

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6237334号
(P6237334)

(45) 発行日 平成29年11月29日(2017.11.29)

(24) 登録日 平成29年11月10日(2017.11.10)

(51) Int.Cl.

F I

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 2 1 0 D

G 0 6 F 17/30 3 5 0 C

請求項の数 11 (全 26 頁)

(21) 出願番号 特願2014-36700 (P2014-36700)
 (22) 出願日 平成26年2月27日(2014.2.27)
 (65) 公開番号 特開2015-162076 (P2015-162076A)
 (43) 公開日 平成27年9月7日(2015.9.7)
 審査請求日 平成28年11月2日(2016.11.2)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100094525
 弁理士 土井 健二
 (74) 代理人 100094514
 弁理士 林 恒徳
 (72) 発明者 高橋 哲朗
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内
 審査官 樋口 龍弥

最終頁に続く

(54) 【発明の名称】 クエリ生成方法、クエリ生成プログラム、及び、クエリ生成装置

(57) 【特許請求の範囲】

【請求項 1】

処理ユニットが、

入力された検索語に基づいて複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示する第1の工程と、

前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成する第2の工程と、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示する第3の工程と、

を実行することを特徴とするクエリ生成方法。

【請求項 2】

請求項 1 において、

前記第2の工程は、前記抽出した文字列毎に、前記排除文書集合内の当該文字列の頻度に応じて、前記検索された複数の文書における前記排除文書集合内での前記出現分布値を加算して前記スコアを算出するクエリ生成方法。

【請求項 3】

10

20

請求項 2 において、

前記第 2 の工程は、前記抽出した文字列毎に、さらに、前記検索された複数の文書における前記排除文書集合以外の文書集合内での前記出現分布値を前記スコアから減算して、前記スコアを算出するクエリ生成方法。

【請求項 4】

請求項 1 乃至 3 のいずれかにおいて、

前記第 3 の工程は、前記文字列リスト内の文字列数に比例する値を示すオブジェクトを前記表示ユニットに表示し、前記オブジェクトに対する前記ユーザの入力に比例する、前記文字列リストの上位数分の前記文字列を選択するクエリ生成方法。

【請求項 5】

請求項 4 において、

前記第 3 の工程は、前記オブジェクトに対する入力に応じて選択される前記文字列を前記表示ユニットに更に表示するクエリ生成方法。

【請求項 6】

請求項 4 または 5 において、

前記オブジェクトは、スライドバーであるクエリ生成方法。

【請求項 7】

請求項 1 乃至 6 のいずれかにおいて、

前記第 2 の工程は、前記文書集合の特徴を示す特徴情報を前記表示ユニットに更に表示するクエリ生成方法。

【請求項 8】

請求項 7 において、

前記特徴情報は、前記文書集合の文書に含まれる前記検索語の使用文字列、前記文書集合の文書に含まれる頻出文字列、前記文書集合の文書に含まれる文字列であって前記検索された複数の文書における前記書集合内での出現分布率が高い文字列、のうち少なくともいずれかであるクエリ生成方法。

【請求項 9】

請求項 1 乃至 8 のいずれかにおいて、

前記第 2 の工程は、前記排除文書集合の指定を受け付けるオブジェクトを更に前記表示ユニットに表示するクエリ生成方法。

【請求項 10】

入力された検索語に基づいて複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示し、

前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成し、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示する、

処理をコンピュータに実行させるクエリ生成プログラム。

【請求項 11】

処理ユニットと、

複数の文書を記憶する記憶装置と、

表示装置と、を有し、

入力された検索語に基づいて前記複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示し、前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数

10

20

30

40

50

の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成し、前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示するクエリ生成装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、クエリ生成方法、クエリ生成プログラム、及び、クエリ生成装置に関する。

【背景技術】

【0002】

マーケティング等の目的のために、ソーシャルメディアから情報を得るというニーズが増えている。例えば、大量にあるソーシャルメディアの文書（記事）のうち、一部の文書の集合がマーケティング等の分析の対象となる。

【0003】

分析の対象となる文書集合を選択するために、文書検索技術が使われる。例えば、マーケティング担当者は、所定の文書を検索するためクエリ（検索語、検索条件）をソーシャルメディアの文書を格納するデータベースに指定することによって、特定の条件に合致する文書を検索する。これにより、マーケティング担当者は、分析の対象となる文書集合を抽出することができる。マーケティング担当者は、分析の対象となる文書集合を選択するために適切なクエリを設定して発行する。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2008-77137号公報

【特許文献2】特開2006-251935号公報

【特許文献3】特開2012-84029号公報

【特許文献4】特開平11-272709号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、マーケティング担当者の所望の文書集合を選択するためのクエリの生成は容易ではない。例えば、「セブンイレブン（登録商標）」の省略形である「セブン」をクエリとする場合、「ウルトラセブン（登録商標）」や「セブンスター（登録商標）」、映画の「セブン」に関する文書も合わせて抽出される。抽出された文書集合に基づいて分析を行う場合、「セブンイレブン」に関する文書以外の文書が含まれることにより、分析の精度が下がる。したがって、抽出された文書集合内の「セブンイレブン」以外に関する文書は、少ないことが望ましい。

【0006】

1つの側面は、本発明は、ユーザの目的に適合する文書を抽出するクエリを生成するクエリ生成方法、クエリ生成プログラム、及び、クエリ生成装置を提供する。

【課題を解決するための手段】

【0007】

第1の側面は、処理ユニットが、入力された検索語に基づいて複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示する第1の工程と、前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成する第2の工程と、前記文字列リスト内の文字列数に比例する入力に応じて、前

10

20

30

40

50

記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示する第3の工程と、を実行する。

【発明の効果】

【0008】

第1の側面によれば、ユーザの選択により、ユーザの目的に適合する文書を抽出するクエリを生成する。

【図面の簡単な説明】

【0009】

【図1】本実施形態における文書検索システムの構成を示す図である。

10

【図2】図1に示す本実施の形態例における検索サーバのハードウェア構成を説明する図である。

【図3】図1、図2に示す検索サーバのソフトウェアブロック図である。

【図4】クラスタを用いるクエリの基本的な生成処理を説明する第1の図である。

【図5】クラスタを用いたクエリの基本的な生成処理を説明する第2の図である。

【図6】本実施の形態例におけるクエリ生成処理を説明するフローチャート図である。

【図7】本実施の形態例におけるクラスタの生成処理を説明する図である。

【図8】クラスタの指定を受け付けるクライアント装置の画面の一例を示す図である。

【図9】キーワードリストの生成処理を説明する図である。

【図10】キーワードリストの生成処理を説明するフローチャート図である。

20

【図11】キーワード毎のtfidf値を例示する図である。

【図12】キーワード生成処理の具体例を説明する第1の図である。

【図13】キーワード生成処理の具体例を説明する第2の図である。

【図14】キーワード生成処理の具体例を説明する第3の図である。

【図15】キーワード生成処理の具体例を説明する第4の図である。

【図16】具体例におけるキーワードリストを示す図である。

【図17】本実施の形態例におけるキーワードの選択処理を説明する図である。

【図18】スライダーを表示するクライアント装置の画面の一例を示す図である。

【図19】排除される文書の割合を表示するクライアント装置の画面の一例を示す図である。

30

【発明を実施するための形態】

【0010】

以下、図面にしたがって本発明の実施の形態について説明する。ただし、本発明の技術的範囲はこれらの実施の形態に限定されず、特許請求の範囲に記載された事項とその均等物まで及ぶものである。

【0011】

〔文書検索システム〕

図1は、本実施形態における文書検索システムの構成を示す図である。図1に示すように、本実施形態における文書検索システムは、クライアント装置80と検索サーバ（クエリ生成装置）10とを有する。クライアント装置80は、検索サーバ10と通信ネットワーク50を介して接続する。クライアント装置80は、例えば、パーソナルコンピュータ等である。なお、図1の例では、1台のクライアント装置80を図示しているが、検索サーバ10は、複数のクライアント装置80と接続してもよい。

40

【0012】

クライアント装置80は、Webページを閲覧するブラウザを介して、検索サーバ10に、検索語等の検索条件を入力する。検索サーバ10は、クライアント装置80から入力された検索条件に基づいてクエリを生成し、文書集合からクエリに対応する文書を抽出する。そして、検索サーバ10は、抽出した文書をクライアント装置80に送信する。本実施の形態例における検索サーバ10は、クライアント装置80からの検索条件や入力情報に基づいて、適切なクエリを生成する。クライアント装置80は、例えば、CPU（Cent

50

ral Processing Unit)、メモリ、表示ユニット、入力ユニット等を有する(不図示)。

【0013】

[検索サーバの構成]

図2は、図1に示す本実施の形態例における検索サーバ10のハードウェア構成を説明する図である。図2に示す検索サーバ10は、例えば、CPU(Central Processing Unit)101、RAM(Random Access Memory)201や不揮発性メモリ202等を備えるメモリ102、通信インタフェース部103を有する。各部は、バス104を介して相互に接続する。

【0014】

CPU101は、バス104を介してメモリ102等と接続すると共に、検索サーバ10全体の制御を行う。メモリ102のRAM201は、CPU101が処理を行うデータ等を記憶する。メモリ102の不揮発性メモリ202は、CPU101が実行するOS(Operating System)のプログラムを格納する領域(図示せず)や、本実施の形態例におけるクエリ生成プログラムを格納する領域210を備える。また、不揮発性メモリ202は、文書集合を格納する領域20を有する。文書集合を格納する領域(以下、文書集合と称する)20は、例えば、ソーシャルネットワークの記事の集合である。記事は、例えば、ブログの記事やコメント等である。不揮発性メモリ202は、HDD(Hard disk drive)、不揮発性半導体メモリ等によって構成される。

【0015】

クエリ生成プログラム領域210のクエリ生成プログラム(以下、クエリ生成プログラム210と称する)は、CPU101の実行によって、本実施の形態例におけるクエリ生成処理を実現する。また、通信インタフェース部103は、ネットワーク50を介して、クライアント装置80等の通信機器との間でデータの送受信を制御する。

【0016】

図3は、図1、図2に示す検索サーバ10のソフトウェアブロック図である。検索サーバ10のクエリ生成プログラム210(図2)は、例えば、文書検索モジュール(以下、文書検索部と称する)11、クラスタリングモジュール(以下、クラスタリング部と称する)12、キーワード生成モジュール(以下、キーワード生成部と称する)13を有する。

【0017】

文書検索部11は、ユーザが入力した検索語をクライアント装置80から受信し(a1)、検索語に基づいて検索サーバ10が格納する文書集合20を検索して、複数の文書を抽出する。なお、図2、図3の例では、検索対象となる文書集合20は、検索サーバ10に格納される。しかしながら、この例に限定されるものではない。文書集合20は、ネットワーク50を介して接続する1つまたは複数の他のサーバに格納されていてもよい。クラスタリング部12は、文書検索部11が検索して抽出した複数の文書を、文書の類似度に基づいて複数の文書集合(以下、クラスタと称する)に分類する。

【0018】

キーワード生成部13は、クラスタ指定受付モジュール(以下、クラスタ指定受付部と称する)31と、キーワードリスト生成モジュール(以下、キーワードリスト生成部と称する)32と、キーワード選択モジュール(以下、キーワード選択部と称する)33を有する。

【0019】

クラスタ指定受付部31は、複数のクラスタを識別する情報をクライアント装置80の表示ユニットに表示させるとともに、複数のクラスタのうち、検索結果から排除すべき排除対象のクラスタ(以下、排除クラスタと称する)、及び、検索結果として選択すべき選択対象のクラスタ(以下、選択クラスタと称する)の指定をユーザから受け付ける(a2)。キーワードリスト生成部32は、排除クラスタ、選択クラスタに基づいて、クエリ23における検索語の排除条件となるキーワードのリスト22を作成する。

【0020】

10

20

30

40

50

キーワード選択部 33 は、キーワードリスト 22 内のキーワード数に比例する値を示すスライダーをクライアント装置 80 の表示ユニットに表示する。また、キーワード選択部 33 は、ユーザによるスライダーの操作を受け付けるとともに (a3)、スライダーの値に対応するキーワードを含む文書の割合をクラスタ毎に表示する。スライダーの値に対応するキーワードがクエリ 23 の排除条件となる。

【0021】

次に、本実施の形態例におけるクエリ生成処理を説明する前に、クラスタを用いたクエリ 23 の基本的な生成処理を図に基づいて説明する。

【0022】

[クラスタを用いたクエリの生成]

図 4 は、クラスタを用いたクエリ 23 の基本的な生成処理を説明する第 1 の図である。本実施の形態例では、ユーザが、コンビニエンスストアの「セブンイレブン」に関する文書を検索して抽出する場合を前提とする。「セブンイレブン」は、「セブン」と省略して用いられることが多い。したがって、ユーザは、例えば、「セブン」をクエリ 23 の検索語として入力する。この結果、文書検索部 11 は、文書内に検索語「セブン」を含む複数の文書 20a を取得する。

【0023】

ただし、検索語「セブン」にしたがって検索された複数の文書 20a は、「セブンイレブン」に関する記事に加えて、「ウルトラセブン」や「セブンスター」、映画の「セブン」に関する記事も含む。検索された複数の文書 20a を対象として、マーケティング等の分析処理を行う場合、「セブンイレブン」に関する記事以外の記事が含まることにより、分析の精度が下がる。したがって、検索結果から、「セブンイレブン」に関する記事以外の記事が排除されることが望ましい。そこで、クラスタリング部 12 は、検索語「セブン」にしたがって検索された複数の文書 20a を、複数のクラスタに分類する。

【0024】

図 5 は、クラスタを用いたクエリ 23 の基本的な生成処理を説明する第 2 の図である。クラスタリング部 12 は、検索語「セブン」に基づいて抽出された複数の文書 20a を、類似性のある文書同士をまとめて、複数のクラスタを生成する。図 5 の例では、クラスタリング部 12 は、例えば、複数の文書 20a 内の各文書を、いずれかのクラスタに分類する。

【0025】

図 5 の例において、クラスタリング部 12 は、複数の文書 20a を、複数のクラスタ C11 ~ C14 に分類する。クラスタ C11 は、例えば、セブンイレブンの話題に関する文書を多く含む。また、クラスタ C12 は、ウルトラマンのセブンの話題に関する文書を、クラスタ C13 は煙草のマイルドセブンの話題に関する文書を多く含む。クラスタ C14 は、セブンイレブン、ウルトラマン、煙草以外の話題に関する文書を含む。

【0026】

次に、クラスタリング部 12 は、例えば、クラスタ C11 ~ C14 に基づいて、各クラスタを代表する語句 (以下、代表語と称する) Ck を抽出する。代表語 Ck とは、例えば、対象のクラスタにより多く頻出し、対象のクラスタ以外のクラスタにはほとんど含まれない単語である。例えば、クラスタリング部 12 は、対象のクラスタへの出現頻度が高く、かつ、対象外のクラスタへの出現頻度が少ない単語を代表語として抽出する。図 5 の例において、クラスタ C11 の代表語は、「コンビニ」「アイス」等である。また、クラスタ C12 の代表語は、「ウルトラマン」「フィギュア」等であって、クラスタ C13 の代表語は、「マイルドセブン」「煙草」等である。ユーザは、代表語を、検索語と組み合わせる排除条件 (NOT) の候補として使用する。

【0027】

例えば、ユーザは、代表語「フィギュア」を用いて、クエリ 23 「セブン and NOT フィギュア」を生成する。これにより、ユーザは、検索語「セブン」を含み、かつ、代表語「フィギュア」を含まない文書の集合を抽出できる。つまり、ユーザは、ウルトラマ

10

20

30

40

50

ンの話題を示すクラスタC 1 2 が含む文書の多くを検索結果から排除することができる。しかしながら、クエリ2 3「セブン a n d N O T フィギュア」に基づいて文書集合を検索すると、例えば、「セブンイレブンでフィギュアを買った」等の文章を含む文書が検索結果から排除されてしまう。つまり、検索によって抽出したい文書についても、検索結果から排除されてしまう。

【 0 0 2 8 】

このように、クラスタの代表語を用いることによって、検索語による検索結果から不要な文書を排除するキーワード候補が取得可能になるものの、検索によって抽出したいユーザ所望の文書についても検索結果から排除されてしまうことがある。したがって、代表語を用いてクエリ2 3 を生成する場合であっても、所望の文書を抽出できるとは限らない。また、検索結果を分類した複数のクラスタから所望のクラスタを選択して文書を抽出する場合であっても、クラスタリングの精度が完全ではないことから、不要な文書が抽出されてしまう場合がある。

10

【 0 0 2 9 】

また、クエリ2 3 による文書検索では、根本的には、ユーザが所望する完全な文書集合を取得することは困難である。つまり、クエリ2 3 による文書検索によると、不要な文書を完全に排除することや、抽出したい文書を完全に選択することは困難である。したがって、ユーザは、所望の文書集合とできるだけ近い文書集合を抽出可能にするクエリ2 3 を生成する。ただし、ユーザは、所望する完全な文書の集合の内容を、予め検知していない場合がある。したがって、ユーザは、クラスタの代表語を用いて試行錯誤を重ねながら、所望の文書集合と近いと思われる文書集合を抽出可能なクエリ2 3 を生成する。

20

【 0 0 3 0 】

しかしながら、各クラスタの代表語を組み合わせたととしても、所望の文書集合と近い文書集合を抽出するクエリ2 3 を生成することは容易ではない。また、ユーザは所望の文書集合の内容を予め検知しているわけではないため、試行錯誤を重ねたととしても最適な文書集合が抽出されるとは限らない。

【 0 0 3 1 】

本実施の形態例における検索サーバ1 0 は、複数のクラスタを識別する情報を表示して、ユーザに、検索結果から排除すべき文書を多く有する排除クラスタを指定させる。そして、検索サーバ1 0 は、排除クラスタ内のキーワードを抽出し、抽出したキーワード毎の、検索された複数の文書における排除クラスタ内での出現分布率を示すスコアを計算し、スコアの降順にソートされたキーワードのキーワードリスト2 2 を生成する。そして、検索サーバ1 0 は、キーワードリスト2 2 内のキーワード数に比例する値を示すスライドバー等の操作オブジェクトを表示して、ユーザのスライドバー等の入力に応じて、キーワードリスト2 2 の上位数分のキーワードを排除条件のクエリ2 3 の候補として選択し、選択したキーワードを含む文書の割合をクラスタ毎に計算して表示する。

30

【 0 0 3 2 】

即ち、本実施の形態例における検索サーバ1 0 は、検索語による検索結果を分類した複数のクラスタから、排除すべき文書を多く有する排除クラスタをユーザに指定させる。検索サーバ1 0 は、排除クラスタの指定に基づいて、排除クラスタ内の文書をより多く排除可能であって、排除クラスタ以外のクラスタから排除される文書量を抑えるキーワードをその有効性の順に有するリストを生成する。したがって、ユーザは、排除クラスタを指定するだけで、検索サーバ1 0 に、クエリ2 3 の排除条件のキーワード候補（キーワードリスト2 2 ）を生成させることができる。

40

【 0 0 3 3 】

そして、検索サーバ1 0 は、キーワードのリスト内のキーワード数に比例する値を示すスライドバーとともに、スライドバーの値に対応する上位数分のキーワードを含む文書の割合をクラスタ毎に表示する。これにより、ユーザは、スライドバーの値に応じて、排除クラスタ内で排除される文書の割合と、排除クラスタ以外のクラスタ内で排除される文書の割合とのバランスを確認しながら、スライドバーの値を指定することができる。つまり

50

、ユーザは、排除すべきクラスタから排除される文書の割合と、選択したいクラスタから排除される文書の割合とのバランスを確認できることにより、最適なバランスを選択することができる。

【0034】

検索サーバ10は、スライダーの値が示す上位数分のキーワードを、クエリ23の排除条件として選択しクエリ23を生成する。したがって、ユーザは、キーワード自体を意識することなく、最適なクエリ23を取得することができる。

【0035】

このように、ユーザは、所望する完全な文書の集合の内容を予め検知していなくても、排除クラスタの指定とスライダーによる指定とを行うだけで、ユーザが所望する文書集合を抽出可能にするクエリ23を生成させることが可能になる。つまり、本実施の形態例における検索サーバ10は、試行錯誤を重ねることなく、簡易な操作にしたがって、ユーザの意図を反映させた排除条件のキーワードの絞り込みを可能にする。また、ユーザは、所望の文書集合に近い文書集合を確実に抽出可能になる。

【0036】

次に、本実施の形態例におけるクエリ生成処理をフローチャート図に基づいて説明する。

【0037】

[フローチャート]

図6は、本実施の形態例におけるクエリ生成処理を説明するフローチャート図である。初めに、検索サーバ10の文書検索部11は、クライアント装置80から文書集合20の検索を行うための検索語を受信する(S11)。なお、検索語は、複数のキーワードによる組み合わせ(キーワード集合)であってもよい。次に、文書検索部11は、検索語に基づいてクエリ23を生成し、文書集合20の検索処理を行う(S12)。文書検索部11は、検索処理の結果、クエリ23が示す条件に合致する複数の文書20aを取得する。

【0038】

次に、検索サーバ10のクラスタリング部12は、検索によって取得した複数の文書20aを、文書の類似性に基づいて複数のクラスタに分類する(S12)。クラスタリング部12は、例えば、k-means、ワン・パスクラスタリング等の公知の技術を用いて、クラスタを生成する。例えば、本実施の形態例におけるクラスタリング部12は、論文「Criterion Functions for Document Clustering Experiments and Analysis (文献名: Technical Report CS Dept. 01-40, Univ. Minnesota, Ying Zhao and George Karypis 2001年)」に記述されるクラスタリングの技術に基づいて、クラスタを生成する。

【0039】

次に、キーワード生成部13のクラスタ指定受付部31は、クライアント装置80の表示ユニットに、生成した複数のクラスタを識別する情報を表示させる(S14)。複数のクラスタを識別する情報は、例えば、クラスタIDやアイコン等である。クラスタ指定受付部31は、例えば、複数のクラスタを識別する情報をクライアント装置80にダウンロードさせ、クライアント装置80で動作するウェブブラウザを介して表示させる。また、クラスタ指定受付部31は、さらに、複数のクラスタのうち、排除クラスタ、選択クラスタに対する指定を受け付けるオブジェクトを表示する(S15)。オブジェクトは、例えば、ラジオボタンである。

【0040】

また、クラスタ指定受付部31は、選択クラスタ及び排除クラスタのユーザによる指定を可能にするために、各クラスタの特徴を表す特徴情報を表示する(S16)。ユーザは、特徴情報を参照することによって、クラスタが有する話題を識別することができる。ユーザは、クラスタが有する話題に基づいて、クラスタが有する文書が検索結果から排除されることが望ましいか否か、クラスタが有する文書が検索結果として選択されることが望ましいか否かを判断可能になる。

【0041】

10

20

30

40

50

特徴情報とは、例えば、クラスタの文書に含まれる検索語を使用する文字列の一部、クラスタの文書に含まれる頻出語、クラスタの文書に含まれる単語であって検索された複数の文書に対する文書集合への出現比率が高い単語（代表語）等である。ただし、特徴情報は、この例に限定されるものではなく、クラスタが有する文書の主題が識別可能になる情報であればいずれの情報であってもよい。

【 0 0 4 2 】

次に、クラスタ指定受付部 3 1 は、ユーザによる選択クラスタ、排除クラスタの指定を受け付ける（S 1 7）。ユーザは、少なくとも、排除クラスタを 1 つ指定する。また、選択クラスタは、必ずしも指定されなくてもよい。ユーザは、排除すべき文書を多く有するクラスタを排除クラスタに指定する。また、ユーザは、検索結果として選択されることが望ましい文書を多く有するクラスタを選択クラスタに指定する。

10

【 0 0 4 3 】

ユーザによるクラスタの指定を受け付けると（S 1 8 の Y E S）、キーワード生成部 1 3 のキーワードリスト生成部 3 2 は、キーワードリスト 2 2 を生成する（S 1 9）。キーワードリスト 2 2 は、クエリ 2 3 の排除条件の候補となる複数のキーワードを有する。また、キーワードリスト 2 2 は、検索された文書に対する排除クラスタへの出現分布率が高いときにより大きい値を有するスコアの降順に、キーワードを有する。キーワードリスト 2 2 の詳細については、別の図にしたがって後述する。

【 0 0 4 4 】

このように、ユーザは、複数のクラスタから排除クラスタを指定するだけで、クラスタに対する指定項目（排除、選択）を反映させた、クエリ 2 3 の排除条件のキーワード候補を検索サーバ 1 0 に生成させることができる。したがって、ユーザは、クエリ 2 3 の排除条件のキーワードの候補を考える必要がない。

20

【 0 0 4 5 】

次に、キーワード生成部 1 3 のキーワード選択部 3 3 は、クライアント装置 8 0 にスライダーを表示させ（S 2 0）、ユーザによるスライダーに対する操作を受け付ける（S 2 1）。スライダーの値は、キーワードリスト 2 2 内のキーワード数に比例する。つまり、ユーザは、スライダーの値を変更させることによって、クエリ 2 3 の排除条件となるキーワードの数を変動させることができる。

【 0 0 4 6 】

30

そして、キーワード選択部 3 3 は、スライダーの値に応じた、キーワードリスト 2 2 の上位数分のキーワードをクエリ 2 3 の排除条件の候補として選択し、選択したキーワードを含む文書の割合をクラスタ毎に計算する（S 2 2）。選択したキーワードを含む文書の割合とは、選択したキーワードが排除条件として適用されることによって、排除される文書の割合を表す。そして、キーワード選択部 3 3 は、クラスタ毎に、選択したキーワードを含む文書の割合を表示する（S 2 3）。

【 0 0 4 7 】

キーワード選択部 3 3 は、スライダーが示す値が更新される度に、クラスタ毎の排除される文書の割合を計算し直して、表示する（S 2 2、S 2 3）。スライダーの値が確定すると（S 2 4 の Y E S）、キーワード選択部 3 3 は、検索語と、排除条件とする上位数分のキーワードとの組み合わせによって、クエリ 2 3 を生成する（S 2 5）。そして、例えば、キーワード選択部 3 3 は、検索語と排除条件である上位数分のキーワードとに基づいたクエリ 2 3 をクライアント装置 8 0 の表示ユニットや、検索サーバ 1 0 のメモリ 1 0 2 に出力する。または、キーワード選択部 3 3 は、検索語と排除条件である上位数分のキーワードとに基づいたクエリ 2 3 にしたがって文書集合 2 0 を検索し直し、文書集合を抽出する。

40

【 0 0 4 8 】

クエリ 2 3 の排除条件の候補として選択するキーワードの数が変動することによって、各クラスタから排除される文書量も変動する。このとき、排除クラスタから排除される文書の量と、選択クラスタから抽出される文書の量とは、トレードオフの関係にある。具体

50

的に、キーワードの数を増加させて排除クラスタから排除される文書量を増加させるようにする、選択クラスタから抽出される文書量も低減する傾向にある。一方、キーワードの数を低減させて選択クラスタから抽出される文書量を増加させようすると、排除クラスタから排除される文書量も低減してしまう傾向にある。

【 0 0 4 9 】

排除クラスタから排除される文書量と、選択クラスタから抽出される文書量との望ましいバランスは、検索ケースによって異なる。例えば、排除クラスタの文書を可能な限り検索結果から排除したい場合、選択クラスタから抽出される文書量が低下したとしても、排除クラスタから排除される文書量が多い方が望ましい。したがって、排除クラスタの文書を可能な限り検索結果から排除したい場合、キーワードの数が多い方が、ユーザが所望する文書集合に近い文書集合を抽出可能になり易い。

10

【 0 0 5 0 】

一方、選択クラスタの文書を可能な限り検索結果として抽出したい場合、排除クラスタから排除される文書量が少なかったとしても、選択クラスタから抽出される文書量が多い方が望ましい。したがって、選択クラスタの文書を可能な限り検索結果として抽出したい場合、キーワードの数が少ない方が、ユーザが所望する文書集合に近い文書集合を抽出可能になり易い。

【 0 0 5 1 】

このように、排除クラスタから排除される文書量と選択クラスタから抽出される文書量との望ましいバランス数が検索ケースによって異なるところ、ユーザは、排除クラスタから排除される文書量と選択クラスタから抽出される文書量とのバランスを確認しながら、スライドバーの値を選択することができる。これにより、ユーザは、意図に沿った文書集合を抽出可能にするクエリ 2 3 の排除条件のキーワード数を選択可能になる。また、ユーザは、スライドバー等のオブジェクトを操作するだけで、キーワード自体を意識することなくクエリ 2 3 を生成させることが可能になる。

20

【 0 0 5 2 】

次に、図 6 で説明したフローチャート図の処理を具体例に対応させて説明する。

【 0 0 5 3 】

図 7 は、本実施の形態例におけるクラスタの生成処理を説明する図である。本実施の形態例における検索サーバ 1 0 の文書検索部 1 1 は、図 4 で説明した処理と同様にして、受け付けた検索語「セブン」に基づいて、文書集合 2 0 を検索する（図 6 の S 1 1、S 1 2）。この結果、文書検索部 1 1 は、文書内に検索語「セブン」を含む複数の文書 2 0 a を取得する。次に、クラスタリング部 1 2 は、検索語「セブン」にしたがって検索された複数の文書 2 0 a を、複数のクラスタ C 1 1 ~ C 1 4 に分類する（S 1 3）。クラスタ C 1 1 ~ C 1 4 は、図 4 で説明したとおりである。

30

【 0 0 5 4 】

次に、クラスタ指定受付部 3 1 は、クライアント装置 8 0 の表示ユニットに、複数のクラスタを識別する情報を表示させるとともに（図 6 の S 1 4）、ラジオボタン等のオブジェクトを表示して、排除クラスタ及び選択クラスタへの指定を受け付ける（S 1 5）。このとき、クラスタ指定受付部 3 1 は、選択クラスタ及び排除クラスタのユーザによる指定を可能にするために、各クラスタの特徴を表す特徴情報を表示する（S 1 6）。

40

【 0 0 5 5 】

図 8 は、クラスタの指定を受け付けるクライアント装置 8 0 の表示ユニットが表示する画面の一例を示す図である。図 8 は、クラスタ毎に、クラスタを識別する文書マークで示されるアイコンに加えて、クラスタの指定を受け付けるボタンメニュー R 1 ~ R 4 と、クラスタの特徴情報であるクラスタ内の検索語の使用例 T 1 ~ T 4 とを有する。図 8 の例において、各クラスタ C 1 1 ~ C 1 4 を識別する情報は、文書マークで示されるアイコンである。

【 0 0 5 6 】

また、ボタンメニュー R 1 ~ R 4 は、「選択」「排除」「その他」のいずれかを指定さ

50

せるラジオボタンである。例えば、初め、全てのクラスタは、「その他」に指定される。「その他」に指定されるクラスタの文書は、キーワードリスト 22 の生成処理に使用されない。処理の詳細については後述するが、キーワード生成部 13 は、排除クラスタ、選択クラスタに基づいて、キーワードリスト 22 を生成する。クラスタの話題が識別できない場合、ユーザは、例えば、クラスタを「その他」のクラスタとする。

【0057】

図 8 の例において、クラスタ内の検索語「セブン」の使用例 T1 ~ T4 は、クラスタ内の文書のうち、一部の文書における検索語の使用部分の文字列である。例えば、クラスタ C11 の検索語の使用部分の文字列 T1 は、「この夏 6、7、8、9 月、私がよく購入した商品を発売したと思います。第 1 位 昆布とかつおのうま味・おでん [セブンイレブン]」である。また、図 8 に示すように、検索語の使用部分の文字列のうち検索語「セブン」は、例えば、太字等によって強調され表示される。なお、検索語は、例えば、斜体、下線等によって強調されて表示されてもよい。ユーザは、クラスタ内の検索語の使用例を参照することによって、クラスタ内の文書で検索語「セブン」がどのように参照されているかを確認可能になり、クラスタが有する話題を識別することができる。

【0058】

なお、クラスタ内の検索語の使用例 T1 ~ T4 は、ユーザがクラスタ C11 ~ C14 それぞれからランダムに選択した文書内での検索語の使用例であってもよい。また、図 6 のフローチャート図で前述したとおり、特徴情報は、検索語の使用例 T1 ~ T4 の他に、例えば、クラスタ内の頻出語や代表語等であってもよい。

【0059】

図 8 の例において、クラスタ C12 の検索語の使用例 T2 によるとクラスタ C12 が「ウルトラマン」に関する話題を有することが識別可能になる。本実施の形態例では、「セブンイレブン」に関する文書を抽出することを目的とする。したがって、「ウルトラマン」の話題を有する文書は、検索結果として抽出したい文書に当たらない可能性があることから、ユーザは、例えば、クラスタ C12 を排除クラスタとしてボタンメニュー R2 に指定する。また、クラスタ C13 の検索語の使用例 T3 によると、クラスタ C13 が「タバコ」に関する話題を有することが識別可能になる。同様に、タバコの話題を有する文書は、検索結果として抽出したい文書に当たらない可能性があることから、ユーザは、例えば、クラスタ C13 を排除クラスタとしてボタンメニュー R3 に指定する。

【0060】

図 8 の例では、排除クラスタに加えて、選択クラスタが指定される。クラスタ C11 の検索語の使用例 T1 によると、クラスタ C11 が「セブンイレブン」に関する話題を有することが識別可能になる。したがって、ユーザは、例えば、クラスタ C11 を選択クラスタとしてボタンメニュー R1 に指定する。また、クラスタ C14 は、特定の話題を有していない。したがって、ユーザは、例えば、クラスタ C14 を、その他のクラスタとする (R4)。

【0061】

図 9 は、キーワードリスト 22 の生成処理を説明する図である。キーワードリスト生成部 32 は、図 8 の画面において指定された排除クラスタ C12、C13、及び、選択クラスタ C11 が有する文書を入力として、クエリ 23 の排除条件の候補となるスコア付きのキーワードのリストを生成する (図 6 の S18、S19)。排除クラスタのみが指定される場合、キーワードリスト生成部 32 は、排除クラスタ内の文書に含まれるキーワードを抽出し、キーワードの複数のクラスタにおける排除クラスタ内での出現分布率に基づいてスコアを算出する。この場合、スコアは、複数のクラスタにおける排除クラスタ内での出現分布率が高いときにより大きい値となる。

【0062】

また、排除クラスタに加えて選択クラスタが指定される場合、キーワードリスト生成部 32 は、排除クラスタ及び選択クラスタ内の文書に含まれるキーワードを抽出し、キーワードの複数のクラスタにおける排除クラスタ内での出現分布率に基づいてスコアを算出す

る。この場合、スコアは、複数のクラスタにおける排除クラスタ内での出現分布率が高く、かつ、選択クラスタ内での出現分布率が低いときにより大きい値となる。

【 0 0 6 3 】

ここで、キーワードリスト 2 2 の生成処理をより具体的に説明する。

【 0 0 6 4 】

[キーワードリストの生成]

図 1 0 は、キーワードリスト 2 2 の生成処理を説明するフローチャート図である。まず、キーワードリスト生成部 3 2 は、排除クラスタ及び選択クラスタ内の文書に含まれるキーワードを複数抽出する (S 3 1)。そして、キーワードリスト生成部 3 2 は、抽出したキーワードのスコアの値を 0 に初期化する (S 3 2)。

10

【 0 0 6 5 】

なお、各クラスタは、クラスタ内に含まれる単語毎に $t f i d f$ (term frequency inverse document frequency) 値を有する。 $t f i d f$ 値は、例えば、単語が、クラスタ内で特徴的である度合いを識別するための指標である。この例では、 $t f i d f$ 値は、対象のクラスタにより偏って出現する度合いを表す。 $t f i d f$ 値は、単語がクラスタ内に出現する回数「 $t f$ 」と、全文書において当該単語が出現する文書数「 $d f$ 」とに基づいて、計算式「 $t f i d f = t f / \log (d f / N)$ 」にしたがって算出される。計算式内の値「 N 」は全文書数を表す。

【 0 0 6 6 】

$t f i d f$ 値は、例えば、対象のクラスタへの出現頻度が高く、かつ、全クラスタ内への出現率が少ない場合により大きい値を有する。言い換えると、対象のクラスタへの出現頻度が高くて全クラスタ内での出現率が高い場合は、いずれのクラスタにも出現することを示すため、 $t f i d f$ 値は大きな値にはならない。単語毎の $t f i d f$ 値の例については、図 1 1 で例示する。

20

【 0 0 6 7 】

そして、キーワードリスト生成部 3 2 は、選択クラスタとして指定されたクラスタがあるか否かを判定する (S 3 3)。選択クラスタがある場合 (S 3 3 の Y E S)、キーワードリスト生成部 3 2 は、抽出したキーワードにしたがって、選択クラスタ内の各文書を検索する (S 3 4)。文書にキーワードが含まれる場合 (S 3 5 の Y E S)、キーワードリスト生成部 3 2 は、キーワードの $t f i d f$ 値をスコアから減算する (S 3 6)。つまり、キーワードが選択クラスタ内に出現する場合、排除条件のキーワードとして不適切である可能性が高いため、スコアの値は減算される。一方、キーワードが含まれない場合 (S 3 5 の N O)、キーワードリスト生成部 3 2 は、 $t f i d f$ 値をスコアから減算しない。キーワードリスト生成部 3 2 は、選択クラスタ内の全ての文書について (S 3 7 の Y E S)、工程 S 3 4 ~ S 3 6 の処理を行う。

30

【 0 0 6 8 】

次に、キーワードリスト生成部 3 2 は、抽出したキーワードにしたがって、排除クラスタ内の各文書を検索する (S 3 8)。文書にキーワードが含まれる場合 (S 3 9 の Y E S)、キーワードリスト生成部 3 2 は、キーワードの $t f i d f$ 値をスコアに加算する (S 4 0)。つまり、キーワードが排除クラスタ内に出現する場合、排除条件のキーワードとして適切である可能性が高いため、スコアの値は加算される。一方、キーワードが含まれない場合 (S 3 9 の N O)、キーワードリスト生成部 3 2 は、 $t f i d f$ 値をスコアに加算しない。

40

【 0 0 6 9 】

キーワードリスト生成部 3 2 は、排除クラスタ内の全ての文書について (S 4 1 の Y E S)、工程 S 3 8 ~ S 4 0 の処理を行う。そして、キーワードリスト生成部 3 2 は、スコアの降順にキーワードをソートする (S 4 2)。これにより、キーワードリスト生成部 3 2 は、スコアの降順にキーワードを有するキーワードのリスト 2 2 を生成する。なお、前述したように、キーワード生成部 3 2 は、「その他」に指定されるクラスタ内の文書を加味することなく、キーワードリスト 2 2 を生成する。

50

【 0 0 7 0 】

次に、キーワードリスト 2 2 の生成処理を具体例に対応させて説明する。

【 0 0 7 1 】

図 1 1 は、キーワード毎の $t f i d f$ 値を例示する図である。図 1 1 の表 H 1 は、クラスタ ID、文書 ID、単語、 $t f i d f$ 値を有する。クラスタ ID は、クラスタを識別する情報であって、文書 ID は、文書を識別する情報である。単語は、クラスタ内の文書に含まれる単語である。図 1 1 の表 H 1 において、例えば、ID「1」のクラスタ C 1 は、文書 doc 1、doc 2 等を有する。また、例えば、ID「2」のクラスタ C 2 は、文書 doc 3 0 3 等を有する。なお、この例では、文書の一部の情報を表しているが、実際には、各クラスタは多数の文書を有する。

10

【 0 0 7 2 】

図 1 1 の例において、文書 doc 1 は、例えば、単語「コンビニ」「おにぎり」等を含む。この例では、一部の単語を表しているが、実際には、各文書は多数の単語を有する。単語「コンビニ」の $t f i d f$ 値は「42.7」であって、単語「おにぎり」の $t f i d f$ 値は「40.3」である。また、文書 doc 2 は、例えば、単語「四国」「コンビニ」等を含み、単語「四国」の $t f i d f$ 値は「58.7」であって、単語「コンビニ」の $t f i d f$ 値は「42.7」である。つまり、単語「四国」は、単語「おにぎり」「コンビニ」よりも、ID「1」のクラスタ C 1 により偏って出現することを意味する。

【 0 0 7 3 】

また、表 H 1 において、文書 doc 3 0 3 は、例えば、単語「コンビニ」「マイルドセブン」「煙草」等を含む。この例では、一部の単語を表しているが、実際には、各文書は多数の単語を有する。単語「コンビニ」の $t f i d f$ 値は「38.1」、単語「マイルドセブン」の $t f i d f$ 値は「37.8」、単語「煙草」の $t f i d f$ 値は「33.6」である。ID「1」のクラスタ C 1 内の文書 doc 1、doc 2、ID「2」のクラスタ C 2 内の文書 doc 3 0 3 はいずれも、単語「コンビニ」の $t f i d f$ 値を有する。これは、単語「コンビニ」が、文書 doc 1、doc 2、doc 3 0 3 のいずれにも含まれることを示す。また、ID「2」のクラスタの単語「コンビニ」の $t f i d f$ 値は、ID「1」のクラスタ C 1 の単語「コンビニ」の $t f i d f$ 値より小さい。これは、単語「コンビニ」が、ID「2」のクラスタ C 2 内の文書 doc 3 0 3 よりも、ID「1」のクラスタ C 1 内の文書 doc 1、doc 2 に、より偏って出現することを示す。

20

30

【 0 0 7 4 】

図 1 2 は、キーワード生成処理の具体例を説明する第 1 の図である。図 1 2 の表 H 1 は、図 1 1 の表 H 1 と同一である。また、この例において、ID「1」のクラスタ C 1 は選択クラスタ、ID「2」のクラスタ C 2 は排除クラスタに該当する。

【 0 0 7 5 】

キーワードリスト生成部 3 2 は、ID「1」の選択クラスタ C 1、及び、ID「2」の排除クラスタ C 2 内の文書から、例えば、キーワード「コンビニ」「おにぎり」「四国」「マイルドセブン」「煙草」等を抽出する（図 9 の S 3 1）。そして、キーワードリスト生成部 3 2 は、抽出した各キーワードのスコアを 0 に初期化したキーワードリスト 2 2 を生成する（S 3 2）。具体例において、選択クラスタ（ID「1」）が存在することから（S 3 3 の YES）、キーワードリスト生成部 3 2 は、キーワード「コンビニ」「おにぎり」「四国」「マイルドセブン」「煙草」にしたがって、ID「1」の選択クラスタ C 1 内の文書を検索する（S 3 4）。

40

【 0 0 7 6 】

具体例において、ID「1」の選択クラスタ C 1 内の文書 doc 1 は、キーワード「コンビニ」「おにぎり」を含む（S 3 5 の YES）。したがって、キーワードリスト生成部 3 2 は、キーワード「コンビニ」の $t f i d f$ 値「42.7」（Y 1 1）、キーワード「おにぎり」の $t f i d f$ 値「40.3」（Y 1 2）をそれぞれスコアから減算する（S 3 6）。したがって、図 1 2 のキーワードリスト 2 2 - 1 におけるキーワード「コンビニ」のスコアは値「- 42.7」（Y 1 3）、キーワード「おにぎり」のスコアは値「- 40

50

、 3」(Y14)となる。

【0077】

図13は、キーワード生成処理の具体例を説明する第2の図である。図13の表H1は、図11の表H1と同一である。具体例において、ID「1」の選択クラスタC1内の文書doc2は、キーワード「四国」「コンビニ」を含む(S35のYES)。したがって、キーワードリスト生成部32は、キーワード「四国」のtfidf値「58.4」(Y21)、キーワード「コンビニ」のtfidf値「42.7」(Y22)をそれぞれスコアから減算する(S36)。したがって、図13のキーワードリスト22-2におけるキーワード「コンビニ」のスコアは値「-85.4(=-42.7-42.7)」(Y23)、キーワード「四国」のスコアは値「-58.4」(Y24)となる。

10

【0078】

図14は、キーワード生成処理の具体例を説明する第3の図である。図14の表H1は、図11の表H1と同一である。次に、キーワードリスト生成部32は、抽出したキーワード「コンビニ」「おにぎり」「四国」「マイルドセブン」「煙草」にしたがって、ID「2」の排除クラスタC2内の各文書を検索する(S38)。具体例において、ID「2」の排除クラスタC2内の文書doc303は、キーワード「コンビニ」「マイルドセブン」「煙草」を含む(S35のYES)。したがって、キーワードリスト生成部32は、キーワード「コンビニ」のtfidf値「38.1」(Y31)をスコアに加算する(S36)。したがって、図14のキーワードリスト22-3におけるキーワード「コンビニ」のスコアは値「(-47.3=-85.4+38.1)」(Y32)となる。

20

【0079】

図15は、キーワード生成処理の具体例を説明する第4の図である。図15の表H1は、図11の表H1と同一である。次に、キーワードリスト生成部32は、キーワード「マイルドセブン」のtfidf値「37.8」(Y41)をスコアに加算するとともに、キーワード「煙草」のtfidf値「33.6」(Y42)をスコアに加算する(S36)。したがって、図13のキーワードリスト22-4におけるキーワード「マイルドセブン」のスコアは値「37.8」(Y43)、キーワード「煙草」のスコアは値「33.6」(Y44)となる。

【0080】

図16は、具体例において生成されるキーワードリスト22-5を示す図である。図16のキーワードリスト22-5は、スコアの降順にキーワードを有する。図16のキーワードリスト22-5において、最もスコアの高いキーワードは「マイルドセブン」である。これは、キーワード「マイルドセブン」が、クエリ23の排除条件として有効性が高いことを示す。次に、排除条件として有効性が高いキーワードは、「煙草」である。

30

【0081】

このように、キーワードリスト生成部32は、ユーザから指定された排除クラスタ、選択クラスタに基づいて、排除クラスタ内への出現分布率が高く、かつ、選択クラスタ内での出現分布率が低いときにより大きい値を有するスコアの降順にキーワードを有するキーワードリスト22を生成することができる。なお、図11～図16の具体例では、排除クラスタに加えて選択クラスタが指定される場合を例示しているが、排除クラスタのみが指定される場合、キーワードリスト生成部32は、排除クラスタへの出現分布率が高いときにより大きい値を有するスコアの降順にキーワードを有するキーワードリスト22を生成する。

40

【0082】

図17は、本実施の形態例におけるキーワードの選択処理を説明する図である。図16のようなキーワードリスト22を生成すると、キーワード選択部33は、キーワードリスト22内のキーワード数の比例した値を示すスライドバー等のオブジェクトを、クライアント装置80に表示させる(図6のS20)。そして、キーワード選択部33は、ユーザのスライドバーに対する操作にしたがって、キーワードリスト22の上位数分のキーワード22aを排除条件のクエリ23候補として選択する(S22)。キーワードリスト22

50

は、スコアの降順にキーワードを有する。したがって、キーワードリスト 22 の上位から順に排除条件とするキーワードが選択されることによって、排除条件とするキーワードを効率的に選択することが可能になる。

【0083】

図 18 は、スライダー S B を表示するクライアント装置 80 の表示ユニットが表示する画面の一例を示す図である。図 18 は、クラスタを識別する情報とクラスタの指定を示すボタンメニュー R 1 ~ R 4 とに加えて、スライダー S B を有する。なお、キーワード選択部 33 は、例えば、スライダー S B の代わりに、キーワードリスト 22 内のキーワード数に比例する複数の項目（例えば、高、中、低等）を表示するドロップダウンリストや、ボタン等を表示して、ユーザに選択させてもよい。

10

【0084】

また、図 18 の例において、例えば、スライダー S B の左端の値に対応するキーワードの数は 0 個である。一方、スライダー S B の右端の値に対応するキーワードの数は、例えば、排除クラスタ内の文書がすべて検索結果から排除される上位数分のキーワード数に対応する。この場合、スライダー S B の値が右端に設定される場合のキーワードの数は、検索ケースによって異なる。ただし、スライダー S B の右端の値に対応するキーワードの数は、所定の値に予め設定されていてもよいし、スコアが所定値以上のキーワードの数であってもよい。

【0085】

ユーザは、例えば、図 18 に示すスライダー S B のノブ（つまみ）p p の位置を変化させることによって、スライダー S B の示す値を変更させる。図 18 の例において、スライダー S B のノブ p p を右方向に変更させた場合、スライダー S B の値に応じて選択されるキーワード数が増加する。

20

【0086】

図 19 は、排除される文書の割合を表示するクライアント装置 80 の表示ユニットが表示する画面の一例を示す図である。図 19 は、スライダー S B に加えて、クラスタ毎に、スライダー S B によって示される排除条件のキーワードが適用された場合に排除される文書の割合を示す棒グラフ E B 1 ~ E B 4 を表示する。図 19 の例において、棒グラフ E B 1 ~ E B 4 に加えて、クラスタが排除クラスタであるか、選択クラスタであるかが識別可能に表示されることにより、ユーザは、排除クラスタの文書がどの程度排除され、選択クラスタの文書がどの程度、排除されずに残るかを確認することができる。

30

【0087】

したがって、ユーザは、スライダーの値に応じて、排除すべきクラスタから排除される文書の割合と、選択したいクラスタから排除される文書の割合とのバランスを確認できる。ユーザは、排除クラスタ内で排除される文書の割合と、排除クラスタ以外のクラスタ内で排除される文書の割合とのバランスを確認しながら、最適なバランスを実現するスライダーの値を指定することができる。したがって、ユーザは、スライダーの値に対応する、所望の文書に近い文書集合を抽出可能な排除条件のキーワードを取得することができる。

【0088】

なお、ユーザは、その他に指定されるクラスタの文書がどの程度、排除されるかを検知することによって、その他に指定されるクラスタとして指定されるクラスタが有する特徴を識別することが可能になる。クラスタの文書が有する特徴を識別することが可能になることによって、ユーザは、その他に指定されるクラスタを、例えば、排除クラスタや選択クラスタに指定し直すことが可能になる。これにより、ユーザは、クラスタの指定と、スライダー等のオブジェクト操作によるキーワード数の調整とを繰り返すことによって、所望の文書により近い文書の集合を抽出することができる。

40

【0089】

以上のように、本実施の形態例におけるクエリ生成方法は、処理ユニットが、入力された検索語に基づいて複数の文書を検索し、検索された複数の文書を類似度にしたがって複

50

数の文書集合に分類し、複数の文書集合を識別する情報を表示ユニットに表示する第1の工程を有する。また、本実施の形態例におけるクエリ生成方法は、表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、抽出した文字列毎の、検索された複数の文書における排除文書集合内での出現分布率を示すスコアを計算し、スコアの降順にソートされた文字列の文字列リストを生成する第2の工程を有する。また、本実施の形態例におけるクエリ生成方法は、文字列リスト内の文字列数に比例する入力に応じて、文字列リストの上位数分の文字列を排除条件のクエリ23の候補として選択し、選択した文字列を含む文書の割合を文書集合毎に計算し表示ユニットに表示する第3の工程を有する。

【0090】

10

したがって、ユーザは、排除クラスタを指定するだけで、検索サーバ10に、クエリ23の排除条件のキーワード候補(キーワードリスト22)を生成させることができる。また、ユーザは、スライダーの値に応じて、排除クラスタ内で排除される文書の割合と、排除クラスタ以外のクラスタ内で排除される文書の割合とのバランスを確認しながら、スライダーの値を指定することができる。つまり、ユーザは、排除すべきクラスタから排除される文書の割合と、選択したいクラスタから排除される文書の割合とのバランスを確認できることにより、最適なバランスを選択することができる。また、ユーザは、キーワード自体を意識することなく、最適なクエリ23を取得することができる。

【0091】

したがって、ユーザは、所望する完全な文書の集合の内容を予め検知していなくても、排除クラスタの指定とスライダーによる指定とを行うだけで、ユーザが所望する文書集合を抽出可能にするクエリ23を生成させることが可能になる。つまり、本実施の形態例における検索サーバ10は、試行錯誤を重ねることなく、簡易な操作にしたがって、ユーザの意図を反映させた排除条件のキーワードの絞り込みを可能にする。また、ユーザは、所望の文書集合に近い文書集合を確実に抽出可能になる。

20

【0092】

また、本実施の形態例におけるクエリ生成方法によると、第2の工程は、抽出した文字列毎に、排除文書集合内の当該文字列の頻度に応じて、検索された複数の文書における排除文書集合内での出現分布値を加算してスコアを算出する。これにより、全てのクラスタにおける排除クラスタ内での出現分布率が高いときにより大きい値を有するスコアを算出可能になる。

30

【0093】

また、本実施の形態例におけるクエリ生成方法によると、第2の工程は、抽出した文字列毎に、さらに、検索された複数の文書における排除文書集合以外の文書集合内での出現分布値をスコアから減算して、スコアを算出する。これにより、全てのクラスタにおける排除クラスタ内での出現分布率が高く、かつ、選択クラスタ内での出現分布率が低いときにより大きい値を有するスコアを算出可能になる。

【0094】

また、本実施の形態例におけるクエリ生成方法によると、第3の工程は、文字列リスト内の文字列数に比例する値を示すオブジェクトを表示ユニットに表示し、オブジェクトに対するユーザの入力に比例する、文字列リストの上位数分の文字列を選択する。これにより、ユーザは、オブジェクトを操作することによって、クエリの排除条件として選択するキーワード自体を意識することなく、キーワードの数を指定することができる。

40

【0095】

また、本実施の形態例におけるクエリ生成方法によると、第2の工程は、文書集合の特徴を示す特徴情報を表示ユニットに更に表示する。また、特徴情報は、文書集合の文書に含まれる検索語の使用文字列、文書集合の文書に含まれる頻出文字列、文書集合の文書に含まれる文字列であって検索された複数の文書における書集合内での出現分布率が高い文字列、のうち少なくともいずれかである。ユーザは、クラスタの特徴情報を参照することにより、クラスタが有する話題を識別可能になり、複数のクラスタから検索結果から排除

50

すべきクラスタ、及び、検索結果に残したいクラスタを指定することができる。

【 0 0 9 6 】

また、本実施の形態例におけるクエリ生成方法によると、第 2 の工程は、排除文書集合の指定を受け付けるオブジェクトを更に表示ユニットに表示する。これにより、ユーザは、オブジェクトを操作することによって、複数のクラスタのうち、検索結果から排除すべきクラスタを簡易に指定することができる。

【 0 0 9 7 】

[他の実施の形態例]

なお、本実施の形態例における検索サーバ 1 0 は、図 1 8、図 1 9 に示すスライダー等のオブジェクトに加えて、ユーザのオブジェクト操作による入力に応じて選択されるキーワードを表示してもよい。例えば、検索サーバ 1 0 のキーワード選択部 3 3 は、クライアント装置 8 0 の表示ユニットに、スライダーの値に応じて選択される排除条件のキーワードの一覧を更に表示する。

10

【 0 0 9 8 】

これにより、ユーザは、スライダー等のオブジェクトの値の変化に応じて選択されるキーワードと、当該キーワードを含むクラスタ毎の文書量とを同時に把握しながら、オブジェクトの値を調整することができる。したがって、クエリの排除条件となるキーワード自体を把握しながらオブジェクトを操作して、キーワード数を選択したいユーザにとって利便性が高い。

【 0 0 9 9 】

20

なお、上記の例では、検索サーバ 1 0 が、複数のクラスタを生成し、排除クラスタ、選択クラスタのユーザによる指定に基づいて、キーワードリスト 2 2 を生成し、ユーザのスライダー等のオブジェクトの操作に応じて選択されるキーワードを含む文書の量をクラスタ毎に表示する。ただし、例えば、検索サーバ 1 0 が検索対象となる文書集合 2 0 を格納し、クライアント装置 8 0 が、文書集合 2 0 の検索結果に基づいて、クラスタの生成、キーワードリスト 2 2 の生成、及び、ユーザのオブジェクトの操作に応じて選択されるキーワードを含む文書の量の表示を行ってもよい。

【 0 1 0 0 】

以上の実施の形態をまとめると、次の付記のとおりである。

【 0 1 0 1 】

30

(付記 1)

処理ユニットが、

入力された検索語に基づいて複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示する第 1 の工程と、

前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成する第 2 の工程と、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示する第 3 の工程と、

40

を実行することを特徴とするクエリ生成方法。

【 0 1 0 2 】

(付記 2)

付記 1 において、

前記第 2 の工程は、前記抽出した文字列毎に、前記排除文書集合内の当該文字列の頻度に応じて、前記検索された複数の文書における前記排除文書集合内での前記出現分布値を加算して前記スコアを算出するクエリ生成方法。

【 0 1 0 3 】

50

(付記 3)

付記 2 において、

前記第 2 の工程は、前記抽出した文字列毎に、さらに、前記検索された複数の文書における前記排除文書集合以外の文書集合内での前記出現分布値を前記スコアから減算して、前記スコアを算出するクエリ生成方法。

【 0 1 0 4 】

(付記 4)

付記 1 乃至 3 のいずれかにおいて、

前記第 3 の工程は、前記文字列リスト内の文字列数に比例する値を示すオブジェクトを前記表示ユニットに表示し、前記オブジェクトに対する前記ユーザの入力に比例する、前記文字列リストの上位数分の前記文字列を選択するクエリ生成方法。

10

【 0 1 0 5 】

(付記 5)

付記 4 において、

前記第 3 の工程は、前記オブジェクトに対する入力に応じて選択される前記文字列を前記表示ユニットに更に表示するクエリ生成方法。

【 0 1 0 6 】

(付記 6)

付記 4 または 5 において、

前記オブジェクトは、スライドバーであるクエリ生成方法。

20

【 0 1 0 7 】

(付記 7)

付記 1 乃至 6 のいずれかにおいて、

前記第 2 の工程は、前記文書集合の特徴を示す特徴情報を前記表示ユニットに更に表示するクエリ生成方法。

【 0 1 0 8 】

(付記 8)

付記 7 において、

前記特徴情報は、前記文書集合の文書に含まれる前記検索語の使用文字列、前記文書集合の文書に含まれる頻出文字列、前記文書集合の文書に含まれる文字列であって前記検索された複数の文書における前記書集合内での出現分布率が高い文字列、のうち少なくともいずれかであるクエリ生成方法。

30

【 0 1 0 9 】

(付記 9)

付記 1 乃至 8 のいずれかにおいて、

前記第 2 の工程は、前記排除文書集合の指定を受け付けるオブジェクトを更に前記表示ユニットに表示するクエリ生成方法。

【 0 1 1 0 】

(付記 1 0)

入力された検索語に基づいて複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示し、

40

前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成し、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示する、

処理をコンピュータに実行させるクエリ生成プログラム。

50

【 0 1 1 1 】

(付記 1 1)

付記 1 0 において、

前記抽出した文字列毎に、前記排除文書集合内の当該文字列の頻度に応じて、前記検索された複数の文書における前記排除文書集合内での前記出現分布値を加算して前記スコアを算出するクエリ生成プログラム。

【 0 1 1 2 】

(付記 1 2)

付記 1 1 において、

前記抽出した文字列毎に、さらに、前記検索された複数の文書における前記排除文書集合以外の文書集合内での前記出現分布値を前記スコアから減算して、前記スコアを算出するクエリ生成プログラム。

10

【 0 1 1 3 】

(付記 1 3)

付記 1 0 乃至 1 2 のいずれかにおいて、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リスト内の文字列数に比例する値を示すオブジェクトを前記表示ユニットに表示し、前記オブジェクトに対する前記ユーザの入力に比例する、前記文字列リストの上位数分の前記文字列を選択するクエリ生成プログラム。

20

【 0 1 1 4 】

(付記 1 4)

付記 1 3 において、

前記オブジェクトに対する入力に応じて選択される前記文字列を前記表示ユニットに更に表示するクエリ生成プログラム。

【 0 1 1 5 】

(付記 1 5)

付記 1 3 または 1 4 において、

前記オブジェクトは、スライダーであるクエリ生成プログラム。

【 0 1 1 6 】

(付記 1 6)

処理ユニットと、

複数の文書を記憶する記憶装置と、

表示装置と、を有し、

入力された検索語に基づいて前記複数の文書を検索し、前記検索された複数の文書を類似度にしたがって複数の文書集合に分類し、前記複数の文書集合を識別する情報を表示ユニットに表示し、前記表示された複数の文書集合のうち排除すべき文書集合として指定された、排除文書集合内の文字列を抽出し、前記抽出した文字列毎の、前記検索された複数の文書における前記排除文書集合内での出現分布率を示すスコアを計算し、前記スコアの降順にソートされた前記文字列の文字列リストを生成し、前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リストの上位数分の前記文字列を排除条件のクエリの候補として選択し、前記選択した文字列を含む文書の割合を文書集合毎に計算し前記表示ユニットに表示するクエリ生成装置。

30

40

【 0 1 1 7 】

(付記 1 7)

付記 1 6 において、

前記抽出した文字列毎に、前記排除文書集合内の当該文字列の頻度に応じて、前記検索された複数の文書における前記排除文書集合内での前記出現分布値を加算して前記スコアを算出するクエリ生成方法。

【 0 1 1 8 】

(付記 1 8)

50

付記 17 において、

前記抽出した文字列毎に、さらに、前記検索された複数の文書における前記排除文書集合以外の文書集合内での前記出現分布値を前記スコアから減算して、前記スコアを算出するクエリ生成方法。

【 0 1 1 9 】

(付記 1 9)

付記 16 乃至 18 のいずれかにおいて、

前記文字列リスト内の文字列数に比例する入力に応じて、前記文字列リスト内の文字列数に比例する値を示すオブジェクトを前記表示ユニットに表示し、前記オブジェクトに対する前記ユーザの入力に比例する、前記文字列リストの上位数分の前記文字列を選択するクエリ生成方法。

【 0 1 2 0 】

(付記 2 0)

付記 19 において、

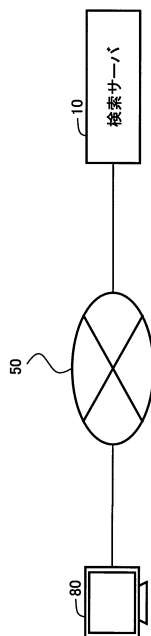
前記オブジェクトに対する入力に応じて選択される前記文字列を前記表示ユニットに更に表示するクエリ生成方法。

【 符号の説明 】

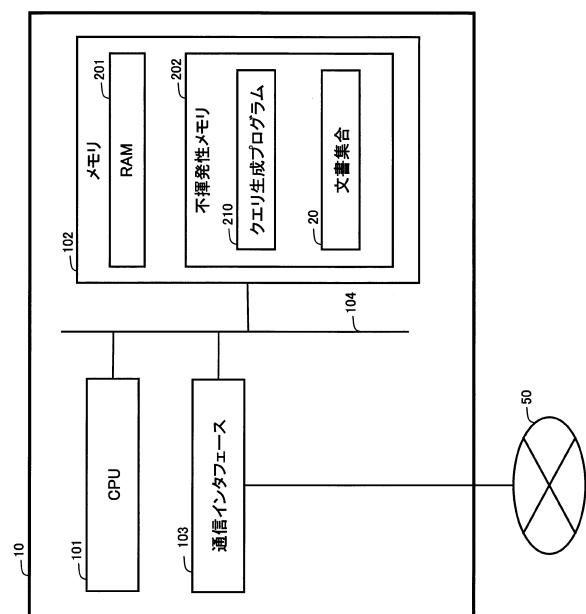
【 0 1 2 1 】

10 : 検索サーバ、80 : クライアント装置、210 : クエリ生成プログラム、C1 ~ C4 : クラスタ、11 : 文書検索部、12 : クラスタリング部、13 : キーワード生成部、31 : クラスタ指定受付部、32 : キーワードリスト生成部、33 : キーワード選択部

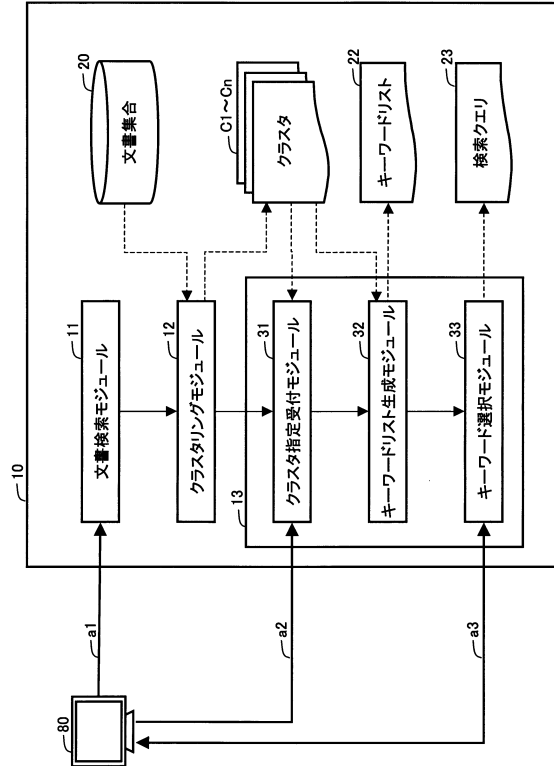
【 図 1 】



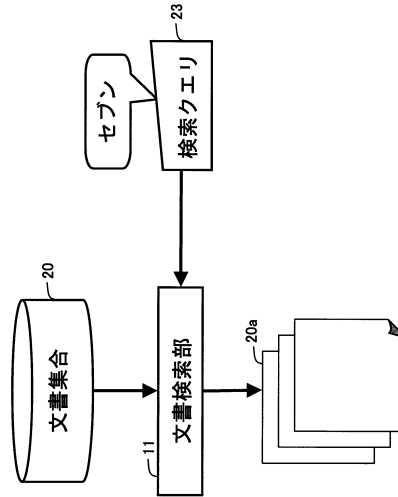
【 図 2 】



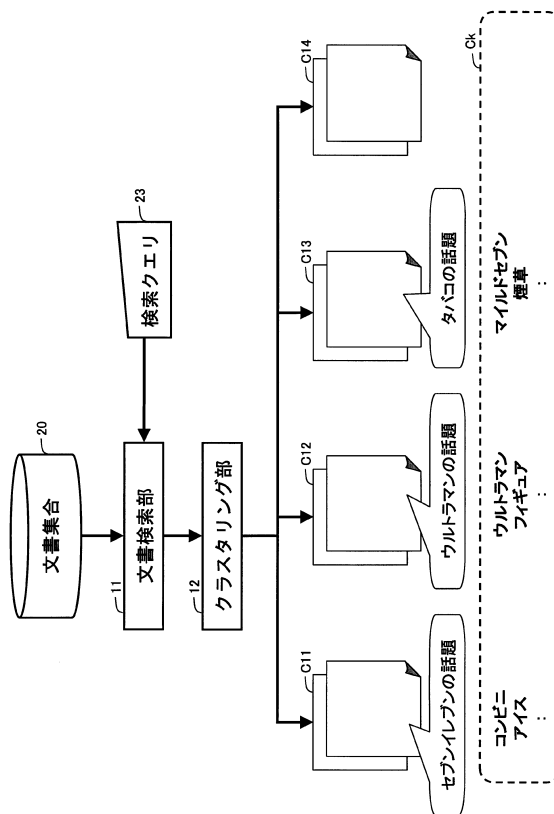
【図 3】



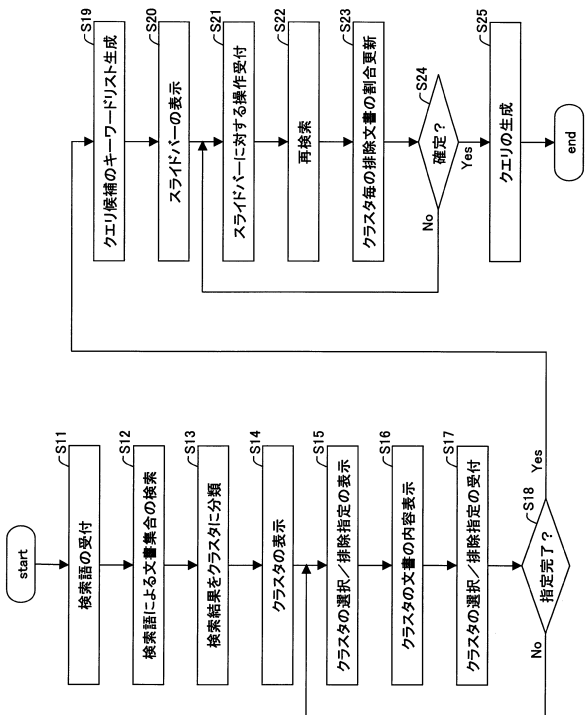
【図 4】



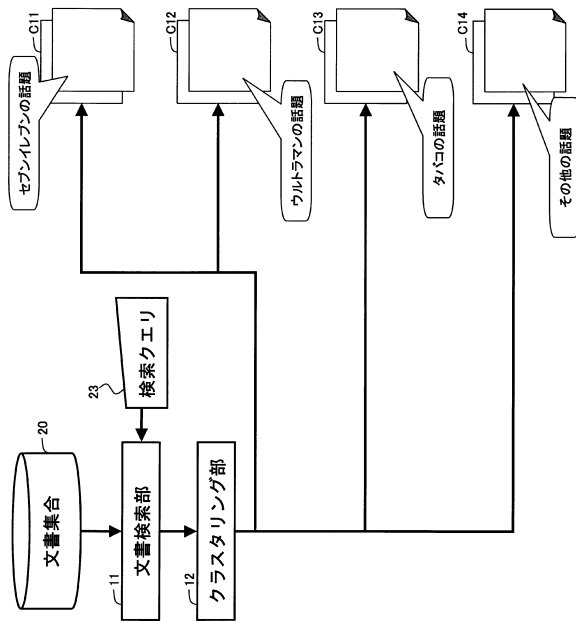
【図 5】



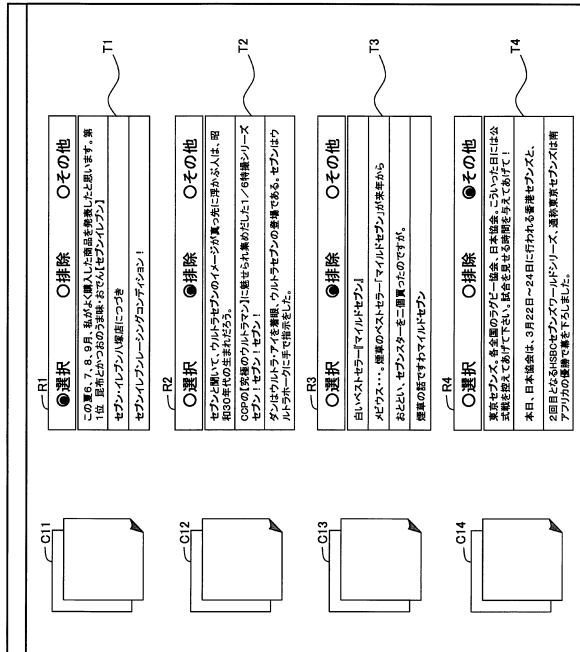
【図 6】



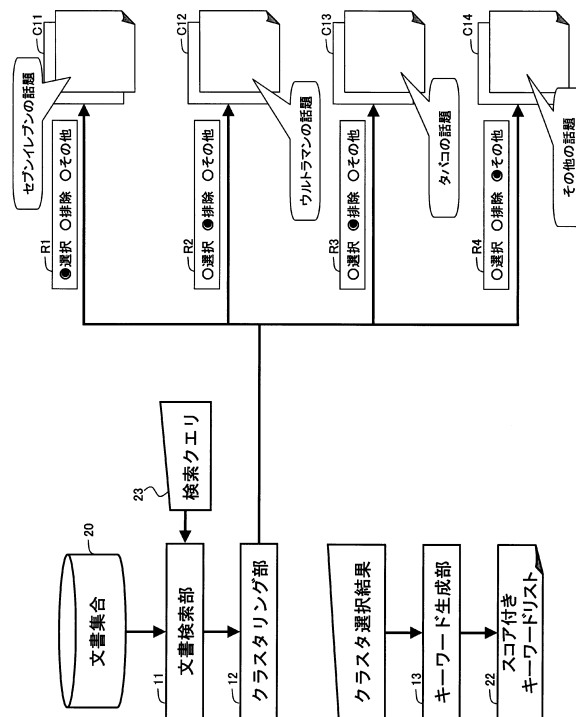
【圖 7】



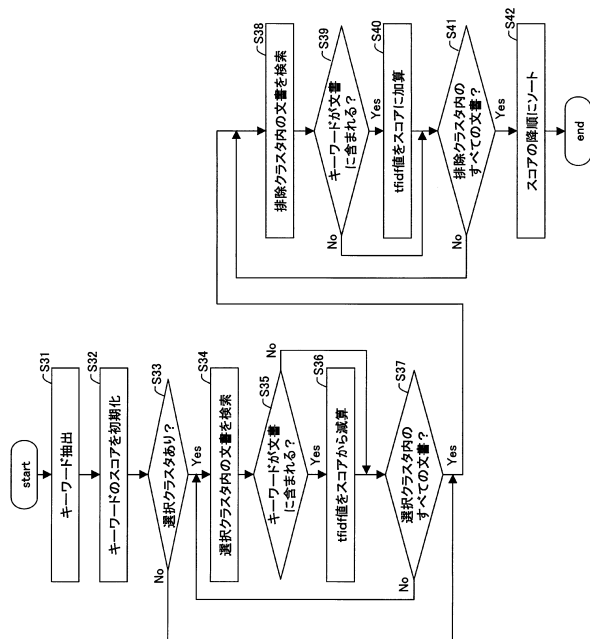
【 図 8 】



【 図 9 】



【 図 1 0 】



【図 1 1】

クラスID	文書ID	単語	tfidf値
1	doc1	コンビニ	42.7
1	doc1	おにぎり	40.3
1	:	:	:
1	doc2	四国	58.4
1	doc2	コンビニ	42.7
2	doc303	コンビニ	38.1
2	doc303	マイルドセブン	37.8
2	doc303	煙草	33.6
:	:	:	:

【図 1 2】

クラスID	文書ID	単語	tfidf値
1	doc1	コンビニ	42.7
1	doc1	おにぎり	40.3
1	:	:	:
1	doc2	四国	58.4
1	doc2	コンビニ	42.7
2	doc303	コンビニ	38.1
2	doc303	マイルドセブン	37.8
2	doc303	煙草	33.6
:	:	:	:

キーワード	スコア
コンビニ	-42.7
おにぎり	-40.3
四国	0
マイルドセブン	0
煙草	0
:	:

【図 1 3】

クラスID	文書ID	単語	tfidf値
1	doc1	コンビニ	42.7
1	doc1	おにぎり	40.3
1	:	:	:
1	doc2	四国	58.4
1	doc2	コンビニ	42.7
2	doc303	コンビニ	38.1
2	doc303	マイルドセブン	37.8
2	doc303	煙草	33.6
:	:	:	:

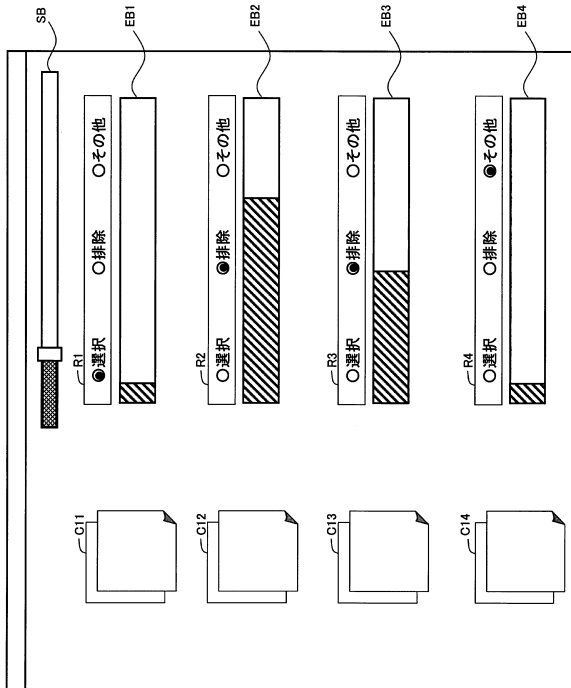
キーワード	スコア
コンビニ	-85.4
おにぎり	-40.3
四国	-58.4
マイルドセブン	0
煙草	0
:	:

【図 1 4】

クラスID	文書ID	単語	tfidf値
1	doc1	コンビニ	42.7
1	doc1	おにぎり	40.3
1	:	:	:
1	doc2	四国	58.4
1	doc2	コンビニ	42.7
2	doc303	コンビニ	38.1
2	doc303	マイルドセブン	37.8
2	doc303	煙草	33.6
:	:	:	:

キーワード	スコア
コンビニ	-47.3
おにぎり	-40.3
四国	-58.4
マイルドセブン	0
煙草	0
:	:

【図 19】



フロントページの続き

(56)参考文献 特開2004-341753(JP,A)
特開2007-172616(JP,A)
特開2007-310734(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 17/30