

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5605730号  
(P5605730)

(45) 発行日 平成26年10月15日(2014.10.15)

(24) 登録日 平成26年9月5日(2014.9.5)

(51) Int. Cl. F 1  
**G 0 6 F 17/30 (2006.01)**  
 G 0 6 F 17/30 3 5 0 C  
 G 0 6 F 17/30 2 1 0 D

請求項の数 6 (全 20 頁)

(21) 出願番号	特願2011-32415 (P2011-32415)	(73) 特許権者	801000027
(22) 出願日	平成23年2月17日 (2011.2.17)		学校法人明治大学
(65) 公開番号	特開2012-173800 (P2012-173800A)		東京都千代田区神田駿河台 1-1
(43) 公開日	平成24年9月10日 (2012.9.10)	(74) 代理人	100064908
審査請求日	平成25年11月26日 (2013.11.26)		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

最終頁に続く

(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) 【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乗じることにより、発見性指数を算出する発見性指数算出部と、

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と、

前記算出された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】

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1に記載の抽出装置。

【請求項3】

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度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2に記載の抽出装置。

【請求項4】

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関連性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1から請求項3のいずれか1項に記載の抽出装置。

【請求項5】

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乗じることにより、発見性指数を算出する発見性指数算出手順と、

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手順と、

前記算出された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6】

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乗じることにより、発見性指数を算出する発見性指数算出ステップと、

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

10

20

30

40

50

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ステップと、

前記算出された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乗じることにより、発見性指数を算出する発見性指数算出部と、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

10

20

30

40

50

性指数算出部と、前記算出された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クラスタ組抽出部と、を備えることを特徴とする。

【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

10

【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。

20

【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関連性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。

【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乗じることにより、発見性指数を算出する発見性指数算出手順と、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手順と、前記算出された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クラスタ組抽出手順と、を有することを特徴とする。

30

40

【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乗じることにより、発見性指数を算出する発見性指数算出ステップと、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

50

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ステップと、前記算出された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。

【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。

【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C1、C2をつなぐ役目をする概念Cnが存在する。抽出装置100は、C1、C2、Cn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。

【0019】

抽出装置100は、概念C1と概念C2の直接の関連度は低いが、概念Cnを経由したC1-Cn-C2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(C1)と関連がある所定の雑誌の読者だとすると、概念C1、概念Cn、概念C2の組み合わせとして、ゴルフクラブ(C1)、口紅(C2)、Cn(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Cn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(C1)と口紅(C2)の組み合わせを抽出する価値は高い。

【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C1、C2、Cn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、

10

20

30

40

50

ターゲットと関連している概念（例えば、概念C1）と、接続概念Cnを介して関連している概念C2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C1）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Cn）のための口紅（概念C2）特集を提供することができる。

【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。

10

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（活性化度上昇期待値算出部）122とを備える。

【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。

20

【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 $t_f$ 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。

【0025】

重要度算出部101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101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102のワードベクトルテーブルT1に記憶させる。

30

【0026】

図2は、重要度記憶部102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。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（ $W\_1$ 、 $W\_2$ 、 $W\_3$ 、...、 $W\_30$ ）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。

【0027】

40

図1に戻って、重要度算出部101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103に出力する。

クラスタ生成部103は、重要度算出部101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。

【0028】

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。

【0029】

50

そこで、クラスタ生成部 103 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  $Mem_C(w)$  を有する。ここで、所属度  $Mem_C(w)$  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。

【0030】

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 103 は、一例として、 $k$ -means 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 103 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。

10

【0031】

【数 1】

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

【0032】

但し、以下の条件式 (2) を満たすものとする。

【0033】

【数 2】

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

20

【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。

【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、 $k$  は事前に与えられるものとする。

30

【0036】

【数 3】

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

40

【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、 $x_i = (x_{i1},$

50

$x_{i2}$  )、 $K$ 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。

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。

10

【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、 $xy$ 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、 $x$ 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 $x$ 印で示された単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。

20

【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。

【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。

30

【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。

【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。

40

【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

50



当該単語の位置を示す情報とを関連付けて記憶させる。

【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。

【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ることの意外度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。

【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CN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CNをC\_\_1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CNは、C\_\_1からC\_\_Nまでのクラスタを取りうる。

【0053】

$$MIR(C\_i, C\_j) = \max_{C\_N} \{ A(C\_i, C\_N) \times A(C\_N, C\_j) \} \quad (5)$$

【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CN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CN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】

$$MIR(C\_i, C\_j) = \max_{C\_N} \{ A(C\_i, C\_N) + A(C\_N, C\_j) \} \quad (6)$$

【0056】

10

20

30

40

50

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】

$$A(C\_i, C\_j) = x \cdot y / (|x| \times |y|) \quad (7)$$

10

【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。

20

【0061】

【数5】

$$A(C\_i, C\_j) = \frac{|C\_i \cap C\_j|}{|C\_i \cup C\_j|} \quad (8)$$

【0062】

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。

30

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】

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

40

【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】

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

【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。

10

【0067】

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。

20

【0068】

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。

30

【0071】

なお、意外度算出部 112 は、ジャックカード係数の逆数 (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】

$$U(C\_i, C\_j) = \frac{|C\_i \cup C\_j|}{|C\_i \cap C\_j|} \quad (12)$$

40

【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】

50

$$U(C\_i, C\_j) = 1 / MI(C\_i, C\_j) \quad (13)$$

【0075】

意外度算出部112は、クラスタC<sub>i</sub>とクラスタC<sub>j</sub>の全ての組み合わせで、意外度U(C<sub>i</sub>, C<sub>j</sub>)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<sub>i</sub>, C<sub>j</sub>)を示す情報と、その各意外度U(C<sub>i</sub>, C<sub>j</sub>)が算出された際に用いられたクラスタC<sub>i</sub>の識別情報とクラスタC<sub>j</sub>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MIR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<sub>i</sub>, C<sub>j</sub>)以外の第3のクラスタCNを経由した該2つのクラスタ(C<sub>i</sub>, C<sub>j</sub>)間の関連度と、該2つのクラスタ(C<sub>i</sub>, C<sub>j</sub>)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。

【0077】

$$S(C\_i, C\_j) = MIR(C\_i, C\_j) \times U(C\_i, C\_j) \quad (14)$$

【0078】

発見性指標Sは、クラスタC<sub>i</sub>とクラスタC<sub>j</sub>との間でクラスタCNを経由した関連性が必要なこと、また同時にクラスタC<sub>i</sub>とクラスタC<sub>j</sub>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<sub>i</sub>, C<sub>j</sub>)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(CN<sub>(i, j)</sub>)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<sub>i</sub>とクラスタC<sub>j</sub>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<sub>i</sub>を示す情報とクラスタC<sub>j</sub>を示す情報と接続クラスタCN<sub>(i, j)</sub>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<sub>i</sub>を示す情報とクラスタC<sub>j</sub>を示す情報と接続クラスタCN<sub>(i, j)</sub>を示す情報とを受け付ける。

【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<sub>i</sub>, C<sub>j</sub>, CN<sub>(i, j)</sub>)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<sub>i</sub>, C<sub>j</sub>, CN<sub>(i, j)</sub>)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。

【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

10

20

30

40

50

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。

具体的には、例えば、活性度算出部 1 2 1 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 1 2 1 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。

また、関数  $\text{sim}$  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  $\text{tf-idf}$  値が示されている。

【0090】

なお、活性度算出部 1 2 1 は、関数  $\text{sim}$  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 1 2 1 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】

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

【0092】

ここで、 $\text{mem}C\_i(y_q)$  は単語  $y_q$  のクラスタ C\_\_i への所属度である。 $\text{MI}(x_p, y_q)$  は、単語  $x_p$  と単語  $y_q$  との相互情報量である。 $\text{tfidf}(x)$  はワードベクトル中の単語  $x_p$  の  $\text{tf-idf}$  値である。

【0093】

なお、活性度算出部 1 2 1 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、 $\text{tf-idf}$  値の高い一定数の上位単語または  $\text{tf-idf}$  値が所定の値を超えた単語の  $\text{tf-idf}$  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 1 2 1 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。

【0094】

活性度算出部 1 2 1 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 1 2 2 に出力する。

相対力指数算出部 1 2 2 は、活性度算出部 1 2 1 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变化に注目し、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

10

20

30

40

50

れる度合い（活性度上昇期待値）を算出する。

【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<sub>i</sub>)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】

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

【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<sub>u</sub>、下限R<sub>L</sub>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。

【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、発見性指標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<sub>i</sub>、C<sub>j</sub>、CN(i, j))を抽出する。

【0101】

(1)新規発見性指標Sの条件として、クラスタの組C<sub>i</sub>、C<sub>j</sub>、CN(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<sub>i</sub>、クラスタC<sub>j</sub>、クラスタCN(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、

(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<sub>i</sub>、C<sub>j</sub>、CN(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<sub>i</sub>、C<sub>j</sub>、CN(i, j))を、下記の式(20)から算出する。

【0103】

10

20

30

40

50

$$\arg \max \{ a S ( C\_i , C\_j , C N ( i , j ) ) + b N ( C\_i , C\_j , C N ( i , j ) , T ) \} \quad ( 2 0 )$$

【 0 1 0 4 】

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、 $C\_i$ 、 $C\_j$ 、 $C N ( i , j )$ のうちいずれかの相対力指数( $R S I$ )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( $i i$ )を満足していることとする。

【 0 1 0 5 】

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。

10

【 0 1 0 6 】

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 $C\_i$ を示す情報とクラスタ $C\_j$ 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 $T 2$ から読み出し、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。

20

【 0 1 0 7 】

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 $t f - i d f$ 値の算出する(ステップ $S 1 0 1$ )。次に、重要度算出部101は、所定期間毎に、各単語の $t f - i d f$ 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 $S 1 0 2$ )。

【 0 1 0 8 】

30

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 $S 1 0 3$ )。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 $S 1 0 3$  NO)、ステップ $S 1 0 1$ 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 $S 1 0 3$  YES)、クラスタ生成部103は、クラスタを生成する(ステップ $S 1 0 4$ )。

【 0 1 0 9 】

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 $S 1 0 5$ )。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 $S 1 0 6$ )。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 $S 1 0 7$ )。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 $S 1 0 8$ )。以上で、本フローチャートの処理を終了する。

40

【 0 1 1 0 】

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。

【 0 1 1 1 】

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

50

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 M I R を算出する (ステップ S 201)。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 M I R を算出したか否か判定する (ステップ S 202)。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 M I R を算出していない場合 (ステップ S 202 NO)、ステップ S 201 の処理に戻る。

【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 M I R を算出した場合 (ステップ S 202 YES)、意外度算出部 112 は、意外度 U を算出する (ステップ S 203)。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 S 204)。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 S 204 NO)、ステップ S 203 の処理に戻る。

10

【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 S 204 YES)、積算部 113 は、発見性指標を算出する (ステップ S 205)。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 S 206)。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 S 206 NO)、ステップ S 201 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 S 206 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する (ステップ S 207)。

20

【0115】

ステップ S 201 ~ ステップ S 207 までの処理に並行して、抽出装置 100 は、ステップ S 208 ~ ステップ S 215 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<sub>i</sub> の活性度を算出する (ステップ S 208)。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 S 209)。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 S 209 NO)、i を 1 増やし (ステップ S 210)、ステップ S 208 の処理に戻る。

30

【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 S 209 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 S 211)。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 S 211 NO)、k を 1 増やし (ステップ S 212)、ステップ S 208 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 S 211 YES)、相対力指数算出部 122 は、j 番目のクラスタ C<sub>j</sub> の相対力指数 (RSI) を算出する (ステップ S 213)。

【0117】

40

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 S 214)。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 S 214 NO)、j を 1 増やし (ステップ S 215)、ステップ S 213 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 S 214 YES)、抽出装置 100 は、ステップ S 216 の処理に進む。

【0118】

次に、ステップ S 216 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

50



タの組み合わせを抽出する（ステップS 2 1 6）。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。

10

【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。

20

【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よい。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）であってもよい。

30

【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。

40

【符号の説明】

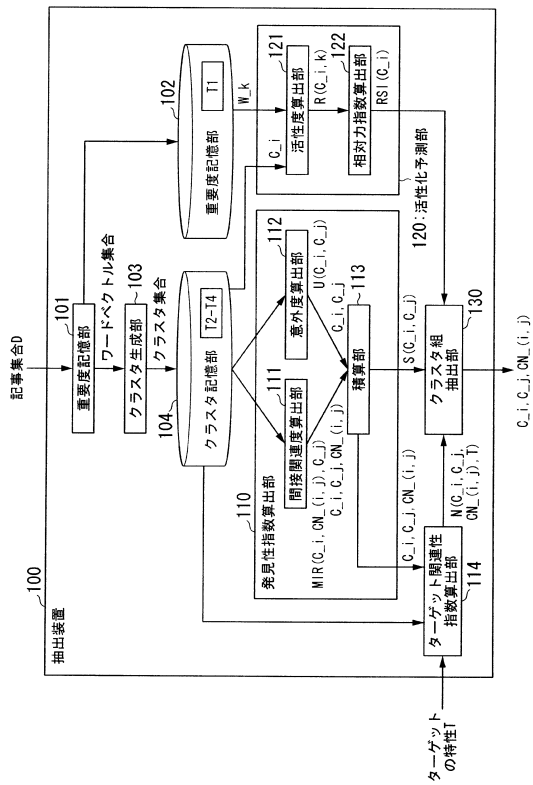
【0124】

- 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

50

- 1 1 2 意外度算出部
- 1 1 3 積算部
- 1 1 4 ターゲット関連性指数算出部
- 1 2 0 活性化予測部
- 1 2 1 活性度算出部
- 1 2 2 相対力指数算出部 (活性度上昇期待値算出部)
- 1 3 0 クラスタ組抽出部

【図1】

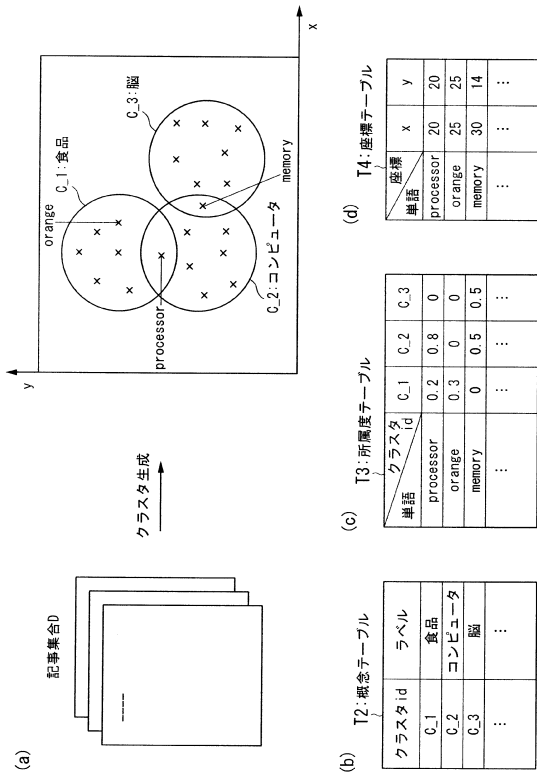


【図2】

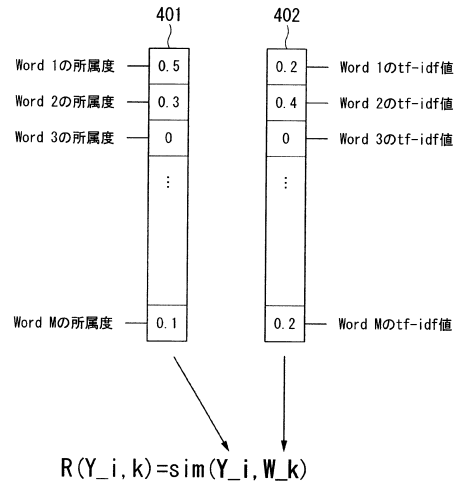
T1: ワードベクトルテーブル

	W_1	W_2	W_3	...	W_30
word1	0.2	0.1	0.4	...	0.2
word2	0.4	0.3	0.1	...	0.2
word3	0.1	0.4	0.5	...	0.1
word4	0.3	0.5	0.3	...	0.3
word5	0.2	0.6	0.7	...	0.2
word6	0.6	0.7	0.4	...	0.6
...	...	...	...	...	...
...	...	...	...	...	...
...	...	...	...	...	...
時間(日)	1	2	3	...	30

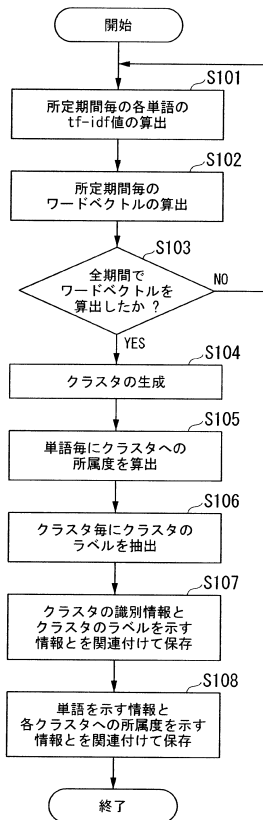
【図3】



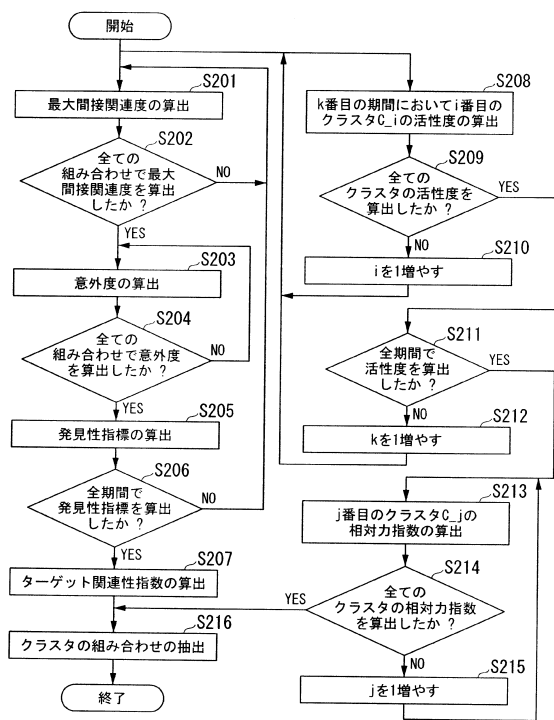
【図4】



【図5】



【図6】



---

フロントページの続き

審査官 加舎 理紅子

- (56)参考文献 国際公開第2010/026900(WO, A1)  
特開平07-143463(JP, A)  
特開2010-061600(JP, A)  
国際公開第2008/004663(WO, A1)  
入江毅 他, 知的判断メカニズムのための概念間の類似度評価モデル, 情報処理学会研究報告書, 1999年 1月12日, Vol. 99, No. 1, p. 93 - 100  
中野俊亮 他, ユーザ対話による意外性を持つキャッチフレーズ作成支援, 第70回(平成20年)全国大会講演論文集(2), 2008年 3月31日, p. 2-201~2-202  
砂山渡 他, 検索語のクラスタリングによるWeb情報の傾向の獲得, 第43回人工知能基礎研究会資料(SIG-FAI-A003), 2000年11月 9日, p. 7 - 11

- (58)調査した分野(Int.Cl., DB名)  
G06F 17/30  
Cini