

(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2012-173800

(P2012-173800A)

(43) 公開日 平成24年9月10日 (2012.9.10)

(51) Int. Cl. F I テーマコード (参考)
G06F 17/30 (2006.01) G06F 17/30 350C 5B075
 G06F 17/30 210D

審査請求 未請求 請求項の数 8 O L (全 20 頁)

(21) 出願番号 特願2011-32415 (P2011-32415)
 (22) 出願日 平成23年2月17日 (2011.2.17)

(71) 出願人 801000027
 学校法人明治大学
 東京都千代田区神田駿河台 1-1
 (74) 代理人 100064908
 弁理士 志賀 正武
 (74) 代理人 100106909
 弁理士 棚井 澄雄
 (74) 代理人 100108578
 弁理士 高橋 詔男
 (74) 代理人 100126882
 弁理士 五十嵐 光永
 (72) 発明者 高木 友博
 神奈川県川崎市多摩区東三田 1-1-1
 学校法人明治大学 生田校舎内
 Fターム(参考) 5B075 NR12 PQ38 PQ75 PR06 QM08

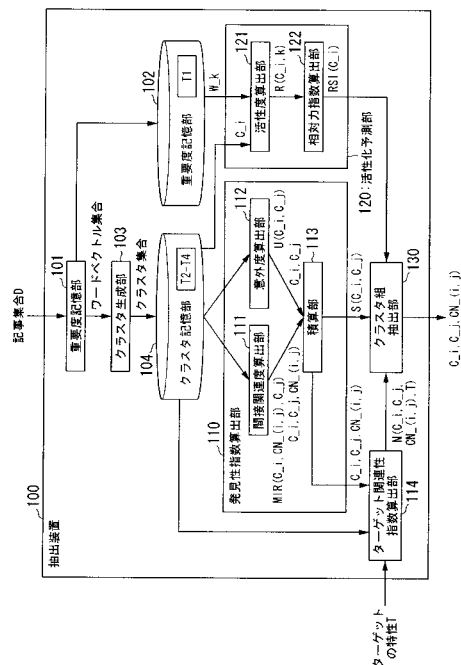
(54) 【発明の名称】 抽出装置、抽出方法および抽出プログラム

(57) 【要約】

【課題】 意外性のある概念の組み合わせを提供する。

【解決手段】 単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部104と、クラスタ記憶部104から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度と、該2つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出部110と、算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出部130を備える。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部と、

前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度と、該2つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出部と、

前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出部と、

を備えることを特徴とする抽出装置。

【請求項 2】

前記クラスタ記憶部から前記クラスタ毎に所属度を示す情報を読み出し、該読み出された所属度を示す情報と、自装置の外部から入力されたターゲットの特性を示す情報とに基づいて、前記対象となる2つのクラスタおよび前記第3のクラスタとターゲットとの関連性を示すターゲット関連性指数を算出するターゲット関連性指数算出部を更に備え、

前記クラスタ組抽出部は、前記発見性指数と前記ターゲット関連性指数とに基づいて、前記クラスタの組み合わせを抽出することを特徴とする請求項1に記載の抽出装置。

【請求項 3】

所定の期間毎に、前記単語を示す情報と該単語の重要度を示す情報とが関連付けられて記憶されている重要度記憶部と、

前記重要度記憶部から所定の期間毎に前記単語の重要度を示す情報を読み出し、前記クラスタ記憶部からクラスタ毎に前記所属度を示す情報を読み出し、該単語の重要度を示す情報と該所属度を示す情報とに基づいて、所定の期間毎に各クラスタの活性化を予測する活性化予測部を更に備え、

前記クラスタ組抽出部は、前記活性化予測部による予測により前記クラスタの組み合わせのうち少なくとも1つのクラスタの活性化が予測された場合、前記発見性指数とターゲット関連性指数とに基づいて、前記クラスタの組み合わせを抽出することを特徴とする請求項2に記載の抽出装置。

【請求項 4】

前記活性化予測部は、

所定の期間毎に、単語が所定のクラスタへ所属している所属度を示す情報と、前記重要度記憶部から読み出された該期間における前記単語の重要度を示す情報とに基づいて、該クラスタの活性化度を算出する活性化度算出部と、

前記算出された活性化度に基づき、各クラスタの活性化の上昇が期待される度合いである活性化上昇期待値を算出する活性化上昇期待値算出部と、

を備え、

前記算出された活性化度と、前記算出された活性化上昇期待値とに基づいて、前記クラスタの活性化を予測することを特徴とする請求項3に記載の抽出装置。

【請求項 5】

前記発見性指数算出部は、

前記クラスタ記憶部から前記所属度が所定値以上の単語の位置を示す情報を3つ以上のクラスタ分読み出し、該読み出された単語の位置を示す情報に基づき、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度を算出する間接関連度算出部と、

前記読み出された単語の位置を示す情報に基づき、前記クラスタの組み合わせの意外度を算出する意外度算出部と、

を備え、

10

20

30

40

50

前記間接関連度と前記意外度とを乗じることにより、前記発見性指数を算出することを特徴とする請求項 1 から請求項 4 のいずれか 1 項に記載の抽出装置。

【請求項 6】

前記発見性指数は、前記間接関連度と前記意外度が高くなるほど高くなり、
前記クラスタ組抽出部は、前記発見性指数と前記ターゲット指数との重み付き和に基づいて、前記クラスタの組み合わせを抽出することを特徴とする請求項 2 から請求項 5 のいずれか 1 項に記載の抽出装置。

【請求項 7】

単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部を備える抽出装置が実行する抽出方法であって、

前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を 3 つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる 2 つのクラスタ以外の第 3 のクラスタを経由した該 2 つのクラスタ間の間接関連度と、該 2 つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出手順と、

前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出手順と、

を有することを特徴とする抽出方法。

【請求項 8】

単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部を備える抽出装置のコンピュータに、

前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を 3 つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる 2 つのクラスタ以外の第 3 のクラスタを経由した該 2 つのクラスタ間の間接関連度と、該 2 つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出ステップと、

前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出ステップと、

を実行させるための抽出プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、抽出装置、抽出方法および抽出プログラムに関する。

【背景技術】

【0002】

現在、既存の単語を組み合わせることによって作られた造語を新商品の名前に用いることが行われている。その造語が流行するかどうかは、その造語を構成する単語の組み合わせによって変わってくるが、世の中には用いる単語の組み合わせの候補がたくさんあるので、どの単語を組み合わせればよいのかは一見ただけでは分からない。また、あらゆる単語の組み合わせに対して造語が流行するか否かを検証することは難しい。

【0003】

その問題に対して、非特許文献 1 では、組み合わせ評価システムが WEB ページ上におけるキーワードの登場回数から、単語の組み合わせの斬新さと大衆に受け入れられる可能性とを推定し、それによって組み合わせの有効度を定めることが示されている。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献 1】西原陽子、砂山渡、谷内田正彦「有効な組み合わせの発見による創造活

10

20

30

40

50

動支援」、電子情報通信学会論文誌 D - I, Vol. J87 - D - I, No. 10, pp. 939 - 949, 2004年10月

【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、非特許文献1における組み合わせ評価システムは、WEBページなどの文章に活字として掲載されているキーワードを抽出することはできるが、その文章には活字として掲載されていないが、その文章の一部あるいは全体から捉えられる概念を抽出することができず、意外性のある概念の組み合わせを提供できないという問題があった。

【0006】

そこで本発明は、上記問題に鑑みてなされたものであり、意外性のある概念の組み合わせを提供することを可能とする抽出装置、抽出方法および抽出プログラムを提供することを課題とする。

【課題を解決するための手段】

【0007】

上記の課題を解決するために、本発明の一態様である抽出装置は、単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部と、前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度と、該2つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出部と、前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出部と、を備えることを特徴とする。

【0008】

上記抽出装置は、前記クラスタ記憶部から前記クラスタ毎に所属度を示す情報を読み出し、該読み出された所属度を示す情報と、自装置の外部から入力されたターゲットの特性を示す情報とに基づいて、前記対象となる2つのクラスタおよび前記第3のクラスタとターゲットとの関連性を示すターゲット関連性指数を算出するターゲット関連性指数算出部を更に備え、前記クラスタ組抽出部は、前記発見性指数と前記ターゲット関連性指数とに基づいて、前記クラスタの組み合わせを抽出することを特徴とするものであってもよい。

【0009】

上記抽出装置は、所定の期間毎に、前記単語を示す情報と該単語の重要度を示す情報とが関連付けられて記憶されている重要度記憶部と、前記重要度記憶部から所定の期間毎に前記単語の重要度を示す情報を読み出し、前記クラスタ記憶部からクラスタ毎に前記所属度を示す情報を読み出し、該単語の重要度を示す情報と該所属度を示す情報とに基づいて、所定の期間毎に各クラスタの活性化を予測する活性化予測部を更に備え、前記クラスタ組抽出部は、前記活性化予測部による予測により前記クラスタの組み合わせのうち少なくとも1つのクラスタの活性化が予測された場合、前記発見性指数とターゲット関連性指数とに基づいて、前記クラスタの組み合わせを抽出することを特徴とするものであってもよい。

【0010】

上記抽出装置の前記活性化予測部は、所定の期間毎に、単語が所定のクラスタへ所属している所属度を示す情報と、前記重要度記憶部から読み出された該期間における前記単語の重要度を示す情報とに基づいて、該クラスタの活性化度を算出する活性化度算出部と、前記算出された活性化度に基づき、各クラスタの活性化の上昇が期待される度合いである活性化上昇期待値を算出する活性化上昇期待値算出部と、を備え、前記算出された活性化度と、前記算出された活性化上昇期待値とに基づいて、前記クラスタの活性化を予測することを特徴とするものであってもよい。

【0011】

10

20

30

40

50

上記抽出装置の前記発見性指数算出部は、前記クラスタ記憶部から前記所属度が所定値以上の単語の位置を示す情報を3つ以上のクラスタ分読み出し、該読み出された単語の位置を示す情報に基づき、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度を算出する間接関連度算出部と、前記読み出された単語の位置を示す情報に基づき、前記クラスタの組み合わせの意外度を算出する意外度算出部と、を備え、前記間接関連度と前記意外度とを乗じることにより、前記発見性指数を算出することを特徴とするものであってもよい。

【0012】

上記抽出装置の前記発見性指数は、前記間接関連度と前記意外度が高くなるほど高くなり、前記クラスタ組抽出部は、前記発見性指数と前記ターゲット指数との重み付き和に基づいて、前記クラスタの組み合わせを抽出することを特徴とするものであってもよい。

10

【0013】

本発明の一態様である抽出方法は、単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部を備える抽出装置が実行する抽出方法であって、前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度と、該2つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出手順と、前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出手順と、を有することを特徴とする。

20

【0014】

本発明の一態様である抽出プログラムは、単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが関連付けられ、前記単語を示す情報と該単語の位置を示す情報とが関連付けられて記憶されているクラスタ記憶部を備える抽出装置のコンピュータに、前記クラスタ記憶部から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度と、該2つのクラスタを組み合わせることの意外度とを反映する発見性指数を算出する発見性指数算出ステップと、前記算出された発見性指数に基づいて、前記クラスタの組み合わせを抽出するクラスタ組抽出ステップと、を実行させるための抽出プログラムである。

30

【発明の効果】

【0015】

本発明によれば、意外性のある概念の組み合わせを提供することができる。

【図面の簡単な説明】

【0016】

【図1】本発明の実施形態における抽出装置のブロック構成図である。

【図2】重要度記憶部に記憶されているワードベクトルテーブルT1の一例である。

【図3】クラスタ生成部による処理を説明するための図である。

40

【図4】活性度の算出方法を説明するための図である。

【図5】本実施形態の抽出装置がクラスタを生成する処理の流れを示したフローチャートである。

【図6】本実施形態の抽出装置がクラスタの組み合わせを抽出する処理の流れを示したフローチャートである。

【発明を実施するための形態】

【0017】

以下、本発明の実施形態について、図面を参照して詳細に説明する。まず、本発明の実施形態における抽出装置100の概要について説明する。抽出装置100は、流行語の重要要素である流行に乗っていることと、新しい驚きがあることとを両立する概念を、その

50

概念を提供する対象であるターゲット（人）の特性と関連性がある、複数の概念を組み合わせる事によって生成する。これにより、抽出装置 100 は、ターゲットの特性に応じて、世間で流行している概念であって、ターゲットにとって意外性がある概念（ヒットコンセプト）を提示することができる。

【0018】

ここで、概念は、データに出現する語の集合として表される。その特殊な場合として 1 語による概念も存在する。

組合せ要素となる概念として、2 つの概念 C 1、C 2 をつなぐ役目をする概念 C n が存在する。抽出装置 100 は、C 1、C 2、C n それぞれの概念を、新聞やウェブ上の時系列データなどから、流行要因（ヒット要因）として定められた特徴を測る測度を測定することで抽出する。

10

【0019】

抽出装置 100 は、概念 C 1 と概念 C 2 の直接の関連度は低いが、概念 C n を経由した C 1 - C n - C 2 の間接関連度は高くなる組合せを抽出する。例えば、抽出装置 100 は、ターゲットがゴルフクラブ（C 1）と関連がある所定の雑誌の読者だとすると、概念 C 1、概念 C n、概念 C 2 の組み合わせとして、ゴルフクラブ（C 1）、口紅（C 2）、C n（プレゼント）を抽出する。一見、ゴルフクラブと口紅の関連度は低いが、プレゼントという概念 C n を経由すると両者の間接関連度は高くなるので、ゴルフクラブ（C 1）と口紅（C 2）の組み合わせを抽出する価値は高い。

20

【0020】

さらに、抽出装置 100 は、それら概念が対象とする期間において活性化傾向にあることと、それらの概念の少なくとも 1 つがターゲットの特性と関連があることも概念の抽出の条件とする。例えば、クリスマス時期において、プレゼントという概念の活性化傾向は強くなり、ターゲットとしての所定の雑誌の読者にとってゴルフクラブの関連度は高い。

【0021】

抽出装置 100 は、上記概念 C 1、C 2、C n の組合せを、ターゲットにとって目新しい概念を示す情報とし出力する。これにより、抽出装置 100 は、ターゲットに対して、ターゲットと関連している概念（例えば、概念 C 1）と、接続概念 C n を介して関連している概念 C 2 を示す情報を提供することができる。これにより、例えば、ターゲットであるゴルフクラブ（概念 C 1）と関連している所定の雑誌の読者に対して、クリスマス時期の流行概念（ヒットコンセプト）として、プレゼント（概念 C n）のための口紅（概念 C 2）特集を提供することができる。

30

【0022】

図 1 は、本発明の実施形態における抽出装置 100 のブロック構成図である。抽出装置 100 は、重要度算出部 101 と、重要度記憶部 102 と、クラスタ生成部 103 と、クラスタ記憶部 104 と、発見性指数算出部 110 と、ターゲット関連性指数算出部 114 と、活性化予測部 120 と、クラスタ組抽出部 130 とを備える。

また、発見性指数算出部 110 は、間接関連度算出部 111 と、意外度算出部 112 と、積算部 113 とを備える。活性化予測部 120 は、活性化度算出部 121 と、相対力指数算出部（活性化上昇期待値算出部）122 とを備える。

40

【0023】

重要度算出部 101 は、自装置の外部から入力された記事集合 D を受け付ける。ここで、入力される記事集合 D は新聞のような世相を表すドキュメントや雑誌のような市場の特性を表すドキュメントの時系列データである。そして、重要度算出部 101 は、記事集合 D から所定期間のドキュメントを一区切りとし、それを時系列順にならべたものを生成する。ここで、一区切りのドキュメントを 1 つのドキュメント、全期間のドキュメントを全ドキュメントと称する。

【0024】

重要度算出部 101 は、各期間における単語の重要度を示す情報を算出する。具体的には、例えば、重要度算出部 101 は、期間毎にドキュメント中に注目語が出現した頻度 t

50

f を、当該ドキュメント中の総単語数で割ることにより、各期間における各語の t f - i d f 値を算出する。ここで、t f - i d f 値とは、情報検索で一般的に語の重要度として使用されている指標である。

【 0 0 2 5 】

重要度算出部 1 0 1 は、この t f - i d f 値を事前に定められた語順に並べたものであるワードベクトルを当該所定期間毎に算出する。このワードベクトルは、各語の t f - i d f 値のリストであり、その期間の特徴を表している。重要度算出部 1 0 1 は、算出したワードベクトルを示す情報を該単語と関連付けて、期間毎に重要度記憶部 1 0 2 のワードベクトルテーブル T 1 に記憶させる。

【 0 0 2 6 】

図 2 は、重要度記憶部 1 0 2 に記憶されているワードベクトルテーブル T 1 の一例である。同図において、上記所定期間を 1 日と定め、1 日毎の単語の t f - i d f 値が予め決められた単語の順番で示されている。また、各列はワードベクトル (W _ 1、W _ 2、W _ 3、...、W _ 3 0) を表している。

このように、このワードベクトルを時系列順に並べることによって、時間順に所定期間毎の記事の特徴が示される。

【 0 0 2 7 】

図 1 に戻って、重要度算出部 1 0 1 は、ワードベクトルの情報の集合 (以下、ワードベクトル集合と称する) をクラスタ生成部 1 0 3 に出力する。

クラスタ生成部 1 0 3 は、重要度算出部 1 0 1 から入力されたワードベクトル集合を用いて、単語を所定のまとまりであるクラスタに分類し、クラスタ毎にラベルを付与する。

【 0 0 2 8 】

本実施形態では、概念は何らかの共通性や関連性によって類似の語の集合で表されると仮定する。ここで言う集合とは、その集合の要素であるかどうかの所属度が 0 または 1 で決まる通常の集合の場合も、要素の所属度を 0 から 1 までの間の任意の値で表すファジィ集合の場合の、両方の可能性がある。

【 0 0 2 9 】

そこで、クラスタ生成部 1 0 3 は、所定のクラスタリング方法に従って、記事集合 D に出現する単語をクラスタリングする。通常 1 つのクラスタには数万の単語が含まれ、それぞれの単語はクラスタに所属する値である所属度 $M e m_{C}(w)$ を有する。ここで、所属度 $M e m_{C}(w)$ は、単語 w がクラスタ C に所属する値を表している。この値は、クラスタが対応している概念に所属する程度を意味する。

【 0 0 3 0 】

クラスタリングにはすでに様々な手法が提案されているが、クラスタ生成部 1 0 3 は、一例として、k - m e a n s 法によって、記事集合 D に出現する単語をクラスタリングする。具体的には、クラスタ生成部 1 0 3 は、下記式 (1) で表される評価値を最小化するクラスタを算出する。ここで、k は事前に与えられるものとする。

【 0 0 3 1 】

【 数 1 】

$$\sum_{k=1}^M \sum_{i=1}^N g_{ik} |x_i - v_k|^2 \quad (1)$$

【 0 0 3 2 】

但し、以下の条件式 (2) を満たすものとする。

【 0 0 3 3 】

10

20

30

40

【数 2】

$$\sum_{k=1}^K g_{ik} = 1, \quad g_{ik} \in \{0,1\} \quad (2)$$

【0034】

ここで、 x_i は i 番目の単語データ (i は 1 から I までの整数) で、 $x_i = (x_{i1}, x_{i2})$ 、 K はクラスタ数、 v_k は k 番目のクラスタの重心 (k は 1 から K までの整数) で、 $v_k = (v_{k1}, v_{k2})$ 、 g_{ik} は i 番目のデータの k 番目のクラスタへの所属度である。

10

【0035】

なお、クラスタ生成部 103 は、 k -means 法を用いたがこれに限らず、fuzzy c-means 法を用いてもよい。その場合、具体的には、クラスタ生成部 103 は、下記式 (3) で表される評価値を最小化するクラスタを算出する。ここで、 k は事前与えられるものとする。

【0036】

【数 3】

$$\sum_{k=1}^K \sum_{i=1}^I (g_{ik})^m |x_i - v_k|^2 \quad (3)$$

20

【0037】

但し、以下の条件式 (4) を満たすものとする。

【0038】

【数 4】

$$\sum_{k=1}^K g_{ik} = 1, \quad g_{ik} \in [0,1] \quad (4)$$

【0039】

ここで、 x_i は i 番目の単語データ (i は 1 から I までの整数) で、 $x_i = (x_{i1}, x_{i2})$ 、 K はクラスタ数、 v_k は k 番目のクラスタの重心 (k は 1 から K までの整数) で、 $v_i = (v_{i1}, v_{i2})$ 、 g_{ik} は i 番目のデータの k 番目のクラスタへの所属度である。

30

このように、クラスタ生成部 103 は、 k -means 法、fuzzy c-means 法のいずれを用いても、要素毎にクラスタに所属する所属度を算出する。

【0040】

クラスタ生成部 103 は、得られたクラスタ 1 つずつに 1 つの概念を割り当てるためにラベルを付与する。具体的には、クラスタ生成部 103 は、クラスタ重心に最も近い語をそのクラスタの代表として、そのクラスタのラベルとする。なお、クラスタ生成部 103 は、クラスタ中の最大の所属度を持つ語をそのクラスタの代表としてそのクラスタのラベルとしてもよい。

40

【0041】

図 3 は、クラスタ生成部 103 による処理を説明するための図である。図 3 (a) は、クラスタ生成部 103 によって生成されるクラスタを説明するための図である。同図において、向かって左側に記事集合 D が示されている。向かって右側には、 x, y の 2 次元平面上にクラスタの 1 例が示されている。

【0042】

その 2 次元平面上で、クラスタの各要素である各単語は、 \times 印で示されている。3 つのクラスタ C_1 、 C_2 、 C_3 が示されており、各クラスタは円内の \times 印で示された単

50

語を含むものとする。クラスタC__1は農産物のラベルが付与されたクラスタであり、その要素にはprocessorとorangeを含む。一方、クラスタC__2はコンピュータのラベルが付与されたクラスタであり、要素にはprocessor、memoryを含む。すなわち、processorは、食品加工機（フードプロセッサ）という意味でクラスタC__1に所属し、コンピュータのプロセッサの意味でクラスタC__2に所属している。

【0043】

クラスタC__3は脳のラベルが付与されたクラスタであり、要素にはmemoryを含む。すなわち、memoryは、コンピュータのメモリという意味でクラスタC__2に所属し、脳の記憶という意味でクラスタC__3に所属している。

10

【0044】

図3(b)は、クラスタ記憶部104に記憶されている概念テーブルT2の1例である。概念テーブルT2には、図3(a)に示されたクラスタを識別する識別情報C__i(iは正の整数)と、図3(a)に示されたクラスタ毎に付与されたラベルを示す情報とが関連付けられている。

【0045】

図3(c)は、クラスタ記憶部104に記憶されている所属度テーブルT3の1例である。所属度テーブルT3には、図3(a)に示された単語を示す情報と、該単語がクラスタに所属している程度である所属度を示す情報とが該クラスタを識別する識別情報C__i毎に関連付けられている。

20

【0046】

図3(d)は、クラスタ記憶部104に記憶されている座標テーブルT4の1例である。座標テーブルT4には、図3(a)に示された単語を示す情報と、該単語の位置を示す情報である座標を示す情報とが関連付けられている。

【0047】

図1に戻って、クラスタ生成部103は、クラスタ識別情報C__i(これ以降、iはクラスタのインデックスを表す1からnまでの正の整数)と、クラスタ毎に付与されたラベルを示す情報とを関連付けてクラスタ記憶部104に記憶させる。また、クラスタ生成部103は、単語を示す情報と、該単語がクラスタに所属している程度である所属度を示す情報とを該クラスタを識別する識別情報C__i毎に関連付けてクラスタ記憶部104に記憶させる。また、クラスタ生成部103は、クラスタ記憶部104に、単語を示す情報と当該単語の位置を示す情報とを関連付けて記憶させる。

30

【0048】

またクラスタ記憶部104には、図3(b)に示されたように、クラスタ生成部103による処理の結果、クラスタを識別する識別情報C__iと、クラスタ毎に付与されたラベルを示す情報とが関連付けられて記憶されている。

またクラスタ記憶部104には、図3(c)に示されたように、クラスタ生成部103による処理の結果、単語を示す情報と該単語がクラスタに所属している程度である所属度を示す情報とが該クラスタ毎に関連付けられて記憶されている。

【0049】

クラスタ記憶部104には、クラスタ生成部103による処理の結果、図3(d)に示されるように、単語を示す情報と、当該単語の位置を示す情報とが関連付けられて記憶されている。ここで、例えば、クラスタ生成部103によるクラスタリングにより2次元平面上に、各単語の位置が割り当てられている場合、当該各単語の位置を示す情報は、2次元平面上における座標を示す情報である。

40

【0050】

発見性指数算出部110は、クラスタ記憶部104から異なるクラスタに関連付けられている所属度を示す情報を所定の数(例えば、3つ)のクラスタ分読み出し、当該読み出された所属度を示す情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した当該2つのクラスタ間の関連度と、該2つのクラスタを組み合わせことの意外度

50

とを反映する発見性指数を算出する。

ここで、発見性指数は2つのクラスタ同士の直接の関連性が低くなるほど高くなり、該2つのクラスタが残りの第3のクラスタと関連性が高くなるほど高くなる。

【0051】

間接関連度算出部111は、クラスタ記憶部104から所属度が所定値以上の単語を示す情報に関連付けられている単語の位置の情報を3つ以上のクラスタ分読み出し、該単語の位置の情報に基づいて、対象となる2つのクラスタ以外の第3のクラスタを経由した該2つのクラスタ間の間接関連度を算出する。一例として、間接関連度算出部111は、対象となる2つのクラスタ以外の第3のクラスタを経由したクラスタ間の関連度のうち最大となる最大間接関連度MIRを算出する。

10

【0052】

具体的には、例えば、間接関連度算出部111は、クラスタC_iとクラスタC_j(これ以降、jはクラスタのインデックスを表す1からnまでの整数)が、接続クラスタC_Nを経由して関連している程度を示す間接関連度のうち、接続クラスタC_NをC₁からC_nまで変化させながら間接関連度を算出し、算出されたn個の間接関連度のうち最大となる最大間接関連度MIRを、下記式(5)を用いて算出する。ここで、接続クラスタC_Nは、C₁からC_Nまでのクラスタを取りうる。

【0053】

$$MIR(C_i, C_j) = \max_{C_N} \{ A(C_i, C_N) \times A(C_N, C_j) \} \quad (5)$$

20

【0054】

ここで、 \max_{C_N} は、引数である右辺の間接関連度が最大となる接続クラスタC_Nを抽出し、そのときの引数の値を出力する関数で、Aは第1の引数と第2の引数の関連度を算出する関数である。

なお、間接関連度算出部111は、クラスタC_iとクラスタC_jが、接続クラスタC_Nを経由して関連している程度を示す最大間接関連度MIRを、下記式(6)を用いて算出してもよい。

【0055】

$$MIR(C_i, C_j) = \max_{C_N} \{ A(C_i, C_N) + A(C_N, C_j) \} \quad (6)$$

30

【0056】

間接関連度算出部111は、式(5)または式(6)の中の関連度Aを、コサイン類似度を用いて算出する。

【0057】

一例として、間接関連度算出部111がコサイン類似度を用いて関連度Aを算出する方法について説明する。

ベクトルxは原点からクラスタC_iの重心へのベクトル、ベクトルyを原点からクラスタC_jの重心へのベクトルである。例えば、間接関連度算出部111は、以下の式(7)に従って、関連度Aを算出する。

【0058】

$$A(C_i, C_j) = x \cdot y / (|x| \times |y|) \quad (7)$$

40

【0059】

ここで、 $x \cdot y$ はベクトルx、yの内積であり、 $(x_1 \times y_1 + x_2 \times y_2 + \dots + x_m \times y_m)$ で表される(mは正の整数)。また、 $|x|$ はベクトルxのノルム $= \sqrt{(x \cdot x)}$ である。式(7)の右辺は、ベクトルx、yのなす角の余弦 \cos を表し、コサイン類似度と呼ばれ、ベクトルの向きの近さ類似性を表す。

【0060】

なお、間接関連度算出部111は、式(5)または式(6)の中の関連度Aを、ジャカード係数または相互情報量などの方法を用いて算出してもよい。

ジャカード係数を用いる場合には、間接関連度算出部111は、C_i、C_jが通常

50

のクラスタの場合、2つのクラスタ C_i 、 C_j のどちらかに出現した単語の出現回数によって関連度 A を算出する。具体的には、間接関連度算出部 111 は、以下の式 (8) に従って関連度 A を算出する。

【0061】

【数5】

$$A(C_i, C_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|} \quad (8)$$

【0062】

10

ここで、 $|C|$ はクラスタ C に含まれる要素 (単語) 数である。この関連度 A が大きいほど、二つのクラスタの類似性は高い。

クラスタ C_i 、クラスタ C_j が *fuzzy c-means* 法で算出されたファジィ集合である場合、間接関連度算出部 111 は、 x_p をクラスタ C_i のワードベクトル x の p 番目要素 (p は 1 から P までの整数)、 y_q をクラスタ C_j のワードベクトル y の q 番目の要素とすると (q は 1 から Q までの整数)、クラスタ C_i 、クラスタ C_j の関連度を次式 (9) で算出する。

【0063】

【数6】

$$A(C_i, C_j) = x \cdot y / (\sum_{p=1}^P x_p + \sum_{q=1}^Q y_q - x \times y) \quad (9)$$

20

【0064】

一方、相互情報量を用いる場合には、間接関連度算出部 111 は、下記の式 (10) に従って、クラスタ C_i 、クラスタ C_j の相互情報量 $MI(C_i, C_j)$ を関連度 A として算出する。ここで、相互情報量は、ある2つの単語が共起する割合によって求められる関連性の指標である。

【0065】

【数7】

$$A(C_i, C_j) = MI(C_i, C_j) = \sum_{p=1}^P \sum_{q=1}^Q [P(x_p, y_q) \times \log\{P(x_p, y_q) / (P(x_p) \times P(y_q))\}] \quad (10)$$

30

【0066】

ここで、 x_p は C_i のワードベクトル x の p 番目の要素、 y_q は C_j のワードベクトル y の q 番目の要素、 $P(x_p, y_q)$ は x_p と y_q の同時出現確率、 $P(x_p)$ 、 $P(y_q)$ は、それぞれ x_p 、 y_q の周辺出現確率である。

【0067】

40

間接関連度算出部 111 は、クラスタ C_i とクラスタ C_j の全ての組み合わせで、最大間接関連度 $MIR(C_i, CN_i, j, C_j)$ を算出する。ここで、 CN_i, j は、クラスタ C_i とクラスタ C_j との間接関連度が最大となる時に選択されたクラスタであり、クラスタ C_i とクラスタ C_j の組み合わせ毎にクラスタ $C_1 \sim C_N$ までの中から選択されたクラスタである。

間接関連度算出部 111 は、算出した全ての最大間接関連度 $MIR(C_i, CN_i, j, C_j)$ を示す情報と、その各最大間接関連度 MIR を算出する際に用いたクラスタ C_i 、 CN_i, j 、 C_j の組み合わせを示す情報とを積算部 113 に出力する。

【0068】

50

意外度算出部 112 は、クラスタ記憶部 104 から所属度が所定値以上の単語を示す情報を 3 つ以上のクラスタ分読み出し、該読み出された単語の位置を示す情報に基づき、クラスタの組み合わせの意外度 U を算出する。具体的には、例えば、意外度算出部 112 は、式 (7) の関連度の式の逆数を意外度として使用し、以下の式に従って、クラスタ C_i とクラスタ C_j 間の意外度 $U(C_i, C_j)$ を算出する。

【0069】

$$U(C_i, C_j) = (|x| \times |y|) / x \cdot y \quad (11)$$

【0070】

ここで、ベクトル x は原点からクラスタ C_i の重心へのベクトル、ベクトル y を原点からクラスタ C_j の重心へのベクトルである。

10

【0071】

なお、意外度算出部 112 は、ジャカード係数の逆数 (式 (7) の右辺の逆数) を用いて、意外度を算出してもよい。その場合、具体的には、意外度算出部 112 は、下記の式 (12) に従って、クラスタ C_i とクラスタ C_j 間の意外度 $U(C_i, C_j)$ を算出する。

【0072】

【数 8】

$$U(C_i, C_j) = \frac{|C_i \cup C_j|}{|C_i \cap C_j|} \quad (12)$$

20

【0073】

ここで、クラスタ C_i とクラスタ C_j の関連性が低いほど、意外度 $U(C_i, C_j)$ は高くなり、両クラスタの組み合わせが意外であることを反映している。

また、意外度算出部 112 は、相互情報量 MI の逆数 (式 (10) の右辺の逆数) を用いて、意外度を算出してもよい。その場合、具体的には、意外度算出部 112 は、下記の式 (13) に従って、クラスタ C_i とクラスタ C_j 間の意外度 $U(C_i, C_j)$ を算出する。

【0074】

$$U(C_i, C_j) = 1 / MI(C_i, C_j) \quad (13)$$

30

【0075】

意外度算出部 112 は、クラスタ C_i とクラスタ C_j の全ての組み合わせで、意外度 $U(C_i, C_j)$ を算出する。

意外度算出部 112 は、算出した全ての意外度 $U(C_i, C_j)$ を示す情報と、その各意外度 $U(C_i, C_j)$ が算出された際に用いられたクラスタ C_i の識別情報とクラスタ C_j の識別情報とを積算部 113 に出力する。

【0076】

続いて、積算部 113 は、最大間接関連度 $MI R$ と意外度 U に基づいて、発見性指数を算出する。具体的には、積算部 113 は、対象となる 2 つのクラスタ (C_i 、 C_j) 以外の第 3 のクラスタ $C N$ を経由した該 2 つのクラスタ (C_i 、 C_j) 間の関連度と、該 2 つのクラスタ (C_i 、 C_j) を組み合わせることの意外度とを反映するクラスタ発見性指標 S を下記式 (14) に従って、算出する。

40

【0077】

$$S(C_i, C_j) = MI R(C_i, C_j) \times U(C_i, C_j) \quad (14)$$

【0078】

発見性指標 S は、クラスタ C_i とクラスタ C_j との間でクラスタ $C N$ を経由した関連性が必要なこと、また同時にクラスタ C_i とクラスタ C_j との組み合わせに新たな意外性が必要なことを両立させるための指標である。すなわち、発見性指標 S は、2 つのクラスタ (C_i 、 C_j) 同士の直接の関連性が低くなるほど高くなり、該 2 つのクラ

50

スタが残りの第3のクラスタ (CN__(i, j)) と関連性が高くなるほど高くなる。

【0079】

積算部113は、クラスタC__iとクラスタC__jの全ての組み合わせで、発見性指標Sを算出し、算出した発見性指標Sを示す情報をクラスタ組抽出部130に出力する。また、積算部113は、クラスタC__iを示す情報とクラスタC__jを示す情報と接続クラスタCN__(i, j)を示す情報とをターゲット関連性指数算出部114に出力する。

【0080】

ターゲット関連性指数算出部114は、自装置の外部から入力されたターゲットの特性(例えば、ターゲットとなる世相、市場、個人の特性)Tを示す情報を受け付ける。また、ターゲット関連性指数算出部114は、積算部113から入力されたクラスタC__iを示す情報とクラスタC__jを示す情報と接続クラスタCN__(i, j)を示す情報とを受け付ける。

10

【0081】

ターゲット関連性指数算出部114は、クラスタ記憶部104からクラスタ(C__i、C__j、CN__(i, j))毎に所属度を示す情報を読み出し、該読み出された所属度を示す情報と、自装置の外部から入力されたターゲットの特性を示す情報Tとに基づいて、前記異なる3つのクラスタ(C__i、C__j、CN__(i, j))とターゲットとの関連性を示すターゲット関連性指数Nを算出する。

【0082】

具体的には、例えば、ターゲット関連性指数算出部114は、下記の式(15)に従って、ターゲット関連性指数Nを算出する。

20

【0083】

$$N(C_i, C_j, CN_i, j, T) = \min(A(C_i, T), A(C_j, T), A(CN_i, j, T)) \quad (15)$$

【0084】

ターゲット関連性指数算出部114は、算出したターゲット関連性指数Nを示す情報をクラスタ組抽出部130に出力する。

【0085】

活性度算出部121は、各期間のワードベクトルを示す情報を重要度記憶部102から読み出し、該読み出された各期間のワードベクトルを示す情報に基づいて、各期間における各クラスタの活性度を算出する。

30

具体的には、例えば、活性度算出部121は、k番目の期間においてi番目のクラスタC__iの活性度をR(C__i, k)とすると、下記の式(16)に従って、活性度を算出する。

【0086】

$$R(C_i, k) = \text{sim}(Y_i, W_k) \quad (16)$$

【0087】

ここで、Y__iはクラスタC__iに所属する単語の所属度から構成される所属度ベクトルであり、W__kは、k番目(kは正の整数)の期間の文書のワードベクトルである。

上記の式(15)は、活性度算出部121は、k番目の期間の文書のワードベクトルW__kと、クラスタC__iを表す所属度ベクトルY__iとの類似度を、そのままそのクラスタC__iの活性度として求めるものである。

40

また、関数simは類似度を表す関数で、コサイン類似度を用いた下記の式(17)で表される。

【0088】

$$\text{sim}(Y_i, W_k) = Y_i \cdot W_k / (|Y_i| \times |W_k|) \quad (17)$$

【0089】

図4は、活性度の算出方法を説明するための図である。同図において、所属度ベクトル401の各要素は、そのクラスタに属する単語(Word 1~Word M)の所属度

50

が示されている（Mは正の整数）。また、k番目の期間の文書のワードベクトル402の各要素は、k番目の期間の文書におけるそのクラスタに属する単語（Word 1 ~ Word M）のtf-idf値が示されている。

【0090】

なお、活性度算出部121は、関数simとしてジャカード係数を用いてもよい。また、活性度算出部121は、下記の式（18）に従って、クラスタC_iの活性度R（C_i）を算出してもよい。

【0091】

【数9】

$$R(C_i) = \sum_{p=1}^P \sum_{q=1}^Q \{ \text{mem}C_i(y_q) \times MI(x_p, y_q) \times \text{tfidf}(x_p) \} \quad (18)$$

10

【0092】

ここで、memC_i(y_q)は単語y_qのクラスタC_iへの所属度である。MI(x_p, y_q)は、単語x_pと単語y_qとの相互情報量である。tfidf(x)はワードベクトル中の単語x_pのtfidf値である。

【0093】

なお、活性度算出部121は、各概念に含まれる語すべてを用いて計算する代わりに、tfidf値の高い一定数の上位単語またはtfidf値が所定の値を超えた単語のtfidf値から構成されるワードベクトルに基づいて活性度を算出してもよい。これにより、活性度算出部121は、計算回数を少なくすることができるので、計算に係る時間を短縮することができる。

20

【0094】

活性度算出部121は、算出した各期間のクラスタC_iの活性度R（C_i, k）を示す情報を相対力指数算出部122に出力する。

相対力指数算出部122は、活性度算出部121から入力された各期間のクラスタC_iの活性度R（C_i, k）に基づいて、それぞれのクラスタの活性度の時間的変化に注目し、世の中一般やターゲット市場さらには個人で、各クラスタの活性度の上昇が期待される度合い（活性度上昇期待値）を算出する。

30

【0095】

具体的には、例えば、相対力指数算出部122は、活性度上昇期待値の一例として、相対力指数RSI（C_i）を算出する。ここで、相対力指数（RSI）とは、過去の値の動きに対する上昇幅の割合を求めたもので、一般にRSI値が30を切ると、上昇傾向になると言われている。相対力指数算出部122は相対力指数（RSI）を算出する際に、例えば1カ月あるいは1日のような所定の長さのサンプリング期間を設けて、そのサンプリング期間内の活性度の上昇値と下降値から、相対力指数（RSI）を算出する。

例えば、相対力指数算出部122は、下記の式（19）に従って、相対力指数（RSI）を算出する。

【0096】

$$RSI = u / (u + d) \times 100 \quad (19)$$

40

【0097】

ここで、uは所定のサンプリング期間の活性度の上昇値の合計、dは所定のサンプリング期間の活性度の下降値の合計である。

なお、相対力指数算出部122は、活性度上昇期待値として相対力指数RSIを用いたが、これに限らず、他の経済指標を用いてもよい。

【0098】

そして、活性化予測部120は、算出された活性度と、算出された活性度上昇期待値とに基づいて、クラスタの活性化を予測する。

具体的には、活性化予測部120は、上記の30という値を一般化して閾値Lとし、上

50

昇を予測する条件を下記の2つとする。1つ目は、(i)過去の一定期間の間に相対力指数(RSI)が閾値Lを下回ったことがあること、2つ目は、(ii)現在の活性値Rが上限 R_u 、下限 R_L の間にあることである。活性化予測部120は、これら2つの条件を満たしたときに、これからのクラスタの活性化を予測し、それ以外の場合、これからクラスタが活性化しないと予測する。

【0099】

活性化予測部120は、予測結果を示す情報をクラスタ組抽出部130に出力する。

クラスタ組抽出部130は、積算部113から発見性指標Sを示す情報を、ターゲット関連性指数算出部114からターゲット関連性指数Nを示す情報を、活性化予測部120から予測結果を示す情報を受け取る。

10

【0100】

クラスタ組抽出部130は、活性化予測部120による予測により前記クラスタの組み合わせのうち少なくとも1つのクラスタの活性化が予測された場合、発見性指数Sとターゲット関連性指数Nとに基づいて、クラスタの組み合わせを抽出する。

具体的には、クラスタ組抽出部130は、下記の3つの条件に基づいて、クラスタの組み合わせ($C_{i,j}$ 、 $C_N(i,j)$)を抽出する。

【0101】

(1)新規発見性指数Sの条件として、クラスタの組 $C_{i,j}$ 、 $C_N(i,j)$ の発見性指標Sが所定の値以上であること、

(2)活性化予測の条件として、クラスタ $C_{i,j}$ 、クラスタ $C_N(i,j)$ のいずれかの相対力指数(RSI)と活性度Rが、それぞれ上述のクラスタの活性化予測条件(i)および(ii)を満足していること、

20

(3)ターゲット関連性指数Nの条件として、クラスタの組 $C_{i,j}$ 、 $C_N(i,j)$ のいずれかが、ターゲットの特性Tと所定の値以上の関連度を持つことである。

【0102】

例えば、クラスタ組抽出部130は、あるターゲットの特性Tが存在した時、特性Tにとっての最適なクラスタの組み合わせ($C_{i,j}$ 、 $C_N(i,j)$)を、下記の式(20)から算出する。

【0103】

$$\text{arg max} \{ a S (C_{i,j}, C_N(i,j), T) + b N (C_{i,j}, C_N(i,j), T) \} \quad (20)$$

30

【0104】

ここで、a、bはS、Nに対する重みを表す係数であり、arg maxは、引数が最大となる値を求める関数である。この式(18)により、クラスタ組抽出部130は、引数の値が最大となるクラスタの組み合わせを抽出することができる。ただし、 $C_{i,j}$ 、 $C_N(i,j)$ のうちいずれかの相対力指数(RSI)と活性度Rが、それぞれクラスタの活性化予測条件(i)および(ii)を満足していることとする。

【0105】

なお、本実施形態では、クラスタ組抽出部130は、一例として、式(20)の引数が最大となるクラスタの組み合わせを1つ抽出したが、これに限ったものではない。クラスタ組抽出部130は、式(20)の引数の値が所定の値以上となる1つ以上のクラスタの組み合わせすべてを抽出してもよい。また、クラスタ組抽出部130は、式(20)の引数の値が高いほうからトップM(Mは正の整数)のクラスタの組み合わせすべてを抽出してもよい。

40

【0106】

そして、クラスタ組抽出部130は、抽出したクラスタの組み合わせを構成するクラスタ $C_{i,j}$ を示す情報とクラスタ $C_N(i,j)$ を示す情報とを自装置の外部に出力する。

なお、クラスタ組抽出部130は、抽出したクラスタの組み合わせを構成する各クラスタに関連付けられたラベルをそれぞれクラスタ記憶部104のテーブルT2から読み出し

50

、読み出した各ラベルを示す情報をヒットコンセプトの組み合わせを示す情報として自装置の外部に出力してもよい。

【0107】

図5は、本実施形態の抽出装置100がクラスタを生成する処理の流れを示したフローチャートである。まず、重要度算出部101は、所定期間毎の一区切りのドキュメント中に掲載された各単語のtf-idf値の算出する(ステップS101)。次に、重要度算出部101は、所定期間毎に、各単語のtf-idf値が予め決められた単語順に並べられたワードベクトルを算出する(ステップS102)。

【0108】

重要度算出部101は、全期間のドキュメントでワードベクトルを算出したか判定する(ステップS103)。重要度算出部101は、全期間のドキュメントでワードベクトルを算出していない場合(ステップS103 NO)、ステップS101の処理に戻る。一方、重要度算出部101が、全期間のドキュメントでワードベクトルを算出した場合(ステップS103 YES)、クラスタ生成部103は、クラスタを生成する(ステップS104)。

10

【0109】

次に、クラスタ生成部103は、単語毎にクラスタへの所属度を算出する(ステップS105)。次に、クラスタ生成部103は、クラスタ毎にクラスタのラベルを抽出する(ステップS106)。次に、クラスタ生成部103は、クラスタの識別情報とクラスタのラベルを示す情報とを関連付けて、クラスタ記憶部104に記憶させる(ステップS107)。次に、クラスタ生成部103は、単語を示す情報と各クラスタへの所属度を示す情報とをクラスタ毎に関連付けてクラスタ記憶部104に記憶させる(ステップS108)。以上で、本フローチャートの処理を終了する。

20

【0110】

以上により、抽出装置100は、記事集合Dから所定期間毎の一区切りのドキュメント中に掲載された各単語の重要度を算出することができる。また、抽出装置100は、記事集合Dからクラスタを生成することができる。

【0111】

図6は、本実施形態の抽出装置100がクラスタの組み合わせを抽出する処理の流れを示したフローチャートである。まず、間接関連度算出部111は、最大間接関連度MIRを算出する(ステップS201)。次に、間接関連度算出部111は、全てのクラスタの組み合わせで最大間接関連度MIRを算出したか否か判定する(ステップS202)。間接関連度算出部111は、全てのクラスタの組み合わせで最大間接関連度MIRを算出していない場合(ステップS202 NO)、ステップS201の処理に戻る。

30

【0112】

一方、間接関連度算出部111が全てのクラスタの組み合わせで最大間接関連度MIRを算出した場合(ステップS202 YES)、意外度算出部112は、意外度Uを算出する(ステップS203)。次に、意外度算出部112は、全てのクラスタの組み合わせで意外度Uを算出したか否か判定する(ステップS204)。意外度算出部112は、全てのクラスタの組み合わせで意外度Uを算出していない場合(ステップS204 NO)、ステップS203の処理に戻る。

40

【0113】

一方、意外度算出部112が全てのクラスタの組み合わせで意外度Uを算出した場合(ステップS204 YES)、積算部113は、発見性指標を算出する(ステップS205)。次に、積算部113は、全期間のドキュメントで発見性指標を算出したか否か判定する(ステップS206)。積算部113は、全期間のドキュメントで発見性指標を算出していない場合(ステップS206 NO)、ステップS201の処理に戻る。

【0114】

一方、積算部113が全期間のドキュメントで発見性指標を算出した場合(ステップS206 YES)、ターゲット関連性指数算出部114は、ターゲット関連性指数を算出

50

する（ステップS207）。

【0115】

ステップS201～ステップS207までの処理に並行して、抽出装置100は、ステップS208～ステップS215までの処理を行う。その際、始めに抽出装置100は、i、j、kを初期化する。次に、処理活性度算出部121は、k番目の期間においてi番目のクラスタC_iの活性度を算出する（ステップS208）。次に、活性度算出部121は、全てのクラスタの活性度を算出したか否か判定する（ステップS209）。活性度算出部121は、全てのクラスタの活性度を算出していない場合（ステップS209 NO）、iを1増やし（ステップS210）、ステップS208の処理に戻る。

【0116】

一方、活性度算出部121が全てのクラスタの活性度を算出した場合（ステップS209 YES）、活性度算出部121は、全期間のドキュメントで活性度を算出したか否か判定する（ステップS211）。活性度算出部121は、全期間のドキュメントで活性度を算出していない場合（ステップS211 NO）、kを1増やし（ステップS212）、ステップS208の処理に戻る。

一方、活性度算出部121が全期間のドキュメントで活性度を算出した場合（ステップS211 YES）、相対力指数算出部122は、j番目のクラスタC_jの相対力指数（RSI）を算出する（ステップS213）。

【0117】

次に、相対力指数算出部122は、全てのクラスタの相対力指数（RSI）を算出したか否か判定する（ステップS214）。相対力指数算出部122は、全てのクラスタの相対力指数（RSI）を算出していない場合（ステップS214 NO）、jを1増やし（ステップS215）、ステップS213の処理に戻る。

一方、相対力指数算出部122が、全てのクラスタの相対力指数（RSI）を算出した場合（ステップS214 YES）、抽出装置100は、ステップS216の処理に進む。

【0118】

次に、ステップS216において、クラスタ組抽出部130は、活性化予測条件を満たす下で、新規発見性指数とターゲット関連性指数とに基づいた評価値が最大になるクラスタの組み合わせを抽出する（ステップS216）。以上で、本フローチャートの処理を終了する。

【0119】

以上により、本実施形態の抽出装置100は、抽出された3つのクラスタのうち少なくとも1つが活性化されていること、抽出された2つのクラスタの組み合わせに意外性があること、その2つのクラスタの組み合わせは直接の関連性は薄い、抽出されたもう1つのクラスタ（第3のクラスタ）を経由すると結び付けられるものであること、そのクラスタの組み合わせを提供する対象であるターゲットの特性と抽出されたクラスタのうち少なくとも1つとが関連性があることという条件下で、クラスタの組み合わせを提供することができる。各クラスタは1つの概念と対応しているので、抽出装置100は、所定の期間において、そのターゲットにとって意外性があり、第3のクラスタに対応する第3の概念を介して結び付けられる概念の組み合わせを提供することができる。

【0120】

また、本実施形態の抽出装置100の各処理を実行するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、当該記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することにより、抽出装置100に係る上述した種々の処理を行ってもよい。

【0121】

なお、ここでいう「コンピュータシステム」とは、OSや周辺機器等のハードウェアを含むものであってもよい。また、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。ま

10

20

30

40

50

た、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM、フラッシュメモリ等の書き込み可能な不揮発性メモリ、CD-ROM等の可搬媒体、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。

【0122】

さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムが送信された場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリ（例えばDRAM（Dynamic Random Access Memory））のように、一定時間プログラムを保持しているものも含むものとする。また、上記プログラムは、このプログラムを記憶装置等に格納したコンピュータシステムから、伝送媒体を介して、あるいは、伝送媒体中の伝送波により他のコンピュータシステムに伝送されてもよい。ここで、プログラムを伝送する「伝送媒体」は、インターネット等のネットワーク（通信網）や電話回線等の通信回線（通信線）のように情報を伝送する機能を有する媒体のことをいう。また、上記プログラムは、前述した機能の一部を実現するためのものであっても良い。さらに、前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるもの、いわゆる差分ファイル（差分プログラム）であっても良い。

10

【0123】

以上、本発明の実施形態について図面を参照して詳述したが、具体的な構成はこの実施形態に限られるものではなく、この発明の要旨を逸脱しない範囲の設計等も含まれる。

20

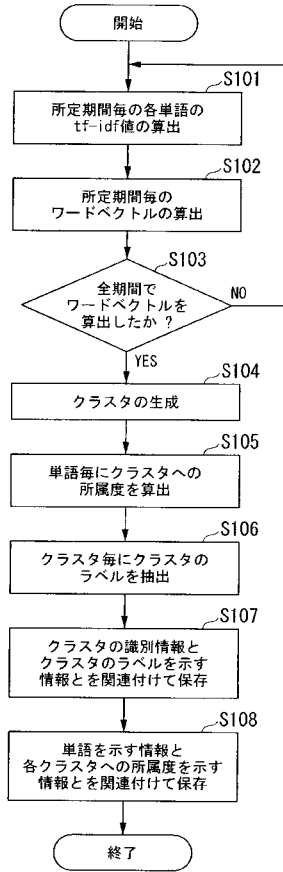
【符号の説明】

【0124】

- 100 抽出装置
- 101 重要度算出部
- 102 重要度記憶部
- 103 クラスタ生成部
- 104 クラスタ記憶部
- 110 発見性指数算出部
- 111 間接関連度算出部
- 112 意外度算出部
- 113 積算部
- 114 ターゲット関連性指数算出部
- 120 活性化予測部
- 121 活性度算出部
- 122 相対力指数算出部（活性度上昇期待値算出部）
- 130 クラスタ組抽出部

30

【 図 5 】



【 図 6 】

