 号公報

【 特許文献 2 】 特開 2 0 0 9 - 7 5 6 0 3 号公報

【 特許文献 3 】 特開 2 0 0 6 - 3 0 9 7 1 8 号公報

【 非特許文献 】

【 0 0 1 0 】

【 非特許文献 1 】 Charikar, M., Similarity estimation techniques from rounding algorithms, In Proceedings of the thirty-fourth annual ACM symposium on Theory of computing, pp.280-288, 2002.

【 非特許文献 2 】 Weiss, Y., Torralba, A., Fergus, R., Spectral hashing, In The Neural Information Processing Systems, Vol.21, pp.1753-1760, 2008. 10

【 非特許文献 3 】 Wang, B., Li, Z., Li, M., Ma, W.-Y., Large-Scale Duplicate Detection for Web Image Search, In Proceedings of IEEE International Conference on Multimedia and Expo, pp.353-356, 2006.

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 1 】

インターネット上に存在するデータの量は拡大の一途をたどっているため、検索サイトに用いられる検索アルゴリズムには更なる精度の向上が期待されている。

そこで本発明者らは、検索サイト用の検索アルゴリズムとして現在利用されている上記 LSH に着目し、その改良、即ち検索精度の向上を検討した。 20

その結果、LSH では、そこで用いられる下記ハッシュ関数

$$h_i(x) = \text{sign}(r_i^T x)$$

において変換ベクトル  $r$  がランダムに決められているため、ベクトルデータをバイナリ符号に変換したとき、もとなるベクトルデータ間の類似関係が十分に維持されないおそれがあると考えた。

非特許文献 3 には、ベクトルデータが含まれるデータベースの特性を当該変換ベクトルに反映させ、もってバイナリ符号に変換する際にベクトルデータ間の類似関係維持を図ろうとする試みが示されている。

しかしながら非特許文献 3 に記載の方法は、主成分分析法を利用しているため、非線形構造をなす多次元空間のデータベースの特性を十分に変換ベクトルに反映させることには無理がある。 30

【 課題を解決するための手段 】

【 0 0 1 2 】

そこで本発明者らは、変換ベクトルにデータベースがなす非線形構造を反映させるべく鋭意検討を重ねた結果、この発明に想到した。

即ち、この発明の第 1 の局面は次のように規定される。

データベースに含まれる第 1 のベクトルデータ  $x$  ( $n$  次元) をバイナリ符号である  $y = [y_1, y_2, \dots, y_d]$ 、ただし、 $n \gg d$  に変換するハッシュ型近似最近傍探索方法において、

下記式 (A) ~ 式 (C) より前記バイナリ符号を得る下記ハッシュ関数  $h(x)$  に適用する変換行列  $r$  を定める方法であって、 40

【 数 6 】

$$h_i(\mathbf{x}) = \text{sign}(\mathbf{r}_i^T \mathbf{x}) \quad \text{式 (A)}$$

ここに、前記バイナリ符号  $y$  は次のように表わされる、

【 数 7 】

$$\mathbf{y} = [y_1, y_2, \dots, y_d]^T \quad \text{式 (B)}$$

【数 8】

$$y_i = (1 + h_i(\mathbf{x}))/2 \quad \text{式 (C)}$$

前記第 1 のベクトルデータ  $\mathbf{x}$  を前記バイナリ符号  $y$  のビット数である前記  $d$  の次元に射影したときの第 2 のベクトルデータ  $\mathbf{h}$  を下記式 (D) で規定したとき、

【数 9】

$$\mathbf{h} = F^T(\mathbf{x} - \bar{\mathbf{x}}) \quad \text{式 (D)}$$

ただし、

10

【数 10】

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}_i$$

下記式 (E) が最小になる変換行列  $F$  を求め、この変換行列  $F$  を式 (A) の変換行列  $r$  とする、

【数 11】

$$\Phi(W) = \sum_i \left\| \mathbf{h}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{h}_j \right\|^2 \quad \text{式 (E)}$$

20

ただし、 $w_{ij}$  は前記データベースの第 1 のベクトルデータ  $\mathbf{x}$  につき、

【数 12】

$$\mathcal{E}(W) = \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j \right\|^2 \quad \text{式 (F)}$$

式 (F) を最小とする重みである、ハッシュ関数の変換行列を定める方法。

【0013】

第 1 の局面に規定の発明の基本原理を以下に説明する。

高次元ベクトルデータをバイナリ符号に変換するアルゴリズムの多くが、次元削減のアルゴリズムに基づいている。次元削減の目的は、高次元空間上でベクトルデータがなす、低次元の部分空間を推定することである。近年、全体では非線形構造を成していても、局所的には通常のユークリッド空間と同じ構造をなす、多様体の性質を利用した非線形次元削減手法が幾つか提案されている。その内、Locally Linear Embedding (LLE) は、局所的な範囲で低次元の線形モデルをあてはめ、それらが滑らかに繋がるように全体の多様体を推定する。この LLE による多様体構造の推定方法を用いて、ベクトルデータがなす非線形な関係をとらえたバイナリ符号を得る。以下、Weissらと同様に、問題の簡単化のため、式(2)の制約条件を緩和する。

30

【0014】

局所的な範囲で低次元の線形モデルをあてはめるため、それぞれのベクトルデータ  $\mathbf{x}_i$  を、近傍のベクトルデータ  $\mathbf{x}_j \in \mathcal{N}_i$  を用いて再構成することを考える。これは、以下の再構成誤差を最小化することで表される。

40

【数 13】

$$\mathcal{E}(W) = \sum_i \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j \right\|^2$$

ここで、 $w_{ij}$  は、再構成する際の重みである。再構成には近傍のベクトルデータを用いるため、

【数 14】

$$j \notin \mathcal{N}_i \Rightarrow w_{ij} = 0,$$

50



重みの大きさの任意性を解決するため  $\sum_j w_{ij} = 1$  とする。各ベクトルデータでの再構成誤差は、

【数 1 5】

$$\begin{aligned} \mathcal{E}_i(W) &= \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j \right\|^2 \\ &= \left\| \mathbf{x}_i \sum_{j \in \mathcal{N}_i} w_{ij} - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j \right\|^2 \\ &= \sum_{j, k \in \mathcal{N}_i} w_{ij} w_{ik} (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_k) \\ &= \sum_{j, k \in \mathcal{N}_i} w_{ij} w_{ik} C_{jk} \end{aligned}$$

10

と表わせる。ここで、 $C_{jk} = (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_k)$  とした。この再構成誤差は、ラグランジュ乗数  $\eta_i$  を用いて

【数 1 6】

$$\mathcal{E}_i(W) = \sum_{j, k \in \mathcal{N}_i} w_{ij} w_{ik} C_{jk} + \eta_i \left( \sum_{j \in \mathcal{N}_i} w_{ij} - 1 \right)$$

20

となる。極致を求めるため、 $w_{ij}$  について偏微分し 0 とおくことで、以下の線形方程式を解く問題となる。

【数 1 7】

$$\sum_{k \in \mathcal{N}_i} C_{jk} w_{ik} = 1$$

以上の重みを求める計算については、de Ridderらや Panらの研究で詳述されている。

【0 0 1 5】

近傍による再構成の重み  $w_{ij}$  で表される、それぞれのベクトルデータにおける局所的な関係を、バイナリ符号に反映させるため、Weissらの最小化問題において、式(1)の目的関数を以下のものに置き換える(式(E))。

30

【数 1 8】

$$\Phi(W) = \sum_i \left\| \mathbf{h}_i - \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{h}_j \right\|^2$$

さらに、新規データのバイナリ符号容易に求められるよう、変換行列  $F \in \mathbb{R}^{n \times d}$  による線形変換を考える(式D)。

【数 1 9】

$$\mathbf{h} = F^T (\mathbf{x} - \bar{\mathbf{x}})$$

40

ここで、式(3)の制約条件を満たすため、平均ベクトル  $\bar{\mathbf{x}}$  で引いた。

【数 2 0】

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}_i$$

【0 0 1 6】

バイナリ符号の集合を  $H = [h_1, \dots, h_m]$  と表すと、式(E)の目的関数は、以下のように表される。

【数 2 1】

$$\begin{aligned}
 \Phi(W) &= \text{tr}[(H - HW)^T(H - HW)] \\
 &= \text{tr}[(H - HW)(H - HW)^T] \\
 &= \text{tr}[H(I - W)(I - W)^T H^T] \\
 &= \text{tr}[F^T X(I - W)(I - W)^T X^T F] \\
 &= \text{tr}[F^T X M X^T F]
 \end{aligned}$$

10

ここで、 $M = (I - W)(I - W)^T$ とした。さらに、式(4)の制約条件は、

【数 2 2】

$$\frac{1}{n} H H^T = \frac{1}{n} F^T X X^T F = I$$

となり、式(6)の目的関数の最小化は、

【数 2 3】

$$\begin{aligned}
 \text{argmin}_F \quad & \text{tr}[F^T X M X^T F] \\
 & F^T X X^T F = nI
 \end{aligned}$$

20

と表わせる。これは、ラグランジュの未定乗数法から、

【数 2 4】

$$\mathcal{L}(F) = F^T (X M X^T) F + \lambda (nI - F^T X X^T F)$$

となるので、Fで偏微分して0としておくことにより、

【数 2 5】

$$\frac{\partial \mathcal{L}}{\partial F} = 2(X M X^T) F - 2\lambda X X^T F = 0$$

30

以下の一般固有値問題へと帰着する。

【数 2 6】

$$X M X^T F = \lambda X X^T F \quad \text{式 (H)}$$

ここで、変換行列 F は、式(H)の一般化固有値問題を解くことで得られる、固有値  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$  に対応する固有関数からなる。

【数 2 7】

$$F = [f_1, f_2, \dots, f_d]$$

40

【0017】

上記において、ベクトルデータ h は、ベクトルデータ x をバイナリ符号 y のビット数である d の次元に射影したものである。変換行列 F を式(A)の変換行列 r へ当てはめる。

上記において、局所的な近傍構造を表す重み  $w_{ij}$  を、それぞれのベクトルデータの近傍から求める。しかし、注目しているベクトルデータに対して、近傍が密集して位置するものもあれば、遠く離れて位置するものもある。そこで、この分布の偏りによる影響を軽減するため、各ベクトルデータでの近傍の分布を表すスコアを定義する。

【数 2 8】

$$d_i = \sum_{j \in \mathcal{N}_i} \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|/\sigma^2)$$

近傍の分布を考慮した変換行列  $F$  は、分布スコアからなる対角行列  $D = \text{diag}[d_1, \dots, d_m]$  により、以下の最小化問題を解くことで得られる。

【数 2 9】

$$\underset{F}{\text{argmin}} \quad \text{tr}[F^T X M X^T F]$$

$$F^T X X^T F = nI$$

10

【0 0 1 8】

これは、ラグランジュの未定乗数法から、

【数 3 0】

$$X D M X^T F = \lambda X X^T F$$

と表わすことができる。さらに、 $X$  がフルランク行列であれば、 $X^T X$  は正則となることから、

【数 3 1】

$$X^T X D M X^T F = \lambda X^T X X^T F$$

$$D M X^T F = \lambda X^T F$$

20

分布スコアは必ず正数であるので、 $S = D^{-1} = \text{diag}[1/d_1, \dots, 1/d_m]$  とおくと、上記の最小化問題は、以下の一般化固有値問題に帰着する。

【数 3 2】

$$X M X^T F = \lambda X S X^T F$$

30

$X$  がフルランク行列ではない場合は、特異値分解を用いて、行列  $X$  のランク数  $l$  と等しい次元の部分空間に射影する。

【数 3 3】

$$X = U R V^T$$

$$\tilde{X} = U^T X = R V^T$$

ここで、 $R$  は要素が特異値  $\eta_1 \geq \dots \geq \eta_l$  の大きさ  $l \times l$  の対角行列であり、 $U$  は大きさ  $n \times l$  の、 $V$  は大きさ  $m \times l$  の直行行列である。 $\tilde{X}$  はフルランク行列となるので、以下の一般化固有値問題を解くことで、変換行列  $F$  を得る。

40

【数 3 4】

$$\tilde{X} M \tilde{X}^T F = \lambda \tilde{X} S \tilde{X}^T F$$

【0 0 1 9】

上記の処理はいわゆるグラフラプラシアン理論を用いて、重み  $w_{ij}$  を正規化したことを意味し、 $n$  次元空間のベクトルデータ  $x$  を  $d$  次元空間への射影するとき、 $n$  次元空間におけるベクトルデータ  $x$  の分布の偏りの影響を軽減できる。よって、ベクトルデータ得られた変

50

換行列 F はデータベースの特性をより正確に反映したものとなる。

【図面の簡単な説明】

【0020】

【図1】図1は、この発明のハッシュ型近似最近傍検索システムの構成を示すブロック図である。

【図2】図2は、20-newsgroupsでビット数を8から64まで変化させた場合の、検索結果上位400件における適合率(Precision)を表したグラフである。

【図3】図3は、20-newsgroupsで検索結果の上位件数を変化させた場合の、ビット数64における再現率(Recall)を表したグラフである。

【図4】図4は、20-newsgroupsにおける再現率と適合率を表したグラフである。

【図5】図5は、CIFAR-10でビット数を8から64まで変化させた場合の、検索結果上位1,000件における適合率(Precision)を表したグラフである。

【図6】図6は、CIFAR-10で検索結果の上位件数を変化させた場合の、ビット数64における再現率(Recall)を表したグラフである。

【図7】図7は、CIFAR-10における再現率と適合率を表したグラフである。

【図8】図8は、CIFAR-10で、SHとNSHによる32ビットのバイナリ符号を用いた、自動車の画像での検索結果上位20件を並べたものである。

【図9】図9は、実施例のハッシュ型近似最近傍探索装置を示すブロック図である。

【発明を実施するための形態】

【0021】

このようにして定められた変換行列を有するハッシュ関数を用いてハッシュ型近似最近傍検索を行なうシステム1を図1に示す。

図1において、訓練データベース3、テストデータベース13、テストデータバイナリ符号データベース19はサーバのメモリ装置の所定の領域が対応される。

訓練データベース3のデータはベクトル化処理部5において所定の方法でベクトル化される。

ハッシュ関数特定部7は、ベクトル化された訓練データベース3のデータ(第1のベクトルデータx)の一部又は全部を用いて、既述の処理を行ない、変換ベクトルFを特定し、もって、第1のベクトルデータxをバイナリ符号化するハッシュ関数を特定する。

【0022】

テストデータベース13のデータは、ベクトル化処理部15において、ベクトル化処理部5と同一の方法でベクトル化される。

バイナリ符号化部17は、ハッシュ関数特定部7で特定されたハッシュ関数を用いてベクトル化されたテストデータベースのデータの一部又は全部をバイナリ符号に変換し、テストデータバイナリ符号データベース19に保存する。

【0023】

クライアントPCが検索対象特定部20に対応し、この検索対象特定部20においてユーザが検索対象を特定する。特定された検索対象はベクトル化処理部21においてベクトル化処理される。このベクトル化処理の方法は、既述のベクトル化処理部5において、訓練データベース3のデータを第1のベクトルデータに変換したベクトル化処理法と同一である。

このようにしてベクトル化された検索対象はバイナリ符号化部23において、ハッシュ関数特定部7で特定されたハッシュ関数を用いてバイナリ符号化される。比較部25は、検索対象のバイナリ符号をテストデータバイナリ符号データベース19に保存されているテストデータバイナリ符号と比較し、例えば、検索対象のバイナリ符号に対する距離が所定の閾値以内のものを、近い順に出力する。

【0024】

この発明の探索方法を評価するために、ベンチマークに20-newsgroupsとCIFAR-10を用いて、従来手法との比較実験を行った。従来手法には、Locality-Sensitive Hashing (LSH)、Kernelized Locality-Sensitive Hashing (KLSH)、Spectral Hashing (SH)を選択し

10

20

30

40

50

た。この内、アルゴリズム中に乱数生成を含むLSHとKLSHについては、5回実行して平均をとった。

20-newsgroupsは、Usenet newsgroupから取得した18,845件のニュースグループの文書からなる。各文書は、異なる20個のニュースグループのいずれかに分類され、11,314件が訓練データセットとして、7,531件がテストデータとして与えられる。各アルゴリズムの訓練には、訓練データセットからランダムに選択した5,000件を用いた。実験では、前処理として単語のステミングとストップワードの除去を行ったのち、文書頻度の大きいほうから2,000語を選択し、tf-idfにより重み付けを行った文書ベクトルを作成した。本発明のパラメータは、事前実験により求めた最適値、近傍数  $k = 205$ 、分布スコアのガウスカーネル幅  $\sigma = 4.0$ を用いた。

10

#### 【0025】

図1は、20-newsgroupsでビット数を8から64まで変化させた場合の、検索結果上位400件における適合率(Precision)を表したグラフである。全てのビット数で、本発明(NSH)が最も大きな適合率となっている。図3は、20-newsgroupsで検索結果の上位件数を変化させた場合の、ビット数64における再現率(Recall)を表したグラフである。全ての上位件数で、本発明(NSH)が最も大きな再現率となっているのがわかる。図4は、20-newsgroupsにおける再現率と適合率を表したグラフである。曲線が右上に伸びるほど検索精度が高いと言える。本発明(NSH)が、従来手法と比較して、正確性・網羅性ともに高いことがわかる。

20

#### 【0026】

20-newsgroupsは、20個のカテゴリに分類されているが、それらは、コンピュータの話題やスポーツの話題など、大きく分類することもできる。これらは高次元の文書ベクトル空間で、類似した話題のものが集まりながらも、粗密をなして複雑な構造を成していると予想する。ベクトルデータの非線形構造を推定し、分布の偏りを重みにより考慮する本発明(NSH)は、従来手法と比較して、この複雑な文書ベクトル空間を正しくとらえられたと考える。

#### 【0027】

CIFAR-10は、飛行機や自動車、イヌなどの10種類のラベルが付与されクラス分けされた、60,000個の  $32 \times 32$ のカラー画像が含まれており、50,000個が訓練データセットとして、10,000個がテストデータセットとして与えられる。各アルゴリズムの訓練には、訓練データセットからランダムに選択した5,000個を用いた。実験では、RGBごとに  $4 \times 4$ 領域6方向4スケールで  $3 \times 4 \times 4 \times 6 \times 4 = 1$ 、152次元のGIST特徴ベクトルを抽出し、それをバイナリ符号に変換した。本発明(NSH)のパラメータは、事前実験により求めた最適値、近傍数  $k = 90$ 、分布スコアのガウスカーネル幅  $\sigma = 0.5$ を用いた。

30

#### 【0028】

図5は、CIFAR-10でビット数を8から64まで変化させた場合の、検索結果上位1,000件における適合率(Precision)を表したグラフである。全てのビット数で、本発明(NSH)が最も大きな適合率となっている。図6は、CIFAR-10で検索結果の上位件数を変化させた場合の、ビット数64における再現率(Recall)を表したグラフである。全ての上位件数で、本発明(NSH)が最も大きな再現率となっているのがわかる。図7は、CIFAR-10における再現率と適合率を表したグラフである。本発明(NSH)が、従来手法と比較して、正確性・網羅性ともに高いことがわかる。

40

#### 【0029】

CIFAR-10には、飛行機や自動車などからなる乗り物画像と、イヌやシカなどからなる生物画像に、大きく分けられる。しかし、画像全体では特定の色が多く使用されている、背景が類似しているなどの要因で、特徴空間上では、乗り物画像と生物画像がはっきりと区別されることなく、粗密をなして分布していると予想する。分布の偏りを重みにより抑えつつ、ベクトルデータがなす非線形な構造を推定するNSHは、この特徴空間上の複雑なデータ関係を、従来手法と比較して正しくとらえられたと考える。

#### 【0030】

50

図 8 は、CIFAR-10で、SHとNSHによる 32ビットのバイナリ符号を用いた、自動車の画像での検索結果上位20件を並べたものである。左上端の画像が検索質問画像であり、正解の画像は緑色の枠で、不正解の画像は赤色の枠で囲んだ。SHと比較して、本件発明（NSH）では適合画像が多く、検索精度が高いことがわかる。

【 0 0 3 1 】

本件発明（NSH）は、特徴空間上の局所的な近傍構造に基づいて、ベクトルデータがなす非線形構造をとらえた短いバイナリ符号を得る。文書データベンチマーク20-newsgroupsと画像データベンチマーク CIFAR-10を用いた比較実験から、Spectral Hashingなどの従来手法よりも高い検索精度を得られることが確認できた。

【 0 0 3 2 】

実施例のハッシュ型近似近最傍探索装置 1 0 0 を図 9 に示す。

この装置 1 0 0 は、検索サーバ 1 0 1 とクライアント端末 2 0 0 がネットワーク N 1 を介して接続されている。任意の数のクライアント端末 2 0 0 をネットワーク N 1 へ接続可能で、ネットワーク N 1 はインターネット 3 0 0 へ開いていてもよい。

【 0 0 3 3 】

検索サーバ 1 0 1 はデータ保存部 1 1 0 とデータ処理部 1 2 0 とを備える。データ保存部 1 1 0 は原データ取得・保存部 1 1 1 を備える。この原データ取得・保存部 1 1 1 は被検索対象となるデータを、インターネット 3 0 0 を介して外部のデータベースから取得し、保存する。外部のデータベースは特定のデータが整理保存された商用、または検索用のデータベースはもとより、法人若しくは個人が運営するホームページ、ブログ、ツイッター等も含まれるものとする。

訓練データ抽出・更新部 1 1 3 は、原データ取得・保存部 1 1 1 に保存されたデータから無作為にデータを抽出し、訓練データとする。抽出するデータ数は特に限定されるものではないが、既述の試験例に準じて 5 0 0 0 個程度とすることが好ましい。なお、この訓練データは周期的に、若しくは任意のタイミングで更新することが好ましい。

【 0 0 3 4 】

ベクトル化処理部 1 2 1 及びハッシュ関数特定部 1 2 3 の動作は、図 1 のベクトル化処理部 5 及びハッシュ関数特定部 7 と同じである。即ち、訓練データ抽出・更新部 1 2 1 で抽出されたデータを特徴量次元削減法等の汎用的な手法（第 1 のベクトル化方法）でベクトル化処理し、ハッシュ関数特定部 1 2 3 において、本発明の手法に従いハッシュ関数（第 1 のハッシュ関数）を特定する。

ベクトル化処理部 1 2 1 は原データ取得・保存部 1 1 1 のデータの全部をまたは所定のルールで選択されたその一部を訓練データと同様にベクトル化処理する。ベクトル化処理され原データは、バイナリ符号化部 1 2 5 においてバイナリ符号化データに次元削減される（前処理 1 3 0）。このとき、ハッシュ関数特定部 1 2 3 において特定されたハッシュ関数が利用される。バイナリ符号化されたデータはバイナリ符号化データ保存部 1 1 5 に保存される。

【 0 0 3 5 】

クライアント端末 2 0 0 において、検索対象となるデータが入力部 2 1 0 で指定される。指定されたデータはベクトル化処理部 2 2 1 において、検索サーバ 1 0 1 のベクトル化処理部 1 2 1 と同一の方法によりベクトル化される。バイナリ符号化部 2 2 3 には、検索サーバ 1 0 1 のハッシュ関数特定部 1 2 3 で特定されたハッシュ関数が提供され、ベクトル化された検索対象データをバイナリ符号化する。このようにしてバイナリ符号化された検索対象データは検索サーバ 1 0 1 の比較部 1 2 7 へ送られる。比較部 1 2 7 はバイナリ符号化された検索対象データとバイナリ符号化データ保存部 1 1 5 に保存されているバイナリ符号化された原データとを比較し、所定のルールに従い近似するデータを抽出する。

【 0 0 3 6 】

比較部 1 2 7 で抽出されたデータはクライアント端末 2 0 0 の出力部 2 3 0 へ送られ、ここで出力される。なお、出力部 2 3 0 はバイナリ符号化されたデータをデコード（逆ハッシュ関数処理、逆ベクトル処理）し、原データの状態で表示できる。

10

20

30

40

50

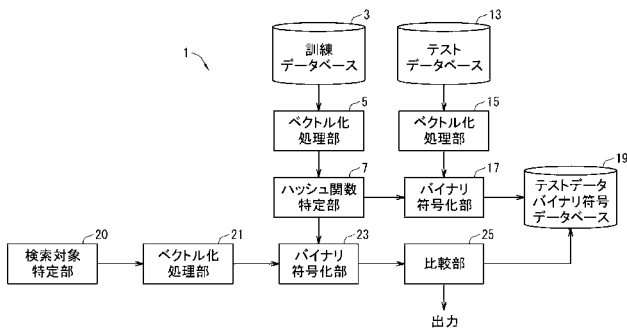
この比較部を端末側に配置することができる。また、比較部を、検索サーバ及び端末と独立して設置することもできる。

【符号の説明】

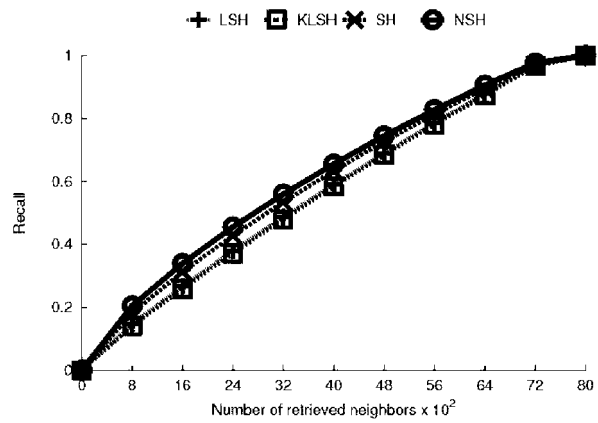
【0037】

- 1 ハッシュ型近似最近傍検索システム
- 3、13 データベース
- 5、15、21 ベクトル化処理部
- 7 ハッシュ関数特定部
- 17、23 バイナリ符号化部
- 25 比較部

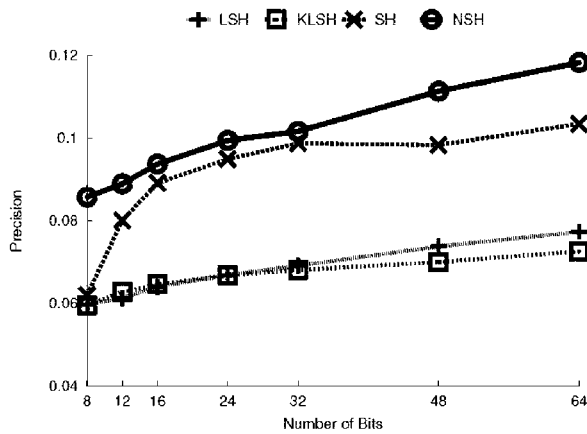
【図1】



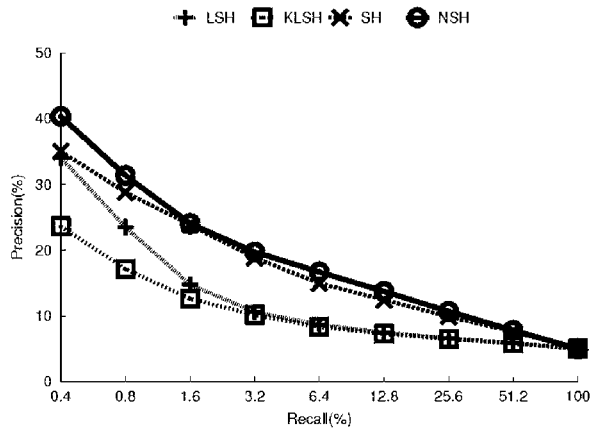
【図3】



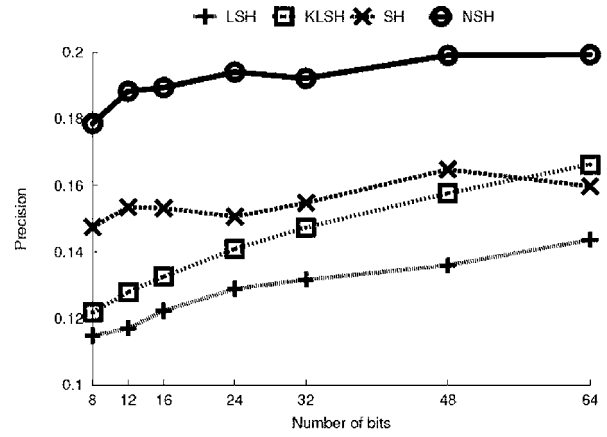
【図2】



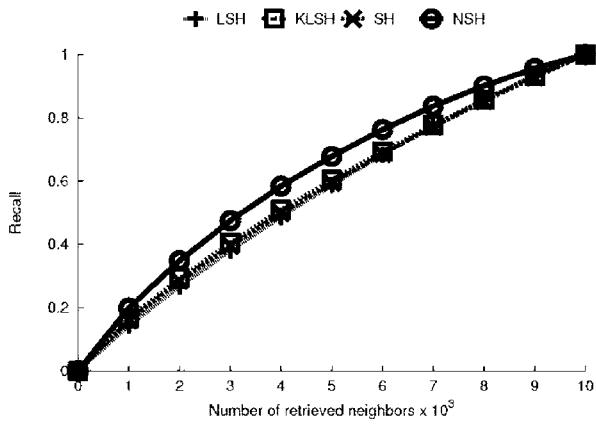
【 図 4 】



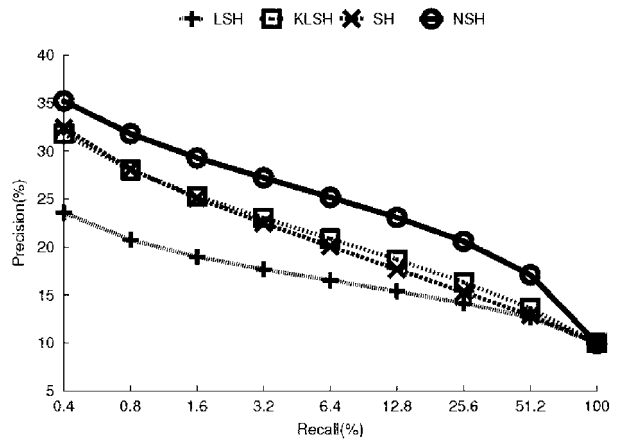
【 図 5 】



【 図 6 】

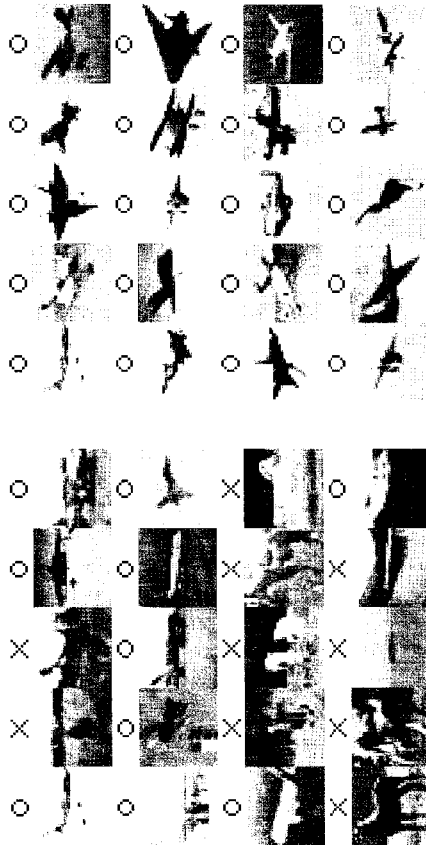


【 図 7 】

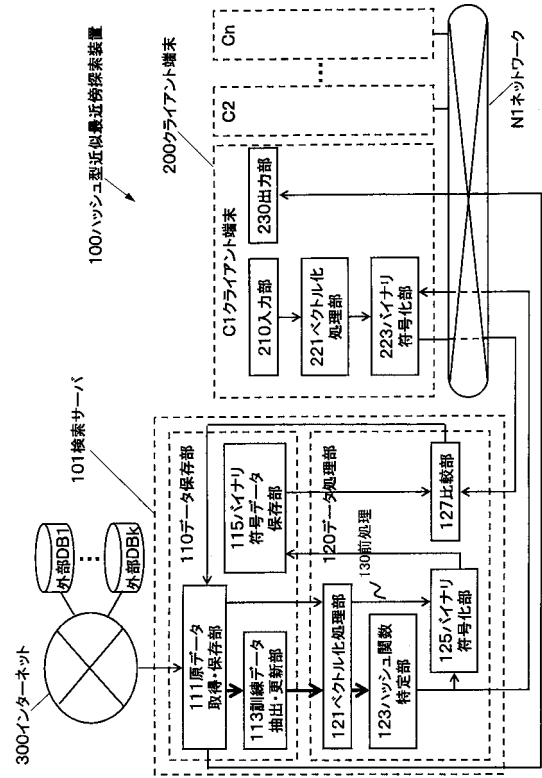




【 図 8 】



【 図 9 】



## 【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/JP2011/078702
<b>A. CLASSIFICATION OF SUBJECT MATTER</b> G06F17/30 (2006.01) i  According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) G06F17/30  Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2012 Kokai Jitsuyo Shinan Koho 1971-2012 Toroku Jitsuyo Shinan Koho 1994-2012  Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) JSTPlus (JDreamII)		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2003-141160 A (International Business Machines Corp.), 16 May 2003 (16.05.2003), entire text; all drawings & US 2003/0159106 A1	1-11
A	JP 2010-256951 A (Dehenken Ltd.), 11 November 2010 (11.11.2010), entire text; all drawings (Family: none)	1-11
A	JP 2010-277522 A (Nippon Telegraph and Telephone Corp.), 09 December 2010 (09.12.2010), entire text; all drawings (Family: none)	1-11
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 10 February, 2012 (10.02.12)		Date of mailing of the international search report 21 February, 2012 (21.02.12)
Name and mailing address of the ISA/ Japanese Patent Office		Authorized officer
Facsimile No.		Telephone No.

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2011/078702

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	Shuhei FURUKAWA et al., "Document Similarity Search and Reranking using Adjacency Information", The First Forum on Data Engineering and Information Management DEIM Forum 2009 Ronbunshu, 09 May 2009 (09.05.2009), no.A9-2, [online] URL: <a href="http://db-event.jpn.org/deim2009/proceedings/files/A9-2.pdf">http://db-event.jpn.org/deim2009/proceedings/files/A9-2.pdf</a>	1-11
A	Masaru NAKANO et al., "Linearized Diffusion Maps and its Application to Documents", Dai 72 Kai (Heisei 22 Nen) Zenkoku Taikai Koen Ronbunshu, 08 March 2010 (08.03.2010), vol.2, pages 2-465 to 2-466	1-11

国際調査報告		国際出願番号 PCT/J P 2011/078702									
A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G06F17/30(2006.01)i											
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G06F17/30											
最小限資料以外の資料で調査を行った分野に含まれるもの <table border="0"> <tr> <td>日本国実用新案公報</td> <td>1922-1996年</td> </tr> <tr> <td>日本国公開実用新案公報</td> <td>1971-2012年</td> </tr> <tr> <td>日本国実用新案登録公報</td> <td>1996-2012年</td> </tr> <tr> <td>日本国登録実用新案公報</td> <td>1994-2012年</td> </tr> </table>				日本国実用新案公報	1922-1996年	日本国公開実用新案公報	1971-2012年	日本国実用新案登録公報	1996-2012年	日本国登録実用新案公報	1994-2012年
日本国実用新案公報	1922-1996年										
日本国公開実用新案公報	1971-2012年										
日本国実用新案登録公報	1996-2012年										
日本国登録実用新案公報	1994-2012年										
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語) JSTPlus(JDreamII)											
C. 関連すると認められる文献											
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号									
A	JP 2003-141160 A (インターナショナル・ビジネス・マシーンズ・コーポレーション) 2003.05.16, 全文、全図 & US 2003/0159106 A1	1-11									
A	JP 2010-256951 A (株式会社データ変換研究所) 2010.11.11, 全文、全図 (ファミリーなし)	1-11									
A	JP 2010-277522 A (日本電信電話株式会社) 2010.12.09, 全文、全図 (ファミリーなし)	1-11									
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。											
* 引用文献のカテゴリー		の日の後に公表された文献									
「A」特に関連のある文献ではなく、一般的技術水準を示すもの		「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの									
「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの		「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの									
「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)		「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの									
「O」口頭による開示、使用、展示等に言及する文献		「&」同一パテントファミリー文献									
「P」国際出願日前で、かつ優先権の主張の基礎となる出願											
国際調査を完了した日 10.02.2012		国際調査報告の発送日 21.02.2012									
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号		特許庁審査官 (権限のある職員) 野崎 大進	5M   9289								
		電話番号 03-3581-1101	内線 3599								

国際調査報告

国際出願番号 PCT/JP2011/078702

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	古川 修平 他, 隣接情報を用いた類似文書検索とリランキング, 第1回データ工学と情報マネジメントに関するフォーラム DEIM フォーラム 2009 論文集, 2009.05.09, No.A9-2, [online]URL: <a href="http://db-event.jpn.org/deim2009/proceedings/files/A9-2.pdf">http://db-event.jpn.org/deim2009/proceedings/files/A9-2.pdf</a>	1-11
A	仲野 将 他, 線形化拡散写像手法の提案とその文書データへの適用, 第72回(平成22年)全国大会講演論文集, 2010.03.08, Vol.2, pp.2-465~2-466.	1-11

---

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, T J, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, R O, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, H U, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI , NO, NZ, OM, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN

(注) この公表は、国際事務局(WIPO)により国際公開された公報を基に作成したものである。なおこの公表に係る日本語特許出願(日本語実用新案登録出願)の国際公開の効果は、特許法第184条の10第1項(実用新案法第48条の13第2項)により生ずるものであり、本掲載とは関係ありません。