

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2014-120140

(P2014-120140A)

(43) 公開日 平成26年6月30日(2014.6.30)

(51) Int.Cl.  
G06F 17/30 (2006.01)

F I  
G06F 17/30 210D

テーマコード (参考)

審査請求 未請求 請求項の数 8 O L (全 22 頁)

(21) 出願番号 特願2012-277491 (P2012-277491)  
(22) 出願日 平成24年12月19日 (2012.12.19)

(71) 出願人 000005223  
富士通株式会社  
神奈川県川崎市中原区上小田中4丁目1番1号  
(74) 代理人 100108187  
弁理士 横山 淳一  
(72) 発明者 岩倉 友哉  
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

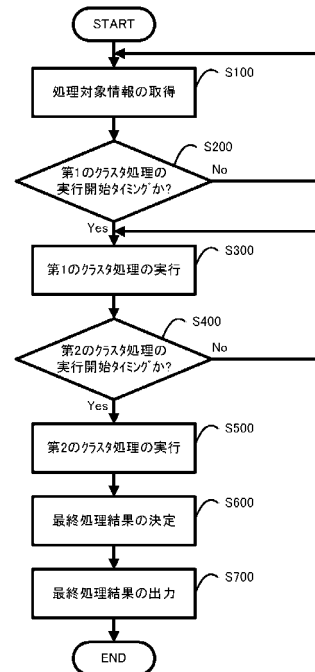
(54) 【発明の名称】 クラスタ処理方法、クラスタ処理装置およびプログラム

(57) 【要約】

【課題】 開示の技術は、高速なクラスタ処理を実現可能とする。

【解決手段】 本開示の技術における解決手段の一観点は、第2のクラスタ手法と異なりかつ該第2のクラスタ手法よりも高速な第1のクラスタ手法によって対象情報に対して第1の分類処理を実行し、前記第1の分類処理の結果に基づいて前記第2のクラスタ手法によって第2の分類処理を実行し、前記第1の分類処理の結果と前記第2の分類処理の結果とに基づいて前記対象情報の分類を決定する、クラスタ処理方法、クラスタ処理装置およびプログラムである。

【選択図】 図4



**【特許請求の範囲】****【請求項 1】**

コンピュータが、

第 2 のクラスタ手法と異なりかつ該第 2 のクラスタ手法よりも高速な第 1 のクラスタ手法によって対象情報に対して第 1 の分類処理を実行し、

前記第 1 の分類処理の結果に基づいて前記第 2 のクラスタ手法によって第 2 の分類処理を実行し、

前記第 1 の分類処理の結果と前記第 2 の分類処理の結果とに基づいて前記対象情報の分類を決定する、

ことを特徴とするクラスタ処理方法。

10

**【請求項 2】**

前記対象情報の分類を決定することは、前記第 1 の分類処理の結果と前記第 2 の分類処理の結果間で分類の変更がなかった前記第 1 の分類処理の結果または前記第 2 の分類処理の結果に基づいて前記対象情報の分類を決定する、

ことを特徴とする請求項 1 記載のクラスタ処理方法。

**【請求項 3】**

前記第 1 の分類処理の結果は、前記対象情報に関連した分類個数情報を含み、

前記第 2 の分類処理は、前記分類個数情報を前記第 2 のクラスタ手法における初期情報として実行される、

ことを特徴とする請求項 1 または 2 記載のクラスタ処理方法。

20

**【請求項 4】**

前記第 1 の分類処理の結果は、前記対象情報に含まれるキーワードを基にした前記対象情報の特徴情報を含み、

前記第 2 の分類処理は、前記対象情報の特徴情報を前記第 2 のクラスタ手法における前記初期情報として実行される、

ことを特徴とする請求項 3 記載のクラスタ処理方法。

**【請求項 5】**

前記第 1 の分類処理は、前記対象情報を用いた前記キーワード別の集合を作成し、異なる前記キーワードに対応する前記集合を異なる分類とする、

ことを特徴とする請求項 4 記載のクラスタ処理方法。

30

**【請求項 6】**

前記第 2 のクラスタ手法は、処理対象情報の類似性を基に処理対象情報を分類する手法である、

ことを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載のクラスタ処理方法。

**【請求項 7】**

第 2 のクラスタ手法と異なりかつ該第 2 のクラスタ手法よりも高速な第 1 のクラスタ手法によって対象情報に対して第 1 の分類処理を実行する第 1 処理部と、

前記第 1 の分類処理の結果に基づいて前記第 2 のクラスタ手法によって第 2 の分類処理を実行する第 2 処理部と、

前記第 1 の分類処理の結果と前記第 2 の分類処理の結果とに基づいて前記対象情報の分類を決定する決定部と、

を有することを特徴とするクラスタ処理装置。

40

**【請求項 8】**

コンピュータに、

第 2 のクラスタ手法と異なりかつ該第 2 のクラスタ手法よりも高速な第 1 のクラスタ手法によって対象情報に対して第 1 の分類処理を実行し、

前記第 1 の分類処理の結果に基づいて前記第 2 のクラスタ手法によって第 2 の分類処理を実行し、

前記第 1 の分類処理の結果と前記第 2 の分類処理の結果とに基づいて前記対象情報の分類を決定する、

50

処理を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】

本発明は、クラスタ処理方法、クラスタ処理装置およびプログラムに関する。

【背景技術】

【0002】

クラスタリング技術に関する先行技術として、文書類似ベクトルを用いて対象文書とクラスタ重心との距離を算出し、さらに同一の対象文書に対して一回目の分類に利用した文書類似ベクトルの次元数を増加させて二回目の分類を行い、安定クラスタの文書を対象から除いて次の対象文書を選定して分類試行を繰り返す技術が知られている。(例えば、特許文献1参照)

10

また、他の先行技術として、データ集合を部分クラスタの集合に変換し(大分類し)、部分クラスタの集合をクラスタリングするにあたり部分クラスタの局所的な密度に関する属性を考慮して詳細分類を行う技術が知られている。(例えば、特許文献2参照)

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開2002-183171号公報

【特許文献2】特開2010-134632号公報

20

【発明の概要】

【発明が解決しようとする課題】

【0004】

前記先行技術は、不安定なクラスタの文書や部分クラスタを詳細分類することで、高精度なクラスタリング結果を得ることができる。

【0005】

しかしながら、前記先行技術では、このような詳細分類が繰り返されるためにクラスタ処理に多くの時間を要するという問題がある。特に、大規模データを処理対象にすると、そのデータ量の多さから必然的に不安定なクラスタの数も多くなる。そのため、前記先行技術により大規模データを対象にしたクラスタ処理においては、更に多くの処理時間が必要となる。

30

【0006】

本願は、高速なクラスタ処理を実現可能とするクラスタ処理方法、クラスタ処理装置およびプログラム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】

【0007】

上記課題を解決するために、本実施例に開示のクラスタ処理方法は、コンピュータが、第2のクラスタ手法と異なりかつ該第2のクラスタ手法よりも高速な第1のクラスタ手法によって対象情報に対して第1の分類処理を実行し、前記第1の分類処理の結果に基づいて前記第2のクラスタ手法によって第2の分類処理を実行し、前記第1の分類処理の結果と前記第2の分類処理の結果とに基づいて前記対象情報の分類を決定する。

40

【発明の効果】

【0008】

本実施例の一観点によれば、高速なクラスタ処理が実現される。

【図面の簡単な説明】

【0009】

【図1】コンピュータシステムを示す。

【図2】コンピュータのハードウェアを示す。

【図3】コンピュータの機能ブロックを示す。

【図4】処理全体のフローチャートを示す。

50

【図5】処理例を示す。

【図6】取得処理の一例のフローチャートを示す。

【図7】第1クラスタ処理の一例のフローチャートを示す。

【図8】第1クラスタ処理におけるキーワード抽出処理の一例のフローチャートを示す。

【図9】第1クラスタ処理における集合作成処理の一例のフローチャートを示す。

【図10】第1クラスタ処理におけるベクトル情報作成処理の一例のフローチャートを示す。

【図11】第2クラスタ処理の一例のフローチャートを示す。

【図12】決定処理の一例のフローチャートを示す。

【発明を実施するための形態】

10

【0010】

以下、図面を参照して開示の技術の実施形態の一例を詳細に説明する。

【0011】

図1は、コンピュータシステム1を示す。コンピュータシステム1は、例えば、コンピュータ100、ネットワーク200、サーバ300、ストレージシステム320、コンピュータ400、コンピュータ500、NAS(Network Attached Storage)600を含む。本例において、コンピュータ100、サーバ300、コンピュータ400、コンピュータ500、NAS600は、ネットワーク200にそれぞれ接続されている。

【0012】

コンピュータ100は、その詳細を後述する装置であり、クラスタ処理を実行する。コンピュータ100は、例えば、サーバ、ワークステーション、パーソナルコンピュータ、インターネットアプライアンス、ゲーム機などである。コンピュータ100として、クラスタ処理にかかる処理負荷やその処理結果の提供形態などに応じて、適宜の装置が選択されればよい。

20

【0013】

ネットワーク200は、LANやインターネット等であり、それに接続される装置間でのデータ通信を可能とする。

【0014】

サーバ300は、ソーシャルメディアを実現する装置である。ソーシャルメディアは、例えば、電子掲示板、ブログ、ウィキ、ツイッター、ポッドキャスト、ソーシャルブックマーク、ソーシャル・ネットワーキング・サービス、画像や動画の共有サイト、通販サイトのカスタマーレビューなどである。サーバ300は、ネットワーク200を介して受信したデータ登録要求やデータ送信要求に応じて、それら要求で示される処理を実行し、また、それら要求を発行した他の装置に要求対象データの送信等を実行する。サーバ300は、ソーシャルメディアを実現するためのデータを格納するストレージシステム320を備えてもよい。ストレージシステム320は、例えば、RAID(Redundant Arrays of Inexpensive Disks)システムである。サーバ300は、ソーシャルメディア上の大規模データを扱う。なお、サーバ300は、ワークステーション、パーソナルコンピュータなど他のコンピュータであってもよい。

30

40

【0015】

コンピュータ100は、上記のようにサーバ300が提供するソーシャルメディア上の大規模データを対象にクラスタ処理を実行し、ソーシャルメディア上で盛り上がっている話題(内容やキーワード)を抽出する。例えば、コンピュータ100は、ソーシャルメディアの情報である例えばテキストをその内容の類似性によってまとめあげることによって、ソーシャルメディア上の情報の中から盛り上がっている話題ごとに情報をクラスタリングする。

【0016】

コンピュータ400は、サーバ300により実現されるソーシャルメディアを利用し、そのサービスを受ける装置である。例えば、コンピュータ400がネットワーク200を

50

介してサーバ300に前述のデータ登録要求やデータ送信要求を発行する。そのような要求に対して、サーバ300は、要求にて示される処理を実行する。

【0017】

コンピュータ500は、例えば、コンピュータ100により実行されたクラスタ処理の処理結果を取得して表示する。また、コンピュータ500は、コンピュータ100の動作設定を管理する装置であってもよい。

【0018】

NAS600は、ネットワーク200に接続されたストレージシステムである。例えば、前述のサーバ300は、直接または他のネットワークを介して間接にネットワーク200に接続されたNAS600をデータの格納先として利用してもよい。この場合、サーバ300は、ネットワーク200を介して、データのライト要求やリード要求をNAS600に発行することになる。NAS600は、他装置、例えばサーバ300からのこのような要求を受信し、その要求で示される処理を実行する。NAS600は、その要求がデータのライトを示す場合は指定されたデータを自身が有する記憶装置（例えばハードディスク）に記録し、要求がデータのリードを示す場合は指定されたデータを記憶装置から読み出して要求の発行元装置にそれを送信する。NAS500は、RAIDシステムであってもよい。

10

【0019】

図2は、コンピュータ100のハードウェアを示す。この図2に示されるハードウェアは、コンピュータ100を構成するハードウェアの一例であり、少なくとも本実施形態に記載される処理の実行に必要なハードウェア構成を備えていればよい。

20

【0020】

図2に示されるコンピュータ100は、例えば、プロセッサ10、RAM(Random Access Memory)20、ドライブ装置30、記憶媒体32、入力インターフェース(I/F)40、入力デバイス42、出力インターフェース(I/F)50、出力デバイス52、通信インターフェース(I/F)60およびバス70などを含む。それぞれのハードウェア構成は、バス70を介して接続されている。

【0021】

プロセッサ10は、例えば、CPU(Central Processing Unit)、MPU(Micro-Processing Unit)、DSP(Digital Signal Processor)などの処理回路である。

30

【0022】

RAM20は読み書き可能なメモリ装置であって、例えば、SRAM(Static RAM)やDRAM(Dynamic RAM)などの半導体メモリである。なお、RAMではなく、フラッシュメモリなどであってもよい。

【0023】

ドライブ装置30は、記憶媒体32にアクセスする。記憶媒体32は、データを記憶している。ドライブ装置30は、記憶媒体32へのデータのライト、記憶媒体32からのデータのリードの少なくともいずれか一方を行う。記憶媒体32は、例えば、ハードディスク、SSD(Solid State Drive)などのフラッシュメモリ、CD(Compact Disc)、DVD(Digital Versatile Disc)、ブルーレイディスクなどである。コンピュータ100は、記憶媒体32の種類に応じたドライブ装置30を備えればよい。

40

【0024】

入力インターフェース40は、入力デバイス42が接続され、入力デバイス42から受信した入力信号をプロセッサ10に伝達する回路である。入力デバイス42は、利用者による操作に応じた入力信号を出力する装置である。入力デバイス42は、例えば、キーボードやコンピュータ100に設置されたボタンなどのキー装置や、マウスやタッチパネルなどのポインティングデバイスである。

【0025】

50

出力インターフェース 50 は、出力デバイス 52 が接続され、出力デバイス 52 に、プロセッサ 100 の指示に応じた出力を実行させる回路である。出力デバイス 52 は、コンピュータ 100 の制御に応じて情報を出力する装置である。出力デバイス 52 は、例えば、ディスプレイなどの画像出力装置（表示デバイス）や、スピーカーなどの音声出力装置、プリンタなどである。また、例えば、タッチスクリーンなどの入出力装置が、入力デバイス 42 および出力デバイス 52 として用いられる。また、入力デバイス 42 及び出力デバイス 52 は、コンピュータ 100 と一体であってもよいし、コンピュータ 100 に含まれず、例えば、コンピュータ 100 に外部から接続する装置であってもよい。また、例えば、コンピュータ 100 が通常動作時にこのコンピュータ 100 を利用するユーザに対して情報提供するための出力デバイス 52 を必要としないサーバのようなコンピュータであれば、入力デバイス 42 および出力デバイス 52 を不要とする構成としてもよい。また、そのような場合、入力インターフェース 40 および出力インターフェース 50 も不要とする構成としてもよい。

#### 【0026】

通信インターフェース 60 は、ネットワーク 200 に接続され、ネットワーク 200 を介した通信の制御を行なう回路である。通信インターフェース 60 は、例えば、Ethernet（登録商標）カードに代表されるネットワークインターフェースカード（NIC）などである。

#### 【0027】

例えば、プロセッサ 10 は、記憶媒体 32 からプログラムをリードし、それを RAM 20 にロードする。また、プロセッサ 10 は、プログラムで利用されるデータを RAM 20 にロードする。このように、RAM 20 は、プロセッサ 10 によるプログラム実行のためのワークエリアとして用いられる。プログラムは、オペレーティング・システム（OS）やアプリケーションプログラムなどであり、所定の処理手順を実行する命令を含む。例えば、このようなプログラムの動作により、本実施形態の処理機能がコンピュータ 100 上に実現される。

#### 【0028】

図 3 は、コンピュータ 100 の機能ブロックを示す。各機能ブロックは、前述のプログラムの動作によってプロセッサ 10 により実現される。

#### 【0029】

コンピュータ 100 上で実現される機能ブロックは、例えば、制御部 110、記憶部 120、取得部 130、第 1 クラスタ処理部 140、第 2 クラスタ処理部 150、決定部 160 および出力部 170 を含む。

#### 【0030】

制御部 110 は、本実施形態の処理機能を実現するため、各機能ブロックを制御する。記憶部 120 は、各種データを記憶する。なお、記憶部 120 は、その機能を実現するため、図 2 における RAM 20 または記憶媒体 32 も含む。

#### 【0031】

取得部 130 は、処理対象となるデータ、すなわち、本実施形態においてはソーシャルメディア上の情報を取得する。以後、取得部 130 によって取得された情報を処理対象情報と呼ぶことがある。また、取得部 130 は、取得した処理対象情報を記憶部 120 に格納する。本実施形態において取得部 130 が取得するソーシャルメディア上の情報は、例えばテキスト情報である。また、処理対象情報は、複数のファイルであってもいいし、複数個の文章を含む 1 つのファイルであってもよい。取得部 130 は、図 2 における通信インターフェース 60 を含んでもよい。

#### 【0032】

第 1 クラスタ処理部 140 は、記憶部 120 に格納された処理対象情報を対象に、第 1 のクラスタ処理を実行する。本実施形態において、第 1 のクラスタ処理は、例えば、キーワード抽出処理である。このキーワード抽出処理は、処理対象情報に含まれるキーワードを抽出する。抽出されたキーワードは、処理対象情報の分類処理に利用される。この処理

の詳細は後述する。第1クラスタ処理部140は、第1のクラスタ処理の結果を記憶部120に格納する。また、第1クラスタ処理部140は、記憶部120をワークスペースとして、各種データを記憶部120に一時格納し、それらデータを利用して第1のクラスタ処理を実行する。

#### 【0033】

第2クラスタ処理部150は、記憶部120に格納されている第1クラスタ処理部140による処理結果を入力情報および初期設定情報として使用し、処理対象情報を対象にした第2のクラスタ処理を実行する。本実施形態において、第2のクラスタ処理は、例えば、K-means (K-平均法) である。このK-meansは、処理対象となる情報の内容を考慮し、類似する情報を分類 (クラスタリング) する。また、後述するように、一般的なK-meansでは、最終解を得るために設定情報を何回か変更して処理試行を行う。これらが理由で、K-meansは、最終の処理結果を得るまでに処理時間がかかる。

10

#### 【0034】

これに対し、第1クラスタ処理部140が実行する前述の第1のクラスタ処理の一例として示したキーワード抽出処理は、K-meansのように処理対象となる情報間の類似性を考慮せず、情報内に含まれるキーワードをピックアップし、キーワードとそれを含む情報との対応付けにより分類処理を行う。このため、同じ情報の分類処理において、キーワード抽出処理を一例にする第1のクラスタ処理は、K-meansを一例にする第2のクラスタ処理より高速である。

20

#### 【0035】

第2クラスタ処理部150は、第2のクラスタ処理の結果を記憶部120に格納する。また、第2クラスタ処理部150は、記憶部120をワークスペースとして、各種データを記憶部120に一時格納し、それらデータを利用して第2のクラスタ処理を実行する。

#### 【0036】

決定部160は、記憶部120に格納されている第1クラスタ処理部140および第2クラスタ処理部150によるそれぞれの処理結果を基に、処理対象情報のクラスタ処理結果を決定する。本実施形態において、決定部160は、第2クラスタ処理部150によるクラスタ処理結果のうち第1クラスタ処理部140によるクラスタ処理結果からクラスタが変更されたものを除外し、それによる第2クラスタ処理部150によるクラスタ処理結果を最終クラスタ処理結果と決定する。なお、決定部160は、第1クラスタ処理部140による処理結果のうち第2クラスタ処理部150によるクラスタ処理結果からクラスタが変更されたものを除外し、それによる第1クラスタ処理部140によるクラスタ処理結果を最終クラスタ処理結果と決定してもよい。つまり、決定部160は、第1クラスタ処理部140の処理結果と第2クラスタ処理部150による処理結果間でクラスタが一致する情報を最終クラスタ処理結果と決定すればよい。決定部160は、最終クラスタ処理結果を記憶部120に格納する。

30

#### 【0037】

出力部170は、記憶部120に格納されている最終クラスタ処理結果を出力する。出力部170は、例えば、図2における入力デバイス42および出力デバイス52を使ってコンピュータ100を使用しているユーザからの要求に応じて、最終クラスタ処理結果を出力デバイス52に出力してもよい。また、出力部170は、例えば、図2における通信インターフェース60を介して外部から受信した要求に応じて、要求元の装置に最終クラスタ処理結果を通信インターフェース60を介して出力してもよい。

40

#### 【0038】

次に、前述した図3に示す各機能ブロックによる処理フローを説明する。図4は、処理全体のフローチャートを示す。なお、各機能ブロックは、制御部110により所定のタイミングで動作有効に設定されるものとする。この所定のタイミングは、本処理の前、処理を開始すべき時点などである。また、本説明においては、制御部110による制御内容を割愛または各機能ブロックによる処理内容の一部として説明している。制御部110と他

50

の各機能ブロックのそれぞれの処理機能は、実施の形態に応じて設定・配分されればよい。

【0039】

まず、取得部130は、処理対象情報を取得する(S100)。取得部130は、取得した処理対象情報を記憶部120に格納する。

【0040】

第1クラスタ処理部140は、第1のクラスタ処理の実行開始タイミングが否かを判定する(S200)。第1クラスタ処理部140は、実行開始タイミングであると判定すると(S200; Yes)、処理対象情報を対象に第1のクラスタ処理を開始する。一方、実行開始タイミングでないと判定すると(S200; No)、第1クラスタ処理部140は次の判定タイミングまで判定処理を待つ。例えば、第1クラスタ処理部140は、記憶部120に格納された処理対象情報が第1のクラスタ処理を行うために必要な情報量であるか否かを判定することによって、第1のクラスタ処理の実行開始タイミングが否かを判定してもよい。

10

【0041】

第1クラスタ処理部140は、実行開始タイミングであると判定した場合、記憶部120に格納されている処理対象情報を対象に第1のクラスタ処理を実行する(S300)。第1クラスタ処理部140は、第1のクラスタ処理の結果を記憶部120に格納する。

【0042】

第2クラスタ処理部150は、第2のクラスタ処理の実行開始タイミングが否かを判定する(S400)。第2クラスタ処理部150は、実行開始タイミングであると判定すると(S400; Yes)、記憶部120に格納されている処理対象情報を対象に第2のクラスタ処理を開始する。一方、実行開始タイミングでないと判定すると(S400; No)、第2クラスタ処理部150は次の判定タイミングまで判定処理を待つ。例えば、第2クラスタ処理部150は、第1のクラスタ処理が行われたか否かを判定することによって、第2のクラスタ処理の実行開始タイミングが否かを判定してもよい。

20

【0043】

第2クラスタ処理部150は、実行開始タイミングであると判定した場合、記憶部120に格納されている第1のクラスタ処理の結果を使って、第2のクラスタ処理を実行する(S500)。第2クラスタ処理部150は、第2のクラスタ処理の結果を記憶部120に格納する。第2クラスタ処理部150は、処理対象情報を対象にした第2のクラスタ処理を実行する場合、例えば、第1のクラスタ処理の結果を第2のクラスタ処理の入力情報および初期設定情報として使用する。

30

【0044】

決定部160は、記憶部120に格納された第1のクラスタ処理の結果および第2のクラスタ処理の結果に基づいて、処理対象情報に対する最終のクラスタ処理結果を決定する(S600)。決定部160は、最終のクラスタ処理結果を記憶部120に格納する。例えば、決定部160は、第2のクラスタ処理の結果のうち第1のクラスタ処理の結果からクラスタが変更されたものを除外し、それによる第2のクラスタ処理の結果を最終クラスタ処理結果として決定してもよい。

40

【0045】

出力部170は、決定部160により記憶部120に格納された最終のクラスタ処理結果を出力する(S700)。出力部170は、最終のクラスタ処理結果の要求元に応じた形態にて最終のクラスタ処理結果を出力する。例えば、出力部170は、前述の通信インターフェース60を介してコンピュータ500から処理要求を受信していた場合、通信インターフェース60を介してコンピュータ500に最終のクラスタ処理結果を送信するよう処理を実行してもよい。また、例えば、コンピュータ100の入力デバイス42を使って処理要求が指示されていた場合、出力インターフェース50を介して出力デバイス52に最終のクラスタ処理結果を出力するよう処理を実行してもよい。この出力処理により、本処理は終了する。なお、本処理は、停止指示を受けるまでS100~S700間の処理

50

が繰り返し実行されてもよい。

【0046】

以上、本処理の説明をしたが、次のような形態で処理が行われるようにしてもよい。

【0047】

例えば、第1クラスタ処理部140は、実行開始タイミングか否かを判定する際に、第1クラスタ処理部140が記憶部120に格納される処理対象情報を監視しそのタイミング判定を行ってもよい。また、実行開始タイミングに到達した旨を取得部130が第1クラスタ処理部140に通知するように構成し、それによって、第1クラスタ処理部140が実行開始タイミングを判定するようにしてもよい。

【0048】

また、例えば、第1クラスタ処理部140および第2クラスタ処理150それぞれは別プロセスとして並列動作するようにしてもよく、第1クラスタ処理部140および第2クラスタ処理150によるそれぞれの実行開始タイミングの判定が並列に行われるようにしてもよい。

【0049】

次に、以上説明した各処理の一例を説明する。以下に説明する各処理は一例であり、本発明はこの処理内容に限定されるものではない。なお、以下説明においては、図2におけるサーバ300を処理対象情報の提供元、コンピュータ500をコンピュータ100の管理装置としている。また、図5を適宜参照して、各処理に関する一例を説明する。

【0050】

図6は、取得処理の一例のフローチャートを示す。図6は、図3における取得部130により実施される処理であって、処理対象情報の取得方法の一例を示す。

【0051】

まず、取得処理の開始が指示されると、取得部130は、コンピュータ500により設定された情報に基づいて、サーバ300にアクセスする(S110)。コンピュータ500により設定される情報は、例えば、処理対象情報の提供元の識別情報やアドレス(本例ではサーバ300を示す)、提供元へのアクセス形態などである。

【0052】

取得部130は、サーバ300と正常に接続されると、サーバ300に処理対象情報となる情報を要求する(S120)。

【0053】

その要求に応答してサーバ300から情報が送信されると、取得部130は、その情報を受信し(S130)、記憶部120に格納する(S140)。このようにして、処理対象情報が取得される。

【0054】

なお、この情報の取得処理は、停止指示があるまで繰り返し実行されるものであってもよい。例えば、サーバ300が、ソーシャルメディアの一例であり投稿サービスであるツイッターのサービスを提供している場合、Streaming API(Application Program Interface)を利用すればそれが実現できる。このStreaming APIを利用してサーバ300にアクセスすると、S120による情報の要求以後、サーバ300は、継続して最新の投稿情報(ツイート(登録商標))を要求元の装置(本例ではコンピュータ100)に順次送信する。したがって、取得部130は、継続して最新の投稿情報を順次受信することになる。取得部130は、このようにして順次受信した情報を記憶部120に順次格納すればよい。また、個々の情報は、他の情報と区別可能な識別情報や区切情報を含む。また、個々の情報がツイート(登録商標)である場合、その情報は固有のURL(Uniform Resource Locator)で示される。従って、各ツイート(登録商標)のURLは、各情報を区別する識別情報として利用可能である。他の処理(例えば、第1のクラスタ処理部140による第1のクラスタ処理)では、このような情報を使って、情報の単位を認識できる。なお、取得部130が情報の単位を認識できるように、取得した情報にシリアル番号等の固有の識別

10

20

30

40

50

情報を付与して管理してもいいし、記憶部 120 における各情報の記憶アドレスを管理したテーブルを作成して情報単位を把握できるようにしてもよい。

【0055】

図5の(A)は、前述のようにして取得された処理対象情報の一例を示す。処理対象情報として、「A社PCを購入」、「A社のPC」、「A社の株価」、「A社株価上昇」の各テキスト情報が取得された例である。

【0056】

図7は、第1クラスタ処理の一例のフローチャートを示す。図7は、図3における第1クラスタ処理部140により実施される処理であって、処理対象情報に対する第1のクラスタ処理の一例を示す。なお、本例では、前述のように、第1のクラスタ処理のクラスタ手法をキーワード抽出処理としている。

10

【0057】

前述のように、第1クラスタ処理部140は、実行開始タイミングであるか否かを判定し、実行開始タイミングである場合に本処理の実行を開始する。なお、処理対象情報は複数個の情報を含むものとする。この実行開始タイミングは、例えば、前述のように、第1クラスタ処理部140による第1のクラスタ処理を行うために必要な情報量を取得したことをそのタイミングとすればよい。例えば、ツイート(登録商標)を処理対象情報とする場合は、その必要な情報量を数万件~数十万件に設定すればよい。この必要な情報量は、処理対象情報の内容や取得頻度などを基に適宜変更すればよい。

【0058】

まず、第1クラスタ処理部140は、処理対象情報の中からキーワードを抽出する(S3100)。次に、第1クラスタ処理部140は、抽出したキーワード別に、処理対象情報の各情報による集合を作成する(S3200)。そして、第1クラスタ処理部140は、キーワード数を基に処理対象情報の各情報のベクトル情報を作成する(S3300)。このベクトル情報は、処理対象情報の各情報の特徴を示す。

20

【0059】

図8は、第1クラスタ処理におけるキーワード抽出処理の一例のフローチャートを示す。この図8に示されるキーワード抽出処理は、処理対象情報の中のすべてのキーワードの抽出処理、およびキーワードとクラスタ番号の対応情報の作成処理を含む。

【0060】

まず、第1クラスタ処理部140は、記憶部120から処理対象情報のうちの1つの情報を取得する(S3110)。そして、第1クラスタ処理部140は、取得した情報について形態素解析を行う(S3120)。

30

【0061】

第1クラスタ処理部140は、形態素解析により得られた情報の形態素の中から名詞を判別し、名詞をキーワードの候補として抽出する(S3130)。なお、記憶部120には、確定されたキーワードがクラスタ番号と対応付けて登録される。例えば、キーワードとクラスタ番号の対応情報は、記憶部120において、テーブル形式や配列形式など種々の形態で記憶されうる。また、本例においては、クラスタ番号は、最小番号を1とした正の整数値である。

40

【0062】

第1クラスタ処理部140は、記憶部120に登録されている対応情報を参照して、キーワードの候補のうち記憶部120にキーワードとして登録されていない候補を特定する(S3140)。第1クラスタ処理部140は、特定した候補を新たなキーワードとし、クラスタ番号を更新しつつ(最終のクラスタ番号をインクリメントしつつ)、キーワードとクラスタ番号の新たな対応情報を記憶部120に登録する(S3150)。

【0063】

第1クラスタ処理部140は、処理対象情報の各情報を対象にして、以上説明したS3110~S3150の処理を実行する。その後、第1クラスタ処理部140は、前述の処理によって抽出したキーワードの個数を記憶部120に登録する(S3160)。以後、

50

この記憶部 120 に登録されたキーワードの個数をキーワードの個数情報と称する。例えば、第 1 クラスタ処理部 140 は、キーワードの個数情報として最終のクラスタ番号を記憶部 120 に登録すればよい。

【0064】

以上の処理によって、処理対象情報の中のすべてのキーワードの抽出、およびキーワードとクラスタ番号の対応情報の作成が完了する。

【0065】

図 9 は、第 1 クラスタ処理における集合作成処理の一例のフローチャートを示す。この図 9 に示される集合作成処理は、キーワードごと（クラスタ番号ごと）に、キーワードを含む情報の集合を作成する処理を含む。

10

【0066】

まず、第 1 クラスタ処理部 140 は、記憶部 120 に登録された対応情報の中から 1 つの対応情報を取得する（S3210）。なお、この取得処理が初回であれば、第 1 クラスタ処理部 140 は、クラスタ番号が 1 の対応情報を記憶部 120 から取得する。第 1 クラスタ処理部 140 は、取得した対応情報に含まれるキーワードが処理対象情報の各情報に含まれるか否かを判定する（S3220）。第 1 クラスタ処理部 140 は、処理対象情報の中でキーワードを含んでいる全ての情報をそのキーワードに対応する集合とする（S3230）。第 1 クラスタ処理部 140 は、作成された集合とクラスタ番号を対応付けて記憶部 120 に登録する。

【0067】

第 1 クラスタ処理部 140 は、記憶部 120 に登録されているすべての対応情報を対象にして、以上説明した S3210 ~ S3230 の処理を実行する。これによって、キーワード別に情報がまとめ上げられ、情報の集合が作成される。

20

【0068】

図 5 の（B）は、第 1 クラスタ処理部 140 による前述の処理結果の一例を示す。前述の処理によって記憶部 120 に登録された情報は、キーワード、テキスト集合、クラスタ番号を含む。なお、図 5 においては、前述の集合に対応するテキスト情報群をテキスト集合と称している。本図は、図 8 に示すキーワード抽出処理によって、処理対象情報から、「A 社」、「PC」、「株価」の各キーワードが抽出された例を示す。また、本図は、そのキーワード抽出処理によって、キーワード「A 社」にクラスタ番号「1」、キーワード「PC」にクラスタ番号「2」、キーワード「株価」にクラスタ番号「3」が対応付けられたことを示す。また、本図は、図 9 に示す集合作成処理によって、各キーワードに該キーワードを含むテキスト情報が対応付けて登録されたことを示す。図 5 の（B）においては、第 1 クラスタ処理部 140 の処理によって 8 つのエントリが記憶部 120 に登録されていることを示す。

30

【0069】

図 10 は、第 1 クラスタ処理におけるベクトル情報作成処理の一例のフローチャートを示す。この図 10 に示されるベクトル情報作成処理は、処理対象情報の各情報についてベクトル情報を作成する処理を含む。

【0070】

第 1 クラスタ処理部 140 は、キーワードの個数を特定する（S3310）。第 1 クラスタ処理部 140 は、例えば、図 8 のキーワード抽出処理において記憶部 120 に登録されたキーワードの個数情報を取得すればよい。なお、記憶部 120 にキーワードの個数情報が登録されていない場合、第 1 クラスタ処理部 140 は、図 8 のキーワード抽出処理によって作成され記憶部 120 に登録されている対応情報に含まれるクラスタ番号の中から最大のクラスタ番号を取得することによってキーワードの個数を特定できる。

40

【0071】

第 1 クラスタ処理部 140 は、各キーワードを要素とした処理対象情報の多次元配列を記憶部 120 に作成する（S3320）。例えば、処理対象情報の情報数を  $i$ 、 $n(i)$  を  $i$  番目の処理対象の情報に含まれるキーワード数とすると、その配列  $x_i$  は  $x_i = (x_{i1}, x_{i2}, \dots, x_{in(i)})$

50

$x_{i,1}, x_{i,2}, \dots, x_{i,n(i)}$  )と表現できる。なお、 $x_{i,j}$  ( $1 \leq j \leq n(i)$ ) は、各キーワードに対応する。クラスタ番号を示し、クラスタ番号対応のキーワードを意味する。第1クラスタ処理部140は、処理対象情報の全ての情報について処理した後、各キーワードが出現する情報数を計算し、指定された閾値以上の出現回数のキーワードを選択し、選択されたキーワードの種類数をクラスタの数  $k$  とする。また、第1クラスタ処理部140は、同じキーワードを含む情報の集合を1つのクラスタとする。また、第1クラスタ処理部140は、選択された各キーワードに数字を付与しクラスタ番号とする。これらクラスタ数とクラスタ情報果は、第2クラスタ処理部150にて利用される。

【0072】

第1クラスタ処理部140は、作成した配列  $x_i$  を使い、処理対象情報の各情報のベクトル情報を作成する (S3330)。その後、第1クラスタ処理部140は、この配列  $x_i$  の各要素であるキーワードが対応するベクトルの次元に対して値を埋める。ベクトルの各次元の値としては、例えば、2値情報 (例えば、出現しない場合は値0、出現する場合は値1)、出現頻度、TF・IDFのような重みづけ手法などを用いればよい。この処理を処理対象情報の各情報について実施することで、各情報に対応するベクトル情報が作成される ( $x_i$  のベクトルが完成する)。第1クラスタ処理部140は、作成したベクトル情報を記憶部120に登録する。

10

【0073】

次に、第2クラスタ処理の一例を説明する。

【0074】

図11は、第2クラスタ処理の一例のフローチャートを示す。図11は、図3における第2クラスタ処理部150により実施される処理であって、第1のクラスタ処理の処理結果を用いた第2のクラスタ処理の一例を示す。なお、本例では、前述のように、第2のクラスタ処理のクラスタ手法を  $K$ -means としている。

20

【0075】

前述のように、第2クラスタ処理部150は、第1のクラスタ処理が行われたか否かを判定し、実行開始タイミングである場合に本処理の実行を開始する。

【0076】

まず、第2クラスタ処理部150は、記憶部120に格納されているキーワードの個数情報および各情報のベクトル情報を取得する (S610)。このキーワードの個数情報は、図8に示すキーワード抽出処理において、第1クラスタ処理部140により作成され、記憶部120に登録されたものである。また、各情報のベクトル情報は、図10に示すベクトル情報作成処理において、第1クラスタ処理部140により作成され、記憶部120に登録されたものである。

30

【0077】

第2クラスタ処理部150は、S610にて取得した情報を用いて第2クラスタ処理の一例である  $K$ -means のクラスタ処理における初期設定を行う (S620)。

【0078】

ここで、一般的な  $K$ -means によるクラスタ処理を説明する。

【0079】

例えば、 $K$ -means は以下の式1で示される目的関数を最小化する分割最適化クラスタ処理である。

40

式1:

$$\text{Err}(\{X_i\}) = \sum_i^k \sum_{x \in X_i} \|x - \bar{x}_i\|^2$$

【0080】

50

ここで、 $X$  は、データ集合であり、ベクトルで表現されたデータ  $x$  の集合である。

$$\{X_i\}$$

は  $k$  個のクラスタ、

$$X_i$$

は  $i$  番目のクラスタであり、データ集合の網羅的で互いに疎な部分集合である。また、 $k$  はクラスタ数であり、

10

$$x_i$$

はセントロイド（重心）である。

【0081】

上記を実現するために、一般的な  $K$ -means は、以下に示す処理を行う。

【0082】

まず、初期設定では、データ集合をランダムに  $k$  個のクラスタに分割し、それを初期クラスタとする。

20

【0083】

その後、各クラスタについて以下の式 2 で示されるセントロイドの計算を行う。  
式 2 :

$$x_i = \frac{1}{|x_i|} \sum_{x \in X_i} x$$

30

【0084】

$$|x_i|$$

はクラスタ

$$X_i$$

に含まれるデータ数である。

40

【0085】

続いて、各データ

$$x \in X$$

において、各クラスタのセントロイド

—  
 $x_i$

との距離

$x \in X$

を計算し、距離が最小であるクラスタ

$X_i$

10

を見つけ、データをそのクラスタに割り当てる。このようにして、全てのデータがクラスタに割り当てると、式 2 においてセントロイドを更新する。

【 0 0 8 6 】

前述の各クラスタのセントロイド計算とクラスタへの割り当て処理は、反復数が設定回数に達するまで繰り返される。反復数が設定回数に達した場合、繰り返し処理が終了され、本アルゴリズムの最終処理結果として  $k$  個のクラスタ

$\{X_i\}$

20

が出力される。なお、反復数が設定回数に達する前に、クラスタの割り当てが前回の反復時のクラスタの割り当てから変化がなかった場合に、繰り返し処理が終了されるようにしてもよい。このアルゴリズムによる計算量は、データ数を  $N$ 、反復回数を定数とすると、 $O(Nk)$ となる。

【 0 0 8 7 】

一般的に、 $K$ - $means$ では、初期クラスタを何回か変更して各初期クラスタにて前述のアルゴリズムを実行し、それぞれのアルゴリズム実行において前述の目的関数を最小化する分割を選択する。そして、この数回のアルゴリズムの実行に基づいて、大域最適に近い解の探索が行われる。

30

【 0 0 8 8 】

このように、 $K$ - $means$ を用いた一般的なクラスタ処理は、初期化においてデータ集合をランダムにクラスタ分割する。また、このようなデータ集合をランダムにクラスタ分割する初期クラスタ処理が何度か行われ（つまり、初期クラスタを何度か変更して）、前述のアルゴリズムが実行され解の探索が行われる。このため、計算量が多く、最終的な処理結果を得るために長い処理時間が必要とされる。

【 0 0 8 9 】

また、前述のように計算量はデータ数に比例することから、例えば処理対象情報が大規模データであると、非常に長い処理時間が必要とされる。そのため、処理対象情報を取得し、リアルタイムにクラスタ情報を提供するというような処理が困難となる。

40

【 0 0 9 0 】

本実施形態では、初期設定に第 1 クラスタ処理部 1 4 0 の処理結果を用いる。つまり、第 2 クラスタ処理部 1 5 0 は、前述の目的関数における初期クラスタのクラスタ数  $k$  に、前述で取得されたキーワードの個数情報を設定し、各クラスタにキーワード対応のクラスタ番号を付与する。また、第 2 クラスタ処理部 1 5 0 は、前述の目的関数におけるデータ集合  $X$  に、前述で取得された各情報のベクトル情報を設定する。また、本実施形態では、第 1 クラスタ処理部 1 4 0 で決定したクラスタ情報を基に  $K$ - $means$ のセントロイドの初期値の計算を行う。このセントロイドの初期値の計算は、通常の  $K$ - $means$ と同

50

様に前述の式 2 により行われるが、一般的な K - m e a n s での処理のように初期のセントロイド計算のためのクラスタ作成がランダムで実施されるのではなく、本実施形態では、第 1 クラスタ処理部 1 4 0 で決定されたクラスタ情報を基に計算を行う。これにより、第 1 クラスタ処理部での処理結果が第 2 クラスタ処理部 1 5 0 に引き継がれる。

【 0 0 9 1 】

このように初期設定が完了した後、第 2 クラスタ処理部 1 5 0 は、K - m e a n s によるクラスタ処理を実行する ( S 6 3 0 )。この第 2 クラスタ処理部 1 5 0 における K - m e a n s によるクラスタ処理では、前述の繰り返し処理が実行され、反復数が設定回数に達した場合に繰り返し処理を終了する。なお、前述のように、設定回数に達する前に、今回のクラスタの割り当てが前回の反復によるクラスタの割り当てから変化がなかった場合に、繰り返し処理が終了されるようにしてもよい。

10

【 0 0 9 2 】

前述のように K - m e a n s を用いた一般的なクラスタ処理では何度か初期クラスタを変更して前述のアルゴリズムが実行されるが、本実施形態においては、第 2 クラスタ処理部 1 5 0 は前述のキーワードの個数情報を用いて決定した初期クラスタにて前述のアルゴリズムを 1 度だけ実行する。このため、本実施形態によれば、一般的な K - m e a n s によるクラスタ処理と比較して、計算量が少なくなり処理時間の短縮が図れる。

【 0 0 9 3 】

第 2 クラスタ処理部 1 5 0 は、K - m e a n s による前述のクラスタ処理の結果を記憶部 1 2 0 に登録する ( S 6 4 0 )。このクラスタ処理の結果は前述のキーワードに対応付けられた各情報に対するクラスタ番号であり、第 2 クラスタ処理部 1 5 0 は、この処理結果であるクラスタ番号を各情報に対応づけて記憶部 1 2 0 に登録する。

20

【 0 0 9 4 】

図 5 の ( C ) は、第 2 クラスタ処理部 1 5 0 による前述の処理結果の一例を示す。本図において、テキスト集合、クラスタ番号の情報は、前述の第 1 クラスタ処理部 1 4 0 による処理結果を示す。本図における新クラスタ番号は、第 2 クラスタ処理部 1 5 0 による前述の処理によって得られ、テキスト集合の各テキスト情報に対応付けて登録されたクラスタ番号である。

【 0 0 9 5 】

図 1 2 は、決定処理の一例のフローチャートを示す。図 1 2 は、図 3 における決定部 1 6 0 により実施される処理であって、第 1 クラスタ処理部 1 4 0 と第 2 クラスタ処理部 1 5 0 のそれぞれの処理結果に基づき、処理対象に対する最終のクラスタ処理結果の決定方法の一例を示す。

30

【 0 0 9 6 】

まず、第 2 クラスタ処理部 1 5 0 による処理が終了すると、決定部 1 6 0 は、前述のようにして登録された第 1 クラスタ処理部 1 4 0 と第 2 クラスタ処理部 1 5 0 の各処理結果のうち 1 つの対の情報を記憶部 1 2 0 から取得する ( S 7 1 0 )。この対の情報は、第 1 クラスタ処理部 1 4 0 の処理によって得られたクラスタ番号と第 2 クラスタ処理部 1 5 0 の処理によって得られたクラスタ番号である。

【 0 0 9 7 】

そして、決定部 1 6 0 は、取得した 2 つのクラスタ番号を比較し、第 2 クラスタ処理部 1 5 0 のクラスタ処理により得られたクラスタ番号が第 1 クラスタ処理部 1 4 0 の処理により得られたクラスタ番号から変更されているか判定する ( S 7 2 0 )。クラスタ番号が変更されている場合、決定部 1 6 0 は、それらクラスタ番号と、それらに対応する情報を記憶部 1 2 0 から削除する。

40

【 0 0 9 8 】

以上の処理を図 5 の ( C ) を用いて以下に説明する。

【 0 0 9 9 】

決定部 1 6 0 は、図 5 の ( C ) に示す記憶部 1 2 0 の登録情報のなかから、テキスト情報に対応するクラスタ番号 ( 第 1 クラスタ処理部 1 4 0 の処理結果 ) と新クラスタ番号 (

50

第2クラスタ処理部150)を取得する。例えば、図5の(C)において最初のエン트리であるテキスト情報「A社のPCを購入」を例にすると、決定部160は、そのテキスト情報に対応するクラスタ番号「1」と新クラスタ番号「2」を取得する。この取得された2つのクラスタ番号は異なるため、決定部160は、テキスト情報「A社のPCを購入」、このテキスト情報に対応するクラスタ番号「1」および新クラスタ番号「2」を記憶部120から削除する。

#### 【0100】

決定部160は、記憶部120に登録された第1クラスタ処理部140と第2クラスタ処理部150の各対の情報(図5の(C)における各エン트리)を対象にして、以上説明したS710~S730の処理を実行する。

10

#### 【0101】

すべての情報を対象に処理を実行した後、決定部160は、削除されず記憶部120に残っている各情報(各エン트리)を処理結果として記憶部120の処理結果格納領域に登録する(S740)。なお、すべての情報を対象にした処理の実行後ではなく、決定部160は、クラスタ番号の比較においてクラスタ番号が一致すると判定された際に、クラスタ番号の対およびそれに対応する情報を記憶部120の処理結果格納領域に登録するようにしてもよい。

#### 【0102】

前述のように、決定部160は、第1クラスタ処理部140による処理結果と第2クラスタ処理部150による処理結果との間で、クラスタ番号が不一致となった情報を削除する。このクラスタ番号が不一致となった情報は、処理対象情報についてクラスタを決定する上で不安定な情報といえる。例えば、大規模データは多様な情報内容を含む。このような多様な情報内容のデータを処理対象情報としてクラスタ処理を行う場合、クラスタを決定する上で不安定な情報が非常に多く出現することになる。

20

#### 【0103】

本実施形態は、クラスタを決定する上で不安定な情報を特定し、先行技術のように不安定な情報を詳細分類するのではなく、不安定な情報を削除してクラスタ処理の対象外とする。このような処理によって、本実施形態は高速なクラスタ処理を実現する。本実施形態は、多様な情報内容を含む例えば大規模データを対象にしたクラスタ処理に特に有効である。

30

#### 【0104】

図5の(C)では、8つのエントリのうち上部の4つのエントリにおいて、第1クラスタ処理部140の処理によって得られたクラスタ番号と、第2クラスタ処理部150の処理によって得られた新クラスタ番号とが相違する。したがって、図5の(C)に示す例の場合、この4つのエントリが決定部160によって削除される。図5の(D)は、決定部160の処理によって、クラスタ番号が不一致のエントリが削除された結果を示す。言い換えれば、図5の(D)は、決定部160の処理によって、記憶部120の処理結果格納領域に登録されたエントリを示す。

#### 【0105】

決定部160によって記憶部120の処理結果格納領域に登録された最終クラスタ処理結果は、出力部170によって出力される。出力部170は、記憶部120の処理結果格納領域に登録されている最終クラスタ処理結果を、例えば、そのまま出力してもいいし、クラスタを区別可能にして各情報を出力してもよい。出力部170による出力形態は、要求元の装置や出力先の装置、処理結果の使用形態に応じて適宜変更すればよい。

40

#### 【0106】

図5の(E)は、図5の(A)に示す処理対象情報について本実施形態によるクラスタ処理が実行された結果を示す。図5の(E)は、図5の(D)の情報内容に基づくものである。図5の(D)に示されるように、処理対象情報の各テキスト情報である、「A社PCを購入」、「A社のPC」、「A社の株価」、「A社株価上昇」は、「A社PCを購入」と「A社のPC」に同じクラスタ番号「2」、「A社の株価」と「A社株価上昇」に同

50

じクラスタ番号「3」が付与されている。つまり、図5の(E)に示すように、「A社PCを購入」と「A社のPC」が1つのクラスタ、「A社の株価」と「A社株価上昇」が他のクラスタとなる。

【0107】

以上説明したように、本実施形態では、K-meansを例にした第2のクラスタ処理よりも高速な(短時間で処理結果を得られる)キーワード抽出処理を例にした第1のクラスタ処理にてまず処理対象情報を分類し(クラスタ処理し)、その分類結果を用いて第2のクラスタ処理にて分類した(クラスタ処理した)後、双方の分類結果で相違する情報を除外した分類結果を処理対象情報に対する最終処理結果と決定する。つまり、本実施形態は、処理スピード重視で第1のクラスタ処理にて大雑把な分類を行い、精度重視で第2のクラスタ処理にて分類した結果を用いて第1のクラスタ処理の結果の精度を評価し、精度の悪い情報を排除している。このようにすることで、本実施形態は、高速なクラスタ処理を実現しつつ、クラスタの精度を維持している。

10

【0108】

なお、前述した本実施形態では、第1のクラスタ処理をキーワード抽出処理、第2のクラスタ処理をK-meansとしたが、それぞれのクラスタ処理はこれに限定されるものではない。例えば、第2のクラスタ処理は混合正規分布推定であってもよい。

【0109】

以上の実施形態に関し、更に以下の付記を開示する。

(付記1)

20

コンピュータが、

第2のクラスタ手法と異なりかつ該第2のクラスタ手法よりも高速な第1のクラスタ手法によって対象情報に対して第1の分類処理を実行し、

前記第1の分類処理の結果に基づいて前記第2のクラスタ手法によって第2の分類処理を実行し、

前記第1の分類処理の結果と前記第2の分類処理の結果とに基づいて前記対象情報の分類を決定する、

ことを特徴とするクラスタ処理方法。

(付記2)

前記対象情報の分類を決定することは、前記第1の分類処理の結果と前記第2の分類処理の結果間で分類の変更がなかった前記第1の分類処理の結果または前記第2の分類処理の結果に基づいて前記対象情報の分類を決定する、

ことを特徴とする付記1のクラスタ処理方法。

30

(付記3)

前記第1の分類処理の結果は、前記対象情報に関連した分類個数情報を含み、

前記第2の分類処理は、前記分類個数情報を前記第2のクラスタ手法における初期情報として実行される、

ことを特徴とする付記1または2のクラスタ処理方法。

(付記4)

前記第1の分類処理の結果は、前記対象情報に含まれるキーワードを基にした前記対象情報の特徴情報を含み、

前記第2の分類処理は、前記対象情報の特徴情報を前記第2のクラスタ手法における前記初期情報として実行される、

ことを特徴とする付記3のクラスタ処理方法。

40

(付記5)

前記第1の分類処理は、前記対象情報を用いた前記キーワード別の集合を作成し、異なる前記キーワードに対応する前記集合を異なる分類とする、

ことを特徴とする付記4のクラスタ処理方法。

(付記6)

前記第2のクラスタ手法は、処理対象情報の類似性を基に処理対象情報を分類する手法

50

である、

ことを特徴とする付記 1 乃至 5 のいずれか 1 つのクラスタ処理方法。

(付記 7)

第 2 のクラスタ手法と異なりかつ該第 2 のクラスタ手法よりも高速な第 1 のクラスタ手法によって対象情報に対して第 1 の分類処理を実行する第 1 処理部と、

前記第 1 の分類処理の結果に基づいて前記第 2 のクラスタ手法によって第 2 の分類処理を実行する第 2 処理部と、

前記第 1 の分類処理の結果と前記第 2 の分類処理の結果とに基づいて前記対象情報の分類を決定する決定部と、

を有することを特徴とするクラスタ処理装置。

10

(付記 8)

前記決定部は、前記第 1 の分類処理の結果と前記第 2 の分類処理の結果間で分類の変更がなかった前記第 1 の分類処理の結果または前記第 2 の分類処理の結果に基づいて前記対象情報の分類を決定する、

ことを特徴とする付記 7 のクラスタ処理装置。

(付記 9)

前記第 1 処理部による前記第 1 の分類処理の結果は、前記対象情報に関連した分類個数情報を含み、

前記第 2 処理部は、前記分類個数情報を前記第 2 のクラスタ手法における初期情報として前記第 2 の分類処理を実行する、

20

ことを特徴とする付記 7 または 8 のクラスタ処理装置。

(付記 10)

前記第 1 処理部による前記第 1 の分類処理の結果は、前記対象情報に含まれるキーワードを基にした前記対象情報の特徴情報を含み、

前記第 2 処理部は、前記対象情報の特徴情報を前記第 2 のクラスタ手法における前記初期情報として前記第 2 の分類処理を実行する、

ことを特徴とする付記 9 のクラスタ処理方法。

(付記 11)

前記第 1 処理部は、前記対象情報を用いた前記キーワード別の集合を作成し、異なる前記キーワードに対応する前記集合を異なる分類とする、

30

ことを特徴とする付記 10 のクラスタ処理装置。

(付記 12)

前記第 2 のクラスタ手法は、処理対象情報の類似性を基に処理対象情報を分類する手法である、

ことを特徴とする付記 7 乃至 11 のいずれか 1 つのクラスタ処理装置。

(付記 13)

コンピュータに、

第 2 のクラスタ手法と異なりかつ該第 2 のクラスタ手法よりも高速な第 1 のクラスタ手法によって対象情報に対して第 1 の分類処理を実行し、

前記第 1 の分類処理の結果に基づいて前記第 2 のクラスタ手法によって第 2 の分類処理を実行し、

40

前記第 1 の分類処理の結果と前記第 2 の分類処理の結果とに基づいて前記対象情報の分類を決定する、

処理を実行させるためのプログラム。

(付記 14)

前記対象情報の分類を決定することは、前記第 1 の分類処理の結果と前記第 2 の分類処理の結果間で分類の変更がなかった前記第 1 の分類処理の結果または前記第 2 の分類処理の結果に基づいて前記対象情報の分類を決定する、

ことを特徴とする付記 13 のプログラム。

(付記 15)

50

前記第 1 の分類処理の結果は、前記対象情報に関連した分類個数情報を含み、

前記第 2 の分類処理は、前記分類個数情報を前記第 2 のクラスタ手法における初期情報として実行される、

ことを特徴とする付記 1 3 または 1 4 のプログラム。

(付記 1 6)

前記第 1 の分類処理の結果は、前記対象情報に含まれるキーワードを基にした前記対象情報の特徴情報を含み、

前記第 2 の分類処理は、前記対象情報の特徴情報を前記第 2 のクラスタ手法における前記初期情報として実行される、

ことを特徴とする付記 1 5 のプログラム。

(付記 1 7)

前記第 1 の分類処理は、前記対象情報を用いた前記キーワード別の集合を作成し、異なる前記キーワードに対応する前記集合を異なる分類とする、

ことを特徴とする付記 1 6 のプログラム。

(付記 1 8)

前記第 2 のクラスタ手法は、処理対象情報の類似性を基に処理対象情報を分類する手法である、

ことを特徴とする付記 1 3 乃至 1 7 のいずれか 1 つのプログラム。

【符号の説明】

【0 1 1 0】

1 : コンピュータシステム

1 0 0、4 0 0、5 0 0 : コンピュータ

2 0 0 : ネットワーク

3 0 0 : サーバ

3 2 0 : ストレージシステム

6 0 0 : N A S

1 0 : プロセッサ

2 0 : R A M

3 0 : ドライブ装置

3 2 : 記憶媒体

4 0 : 入力インターフェース

4 2 : 入力デバイス

5 0 : 出力インターフェース

5 2 : 出力デバイス

6 0 : 通信インターフェース

7 0 : バス

1 1 0 : 制御部

1 2 0 : 記憶部

1 3 0 : 取得部

1 4 0 : 第 1 クラスタ処理部

1 5 0 : 第 2 クラスタ処理部

1 6 0 : 決定部

1 7 0 : 出力部

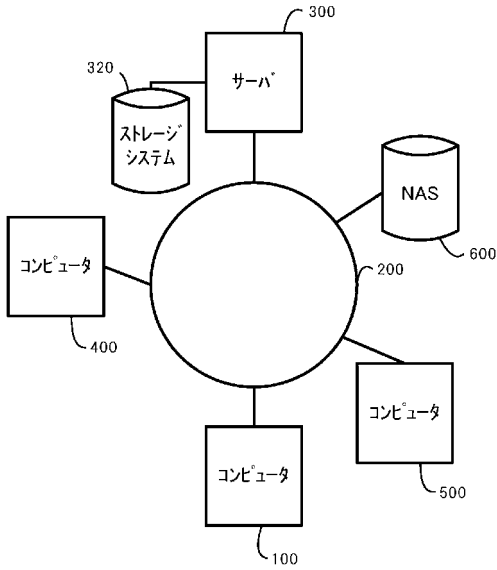
10

20

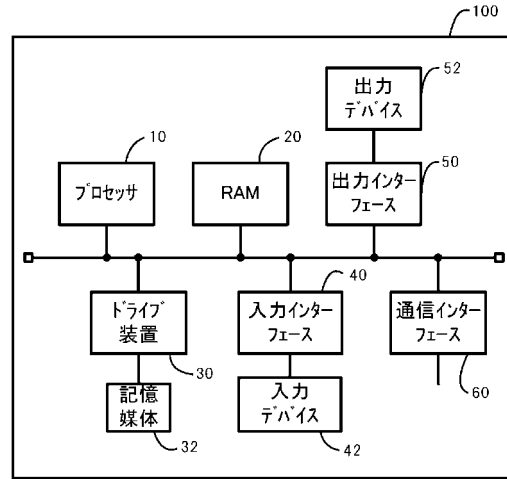
30

40

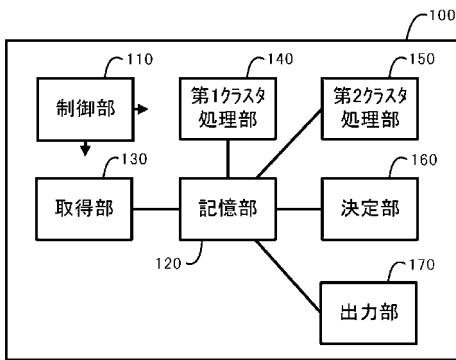
【 図 1 】



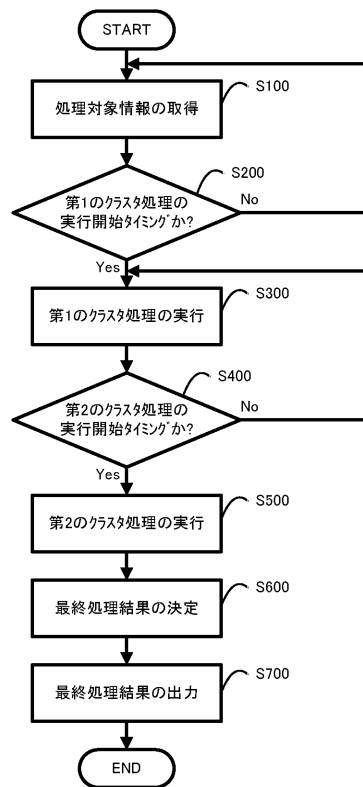
【 図 2 】



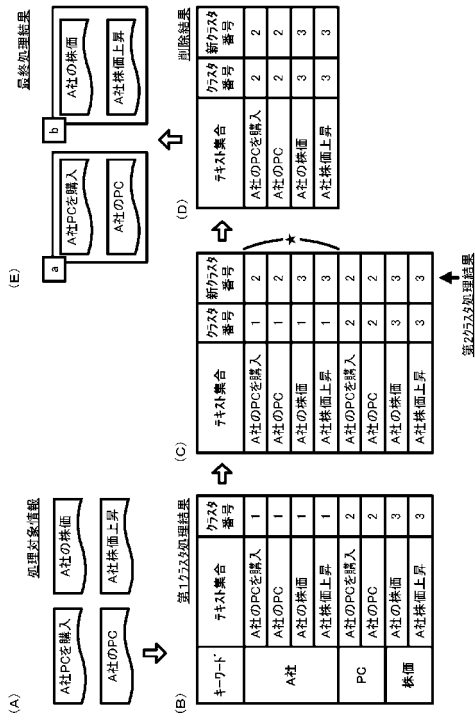
【 図 3 】



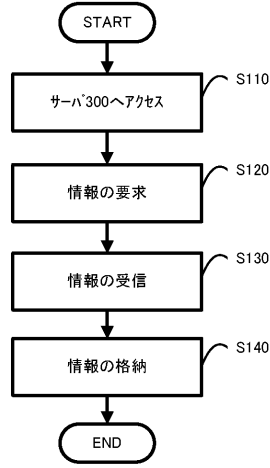
【 図 4 】



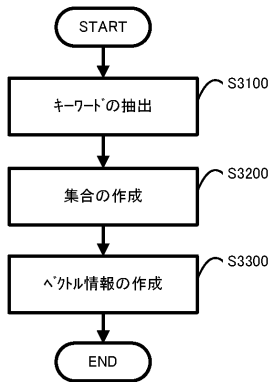
【 図 5 】



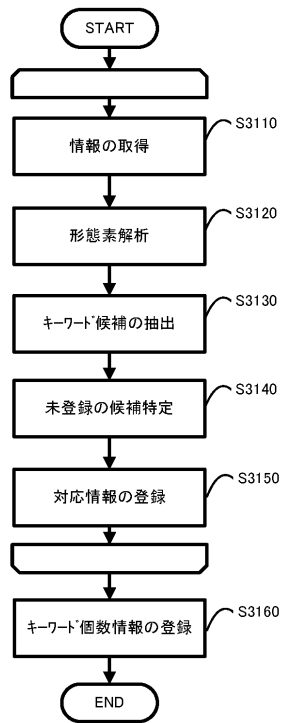
【 図 6 】



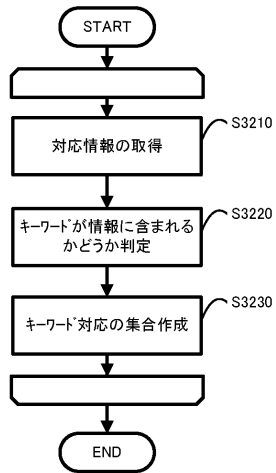
【 図 7 】



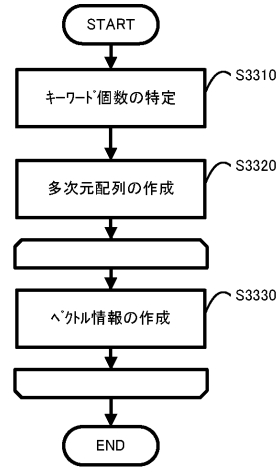
【 図 8 】



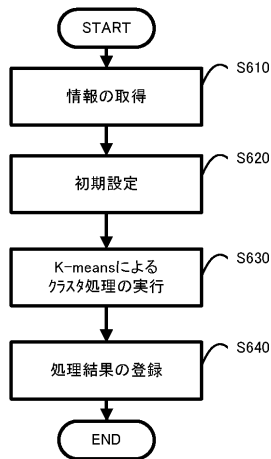
【 図 9 】



【 図 1 0 】



【 図 1 1 】



【 図 1 2 】

