

(19) 日本国特許庁(JP)

再公表特許(A1)

(11) 国際公開番号

W02016/117698

発行日 平成29年11月24日 (2017.11.24)

(43) 国際公開日 平成28年7月28日 (2016.7.28)

(51) Int.Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 350C

テーマコード (参考)

審査請求 未請求 予備審査請求 未請求 (全 29 頁)

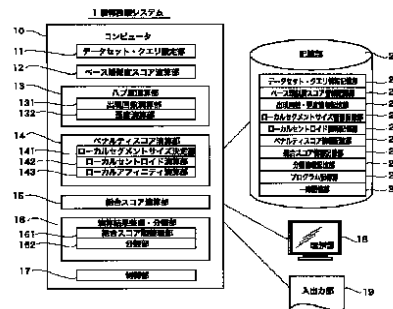
<p>出願番号 特願2016-570732 (P2016-570732)</p> <p>(21) 国際出願番号 PCT/JP2016/051909</p> <p>(22) 国際出願日 平成28年1月22日 (2016.1.22)</p> <p>(31) 優先権主張番号 特願2015-11853 (P2015-11853)</p> <p>(32) 優先日 平成27年1月23日 (2015.1.23)</p> <p>(33) 優先権主張国 日本国(JP)</p>	<p>(71) 出願人 504202472 大学共同利用機関法人情報・システム研究機構 東京都立川市緑町10番3号</p> <p>(74) 代理人 100106909 弁理士 棚井 澄雄</p> <p>(74) 代理人 100188558 弁理士 飯田 雅人</p> <p>(74) 代理人 100161207 弁理士 西澤 和純</p> <p>(74) 代理人 100141139 弁理士 及川 周</p>
---	--

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理方法、プログラム、及び非一時記憶媒体

(57) 【要約】

情報処理装置は、複数のデータにより構成される第1データセット内の処理対象データと、検索対象データとしてのクエリとの間の類似度を算出する第1算出部と、前記第1データセット内の一部のデータにより構成される第2データセットを、前記処理対象データごとに特定する特定部と、前記特定部が特定する前記第2データセットから、前記処理対象データごとの基準を算出する第2算出部と、前記第1算出部が算出する前記類似度と、前記第2算出部が算出する基準とを用いて、前記処理対象データごとにスコアを算出する第3算出部と、を含む。



- 1 Information retrieval system
- 10 Computer
- 11 Data set and query specifying unit
- 12 Base similarity score arithmetic unit
- 13 Base similarity score arithmetic unit
- 14 Penalty score arithmetic unit
- 15 Overall score arithmetic unit
- 16 Arithmetic result arrangement and classification unit
- 17 Control unit
- 18 Display unit
- 19 Input/output unit
- 20 Exchange unit
- 21 Data set and query information storage unit
- 22 Base similarity score information storage unit
- 23 Number of occurrences and arithmetic information storage unit
- 24 Local segment size information storage unit
- 26 Local centroid information storage unit
- 28 Penalty score information storage unit
- 27 Overall score information storage unit
- 28 Classification information storage unit
- 29 Program storage unit
- 30 Temporary storage unit
- 131 Number of occurrences arithmetic unit
- 132 Binary arithmetic unit
- 141 Local segment size determination unit
- 142 Local centroid arithmetic unit
- 143 Local arithmetic arithmetic unit
- 101 Overall score classification unit
- 102 Classification unit

【特許請求の範囲】**【請求項 1】**

複数のデータにより構成される第 1 データセット内の処理対象データと、検索対象データとしてのクエリとの間の類似度を算出する第 1 算出部と、

前記第 1 データセット内の一部のデータにより構成される第 2 データセットを、前記処理対象データごとに特定する特定部と、

前記特定部が特定する前記第 2 データセットから、前記処理対象データごとの基準を算出する第 2 算出部と、

前記第 1 算出部が算出する前記類似度と、前記第 2 算出部が算出する基準とを用いて、前記処理対象データごとにスコアを算出する第 3 算出部と、

を含む情報処理装置。

10

【請求項 2】

前記クエリについて、前記第 1 データセットから第 1 の所定数のデータを、前記スコアに基づいて抽出する抽出部

をさらに備える請求項 1 に記載の情報処理装置。

【請求項 3】

前記特定部は、前記第 1 データセット内のデータと前記処理対象データとの間の類似度に基づいて、前記第 2 データセットを特定する

請求項 1 又は請求項 2 に記載の情報処理装置。

【請求項 4】

前記特定部は、前記第 1 データセット内のデータと前記処理対象データとの間の類似度が高い順に第 2 の所定数のデータを抽出することにより、前記第 2 データセットを特定する

請求項 1 から請求項 3 のいずれか一項に記載の情報処理装置。

20

【請求項 5】

前記第 2 の所定数とは、前記第 1 データセット内の 2 つのデータの全ての組み合わせにおける類似度に基づいて、前記第 1 データセット内の各々のデータに対して当該類似度が高い順に第 1 の所定数のデータを抽出する場合に、前記処理対象データが抽出される回数と、前記処理対象データと当該処理対象データの基準との間の類似度と、の間の相関が最大になる数である

請求項 4 に記載の情報処理装置。

30

【請求項 6】

前記第 2 の所定数とは、前記第 1 データセット内の 2 つのデータの全ての組み合わせにおける類似度に基づいて、前記第 1 データセット内の各々のデータに対して当該類似度が高い順に第 1 の所定数のデータを抽出する場合に、前記処理対象データが抽出される回数に関する分布の歪度が最小になる数である

請求項 4 に記載の情報処理装置。

【請求項 7】

前記第 1 算出部は、内積と、コサインと、距離と、カーネルとのうちの少なくともいずれか 1 つに基づいて前記類似度を算出する

請求項 1 から請求項 6 のいずれか一項に記載の情報処理装置。

40

【請求項 8】

情報処理装置が、

複数のデータにより構成される第 1 データセット内の処理対象データと、検索対象データとしてのクエリとの間の類似度を算出する第 1 ステップと、

前記第 1 データセット内の一部のデータにより構成される第 2 データセットを、前記処理対象データごとに特定する第 2 ステップと、

前記第 2 ステップにおいて特定した前記第 2 データセットから、前記処理対象データごとの基準を算出する第 3 ステップと、

前記第 1 ステップにおいて算出した前記類似度と、前記第 3 ステップにおいて算出した

50

基準とを用いて、前記処理対象データごとにスコアを算出する第4ステップと、
を含む情報処理方法。

【請求項9】

コンピュータに、
複数のデータにより構成される第1データセット内の処理対象データと、検索対象データとしてのクエリとの間の類似度を算出する第1ステップと、
前記第1データセット内の一部のデータにより構成される第2データセットを、前記処理対象データごとに特定する第2ステップと、
前記第1ステップにおいて特定した前記第2データセットから、前記処理対象データごとの基準を算出する第3ステップと、
前記第1ステップにおいて算出した前記類似度と、前記第3ステップにおいて算出した基準とを用いて、前記処理対象データごとにスコアを算出する第4ステップと、
を実行させるためのプログラム。

10

【請求項10】

コンピュータに、
複数のデータにより構成される第1データセット内の処理対象データと、検索対象データとしてのクエリとの間の類似度を算出する第1ステップと、
前記第1データセット内の一部のデータにより構成される第2データセットを、前記処理対象データごとに特定する第2ステップと、
前記第1ステップにおいて特定した前記第2データセットから、前記処理対象データごとの基準を算出する第3ステップと、
前記第1ステップにおいて算出した前記類似度と、前記第3ステップにおいて算出した基準とを用いて、前記処理対象データごとにスコアを算出する第4ステップと、
を実行させるためのプログラムを記録したコンピュータ読み取り可能な非一時記憶媒体

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、情報処理方法、プログラム、及び非一時記憶媒体に関する。詳しくは、高次元又はノ及び大規模データセットに対するk近傍検索で発生するハブを軽減する類似度演算システム及び類似度演算方法に関する。

30

2015年1月23日に、日本に出願された特願2015-11853号に基づき優先権を主張し、その内容をここに援用する。

【背景技術】

【0002】

k近傍法は実装が簡素であるにもかかわらず分類や情報検索で有効であるために多くの分類システムや情報検索システムで用いられている。しかし、データセット内のデータが高次元空間に存在するとみなせる場合（例えばデータが多数の属性を持つベクトルとして表現される場合）、他のデータのk近傍に頻出するデータ（ハブと呼ばれる）が出現し、結果としてk近傍法の性能は低下する。このハブの現象は、Radovanovicら（非特許文献1参照）により、ごく最近発見されたデータの高次元性にまつわる現象である。一方、発明者らは「（グローバル）センタリング」、すなわち、原点をデータセットの平均（グローバルセントロイド）に移動することにより、k近傍法におけるハブの影響を軽減できることを発表した（非特許文献2参照）。ハブはデータセットの平均の近くに位置するデータであり、センタリングはハブを軽減するのに有効である。

40

【先行技術文献】

【非特許文献】

【0003】

【非特許文献1】Radovanovic, M; Nanopoulos, A.; and Ivanovic, M. 2010a, "Hubs in space: Popular

50

nearest neighbors in high-dimensional data." Journal of Machine Learning Research 11:2487-2531, 2010年

【非特許文献2】鈴木郁美、原一夫、新保仁「k近傍法でハブを軽減する類似度尺度」、情報処理学会研究報告、自然言語処理研究会、2012-NL-209、No.11、pp.1-8、2012、2013年

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、次元 d がさほど大きくななくても、データ数 n が大きくなると、例えば $n = 10,000$ 、 $d = 500$ 、異種のハブが存在することが解った。そして、かかるハブはセンタリングでは軽減できないという問題があった。

本発明は、次元数が大きいときだけでなく、データ数が大きいときに出現するハブを軽減できる情報処理装置、情報処理方法、プログラム、及び非一時記憶媒体を提供することを目的とする。

【課題を解決するための手段】

【0005】

本発明の第1の態様に係るハブを軽減する類似度演算システム1は、例えば図7に示すように、検索先となる高次元及び/又は大規模データセット、及び、検索対象データとしてのクエリ q を設定するデータセット・クエリ設定部11と、データセット・クエリ設定部11にて設定されたクエリ q 及びデータセット内の各データ x について、ベースとする類似度尺度を用いて、クエリ q と各データ x との類似度 $Sim(x; q)$ を演算するベース類似度スコア演算部12と、前記設定されたデータセット内の各データ x について、データセット内の他のデータの k 近傍への出現回数 $N_k(x)$ を計算するハブ度演算部13と、ハブ度演算部13で各データ x に対して計算された出現回数 $N_k(x)$ を用い、各データの近傍に入るデータの集合の要素がデータセット全体を一様にカバーせず局在化する度合いに応じて、各データの近傍に入るデータからなるローカルな部分集合の中心としてのローカルセントロイド $c(x)$ を計算すると共に、各データについてローカルセントロイド $c(x)$ を用いた類似度演算により各データのペナルティスコアとして用いるローカルアフィニティ $LocalAffinity(x)$ を計算することによりローカライズドセンタリングの前工程を実行するペナルティスコア演算部14と、ベース類似度スコア演算部12にて計算されたクエリと各データとの類似度 $Sim(x; q)$ とペナルティスコア演算部14にて計算されたローカルアフィニティ $LocalAffinity(x)$ を用い、各データに対して総合スコアを計算することによりローカライズドセンタリングの後工程を実行する総合スコア演算部15とを備える。

【0006】

ここにおいて、高次元及び/又は大規模データセットとは、次元数が大きい(あるいは次元数が大きいとみなせる)及び/又はデータ数が大きいデータセットをいう。例えば、生命情報学分野では、機能未知の配列データの機能を、機能既知の配列データに対する相同性を頼りに予測しようとするが、検索先として使用される機能既知のDNAの塩基配列あるいはタンパク質のアミノ酸配列のデータベースは、過去数十年間に世界中の研究成果を集めた膨大な数の配列データを格納する、大規模なデータセットである。また、DNAの塩基配列あるいはタンパク質のアミノ酸配列の相同性検索アルゴリズムとして広く使われるブラスト(BLAST: Basic Local Alignment Search Tool)は、個々の配列データを特徴付ける属性として、配列内に出現しうるあらゆる部分配列(及びその組み合わせ)を用いるため、配列データを本質的に高次元データとみなして扱う。ここでの高次元及び/又は大規模データセットとは、ハブが発生するあらゆるデータセットを意味する。なぜなら、現在までに報告されているハブは、データセットが高次元である場合と、データセットが大規模である場合に限られて発生しているからである。敢えて数値化するとすれば例えば、高次元は人間の直感で理解が難しくなる4

10

20

30

40

50

次元以上をいい、大規模はデータ数100以上をいい、1000以上が好ましい。なお、後者のハブは、本発明者らにより発見されたハブである。

【0007】

また、ここでのデータセットとは、検索対象データ(クエリ)に対する類似データの検索先となるデータ集合をいう。例えば、生命情報学分野における配列の相同性検索においては、データセットとして用いられるのは、過去に世界中の研究成果として蓄積された塩基配列データ又はアミノ酸配列データからなる配列データ集合群、又はそれらの部分集合である。

また、ここでの類似度尺度とは、2つのデータの類似性を測る尺度として使用できるものすべてを含む。典型的には、内積、コサイン、距離である。内積は2つのベクトルデータのスカラー積であり、コサインは長さ1に規格化された2つのベクトルデータの内積である。さらに、内積の一般化とみなせる(機械学習分野で主に呼ばれるところの各種の)カーネルも含む。距離の典型は、2つのベクトルデータ間のユークリッド距離(l^2 ノルム)であるが、ユークリッド距離を一般化した距離(マンハッタン距離や l^p ノルムなど)も含む。さらに、ドメインの知識を持つ人間が各タスクの目的に応じて適宜定めた類似度スコア計算方法(BLASTなど)が出力する類似度も、ここでの類似度尺度に含まれる。

【0008】

また、「ベースとする類似度尺度」は、ベース類似度スコア演算部にて類似度 $Sim(x; q)$ を計算するとき、及びペナルティスコア演算部にてローカルアフィニティ $LocalAffinity(x)$ を計算するとき使用される類似度尺度である。類似度 $Sim(x; q)$ を計算するとき及びローカルアフィニティ $LocalAffinity(x)$ を計算するとき使用される類似度尺度として、上記全ての類似度尺度を使用できる。また、高次元及び/又は大規模データセットを検索先とする場合、ベース類似度スコア演算部にて計算されるスコアである類似度 $Sim(x; q)$ はハブを発生させ易い。そこで、ベース類似度スコア演算部にて計算されるスコアをそのまま用いず、ハブに成り易い程度に応じて各データに対してペナルティスコアを計算し、ベースとするスコアからペナルティスコアを差し引いた総合スコアを用いる。すなわち、ハブに成り易いデータの総合スコアを小さくすることで、ハブに成り易いデータを他のデータの k 近傍に入りやすくするものである。

【0009】

また、ハブ度演算部13は、データセット内の各データ x に対し、他のデータの k 近傍に出現する回数 $N_k(x)$ を計算する。さらに、 $N_k(x)$ の分布の歪度を計算する。 k 近傍になるとは、類似度が高い方から k 番目以内になることをいう。また、 $N_k(x)$ が大きいほど、データ x がハブであることを意味する。さらに、 $N_k(x)$ の分布の歪度は、 $S_{N_k} = E\left\{\frac{(N_k - \mu_{N_k})^3}{N_k^3}\right\}$ ($E\{\}$ は期待値を計算するオペレータ、 μ_{N_k} と $\sigma_{N_k}^2$ はそれぞれ $N_k(x)$ の分布の平均と分散である)と計算する。 $N_k(x)$ の分布の歪度が大きいほど、データセットにハブが存在することを意味する。

また、大規模なデータセット、すなわち、データ数 n が大きいデータセットにおいては、各データ x についてその近傍に入るデータの集合を考えたとき、その集合がデータセット全体を一様にカバーせず、データセットに局在する部分集合を形成し易い、という事実がある。そこで、ローカルセグメントサイズをパラメタとし、各データに対してその近傍に入るデータをローカルな部分集合として定め、このローカルな部分集合の中心としてローカルセントロイド $c(x)$ を定義する。また、各データの近傍に入るデータの集合の要素がデータセット全体を一様にカバーせず局在化する度合いに応じて、ローカルセントロイド $c(x)$ を計算するのは、ハブ度演算部13で各データ x に対して計算された出現回数 $N_k(x)$ に応じてペナルティを与えるためである。なお、ローカルセントロイドは各データにより異なるが必ず存在するとみなしてよい。したがって、ローカルセントロイドの有無を判断する必要はない。

具体的には、ローカルセントロイド $c(x)$ は例えば、式(2A)で求められる(段

10

20

30

40

50

落 0 0 3 1 参照)。パラメタにより隣接データの範囲が変わるので、ローカルセントロイド $c(x)$ も変わり得る。ローカライズドセンタリングとは、ハブ度と相関が高くなるようなローカルなセントロイドを見つけ出し、ローカルなセントロイドとの類似度をペナルティとして与えることをいう。例えば、式(5)が使用される(段落 0 0 3 7 参照)。類似度尺度にはベースとする類似度尺度を使用できる。ここでは、ローカライズドセンタリングは、先にローカルアフィニティを計算し、その後総合スコアを計算することにより実行される。すなわち、ベース類似度スコア演算部にて計算されるスコア(ベースとするスコア)から、ハブに成り易い程度に応じて各データに対して付与するペナルティスコア(ローカルアフィニティ)を計算し、ベースとするスコアからペナルティスコアを差し引いた総合スコアを計算することにより実行される。

10

また、ローカルアフィニティとは、各データ $x \in D$ について、 x とその近傍(N_N すなわち $-Nearest-Neighbor$) に属するデータ間の平均類似度として定義される(段落 0 0 3 1、式(2)参照)。また、ペナルティスコアとして用いるとは、ローカルセントロイド $c(x)$ 近傍のデータの類似度スコアを低くするために用いることをいい、ここでは、ローカルアフィニティをペナルティスコアとして差し引いている。

【0 0 1 0】

本態様のように構成すると、ローカルセントロイドという新しい概念に基づいて計算されるローカルアフィニティが、データのハブ度と高い相関を持つため、それをペナルティスコアとして用いることにより、次元数が大きいときだけでなく、データ数 n が大きいときに出現するハブをも軽減できる。なお、既存手法において用いられたグローバルセントロイドでは、次元数 d が大きいときに生じるハブは軽減できるがデータ数 n が大きいときに生じるハブを軽減できなかった。

20

【0 0 1 1】

本発明の第 2 の態様に係るハブを軽減する類似度演算システム 1 は、第 1 の態様において、ベースとする類似度尺度として、内積、コサイン、距離、内積の一般化とみなせるカーネル、又は、ドメインの知識を持つ人間が各タスクの目的に応じて適宜定めた類似度スコア計算方法が出力する類似度尺度を用いる。

ここにおいて、内積の一般化とみなせるカーネルとは、機械学習分野で主に呼ばれるところの各種のカーネルをいう。また、ドメインの知識を持つ人間が各タスクの目的に応じて適宜定めた類似度スコア計算方法が出力する類似度には B L A S T などが含まれる。

30

本態様のように構成すると、内積、コサイン、距離、カーネル、B L A S T スコア等は演算が簡素で扱いやすく、また、ローカルアフィニティをペナルティスコアとして用いるローカライズドセンタリングによりハブを軽減できる実績が得られている。なお、ローカライズドセンタリングによりハブが軽減しさえすれば、その他の類似度尺度も使用できる。

【0 0 1 2】

本発明の第 3 の態様に係るハブを軽減する類似度演算システム 1 は、第 1 ないし第 3 のいずれかの態様において、ペナルティスコア演算部 1 4 は、ローカルアフィニティ $LocalAffinity(x)$ を計算する際のパラメタであるローカルセグメントサイズ k を、出現回数 $N_k(x)$ とローカルアフィニティ $LocalAffinity(x)$ との相関が最大になるように定める。

40

【0 0 1 3】

ここにおいて、ローカルアフィニティは、各データ $x \in D$ について、 x とその近傍(k はローカルセグメントサイズであり、ローカルアフィニティを定めるパラメタである) に属するデータ間の平均類似度として定義され、明細書中の式(2)で表される(段落 0 0 3 1 参照)。

本態様のように構成すると、ローカルセグメントサイズ k を適切に定めることができる。

【0 0 1 4】

50

本発明の第4の態様に係るハブを軽減する類似度演算方法は、例えば図8に示すように、検索先となる高次元及び/又は大規模データセット、及び、検索対象データとしてのクエリを設定を行うデータセット・クエリ設定工程(S101)と、データセット・クエリ設定工程(S101)にて設定されたクエリとデータセット内の各データについて、ベースとする類似度尺度を用いてクエリと各データとの類似度を計算するベース類似度スコア演算工程(S102)と、設定されたデータセット内の各データについて、データセット内の他のデータのk近傍への出現回数 $N_k(x)$ を計算するハブ度演算工程(S103)と、ハブ度演算工程(S103)で各データに対して計算された、出現回数 $N_k(x)$ を用い、各データの近傍に入るデータの集合の要素がデータセット全体を一様にカバーせず局在化する度合いに応じて、各データの近傍に入るデータからなるローカルな部分集合の中心としてのローカルセントロイド $c(x)$ を計算すると共に、各データについてローカルセントロイド $c(x)$ を用いた類似度演算により各データのペナルティスコアとして用いるローカルアフィニティ $Local\ Affinity(x)$ を計算することによりローカライズドセンタリングの前工程を実行するペナルティスコア演算工程(S106)と、ベース類似度スコア演算工程(S102)にて計算されたクエリと各データとの類似度 $Sim(x; q)$ とペナルティスコア演算工程(S106)にて計算されたローカルアフィニティ $Local\ Affinity(x)$ を用い、各データに対して総合スコアを計算することによりローカライズドセンタリングの後工程を実行する総合スコア演算工程(S107)とを備える。

10

20

【0015】

ここにおいて、ローカライズドセンタリングは前工程と後工程からなる。

本態様のように構成すると、(次元数 d が大きいときに生じるハブは軽減できるがデータ数 n が大きいときに生じるハブを軽減できない既存手法で用いられたグローバルセントロイドに対し)本発明が新たに導入するローカルセントロイドという概念に基づいて計算されるローカルアフィニティは各データのハブ度と高い相関を持つため、それをペナルティスコアとして用いることにより、次元数が大きいときだけでなく、データ数 n が大きいときに出現するハブをも軽減できる。

【0016】

本発明の第5の態様に係るハブを軽減する類似度演算方法は、第4の態様において、ベースとする類似度尺度として、内積、コサイン、距離、内積の一般化とみなせるカーネル、又は、ドメインの知識を持つ人間が各タスクの目的に応じて適宜定めた類似度スコア計算方法が出力する類似度これらに基づいて作成された類似度尺度を用いる。

30

このように構成すると、内積、コサイン、距離、カーネル、BLASTスコア等は演算が簡素で扱いやすく、また、ローカルアフィニティをペナルティスコアとして用いるローカライズドセンタリングによりハブを軽減できる実績が得られている。なお、ローカライズドセンタリングによりハブが軽減する限りは、その他の類似度尺度も使用できる。

【0017】

本発明の第6の態様に係るハブを軽減する類似度演算方法は、第4又は第5の態様において、ペナルティスコア演算工程(S106)は、ローカルアフィニティ $Local\ Affinity(x)$ を計算する際のパラメタであるローカルセグメントサイズを、出現回数 $N_k(x)$ とローカルアフィニティ $Local\ Affinity(x)$ との相関が最大になるように定める。

40

このように構成すると、ローカルセグメントサイズを適切に定めることができる。

【0018】

本発明の第7の態様に係るプログラムは、第4ないし第6のいずれかの態様におけるハブを軽減する類似度演算方法をコンピュータに実行させるためのプログラムである。

【0019】

本発明の第8の態様に係る記録媒体は、第7の態様におけるプログラムを記録したコンピュータ読み取り可能な記録媒体である。

【発明の効果】

50

【 0 0 2 0 】

本発明の一態様によれば、次元数 d が大きいときに出現するハブのみならず、データ数 n が大きいときに出現するハブを軽減できる情報処理装置、情報処理方法、プログラム、及び非一時記憶媒体を提供できる。

【 図面の簡単な説明 】

【 0 0 2 1 】

【 図 1 】高次元人工データセット ($d = 4000$) におけるハブを説明するための第 1 図である。図 1 は、センタリング (グローバルセンタリング) 前の N_{10} ($k = 10$ のとき、各データが他のデータの k 近傍に出現する回数) の分布を示す図である。

【 図 2 】高次元人工データセット ($d = 4000$) におけるハブを説明するための第 2 図である。図 2 は、センタリング (グローバルセンタリング) 後の N_{10} ($k = 10$ のとき、各データが他のデータの k 近傍に出現する回数) の分布を示す図である。

【 図 3 】高次元人工データセット ($d = 4000$) におけるハブを説明するための第 3 図である。図 3 は、各データの N_{10} とグローバルアフィニティ (データセットの平均への類似度) との関係 (グローバルセンタリング前) を示す散布図である。

【 図 4 】大規模人工データセット ($n = 10000$) におけるハブを説明するための第 1 図である。図 4 は、センタリング (グローバルセンタリング) 前の N_{10} の分布を示す図である。

【 図 5 】大規模人工データセット ($n = 10000$) におけるハブを説明するための第 2 図である。図 5 は、センタリング (グローバルセンタリング) 後の N_{10} の分布を示す図である。

【 図 6 】大規模人工データセット ($n = 10000$) におけるハブを説明するための第 3 図である。図 6 は、各データの N_{10} とグローバルアフィニティ (データセットの平均への類似度) との関係 (グローバルセンタリング前) を示す散布図である。

【 図 7 】データ数 n 及び次元数 d を変数として、生成した様々な人工データセットを用いて作成した、 N_{10} 分布の歪度の等高線第 1 図である。図 7 は、センタリング (グローバルセンタリング) 前の等高線図である。数値が高いほど、当該データセットにハブが多く存在することを示す。なお、等高線図における + マークは、生成したデータセット (データ数 n と次元数 d) に対応する。

【 図 8 】データ数 n 及び次元数 d を変数として、生成した様々な人工データセットを用いて作成した、 N_{10} 分布の歪度の等高線第 2 図である。図 8 は、センタリング (グローバルセンタリング) 後の等高線図である。数値が高いほど、当該データセットにハブが多く存在することを示す。

【 図 9 】図 7、8 と同じ人工データセットを用いたときの、ローカライズドセンタリング後の N_{10} 分布の歪度の等高線図である。

【 図 10 】図 7、8 と同じ人工データセットを用いたときの、グローバル対ローカルアフィニティ比を示す等高線図である。

【 図 11 】図 4 - 6 と同じ大規模人工データセット ($n = 10000$) を用いたときの、各データの N_{10} とローカルアフィニティ (ローカルセントロイドへの類似度) の相関 (ローカライズドセンタリング前) を示す散布図である。

【 図 12 】図 4 - 6 と同じ大規模人工データセット ($n = 10000$) を用いたときの、ローカライズドセンタリング後の N_{10} 分布を示す図である。

【 図 13 】図 1 - 3 と同じ高次元人工データセット ($d = 4000$) を用いたときの、各データの N_{10} とローカルアフィニティ (ローカルセントロイドへの類似度) の相関を示す散布図である。

【 図 14 】実施例 1 におけるハブを軽減する類似度演算システムの構成例を示す図である。

【 図 15 】実施例 1 におけるハブを軽減する類似度演算方法の処理フロー例を示す図である。

【 図 16 】実世界データセットに対する、データ数 n と N_{10} 分布の歪度との関係を示す

10

20

30

40

50

折れ線グラフである。(i)はグローバルセンタリングやローカライズドセンタリングを施す前の折れ線グラフ(比較のベースライン)である。(ii)はグローバルセンタリング後、(iii)はローカライズドセンタリング後の折れ線グラフである。(iv)はグローバル対ローカルアフィニティ比の折れ線グラフである。使用データセットはWebKBである。

【図17】実世界データセットに対する、データ数 n と N_{10} 分布の歪度との関係を示す折れ線グラフである。(i)はグローバルセンタリングやローカライズドセンタリングを施す前の折れ線グラフ(比較のベースライン)である。(ii)はグローバルセンタリング後、(iii)はローカライズドセンタリング後の折れ線グラフである。(iv)はグローバル対ローカルアフィニティ比の折れ線グラフである。使用データセットはReuters-52である。

10

【図18】実世界データセットに対する、データ数 n と N_{10} 分布の歪度との関係を示す折れ線グラフである。(i)はグローバルセンタリングやローカライズドセンタリングを施す前の折れ線グラフ(比較のベースライン)である。(ii)はグローバルセンタリング後、(iii)はローカライズドセンタリング後の折れ線グラフである。(iv)はグローバル対ローカルアフィニティ比の折れ線グラフである。使用データセットはTD2-30である。

【図19】実世界データセットに対する、データ数 n と N_{10} 分布の歪度との関係を示す折れ線グラフである。(i)はグローバルセンタリングやローカライズドセンタリングを施す前の折れ線グラフ(比較のベースライン)である。(ii)はグローバルセンタリング後、(iii)はローカライズドセンタリング後の折れ線グラフである。(iv)はグローバル対ローカルアフィニティ比の折れ線グラフである。使用データセットは20Newsgroupsである。

20

【図20】図16-19で使用したのと同じ実世界データセットを用いて行った、 k 近傍法による多クラス分類の精度、及び、 N_{10} 分布の歪度を示す表形式の第1図である。

【図21】図16-19で使用したのと同じ実世界データセットを用いて行った、 k 近傍法による多クラス分類の精度、及び、 N_{10} 分布の歪度を示す表形式の第2図である。

【図22】図16-19で使用したのと同じ実世界データセットを用いて行った、 k 近傍法による多クラス分類の精度、及び、 N_{10} 分布の歪度を示す表形式の第3図である。

【図23】図16-19で使用したのと同じ実世界データセットを用いて行った、 k 近傍法による多クラス分類の精度、及び、 N_{10} 分布の歪度を示す表形式の第4図である。

30

【発明を実施するための形態】

【0022】

図面を参照して以下に本発明の実施の形態について説明する。なお、各図において、互いに同一又は相当する部分には同一符号を付し、重複した説明は省略する。

【0023】

〔高次元データセット中でのハブ〕

ハブは高次元空間での最近傍検索に係る現象として知られている。 $D \subset R^d$ を d 次元空間上のデータセットとする。 $N_k(x)$ は、データ $x \in D$ が D に含まれる他データの k 近傍に含まれる回数を示すものとする。次元 d が増加すると、 N_k の分布の形状は右に長い尾を引くように変化する。すなわち、少数のデータが予期しない高い N_k 値を持つ。このようなデータをハブという。ここで、テキスト文書集合を模した人工データセット(各データが一つのテキスト文書に対応する)を用いて、ハブの出現を例示する。テキスト文書集合を模したデータセット D を生成するに際し、テキスト文書の表現として通常用いられるバッグオブワーズ「bag-of-words」モデルによって、非負値を要素として持つ d 次元ベクトルを、各テキスト文書を模したベクトルデータとして生成する。 $n = |D|$ を生成するデータ数とするとき、まず n 個のベクトルを生成し、これらをゼロベクトルとして初期化する。次に、各次元 $i = 1, \dots, d$ について、対数正規分布 $\text{LogNormal}(5, 1)$ (5, 1はそれぞれ、平均、分散に係るパラメタ)に従う実数を生成し、端数を丸めた整数 n_i を得る。そして、データセット D の n 個のベクトルから一様分布

40

50

に従って n_i 個のベクトルを選び、選ばれたベクトルの第 i 成分を 0 でない値に置換する。具体的には、 $[0, 1]$ の一様分布に従う値 ($0 \sim 1$ の間の実数) に置換する。最後に、全てのベクトルを長さが 1 になるように正規化する。ここでは、一例として、類似度尺度として内積を使用するが、全てのベクトルの長さが 1 に等しいため、内積はコサインに等しい。なお、上述したように、類似度尺度としては、内積、コサイン、距離等の別の尺度をベースとしてもよく、いずれの尺度を用いるかは、ユーザにより選択可能としてもよい。

【0024】

図 1 - 3 は高次元データセット ($d = 4000$) におけるハブを説明するための図である。図 1、2 はそれぞれセンタリング (グローバルセンタリング) 前後の N_{10} ($k = 10$ のとき、各々のデータが他のデータの k 近傍に出現する回数) の分布を示す図である。図において、ハブ、すなわち、極端に大きな N_{10} 値を持つデータの存在を確認できる。

データセットにどの程度ハブが存在するかを、Radovanovicら (Radovanovic, M.; Nanopoulos, A.; and Ivanovic, M. 2010a, Hubs in space: Popular nearest neighbors in high-dimensional data. Journal of Machine Learning Research 11: 2487 - 2531) に従って、 N_k 分布の歪度 $S_{N_k} = E \{ (N_k - \mu_{N_k})^3 / N_k^3 \}$ によって評価する。ここに $E \{ \}$ は期待値を計算するオペレータ、 μ_{N_k} と N_k はそれぞれ N_k 分布の平均と分散である。歪度は分布の対称性の程度を測るための標準尺度である。その値はガウス分布のような対称な分布では 0 であり、長い右又は左の尾を持つ分布については正又は負である。図 1 の人工データセットにおいては、 $S_{N_{10}}$ は大きな正值 4.45 である。

図 3 は N_{10} とデータセットのグローバルセントロイド (データセットの平均) への類似度 (すなわち、後述するグローバルアフィニティ) との関係 (グローバルセンタリング前) を示す散布図である。図が示すように、それらの間に強い相関が存在する。このように、高次元データセットでは、グローバルセントロイドに類似するデータがハブになる。更なる例としては上記 Radovanovicらを参照されたい。

【0025】

〔ハブ軽減法としてのセンタリング〕

グローバルセンタリング (単にセンタリングともいう) とは、特徴空間の原点をデータセットのグローバルセントロイド (平均) に移動する変換をいう。グローバルセントロイド c は次式で表される。

$$c = (1/D) \sum_{x \in D} x$$

センタリングは、データセットの観測バイアスを除去する古典的技術であるが、最近、発明者らによって、ハブを除く効果を持つことが確認された (Suzuki, I.; Hara, K.; Shimbo, M.; Saerens, M.; and Fukumizu, K. 2013. "Centering similarity measures to reduce hubs." In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, 613 - 623.) センタリングの重要な意義は、データとグローバルセントロイド c 間の類似度 (内積) を一定 (実際にはゼロ) にすることである。具体的には、与えられたクエリ $q \in R^d$ とデータ $x \in D$ の内積

$$\text{Sim}(x, q) = \langle x, q \rangle$$

は、センタリング後、次の様に変化する。

$$\text{Sim}^{\text{CENT}}(x, q) = \langle x - c, q - c \rangle$$

ここで、 $q = c$ と置換すると、任意の $x \in R^d$ に対して、

$$\text{Sim}^{\text{CENT}}(x, c) = 0$$

となる。

10

20

30

40

50

【0026】

このことは、グローバルセントロイドに特別に高い類似度を持つデータはもはや存在しないことを意味する。ここで期待されるのは、図1と図3で使用したデータセット(グローバルセントロイドに類似するデータがハブとなっていた)に対する、センタリングによるハブ軽減の効果である。期待されたとおり、このデータセットについての N_{10} 分布の歪度 $S_{N_{10}}$ は、センタリング前の4.45から、センタリング後に0.27に減少する。実際、センタリング後の図2で見られる N_{10} 分布はほぼ対称である。

【0027】

〔大規模データセット中でのハブ〕

しかしながら、(グローバル)センタリングではハブが軽減されない場合がある。さほど高次元ではないが、データ数が大きい次元数 $d = 500$ 、データ数 $n = 10000$ として前記と同様に人工データセット D を生成し、この人工データセットに出現するハブについて以下に説明する。

10

【0028】

図4-6は上記データセットにおけるハブを説明するための図である。図4、5はそれぞれ(グローバル)センタリング前後の N_{10} ($k = 10$ のとき、各々のデータが他のデータの k 近傍に出現する回数)の分布を示す図である。図6は N_{10} とグローバルセントロイドへの類似度(すなわち、後述するグローバルアフィニティ)との関係(センタリング前)を示す散布図である。

このデータセットでは N_{10} 分布の歪度 $S_{N_{10}}$ はセンタリング前に5.88、センタリング後に5.82である。つまり、センタリングによりハブはさほど軽減されない。図4及び図5に示されるように、 N_{10} 分布の形状はセンタリング前後でほぼ同一である。さらに、図6の散布図に示されるように、このデータセットにおいてはセンタリング法が成功するために必要とされる条件(大きな N_k 値を持つデータ、すなわち、ハブはグローバルセントロイドと高い類似度を持つ)を満たさない。

20

【0029】

データセットを、データ数 n を100~10000、次元数 d を100~4000の範囲で変えて生成し、各データセットについて、 N_{10} 分布の歪度を計算する。各 n と d の組み合わせについて、10個のデータセットを生成し、平均の歪度を計算する。

【0030】

図7、8は、データ数 n 及び次元数 d を変数とした N_{10} 分布の歪度の等高線図である。図7はセンタリング前、図8はセンタリング後の等高線図である。高い値はハブの存在を意味する。+マークは生成したデータセット(のデータ数 n と次元数 d)に対応する。図7を見ると、ハブは左上と右下の領域のデータセットで発生していることが分かる。他方、図8によれば、左上の領域のハブはセンタリング後に消え、右下の領域のハブはセンタリングしても消えず、センタリング前と殆ど同程度残ることが分かる。

30

何故、右下の領域のデータセットでは、(グローバル)センタリングによってハブが軽減されないのかを調査するために、各データの特徴を調べるための2つの量、ローカルアフィニティとグローバルアフィニティを導入する。

【0031】

ローカルアフィニティ $Local\ Affinity(x)$ は、各データ $x \in D$ について、 x とその近傍(NN すなわち $- Nearest - Neighbor$)に属するデータ間の平均類似度として定義され、次のように計算される。

40

$$Local\ Affinity(x) = \frac{1}{|NN(x)|} \sum_{x' \in NN(x)} \langle x, x' \rangle = \langle x, c(x) \rangle \quad (2)$$

ここに、ローカルセグメントサイズ k は、各データ x に対してローカルセントロイドを計算するために使用する隣接データの範囲を定めるパラメタである。

そして、次式が定める $c(x)$ を、データ x のローカルセントロイドと呼ぶ(式において、 kNN ではなく NN が使われることに注意)。

$$c(x) = \frac{1}{|NN(x)|} \sum_{x' \in NN(x)} x' \quad (2A)$$

50

【0032】

グローバルアフィニティ $Global\ Affinity(x)$ は、各データ $x \in D$ について、 x と D 中の全てのデータ間の平均類似度として定義される。

$$Global\ Affinity(x)$$

$$(1/D) \sum_{x' \in D} \langle x, x' \rangle = \langle x, c \rangle \quad (3)$$

ここに、データセットのグローバルセントロイド（データセットの平均）は、

$$c = (1/D) \sum_{x' \in D} x'$$

上式から解るように、 x のグローバルアフィニティは単純に x と c の類似度である。もし、 D が非負値を要素として持つベクトルデータの集合であるなら、全ての $x \in D$ について、

$$Local\ Affinity(x) \geq Global\ Affinity(x)$$

である。ここで、例えばテキスト文書は通常、(tf-idf重み付け)ワードカウントベクトル、すなわち非負値を要素として持つベクトルとして表現されるため、上記不等式が成り立つ。さらに、 $D = \{x\}$ なら、次式が成立する。

$$Local\ Affinity(x) = Global\ Affinity(x)$$

【0033】

次に、グローバル対ローカルアフィニティ比

$$Local\ Affinity(x) / Global\ Affinity(x)$$

について説明する。

各データ x が非負ベクトルとして表されるデータセット D においては、あらゆるデータ $x \in D$ に対して、この比は $[0, 1]$ に入る。この比は、 D 中でデータ x がどの程度ローカル化しているかの指標となる。すなわち、(ある $x \in D$ に対して)比が 0 に近いなら、 x は近傍データと特別に高い類似度を持つ、つまり x はローカル化されている局在データの集合を持っている。比が 1 に関わりなく 1 に近いなら、 x は特別に高い類似度を持つデータを D 中に持たない、つまり x はローカル化されていない。

【0034】

図 9、10 はローカライズドセンタリングに係る図(その 1)である。データ数と次元数を変えた(図 7 - 8 に使用したのと同じ)様々なデータセットを用いて作成された。図 9 はローカライズドセンタリング後の N_{10} 分布の歪度を示す等高線図である。図 10 はグローバル対ローカルアフィニティ比を示す等高線図である。

ローカライズドセンタリングとは、ハブ度と相関が高くなるようなローカルなセントロイドを見つけ出し、ローカルなセントロイドとの類似度をペナルティとして与えることをいう。ローカライズドセンタリングは、例えば、先にローカルアフィニティを計算し、その後総合スコアを計算することにより実行される。

グローバル対ローカルアフィニティ比は個々のデータに対して定義されるので、データセットに対するグローバル対ローカルアフィニティ比として、データセット中の全データのグローバル対ローカルアフィニティ比の平均を用いた(ローカルセグメントサイズは $k=20$ に固定)。これを図 8 と比較すると、(グローバル)センタリング後に歪度が軽減されず大きいままである右下の領域は、図 10 ではグローバル対ローカルアフィニティ比が小さい領域に対応する。このことは、この領域のデータセットにおいてはデータがローカル化していることを示す。このようなデータセットにおいては、ローカルアフィニティがグローバルアフィニティよりハブを軽減するために関連性が高いということを示唆する。

【0035】

図 11、12 はローカライズドセンタリングに係る図(その 2)である。図 4 - 6 で使用した大規模データセット ($n=10000$) について、図 11 は N_{10} とローカルアフィニティの相関(ローカライズドセンタリング前)を示す散布図である。図 12 はローカライズドセンタリングを適用した後の N_{10} の分布を示す図である。

図 11 が示すように、 N_{10} とローカルアフィニティ(すなわち、ローカルセントロイドへの類似度)との間に強い相関が観察された。これは図 6 で示した N_{10} とグローバル

10

20

30

40

50

アフィニティ（すなわちグローバルセントロイドへの類似度）との間の弱い相関と対照的である。そして、図12が示すように、ローカライズドセンタリング後に N_{10} の分布の歪度が小さくなること（ハブが軽減されること）が観察された。

【0036】

図13は、図1-3で用いた高次元データセットについて、 N_{10} と、ローカルアフィニティの相関を示す散布図である。図13を見ると、高次元データセット（ $d = 4000$ ）においても N_{10} はローカルアフィニティと強く相関する。これらの結果は、ローカルアフィニティは、大規模データセットだけでなく高次元データセットにおいても、ハブを軽減するための指標となり得ることを示す。

【0037】

〔提案方法：ローカライズドセンタリング〕

以上の所見から、ローカルアフィニティはハブを軽減するための指標となり得る。そこで、ローカルアフィニティに基づくハブ軽減の新たな方法「ローカライズドセンタリング」を提案する。それは、高次元データセットと大規模データセットの両方に効果がある。

ローカライズドセンタリングを、グローバルセンタリングによる式(1)とのアナロジーで説明する。データセット D に属するデータ x とクエリ q との類似度を計算する上で、 x に依存しない項は、異なる x に対して一定である。よって、式(1)から、次式が導かれる。

$$\text{Sim}^{\text{CENT}}(x; q) \\ \langle x - c, q - c \rangle = \langle x, q \rangle - \langle x, c \rangle + \text{constant} \quad (4)$$

言い換えれば、センタリング法では、もともとの類似度スコア $\langle x, q \rangle$ から、グローバルアフィニティ $\langle x, c \rangle$ が、ペナルティ項として差し引かれる。同様にローカライズドセンタリングと称する提案方法では、もともとの類似度スコアから x のローカルアフィニティをペナルティ項として差し引く。

$$\text{Sim}^{\text{LCENT}}(x; q) = \langle x, q \rangle - \langle x, c(x) \rangle \quad (5)$$

【0038】

式(5)はローカルセグメントサイズを定めるパラメタを含む。パラメタは、例えば、 $N_k(x)$ とローカルアフィニティ $\langle x, c(x) \rangle$ との相関が最大になるように定める。もし $k = D$ であれば、提案方法はグローバルセンタリングと同一になる。式(5)への変換後、任意のデータ $x \in D$ はそのローカルセントロイド $c(x)$ と一定の類似度（実際はゼロ）を持つ。その理由は、式(5)に $q = c(x)$ を代入すると、次式が導かれるからである。

$$\text{Sim}^{\text{LCENT}}(x; c(x)) = 0$$

言い換えれば、変換後にはローカルセントロイドとの類似度が特別高くなるデータは存在しない。ローカルセントロイドと高い類似度を持つデータがハブになるという、前記の所見を考慮すれば、この変換がハブを軽減することが期待される。実際、この期待は図12のデータセットに対しては成立している。

【0039】

ローカライズドセンタリングはデータ数 n と次元数 d を変えた他のデータセットでもハブを軽減する。図9にローカライズドセンタリング後の N_{10} 分布の歪度の等高線図を示す。図より、ローカライズドセンタリングは両タイプのハブ、すなわち、大規模データセットに生じるハブ（等高線図の右下の領域に対応する）と、高次元データセットに生じるハブ（左上の領域に対応する）の、どちらのハブも軽減できることが分かる。高次元データセットに対しても、ローカライズドセンタリング（式(5)）はグローバルセンタリング（式(4)）と同様に効果的である。このようなデータセットでは、図10の左上領域に示したように、グローバル対ローカルアフィニティ比は1に近いからである。このことは、この領域のデータセットの殆どのデータについて、 $\langle x, c(x) \rangle \approx \langle x, c \rangle$ を意味し、したがって、式(5)は式(4)と殆ど同じになる。

式(5)はさらに拡張できる。式(5)の右辺の第2項はペナルティ項として説明できる、すなわち、ペナルティ項は x がどの程度のハブであるかを表し、それに応じて類似度

10

20

30

40

50

スコアを小さくするために使用される。式(5)の拡張は、たとえば式(6)あるいは(7)のように、ペナルティの程度をコントロールするためのパラメタの導入により行う。

$$\text{Sim}^{LCE NT}(x; q) = \langle x, q \rangle - \langle x, c(x) \rangle \quad (6)$$

$$\text{Sim}^{LCE NT}(x; q) = \langle x, q \rangle - \langle x, c(x) \rangle \quad (7)$$

パラメタは、 N_k 分布の歪度がもっとも小さくなるように選択すればよい。

【実施例1】

【0040】

〔システム構成〕

図14に本実施例における類似度演算システム1(情報処理装置の一例)のシステム構成を例示する。類似度演算システム1は、データセット・クエリ設定部11、ベース類似度スコア演算部12、ハブ度演算部13、ペナルティスコア演算部14、総合スコア演算部15、演算結果整理・分類部16、制御部17、表示部18、入出力部19、記憶部20を備える。

10

データセット・クエリ設定部11は検索対象データ(クエリ)とデータセットの設定を行う。本実施例で扱うデータセットは、例えば、テキスト文書データセットや、DNAの塩基配列あるいはタンパク質のアミノ酸配列のデータセット等で、高次元又はノ及び大規模データを有する。操作者が例えば、コンピュータ10の表示部18の設定画面のクエリ入力欄とデータセット入力欄にクエリとデータセット名を入出力部19のマウスやキーボードを用いて入力し、類似度演算システム1が類似度演算に用いるクエリとデータセットであると認識するように、検索対象データ(クエリ)とデータセットを設定する。

20

ベース類似度スコア演算部12(第1算出部の一例)は、データセット・クエリ設定部11にて設定されたクエリと設定されたデータセット内の各データについて、ベースとする類似度尺度を用いてクエリとの類似度を演算する。ベースとする類似度尺度として、内積、コサイン、距離等を使用できる。内積、コサイン、距離と異なり、データが明にベクトル表現されていることを前提としない類似度尺度であっても使用できる。例えば、塩基配列間の類似度尺度であるBLASTスコアや、文字列間の類似度尺度となるストリングカーネル、グラフ間の類似度尺度となるグラフカーネルなどを、ベースとする類似度尺度として使用できる。

【0041】

30

ハブ度演算部13はハブ度として出現回数及び歪度を演算する。すなわち、データセット内の各データxについて、他データのk近傍(ベースとする類似度尺度により決定される)への出現回数 $N_k(x)$ を演算する出現回数演算部131、出現回数の分布の歪度を演算する歪度演算部132を有する。分布の歪度は、 $S_{N_k} = E[(N_k - \mu_{N_k})^3 / N_k^3]$ ($E[\]$ は期待値を計算するオペレータ、 μ_{N_k} と N_k はそれぞれ N_k 分布の平均と分散である)で表される。

【0042】

ペナルティスコア演算部14(特定部、第2算出部の一例)は、各データxのローカルセントロイド $c(x) = (1/N_N(x)) \sum_{N_N(x)} x'$ の計算に用いる隣接データの範囲であるローカルセグメントサイズを定めるパラメタを決定するローカルセグメントサイズ決定部141、各データxのローカルセントロイドを計算するローカルセントロイド演算部142を有する。ローカルセグメントサイズは、例えば、ローカルアフィニティ、すなわち、ローカルセントロイドとの類似度 $\langle x, c(x) \rangle$ と、他データのk近傍に含まれる回数 $N_k(x)$ との相関が最大となるようにパラメタを定める。隣接データの範囲とは出現回数 $N_k(x)$ の大きいデータの近傍データが局在化されている範囲で、ローカルセグメントサイズ内をいう。なお、 α が固定されていれば、ローカルセントロイドは各データに対して1つだけ存在する。なお、ローカルセントロイドは各データにより異なるが必ず存在するとみなしてよい。したがって、ローカルセントロイドの有無を判断する必要はない。このように、ペナルティスコア演算部14は、各データxごとに隣接データを特定し、各データの基準座標ともいえるローカルセントロイドを算出する。

40

50

ペナルティスコア演算部 1 4 は、さらに、各データ x に対し、ローカルセントロイド演算部 1 4 2 で演算されたローカルセントロイド $c(x)$ を用い、ローカルアフィニティ、すなわち、ローカルセントロイドとの類似度 $\langle x, c(x) \rangle$ を、各データ x のペナルティスコアとして演算する、ローカルアフィニティ演算部 1 4 3 を有する。ローカルアフィニティの演算はローカライズドセンタリングの前工程を形成する。

総合スコア演算部 1 5 (第 3 算出部の一例) は、各データ x に対し、ベース類似度スコア演算部 1 2 においてベースとする類似度尺度を用いて演算されたクエリとの類似度から、ローカルアフィニティ演算部 1 4 3 で演算されたペナルティスコアを差し引きすることにより、総合スコアを演算する。このように、総合スコア演算部 1 5 は、ベース類似度スコア演算部 1 2 が演算する類似度と、ペナルティスコア演算部 1 4 が算出するローカルセントロイドとに基づいて、各データ x ごとに総合スコアを算出する。総合スコアの演算はローカライズドセンタリングの後工程を形成する。

10

【 0 0 4 3 】

演算結果整理・分類部 1 6 (抽出部の一例) は、各データ x に対して総合スコア演算部 1 5 で計算した総合スコアを、大きい順に整理する総合スコア順整理部 1 6 1 と、この整理結果を用いて k 近傍法により分類する分類部 1 6 2 を有する。分類は、総合スコアが大きい k 個のデータの分類ラベルを用いた (スコア重み付き) 多数決により行う。つまり、分類部 1 6 2 は、クエリについて、総合スコアが大きい順にデータセットから k 個のデータを抽出し、分類を行う。これにより、類似度演算システム 1 は、クエリに対応する分類ラベルを特定することができる。

20

【 0 0 4 4 】

制御部 1 7 は類似度演算システム 1 及びその各部 1 1 ~ 1 6 を制御して、類似度演算システム 1 としての機能を発揮させる。

以上のデータセット・クエリ設定部 1 1 から制御部 1 7 はコンピュータ 1 0 内に実現可能であり、本実施例でもそのようにしている。ただし、ベース類似度スコア演算部 1 2 では、BLAST のような汎用のアルゴリズムが提供する類似度演算を利用しても良い。この場合でも、ベース類似度スコア演算部 1 2 は演算結果を取り込み、記憶部 2 0 に記憶させる機能を保有する。

表示部 1 8 は、類似度演算システム 1 各部 1 1 ~ 1 6 の入力結果、演算結果、記憶部 2 0 の記憶内容、入出力操作時の案内画面等を表示する。入出力部 1 9 は、キーボード、マウス、マイクロフォン等の入力機器、プリンタ、スピーカー等の出力機器を有し、操作者はこれら入出力機器を用いて適宜入出力できる。

30

【 0 0 4 5 】

記憶部 2 0 は、類似度演算システム 1 各部 1 1 ~ 1 6 の入力結果、演算結果、演算途中の情報、類似度演算システム 1 を動作させるのに必要なプログラム及び情報を記憶する。記憶部 2 0 は、データセット・クエリ情報記憶部 2 1、ベース類似度スコア情報記憶部 2 2、出現回数・歪度情報記憶部 2 3、ローカルセグメントサイズ情報記憶部 2 4、ローカルセントロイド情報記憶部 2 5、ペナルティスコア情報記憶部 2 6、総合スコア情報記憶部 2 7、分類情報記憶部 2 8、プログラム記憶部 2 9 及び一時記憶部 3 0 を有する。

40

【 0 0 4 6 】

データセット・クエリ情報記憶部 2 1 は、検索先データセット、検索対象データ (クエリ) に関して、データセット・クエリ設定部 1 1 による設定結果を記憶する。ベース類似度スコア情報記憶部 2 2 は、ベース類似度として用いる類似度の計算アルゴリズム (内積、コサイン、距離等の計算アルゴリズムや、BLAST のような汎用の類似度計算アルゴリズムを含む) や、ベース類似度スコア演算部 1 2 による類似度演算結果等を記憶する。出現回数・歪度情報記憶部 2 3 は、ハブ度演算部 1 3 での演算結果、すなわち、出現回数演算部 1 3 1 により演算された各データが他のデータの k 近傍に出現する回数 (各データのハブ度)、及び、歪度演算部 1 3 2 により演算された出現回数分布の歪度 (データセットのハブ度) を記憶する。

【 0 0 4 7 】

50

ローカルセグメントサイズ情報記憶部 2 4 は、ローカルセグメントサイズ決定部 1 4 1 で決定されたローカルセグメントサイズ を、ローカルセントロイド情報記憶部 2 5 は、ローカルセントロイド演算部 1 4 2 で演算された各データのローカルセントロイドに関する情報を記憶する。ペナルティスコア情報記憶部 2 6 は、ローカルアフィニティ演算部 1 4 3 で演算された各データのペナルティスコアを記憶する。総合スコア情報記憶部 2 7 は、総合スコア順整理部 1 6 1 による総合スコア順の整理結果等を記憶する。分類情報記憶部 2 8 は、分類部 1 6 2 による k 近傍法を用いたクエリの分類結果を記憶する。プログラム記憶部 2 9 は、類似度演算システム 1 及びその各部 1 1 ~ 1 6 のデータの流れ、演算、分類を実行させるためのプログラム、類似度演算システム 1 及びその各部 1 1 ~ 1 6 を制御するために必要なプログラム、類似度演算システム 1 の起動に必要なプログラム等を記憶する。一時記憶部 3 0 は、演算途中の情報等を一時的に記憶する。

10

【 0 0 4 8 】

〔 処理フロー 〕

図 1 5 に本実施例における類似度演算の処理フロー例を示す。まず、データセット・クエリ設定部 1 1 にて高次元又は / 及び大規模データについて、検索対象クエリと検索先データセットを設定する（データセット・クエリ設定工程：S 1 0 1）。例えば、生命情報学分野における配列の相同性検索では、世界中の研究成果として蓄積された機能既知の DNA の塩基配列あるいはタンパク質のアミノ酸配列のデータセット、及び、検索対象クエリとなる機能未知の塩基配列やアミノ酸配列が設定され、データセット・クエリ情報記憶部 2 1 に記憶される。次に、データセット・クエリ設定工程（S 1 0 1）にて設定されたクエリとデータセット内の各データについて、ベース類似度スコア演算部 1 2 にてベースとする類似度尺度を用いてクエリとデータの類似度を演算する（ベース類似度スコア演算工程：S 1 0 2）。ベースとする類似度尺度として、生命情報学分野における配列の相同性検索では、BLAST アルゴリズムが出力するスコアが用いられるが、一般には内積、距離、コサイン等が使用される。演算されたベース類似度スコアはベース類似度スコア情報記憶部 2 2 に記憶される。

20

【 0 0 4 9 】

次に、ハブ度演算部 1 3 にて、各データのハブ度、及び、データセットのハブ度を演算する（ハブ度演算工程：S 1 0 3）。すなわち、出現回数演算部 1 3 1 では、各データ x のハブ度を、データセット内の他のデータの k 近傍への出現回数 $N_k(x)$ を演算する。歪度演算部 1 3 2 では、 $N_k(x)$ の分布の歪度を演算する。出現回数演算部 1 3 1 及び歪度演算部 1 3 2 により作成された出現回数分布図及び歪度（図 1、2、4、5、1 2 参照）が参照される。ハブが存在すれば、出現回数分布図に大きい $N_k(x)$ を持つデータが現われ、大きい歪度が観測される。出現回数 $N_k(x)$ 及び歪度は、出現回数・歪度情報記憶部 2 3 に記憶される。

30

【 0 0 5 0 】

次に、ローカルセグメントサイズ決定部 1 4 1 にて、ローカルセグメントサイズ を決定し（ローカルセグメントサイズ決定工程：S 1 0 4）、 の値はローカルセグメントサイズ情報記憶部 2 4 に記憶する。そして、決定された を使い、ローカルセントロイド演算部 1 4 2 にて、各データ x のローカルセントロイド $c(x)$ を演算し（ローカルセントロイド演算工程：S 1 0 5）、 $c(x)$ はローカルセントロイド情報記憶部 2 5 に記憶する。さらに、演算された $c(x)$ を使い、ローカルアフィニティ演算部 1 4 3 にて、各データ x のローカルアフィニティを $c(x)$ と x の類似度として計算し（ペナルティスコア演算工程：S 1 0 6）、これを各データのペナルティスコアとしてペナルティスコア情報記憶部 2 6 に記憶する。ローカルアフィニティの演算はローカライズドセンタリングの前工程を形成する。

40

【 0 0 5 1 】

次に、総合スコア演算部 1 5 にて、ベース類似度スコア情報記憶部 2 2 及びペナルティスコア情報記憶部 2 6 に記憶された各データに対するベース類似度スコア及びペナルティスコアを用い、各データに対する総合スコアを演算する（総合スコア演算工程：S 1 0 7

50

）。ここでは、ベース類似度スコアからペナルティスコアを差し引くことにより、総合スコアを得る。総合スコアの演算はローカライズドセンタリングの後工程を形成する。ローカライズドセンタリングは前工程と後工程からなる。データセットのデータは、総合スコア演算工程（S107）での演算結果に基づいて、演算結果整理・分類部16にて整理・分類される（演算結果整理工程：S108）。典型的には、総合スコアの大きい順に整理（ランク付け）され、総合スコア情報記憶部27にランキングが記憶される。さらに、検索対象クエリを分類する場合には、k近傍法を用いて（総合スコアが大きいk個のデータの分類ラベルの投票により）検索対象クエリの分類ラベルの予測を行う（分類工程：S109）。予測結果は、分類情報記憶部28に記憶される。

【0052】

〔実世界データを用いた実験〕

本発明による手法（ローカライズドセンタリング）は、（人工データではなく）実世界データセットに対してハブを軽減する効果を持つかどうか、さらに、ハブを軽減する結果としてタスクの精度を向上するかどうかを調べるために、多クラス分類タスクのベンチマークとして用いられる実世界のテキスト文書データセットに対してローカライズドセンタリングを適用する。使用するテキスト文書データセットは、図16はWebKB、図17はReuters-52、図18はTDT2-30、図19は20Newsgroupsである。実験を通して、テキスト文書の分類タスクにおいて標準的なデータ表現であるtf-idf重み付きbag-of-wordsベクトル表現を用い、標準的な類似度尺度（ベースライン）としてコサインを用いる。

【0053】

〔ハブの軽減〕

最初の実験では、ローカライズドセンタリングが大規模データセットに生じるハブを軽減するかどうかを調べた。データセット全体からn個のデータをランダムに選択し、データ数nのサブセットを生成した。データ数nを100から4000の範囲で動かし、各nについてサブセットの生成を10回行った。各サブセットに対して N_{10} 分布の歪度を計算し、各nについて平均歪度（10回の平均）を計算した。そして、グローバルセンタリング（センタリング）後及びローカライズドセンタリング後の平均歪度を比較した。

【0054】

図16-19は、データ数nと N_{10} 分布の平均歪度との関係を示す図である。（i）はグローバルセンタリング又はローカライズドセンタリングを施す前の折れ線グラフ（比較のベースライン）である。（ii）はグローバルセンタリング後、（iii）はローカライズドセンタリングの折れ線グラフである。使用した4つのデータセットを通して、同じ傾向が観察された。すなわち、データ数nの増加につれて、ベースラインでは N_{10} 分布の平均歪度は増加し、ハブが増大する。一方、グローバルセンタリングによって、平均歪度はやはり増加するが、平均歪度はベースラインよりやや小さい。他方、ローカライズドセンタリングでは、データ数nに関わらず、平均歪度はほぼ同じ小さい値のままで変化しない。まとめると、グローバルセンタリングはデータ数nが小さい（約500より小さい）ときに効果が見られるが、データ数nが大きいときには否である。他方、ローカライズドセンタリングは、データ数nが大きくてもハブの発生を低いレベルに維持する。

さらに、式（2）及び式（3）を用い、各データセットについてグローバル対ローカルアフィニティ比を計算し、図16-19にその折れ線グラフを（iv）として示した。すると、人工データセットで観察されたのと同じ傾向が明確に観察された。すなわち、グローバル対ローカルアフィニティ比が小さいとき、グローバルセンタリングはハブを軽減しないが、ローカライズドセンタリングはハブを軽減する。

【0055】

〔k近傍分類の精度〕

このように、ローカライズドセンタリングによりハブが軽減することは確認がされたが、次の実験では、このハブ軽減がk近傍分類精度を改善するかどうかを調べた。タスクはテストデータとしてのテキスト文書の分類ラベルを訓練データとしてのテキスト文書の分

10

20

30

40

50

類ラベルを用いて予測することである。リーブワンナウト leave - one - out 交差検証法（全データから一つをテストデータとして、残りを訓練データとして用い、訓練データを使ってテストデータの分類ラベルを予測することを、全データに亘ってデータ数だけ繰り返し、予測の正解率を評価値とする）によりパフォーマンスを評価した。

【0056】

ベースラインとなるコサイン (COS)、及び、コサインを変形した5つの類似度尺度を試した。グローバルセンタリング (標準的なセンタリング) (CENT)、ローカライズドセンタリング (LCENT)、コミュートタイムカーネル (CT)、ミューチュアルプロキシミティ (MP) 及ローカルスケーリング (LS) である。コミュートタイムカーネルは、グラフラブリアンに基づくグラフノード間の類似度を定めるために Saerens 10
らにより発表され (Saerens, M.; Fous, F.; Yen, L.; and Dupont, P. 2004. "The principal components analysis of graph, and its relationship to spectral clustering." In Proceedings of the 15th European Conference on Machine Learning (ECML '04), 371-383. 2004年)、後に
発明者らによりハブ軽減に効果があることが示された (Suzuki, I.; Hara, K.; Shimbo, M.; Matsumoto, Y.; and Saerens, M. 2012. "Investigating the effectiveness of Laplacian-based kernels in hub reduction." In Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence. 2012年)。ミューチュアルプロキシミティ (Schnitzer, D.; Flexer, A.; Schedl, M.; and Widmer, G. 2012. "Local and global scaling reduce hubs in space." Journal of Machine Learning Research 13(1): 2871-2902. 2012年) とローカルスケーリング (Zelnik-Manor, L. and Perona, P. 2005. "Self-tuning spectral clustering." In Advances in Neural Information Processing Systems 17. MIT Press. 1601-1608. 2005年) はデータ間の近傍関係を対称化しようと試みるものである。 20

【0057】

図20-23に実験結果を示した。LCENTあるいはLSが精度の項目において最良のパフォーマンスを達成した。CENTはコサインに対して歪度と精度をわずかに改善するのみに留まったが、対照的にLCENTはコサインに対して歪度と精度を大きく改善した。 30

【0058】

以上により、本実施例によれば、次元数が大きいときだけでなく、データ数が大きいときに出現するハブを軽減できる類似度演算システム及び類似度演算方法を提供できる。 40

【実施例2】

【0059】

本実施例では、同一データセットについて検索を行う場合、ペナルティスコア情報記憶部26に記憶されているペナルティスコアを用いることにより、検索処理を早められる。あるいは、前処理としてオフラインに各データに対するペナルティスコアを計算しておき、ペナルティスコア情報記憶部26に記憶しておくことで、検索処理を早められる。その他のシステム構成及び処理フローは実施例1と同様であり、実施例1と同様に次元数が大きいときだけでなく、データ数が大きいときに出現するハブを軽減できる類似度演算システム及び類似度演算方法を提供できるという効果を奏する。

【実施例3】

10

20

30

40

50

【0060】

本実施例では、実施例1で説明した総合スコアを示す式(6)、(7)のペナルティスコアを調整する例について説明する。

$$Sim^{LCENT}(x; q) - \langle x, q \rangle - \langle x, c(x) \rangle \quad (6)$$

$$Sim^{LCENT}(x; q) - \langle x, q \rangle - \langle x, c(x) \rangle \quad (7)$$

式(6)、(7)の右側の第2項はペナルティスコアである。パラメタは N_k 分布の歪度を最も軽減するように調整され得る。その他のシステム構成及び処理フローは実施例1と同様であり、実施例1と同様に次元数が大きいときだけでなく、データ数が大きいときに出現するハブを軽減できる類似度演算システム及び類似度演算方法を提供できるという効果を奏する。

10

【実施例4】

【0061】

本実施例では、ベースとする類似度尺度として、これまで説明に用いてきた内積ではなく、ユークリッド距離(L^2 ノルム)を用いた場合の、ハブ軽減方法について説明する。ここでは、クエリ q と各データ x との類似度 $Sim(x; q)$ をユークリッド距離($\sqrt{\|x - q\|^2}$)の2乗にマイナス符号を付した値として定義する。すなわち、

$$Sim(x; q) = -\|x - q\|^2 \quad (8)$$

である。

これを類似度尺度として用いて発生するハブは、ローカルセントロイド $c(x)$ を用い、次のように類似度スコアを補正することによって、軽減できる。

20

$$Sim^{LCENT}(x; q) - (\|x - q\|^2 + \|x - c(x)\|^2) \quad (9)$$

さらに、類似度スコアを調整する右辺第2項に、次のようにパラメタを付加してもよい。

$$Sim^{LCENT}(x; q) - (\|x - q\|^2 + (\|x - c(x)\|^2)) \quad (10)$$

$$Sim^{LCENT}(x; q) - (\|x - q\|^2 + \|x - c(x)\|^2) \quad (11)$$

30

【実施例5】

【0062】

本実施例では、各データのベクトル表現が与えられず、機械学習分野で呼ばれるところのカーネルやBLASTスコアなどのようにデータ間の類似度だけが、ベースとする類似度尺度として与えられる場合の、ハブ軽減方法について説明する。まず、クエリ q と各データ x との類似度が、BLASTスコアにより定義されているとする。すなわち、

$$Sim(x; q) = \text{BLASTスコア}(x; q) \quad (12)$$

とする。ここで、BLASTスコアを類似度尺度として用いて発生するハブは、次のようにスコアを補正することによって、軽減できる。

40

$$Sim^{LCENT}(x; q) - (1/n) \sum_{x' \in N_N(x)} \text{BLASTスコア}(x; x') \quad (13)$$

上記ではBLASTスコアを例に用いて説明したが、ベースとする類似度尺度としてその他の類似度(カーネルなど)を用いる場合は、上記BLASTスコアの部分を、用いる類似度(カーネルなど)が定めるスコアに置き換えればよい。すなわち、

$$Sim^{LCENT}(x; q) - (1/n) \sum_{x' \in N_N(x)} Sim(x; x') \quad (14)$$

とすれば、任意の類似度尺度 $Sim(x; q)$ に対して、ハブを軽減できる。ここで、 $n = |D|$ でも良く、さらに、実施例3、4と同様に、ペナルティスコアを調整するためのパラメタを導入しても良い。

50

【実施例 6】

【0063】

本実施例では、ローカルセグメントサイズ の値を、データセット内で一定としない場合について説明する。例えば、ローカルセグメントサイズ の値は、データセット内のデータ x ごとに異なっていてよい。この場合、 の値は、データ x の近傍におけるデータの密集度に基づいて定めてもよい。具体的には、データ x のハブ度、すなわち、データセット内の他のデータの k 近傍への出現回数 $N_k(x)$ が大きくなる程、 の値も大きくなるようにしてもよい。これにより、データセット内に複数のデータのクラスター（局所集合）が存在し、各クラスターを構成するデータ数が大きく異なる場合であっても、各データ x に対応するローカルセントロイド $c(x)$ を、データ x が属するクラスターの中心付近に設定することができる。したがって、類似度演算システム 1 は、ハブの発生を抑制し、 k 近傍法の精度を向上させることができる。

10

【0064】

次に、変形例に係る類似度演算システム 1 について説明する。

上述した実施例では、一例として、類似度演算システム 1 がローカルセントロイド $c(x)$ に基づいて、ペナルティスコアを算出する場合について説明したがこれには限られない。類似度演算システム 1 は、ローカルセントロイド $c(x)$ に基づいてボーナススコアを算出してもよい。類似度尺度として内積を用いる場合、内積は大きければ大きい程、2つのデータが類似することを表す。したがって、内積を用いる場合には、ローカルアフィニティ $Local\ Affinity(x)$ はペナルティとして算出される。他方、類似度尺度として、距離を用いる場合、距離は小さければ小さい程、2つのデータが類似することを表す。したがって、距離を用いる場合には、ローカルアフィニティ $Local\ Affinity(x)$ はボーナスとして算出される。このように、ローカルアフィニティ $Local\ Affinity(x)$ に係るペナルティ、ボーナスの概念は、類似度尺度等に応じて異なる。つまり、ローカルアフィニティ $Local\ Affinity(x)$ とは、スコアの補正に用いられるパラメタであり、その物理的な意味付けは、システムの構成に応じて異なるパラメタである。

20

【0065】

また、上述した実施例では、一例として、データセット、クエリを構成するデータとして塩基配列データやアミノ酸配列データを用いる場合について説明したが、これには限られない。例えば、上述した配列データの他にも、マイクロアレイデータ等の生命科学分野のデータが用いられてもよい。また、生命科学分野に限られず、音声データ、画像データ、テキストデータ等の任意の種類データが用いられてよい。

30

【0066】

なお、ローカルセグメントサイズ の値は、任意の方法で設定されてよい。例えば、塩基配列データやアミノ酸配列データは、データ空間において、比較的、データ数の小さなクラスターが大量に発生し、これらがハブ発生の原因となる場合がある。そこで、この場合には、ローカルセグメントサイズ の値は、例えば、10以下の比較的小さな値としてもよい。このように、データ空間においてデータ数の小さなクラスターが大量に発生している場合には、ローカルセグメントサイズ として、相対的に小さな値を設定する。これにより、類似度演算システム 1 は、演算量を低減しつつ、ハブを軽減することができる。これに対して、データ空間においてデータ数の大きなクラスターが多く発生している場合には、ローカルセグメントサイズ の値を、相対的に大きな値としてもよい。これにより、データ x がクラスターの外縁部に位置していたとしても、データ x ごとのローカルセントロイド $c(x)$ をクラスターの中心に近づけることができるため、効率よくハブを軽減することができる。また、ローカルセグメントサイズ の値は、 k の値に基づいて定められてもよい。例えば、ローカルセグメントサイズ の値は、 k の値の2倍程度としてもよい。また、ローカルセグメントサイズ の値は、システムごとに予め定められていてもよい。また、ローカルセグメントサイズ の値は、データ x のハブ度、すなわち、データセット内の他のデータの k 近傍への出現回数 $N_k(x)$ に関する分布の歪度が最小になる

40

50

数としてもよい。

【 0 0 6 7 】

また、上述した実施例では、一例として、k近傍法を用いる場合について説明したが、これには限られない。例えば、k近傍法を用いる代わりにイプシロン近傍法を用いてもよい。この場合、隣接データの範囲は、データセット内のデータのうち、データxとの類似度が所定の値よりも高いデータとしてもよい。

【 0 0 6 8 】

また、本発明の一態様は、以上の実施例のフローチャート等に記載の類似度演算方法を含むコンピュータに実行させるためのプログラムとしても実現可能である。プログラムは情報検索装置の内蔵記憶部に蓄積して使用してもよく、外付けの記憶装置に蓄積して使用してもよく、インターネットからダウンロードして使用しても良い。また、当該プログラムを記録した記録媒体としても実現可能である。

10

【 0 0 6 9 】

以上、本発明の一態様に係る実施の形態について説明したが、実施の形態は以上の例に限られるものではなく、本発明の趣旨を逸脱しない範囲で、種々の変更を加え得ることは明白である。例えば、上述の各実施例において説明した各構成は、特定の機能を発揮するのに不要である場合には、省略することができる。また、上述した各実施例において説明した各構成は、複数の装置に分散して備えられても良いし、1つの構成としてまとめられてもよい。

例えば、データセットがインターネット検索におけるWeb文書集合のように、データセットの次元数は明らかではないがデータ数が膨大であることが明白である場合、本発明の適用によりハブを軽減できる効果が認められるならば、本発明が提供する類似度演算システム及び類似度演算方法は有効である。また、以上の実施例では検索対象がクエリと同一である場合について説明したが、クエリが検索対象と関連性が深いキーワードであっても良い。また、データセットが膨大なときには、データセットが複数の記憶装置に分散して記憶されても良い。また、システム構成では、ローカルセントロイド情報記憶部25とペナルティスコア情報記憶部26をまとめて1つにしても良い。また、処理フローではデータセット・クエリ設定工程(S101)とベース類似度スコア演算工程(S102)とを同時に行っても良い。また、k近傍法のパラメタkは目的、状況に応じて適宜定めることができる。

20

30

【 産業上の利用可能性 】

【 0 0 7 0 】

本発明は類似度演算及び情報分類に利用される。

【 符号の説明 】

【 0 0 7 1 】

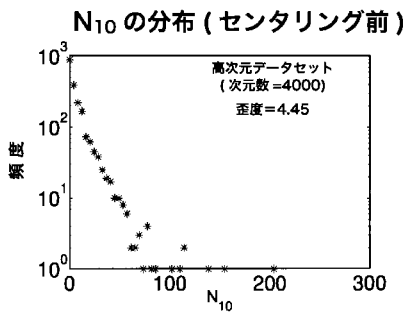
- 1 類似度演算システム
- 10 コンピュータ
- 11 データセット・クエリ設定部
- 12 ベース類似度スコア演算部
- 13 ハブ度演算部
- 14 ペナルティスコア演算部
- 15 総合スコア演算部
- 16 演算結果整理・分類部
- 17 制御部
- 18 表示部
- 19 入出力部
- 20 記憶部
- 21 データセット・クエリ情報記憶部
- 22 ベース類似度スコア情報記憶部
- 23 出現回数・歪度情報記憶部

40

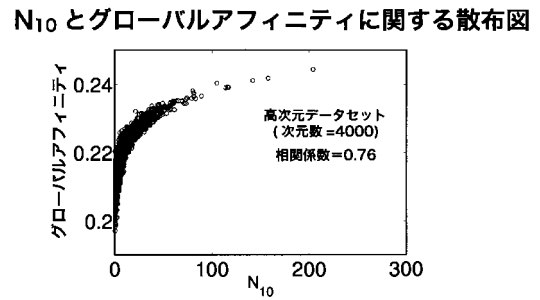
50

- 2 4 ローカルセグメントサイズ情報記憶部
- 2 5 ローカルセントロイド情報記憶部
- 2 6 ペナルティスコア情報記憶部
- 2 7 総合スコア情報記憶部
- 2 8 分類情報記憶部
- 2 9 プログラム記憶部
- 3 0 一時記憶部
- 1 3 1 出現回数演算部
- 1 3 2 歪度演算部
- 1 4 1 ローカルセグメントサイズ決定部
- 1 4 2 ローカルセントロイド演算部
- 1 4 3 ローカルアフィニティ演算部
- 1 6 1 総合スコア順整理部
- 1 6 2 分類部

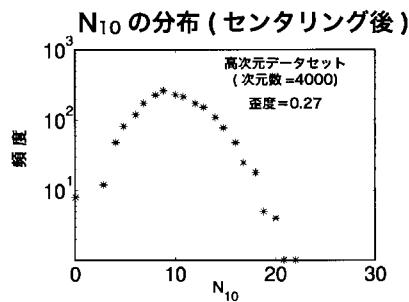
【 図 1 】



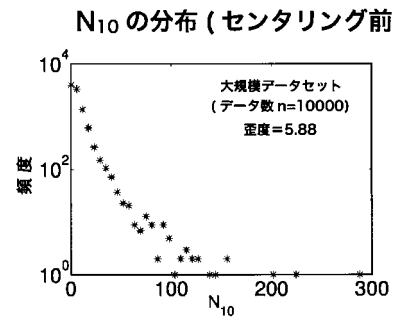
【 図 3 】



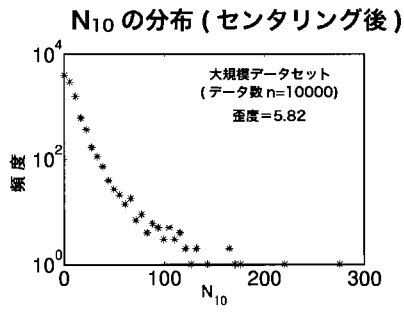
【 図 2 】



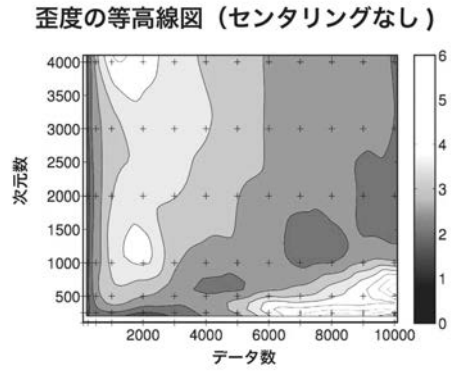
【 図 4 】



【 図 5 】

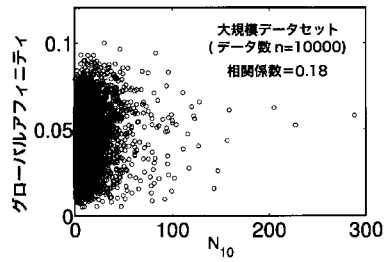


【 図 7 】



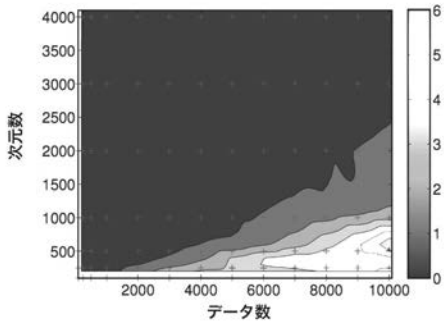
【 図 6 】

N_{10} とグローバルアフィニティに関する散布図



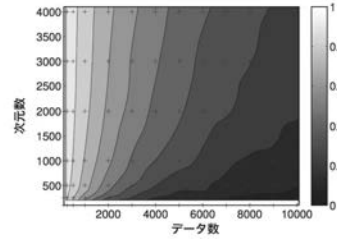
【 図 8 】

歪度の等高線図 (グローバルセンタリング後)



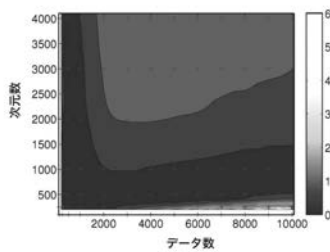
【 図 1 0 】

グローバル対ローカルアフィニティ比



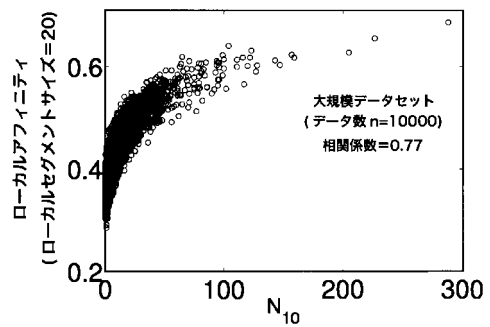
【 図 9 】

歪度の等高線図 (ローカライズドセンタリング後)



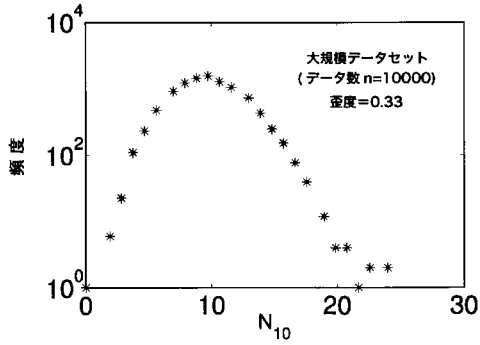
【 図 1 1 】

(a) N_{10} とローカルアフィニティに関する散布図



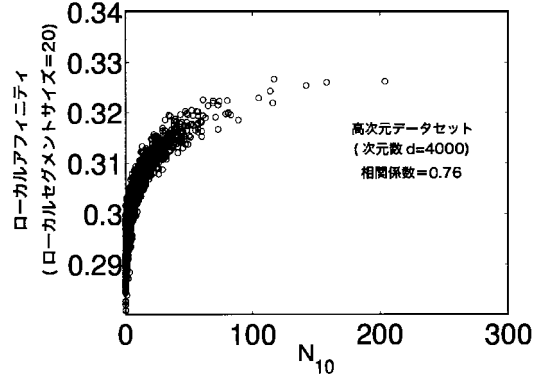
【 図 1 2 】

(b) N_{10} の分布 (ローカライズセンタリング後)

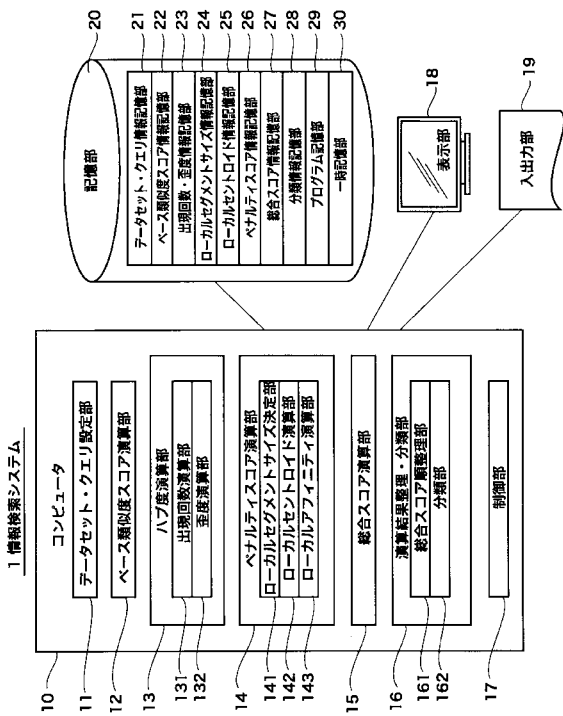


【 図 1 3 】

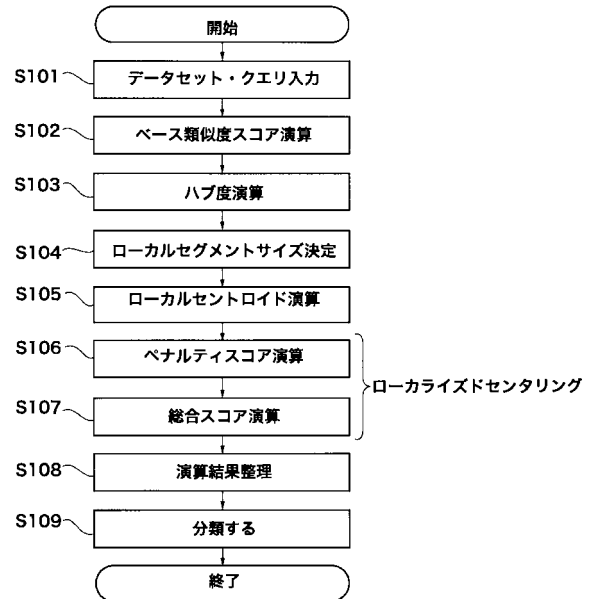
N_{10} とローカルアフィニティに関する散布図



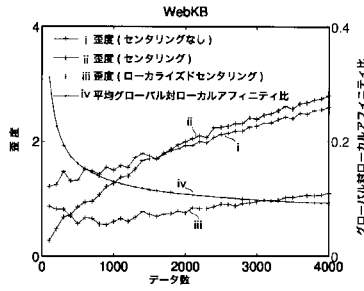
【 図 1 4 】



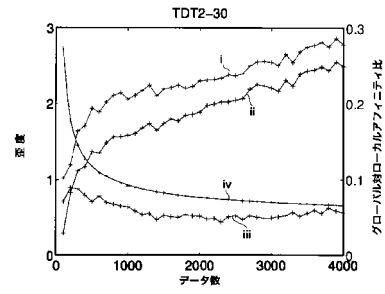
【 図 1 5 】



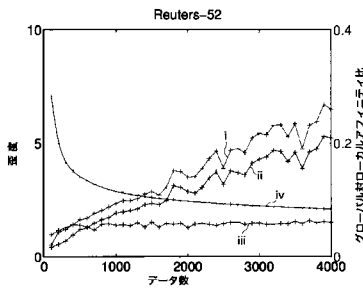
【 図 1 6 】



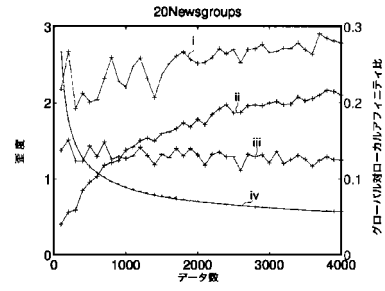
【 図 1 8 】



【 図 1 7 】



【 図 1 9 】



【 図 2 0 】

WebKB (データ数 4168, 分類クラス数 4)												
k	k 近傍法分類の正解率						Nk 分布の密度					
	COS	CENT	LCENT	CT	MP	LS	COS	CENT	LCENT	CT	MP	LS
10	0.753	0.756	0.761	0.743	0.761	0.764	2.64	2.85	1.13	3.52	0.74	1.15
20	0.769	0.766	0.774	0.757	0.771	0.777	2.01	2.12	0.87	2.51	0.69	0.65
30	0.770	0.776	0.780	0.764	0.780	0.786	1.74	1.70	0.73	2.06	0.65	0.47
40	0.777	0.778	0.781	0.772	0.778	0.785	1.63	1.42	0.66	1.78	0.52	0.34
50	0.776	0.780	0.785	0.781	0.783	0.798	1.57	1.25	0.63	1.60	0.39	0.22

【 図 2 3 】

20Newsgroups (データ数 18,820, 分類クラス数 20)												
k	k 近傍法分類の正解率						Nk 分布の密度					
	COS	CENT	LCENT	CT	MP	LS	COS	CENT	LCENT	CT	MP	LS
10	0.859	0.860	0.877	0.838	0.865	0.974	2.99	2.85	0.81	5.53	0.49	0.77
20	0.845	0.845	0.861	0.841	0.852	0.862	2.69	2.12	0.89	4.09	0.49	12.91
30	0.834	0.836	0.849	0.837	0.846	0.855	2.47	1.70	1.05	3.42	0.43	6.71
40	0.831	0.832	0.841	0.833	0.840	0.849	2.37	1.42	1.27	2.97	0.39	4.54
50	0.825	0.827	0.835	0.833	0.835	0.844	2.27	1.25	1.36	2.67	0.36	3.06

【 図 2 1 】

Reuters-52 (データ数 9100, 分類クラス数 5 2)												
k	k 近傍法分類の正解率						Nk 分布の密度					
	COS	CENT	LCENT	CT	MP	LS	COS	CENT	LCENT	CT	MP	LS
10	0.872	0.885	0.901	0.858	0.898	0.895	14.82	11.04	1.76	6.93	0.82	2.36
20	0.876	0.894	0.913	0.885	0.908	0.904	7.92	6.42	1.58	4.92	0.77	1.76
30	0.875	0.896	0.913	0.889	0.909	0.904	5.74	4.64	1.61	4.04	0.78	1.37
40	0.869	0.896	0.913	0.889	0.903	0.901	4.68	3.77	1.63	3.55	0.81	1.08
50	0.866	0.894	0.910	0.892	0.900	0.897	4.02	3.27	1.63	3.22	0.80	0.90

【 図 2 2 】

TDT2-30 (データ数 9394, 分類クラス数 30)												
k	k 近傍法分類の正解率						Nk 分布の密度					
	COS	CENT	LCENT	CT	MP	LS	COS	CENT	LCENT	CT	MP	LS
10	0.964	0.963	0.964	0.953	0.962	0.961	3.18	2.85	0.91	4.20	0.47	0.80
20	0.965	0.963	0.964	0.955	0.962	0.962	2.81	2.12	0.81	3.09	0.38	0.77
30	0.966	0.964	0.966	0.958	0.963	0.963	2.42	1.70	0.45	2.62	0.31	0.71
40	0.965	0.963	0.966	0.959	0.963	0.963	2.16	1.42	0.45	2.34	0.25	0.64
50	0.965	0.963	0.967	0.960	0.963	0.965	1.98	1.25	0.46	2.15	0.22	0.61

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/JP2016/051909
A. CLASSIFICATION OF SUBJECT MATTER G06F17/30(2006.01)i According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED 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 1996-2016 Kokai Jitsuyo Shinan Koho 1971-2016 Toroku Jitsuyo Shinan Koho 1994-2016 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	Ikumi SUZUKI et al., "k Kinboho de Hub o Keigen suru Ruijido Shakudo", IPSJ SIG Notes 2012 (Heisei 24) Nendo (4), 04 January 2013 (04.01.2013) (received date), pages 1 to 8	1-10
A	Ikumi SUZUKI et al., "Effectiveness of Laplacian-based kernels from the hubness point of view", IEICE Technical Report, 02 November 2011 (02.11.2011), vol.111, no.275, pages 257 to 262	1-10
A	Kohei OZAKI et al., "Semi-supervised Word Sense Disambiguation using Degree Bounded Graph Construction", IPSJ SIG Notes Heisei 22 Nendo (4), 04 January 2011 (04.01.2011) (received date), pages 1 to 8	1-10
<input 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 06 April 2016 (06.04.16)		Date of mailing of the international search report 19 April 2016 (19.04.16)
Name and mailing address of the ISA/ Japan Patent Office 3-4-3, Kasumigaseki, Chiyoda-ku, Tokyo 100-8915, Japan		Authorized officer Telephone No.

国際調査報告		国際出願番号 PCT/J P 2 0 1 6 / 0 5 1 9 0 9	
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 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2016年 日本国実用新案登録公報 1996-2016年 日本国登録実用新案公報 1994-2016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)			
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	鈴木 郁美、外2名、k近傍法でハブを軽減する類似度尺度、情報処理学会研究報告 2012 (平成24)年度(4), 2013.01.04 (受入日), pp. 1-8	1-10	
A	鈴木 郁美、外3名、「ハブの出現しやすさ」から見たラプラシアンベースカーネル、電子情報通信学会技術研究報告, 2011.11.02, 第111巻, 第275号, pp. 257-262	1-10	
<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 「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 「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 「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) 「O」口頭による開示、使用、展示等に言及する文献 「P」国際出願日前で、かつ優先権の主張の基礎となる出願		の日の後に公表された文献 「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 「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 「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 「&」同一パテントファミリー文献	
国際調査を完了した日 06.04.2016		国際調査報告の発送日 19.04.2016	
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) 齊藤 貴孝 電話番号 03-3581-1101 内線 3599	5M 4774

国際調査報告		国際出願番号 PCT/JP2016/051909
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	小寄 耕平、外3名、ハブを作らないグラフ構築法を用いた半教師あり語義曖昧性解消、情報処理学会研究報告 平成22年度(4), 2011.01.04 (受入日), pp. 1-8	1-10

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, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, 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, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, 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, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(72)発明者 原 一夫

静岡県三島市谷田 1 1 1 1 番地 大学共同利用機関法人情報・システム研究機構 国立遺伝学研究所内

(72)発明者 鈴木 郁美

静岡県三島市谷田 1 1 1 1 番地 大学共同利用機関法人情報・システム研究機構 国立遺伝学研究所内

(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。