

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5023176号
(P5023176)

(45) 発行日 平成24年9月12日(2012.9.12)

(24) 登録日 平成24年6月22日(2012.6.22)

(51) Int.Cl.

F I

G06F 17/30 (2006.01)
G06Q 10/00 (2012.01)

G06F 17/30 210D
G06F 17/30 170A
G06F 17/30 360Z
G06F 17/30 220Z
G06F 17/60 150

請求項の数 5 (全 34 頁)

(21) 出願番号 特願2010-64821 (P2010-64821)
(22) 出願日 平成22年3月19日(2010.3.19)
(65) 公開番号 特開2011-198111 (P2011-198111A)
(43) 公開日 平成23年10月6日(2011.10.6)
審査請求日 平成22年3月19日(2010.3.19)

(73) 特許権者 000003078
株式会社東芝
東京都港区芝浦一丁目1番1号
(73) 特許権者 301063496
東芝ソリューション株式会社
東京都港区芝浦一丁目1番1号
(74) 代理人 100108855
弁理士 蔵田 昌俊
(74) 代理人 100091351
弁理士 河野 哲
(74) 代理人 100088683
弁理士 中村 誠
(74) 代理人 100109830
弁理士 福原 淑弘

最終頁に続く

(54) 【発明の名称】 特徴語抽出装置及びプログラム

(57) 【特許請求の範囲】

【請求項1】

文書ID及び内容テキスト情報を有する複数の文書を記憶する文書記憶手段と、
カテゴリID毎に1つ以上の文書IDを関連付けて記憶するカテゴリ記憶手段と、
前記文書記憶手段内の文書毎に、当該文書の文書IDと、当該文書の内容テキスト情報から抽出された文書特徴語とを関連付けて記憶する文書特徴語記憶手段と、
前記カテゴリ記憶手段内で関連したカテゴリID及び1つ以上の文書IDと、当該カテゴリIDに関連したカテゴリ特徴語とを関連付けて記憶するカテゴリ特徴語記憶手段と、
前記文書記憶手段内の文書毎に内容テキスト情報を形態素解析し、当該形態素解析の結果から文書特徴語を抽出し、当該抽出した文書特徴語と、当該文書特徴語に対応する文書の文書IDとを関連付けて前記文書特徴語記憶手段に書き込む文書特徴語抽出手段と、
前記複数の文書により構成される文書集合を入力とし、当該文書集合に含まれる文書IDに関連した文書特徴語が、当該文書IDの文書中で出現する文書数を算出する出現文書数算出手段と、

前記出現文書数算出手段により算出された文書数に基づいて、全文書中におけるカテゴリIDに関連した文書に対する当該文書特徴語の特徴度を算出する特徴度算出手段と、

この文書特徴語に当該特徴度を付加したカテゴリ特徴語を作成し、当該作成したカテゴリ特徴語と当該カテゴリ特徴語に関連したカテゴリID及び1つ以上の文書IDとを関連付けて前記カテゴリ特徴語記憶手段に書き込むカテゴリ特徴語作成手段と、

前記カテゴリ記憶手段内のカテゴリID毎に、当該カテゴリIDに関連した文書IDの

個数を含むカテゴリ個数データを提示するカテゴリ個数提示手段と、

前記カテゴリ個数提示手段によるカテゴリ個数データの提示中、いずれかのカテゴリ個数データの選択を受け付けるカテゴリ個数データ選択受付手段と、

前記カテゴリ個数データ選択受付手段による選択を受け付けたカテゴリ個数データのカテゴリIDに関連したカテゴリ特徴語のうち、特徴度が上位のカテゴリ特徴語における文書特徴語をカテゴリ特徴語として提示するカテゴリ特徴語提示手段と、

前記カテゴリ個数提示手段によるカテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象からなる比較対象集合の選択を受け付ける比較対象集合選択受付手段と、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、前記出現文書数算出手段により算出された文書数に基づいて、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として送出するカテゴリ相違特徴語送出手段と、

前記カテゴリ相違特徴語送出手段により送出されたカテゴリ相違特徴語を提示するカテゴリ相違特徴語提示手段と、

を備えたことを特徴とする特徴語抽出装置。

【請求項2】

請求項1に記載の特徴語抽出装置において、

前記出現文書数算出手段は、

前記文書記憶手段内の文書特徴語毎に、前記文書記憶手段の全ての文書中で当該文書特徴語が出現する文書数を算出する全文書中出現文書数算出手段と、

前記カテゴリ記憶手段内のカテゴリID毎に、当該カテゴリIDに関連付けられた文書IDに関連した文書特徴語が当該文書IDの文書中で出現する文書数を算出するカテゴリ文書中出現文書数算出手段と、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、当該各文書IDに関連した全ての文書中で当該文書特徴語が出現する文書数を算出する第1文書数算出手段と、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の比較対象毎に、当該比較対象内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語が当該各文書IDの文書中で出現する文書数を算出する第2文書数算出手段と、

を備えたことを特徴とする特徴語抽出装置。

【請求項3】

請求項1または2に記載の特徴語抽出装置において、

前記選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、前記第1文書数算出手段により算出された文書数と、前記第2文書数算出手段により算出された文書数とに基づいて、前記比較対象集合における当該文書特徴語の共通性を表す度合いの共通特徴度を算出し、共通特徴度が上位の文書特徴語をカテゴリ共通特徴語として送出するカテゴリ共通特徴語送出手段と、

前記カテゴリ共通特徴語送出手段により送出されたカテゴリ共通特徴語を提示するカテゴリ共通特徴語提示手段と、

を更に備えたことを特徴とする特徴語抽出装置。

【請求項4】

請求項1または2に記載の特徴語抽出装置において、

前記カテゴリ個数提示手段によるカテゴリ相違特徴語の提示中、複数のカテゴリ相違特徴語からなる着目語集合の選択を受け付ける着目語集合選択受付手段と、

前記着目語集合選択受付手段による選択を受け付けた着目語集合と、前記カテゴリ特徴語記憶手段内のカテゴリ特徴語との関連度を算出し、関連度の高いカテゴリ特徴語に関連付けられたカテゴリIDに関連したカテゴリ個数データを強調表示する関連カテゴリ提示手段と、

10

20

30

40

50

を更に備えたことを特徴とする特徴語抽出装置。

【請求項 5】

文書記憶手段、カテゴリ記憶手段、文書特徴語記憶手段及びカテゴリ特徴語記憶手段を備えた特徴語抽出装置のプログラムであって、

前記特徴語抽出装置を、

文書 ID 及び内容テキスト情報を有する複数の文書を前記文書記憶手段に書き込む文書書込手段、

カテゴリ ID 毎に 1 つ以上の文書 ID を関連付けて前記カテゴリ記憶手段に書き込むカテゴリ書込手段、

前記文書記憶手段内の文書毎に内容テキスト情報を形態素解析し、当該形態素解析の結果から文書特徴語を抽出し、当該抽出した文書特徴語と、当該文書特徴語に対応する文書の文書 ID とを関連付けて前記文書特徴語記憶手段に書き込む文書特徴語抽出手段、

前記文書記憶手段内の文書特徴語毎に、前記文書記憶手段の全ての文書中で当該文書特徴語が出現する文書数を算出する全文書中出現文書数算出手段、

前記カテゴリ記憶手段内のカテゴリ ID 毎に、当該カテゴリ ID に関連付けられた文書 ID に関連した文書特徴語が当該文書 ID の文書中で出現する文書数を算出するカテゴリ文書中出現文書数算出手段、

前記全文書中出現文書数算出手段により算出された文書数と、前記カテゴリ文書中出現文書数算出手段により算出された文書数とに基づいて、全文書中におけるカテゴリ ID に関連した文書に対する当該文書特徴語の特徴度を算出する特徴度算出手段、

この文書特徴語に当該特徴度を付加したカテゴリ特徴語を作成し、当該作成したカテゴリ特徴語と当該カテゴリ特徴語に関連したカテゴリ ID 及び 1 つ以上の文書 ID とを関連付けて前記カテゴリ特徴語記憶手段に書き込むカテゴリ特徴語作成手段、

前記カテゴリ記憶手段内のカテゴリ ID 毎に、当該カテゴリ ID に関連付けられた文書 ID の個数を含むカテゴリ個数データを提示するカテゴリ個数提示手段、

前記カテゴリ個数データの提示中、いずれかのカテゴリ個数データの選択を受け付けるカテゴリ個数データ選択受付手段、

前記カテゴリ個数データ選択受付手段による選択を受け付けたカテゴリ個数データのカテゴリ ID に関連したカテゴリ特徴語のうち、特徴度が上位のカテゴリ特徴語における文書特徴語をカテゴリ特徴語として提示するカテゴリ特徴語提示手段、

前記カテゴリ個数提示手段によるカテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリ ID の和集合である複数の比較対象からなる比較対象集合の選択を受け付ける比較対象集合選択受付手段、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリ ID に関連付けられた各文書 ID に関連した文書特徴語毎に、当該各文書 ID に関連した全ての文書中で当該文書特徴語が出現する文書数を算出する第 1 文書数算出手段と、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の比較対象毎に、当該比較対象内の各カテゴリ ID に関連付けられた各文書 ID に関連した文書特徴語が当該各文書 ID の文書中で出現する文書数を算出する第 2 文書数算出手段と、

前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリ ID に関連付けられた各文書 ID に関連した文書特徴語毎に、前記第 1 文書数算出手段により算出された文書数と、前記第 2 文書数算出手段により算出された文書数とに基づいて、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として送出するカテゴリ相違特徴語送出手段と、

前記カテゴリ相違特徴語送出手段により送出されたカテゴリ相違特徴語を提示するカテゴリ相違特徴語提示手段、

として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

10

20

30

40

50

【0001】

本発明は、電子化された大量の文書に対し、各文書の内容を特徴づける単語である特徴語を抽出するための特徴語抽出装置及びプログラムに関する。

【背景技術】

【0002】

特許調査やアンケート分析などにおいては、特徴語を抽出し、複数の文書集合の内容や傾向を比較したいというニーズがある。例えば、特許調査においては、自社と競合他社との各年代の特許出願の傾向を比較したいニーズがある。この種の調査の質は、選定する比較範囲や特徴語に影響される。

【0003】

しかしながら、適切な比較範囲や特徴語は、調査の目的や文書集合の内容によって異なる。このため、比較範囲や特徴語の選定には、文書集合の内容に関連した知識や、目的を踏まえた調査自体に関するスキルが必要である。

【0004】

これに対し、特許文献1では適切な分析軸を提示する技術を提案している。特許文献1に記載の技術は、データに含まれる属性毎に単語を抽出し、単語の属性値毎の出現頻度を集計手段で集計し、集計した出現頻度を視認し易いようにユーザに表示するとともに、単語の出現傾向から分析に適した属性を抽出しユーザに提示する。これにより、ユーザに適切な分類軸の選択を支援する。

【0005】

また、特許文献2に記載の技術では、分析に使用する適切な特徴語の選定のために、比較する属性（例えば、作成日）について、属性値毎（例えば、月毎）に抽出される特徴語の共起関係に基づき、各属性において相違点を提示する。これにより、文書集合の内容をより好適に分析可能としている。なお、特許文献2に記載の技術で用いる「共起」については、例えば、特許文献3にまとめられている。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2006-171931号公報

【特許文献2】特開2002-245070号公報

【非特許文献】

【0007】

【非特許文献1】内山将夫，中條清美，山本英子，井佐原均．「英語教育のための分野特徴単語の選定尺度の比較」，自然言語処理，11（3），165-197，2004．

【非特許文献2】岸田和明．「検索実験における評価指標としての平均精度の性質」，情報処理学会論文誌：データベース，第43巻，第SIG2（TOD13）号（2002）

【非特許文献3】相澤彰子．「共起に基づく類似性尺度」，オペレーションズ・リサーチ，2007年11月号，pp.706(20) - 712(26)．

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかしながら、以上のような特許文献1，2に記載の技術では、通常は特に問題ないが、本発明者の検討によれば、以下に述べる点で改良の余地がある。

【0009】

例えば、特許文献1に記載の技術では、ユーザに提示される分析軸が事前に文書データの属性として定義されている必要がある。そのため、提示される分析軸が事前に定義された属性に限られるため、意図する分析を行えない点で改良の余地がある。

【0010】

特許文献2に記載の技術は、各属性値に対する相違点を表すことにより、文書集合の内

10

20

30

40

50

容の明確化を図っている。このため、分析の対象が属性値に縛られ、ユーザが任意の範囲で文書集合を比較できない点で改良の余地がある。

【0011】

また、特許文献2に記載の技術は、膨大な文書集合を比較分析する場合、ユーザが文書集合の中で何に注目すべきかを把握できない場合がある。例えば、注目する「画像認識」の技術を先行調査するために、数千・数万の特許文献を出願人と出願年月(1990年～2008年の各月)でクロス分析する場合を考える。出願人として数10～100社程度の各企業を各行に配置し、出願年月として100個程度の各月を各列に配置した場合、クロス分析のマトリックス全体として1万前後のセルが構成される。

【0012】

これらのセルは、「画像認識」に関連する多数の特許文献が含まれるセルや、「画像認識」に無関係の多数の特許文献が含まれるセルなどがあり、注目する「画像認識」との関連度にはムラがある。企業毎や出願年毎でも同様のことが言える。

【0013】

これに対し、ユーザは、注目する技術に関連が強い企業や出願年に関するセルの文書集合に比較範囲を絞り込むことで、より精緻な調査を行いたいというニーズがある。

【0014】

しかしながら、特許文献1, 2に記載の技術では、注目すべき比較範囲の絞り込みを支援できず、比較範囲を柔軟に変更することもできない。また、注目する技術に関連の強いセルの特徴語を参照すればユーザは意識しなかった関連技術を発見できるが、特許文献1, 2に記載の技術では、文書集合の内容理解までに留まり、新たに注目すべき特徴語の参照を支援することはできない。

【0015】

本発明は上記実情を考慮してなされたもので、事前に定義された属性に限らずに分析軸の候補として特徴語を提示できると共に、注目すべき比較範囲の絞り込みや、注目すべき特徴語の参照を支援し得る特徴語抽出装置及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】

【0016】

本発明の一つの局面は、特徴語抽出装置であって、文書ID及び内容テキスト情報を有する複数の文書を記憶する文書記憶手段と、カテゴリID毎に1つ以上の文書IDに関連付けて記憶するカテゴリ記憶手段と、前記文書記憶手段内の文書毎に、当該文書の文書IDと、当該文書の内容テキスト情報から抽出された文書特徴語とを関連付けて記憶する文書特徴語記憶手段と、前記カテゴリ記憶手段内で関連したカテゴリID及び1つ以上の文書IDと、当該カテゴリIDに関連したカテゴリ特徴語とを関連付けて記憶するカテゴリ特徴語記憶手段と、前記文書記憶手段内の文書毎に内容テキスト情報を形態素解析し、当該形態素解析の結果から文書特徴語を抽出し、当該抽出した文書特徴語と、当該文書特徴語に対応する文書の文書IDとを関連付けて前記文書特徴語記憶手段に書き込む文書特徴語抽出手段と、前記文書記憶手段内の文書特徴語毎に、前記文書記憶手段の全ての文書中で当該文書特徴語が出現する文書数を算出する全文書中出現文書数算出手段と、前記カテゴリ記憶手段内のカテゴリID毎に、当該カテゴリIDに関連付けられた文書IDに関連した文書特徴語が当該文書IDの文書中で出現する文書数を算出するカテゴリ文書中出現文書数算出手段と、前記全文書中出現文書数算出手段により算出された文書数と、前記カテゴリ文書中出現文書数算出手段により算出された文書数とに基づいて、全文書中におけるカテゴリIDに関連した文書に対する当該文書特徴語の特徴度を算出する特徴度算出手段と、この文書特徴語に当該特徴度を付加したカテゴリ特徴語を作成し、当該作成したカテゴリ特徴語と当該カテゴリ特徴語に関連したカテゴリID及び1つ以上の文書IDとを関連付けて前記カテゴリ特徴語記憶手段に書き込むカテゴリ特徴語作成手段と、前記カテゴリ記憶手段内のカテゴリID毎に、当該カテゴリIDに関連付けられた文書IDの個数を含むカテゴリ個数データを提示するカテゴリ個数提示手段と、前記カテゴリ個数データの提示中、いずれかのカテゴリ個数データの選択を受け付けるカテゴリ個数データ選択受

10

20

30

40

50

付手段と、前記カテゴリ個数データ選択受付手段による選択を受け付けたカテゴリ個数データのカテゴリIDに関連したカテゴリ特徴語のうち、特徴度が上位のカテゴリ特徴語における文書特徴語をカテゴリ特徴語として提示するカテゴリ特徴語提示手段と、前記カテゴリ特徴語提示手段によるカテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象からなる比較対象集合の選択を受け付ける比較対象集合選択受付手段と、前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、当該各文書IDに関連した全ての文書中で当該文書特徴語が出現する文書数を算出する第1文書数算出手段と、前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の比較対象毎に、当該比較対象内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語が当該各文書IDの文書中で出現する文書数を算出する第2文書数算出手段と、前記比較対象集合選択受付手段による選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、前記第1文書数算出手段により算出された文書数と、前記第2文書数算出手段により算出された文書数とに基づいて、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として送出するカテゴリ相違特徴語送出手段と、前記カテゴリ相違特徴語送出手段により送出されたカテゴリ相違特徴語を提示するカテゴリ相違特徴語提示手段と、を備えた特徴語抽出装置である。

10

【0017】

なお、本発明の一つの局面は、装置として表現したが、これに限らず、方法、プログラム又はプログラムを記憶したコンピュータ読取り可能な記憶媒体として表現してもよい。

20

【0018】

(作用)

このような本発明の一つの局面においては、カテゴリIDに関連した文書IDの個数を含むカテゴリ個数データの提示中に、選択を受け付けたカテゴリ個数データのカテゴリIDに関連したカテゴリ特徴語のうち、特徴度が上位のカテゴリ特徴語における文書特徴語をカテゴリ特徴語として提示する。

【0019】

また、本発明の一つの局面においては、カテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象からなる比較対象集合の選択を受け付けると、選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として提示する。

30

【0020】

このように、選択したカテゴリ個数データに応じて、カテゴリ特徴語やカテゴリ相違特徴語といった特徴語を提示する構成により、事前に定義された属性に限らずに分析軸の候補として特徴語を提示できると共に、着目すべき比較範囲の絞り込みや、着目すべき特徴語の参照を支援することができる。

【発明の効果】

40

【0021】

以上説明したように本発明によれば、事前に定義された属性に限らずに分析軸の候補として特徴語を提示できると共に、着目すべき比較範囲の絞り込みや、着目すべき特徴語の参照を支援できる。

【図面の簡単な説明】

【0022】

【図1】本発明の一実施形態に係る特徴語抽出装置の構成を示すブロック図である。

【図2】同実施形態における文書記憶部を説明するための模式図である。

【図3】同実施形態におけるカテゴリ記憶部を説明するための模式図である。

【図4】同実施形態における特徴語記憶部を説明するための模式図である。

50

【図5】同実施形態における特徴語抽出部の動作を説明するためのフローチャートである。

【図6】同実施形態におけるカテゴリ特徴語抽出部の動作を説明するためのフローチャートである。

【図7】同実施形態におけるカテゴリ共通特徴語抽出部の動作を説明するためのフローチャートである。

【図8】同実施形態におけるカテゴリ個数データをセルに提示した画面例を示す模式図である。

【図9】同実施形態における比較対象集合を選択した画面例を示す模式図である。

【図10】同実施形態における共通特徴語を表示した画面例を示す模式図である。

10

【図11】同実施形態におけるカテゴリ相違特徴語抽出部の動作を説明するためのフローチャートである。

【図12】同実施形態における共通特徴語と相違特徴語を表示した画面例を示す模式図である。

【図13】同実施形態における関連カテゴリ提示部の動作を説明するためのフローチャートである。

【図14】同実施形態におけるユーザ操作・提示部の動作を説明するためのフローチャートである。

【図15】同実施形態におけるカテゴリ特徴語を表示した画面例を示す模式図である。

【図16】同実施形態における関連カテゴリのセルを強調表示した画面例を示す模式図である。

20

【図17】同実施形態におけるカテゴリ特徴語から着目語を選択したときの画面例を示す模式図である。

【図18】同実施形態における比較対象の絞り込みと特徴語の表示例を示す模式図である。

【図19】同実施形態における着目語の変更と関連カテゴリの表示例を示す模式図である。

【図20】同実施形態におけるクロス分析の画面例を示す模式図である。

【図21】同実施形態における他のクロス分析の画面例を示す模式図である。

【図22】同実施形態における更に他のクロス分析の画面例を示す模式図である。

30

【図23】同実施形態におけるグラフ表示の画面例を示す模式図である。

【発明を実施するための形態】

【0023】

以下、本発明の一実施形態について図面を用いて説明する。なお、以下の装置は、装置毎に、ハードウェア構成、又はハードウェア資源とソフトウェアとの組合せ構成のいずれでも実施可能となっている。組合せ構成のソフトウェアとしては、予めネットワーク又は記憶媒体から対応する装置のコンピュータにインストールされ、対応する装置の機能を実現させるためのプログラムが用いられる。また、以下の説明で用いられる用語と記号の定義は、次の表1及び表2に示す通りである。

【0024】

40

【表 1】

[表 1]

用語と記号	定義
全文書集合 docAll	分析対象の全文書を含む文書集合
カテゴリ cat	複数の所属文書から構成される文書集合。 カテゴリ特徴語の抽出単位。
比較対象 cmp	複数のカテゴリ (cat1, cat2, ...) によって指定される文書集合。 指定されたカテゴリの和集合。 $cmp = cat1 \cup cat2 \cup \dots$ (cat i はカテゴリ)
比較範囲 tgtDocs	複数の比較対象 (比較対象集合) によって指定される文書集合。 複数の比較対象の和集合。 $tgtDocs = cmp1 \cup cmp2 \cup \dots \cup cmpM$ (cmp i は比較対象)
差集合 cmpDocs_i	比較範囲 tgtDocs から比較対象 cmp_i を除いた差集合 $cmpDocs_i = tgtDocs - cmp_i$
比較対象集合 tgtSet	複数の比較対象を要素とする集合。 $tgtSet = \{cmp1, cmp2, \dots\}$ (cmp i は比較対象)
df(t, docSet)	文書集合 docSet 中の単語 t の出現頻度。 docSet としては、以下の文書集合を引数とする。 - 全文書集合 docAll - カテゴリ cat - 比較対象 cmp - 比較範囲 tgtDocs
docSet	文書集合 docSet 中に含まれる文書数 docSet としては、以下の文書集合を引数とする。 - 全文書集合 docAll - カテゴリ cat - 比較対象 cmp - 比較範囲 tgtDocs

10

20

30

【 0 0 2 5 】

【表 2】

[表 2]

用語と記号	定義
特徴度 $score(t, cat)$	全文書集合 $docAll$ における、カテゴリ cat に対する単語 t の単語特徴度
共通特徴度 $com(t, tgtSet)$	比較対象集合 $tgtSet$ に対する、単語 t の共通特徴度。 比較対象 cmp_i の評価値 $eval(t, cmp_i)$ の総和 $com(t, tgtSet) = \sum_i eval(t, cmp_i)$
相違特徴度 $diff(t, cmp)$	比較対象 cmp に対する、単語 t の相違特徴度
関連度 $rel(cat, tgtTerms)$	カテゴリ cat と、着目語集合 $tgtTerms$ との関連度。 平均精度、 又は着目語集合 $tgtTerms$ 内の単語 t の特徴度 $score(t, cat)$ の総和 $rel(cat, tgtTerms) = \sum_t score(t, cat)$
共通特徴語 $comTerms$	比較対象集合 cmp に対して共通点を表す単語の集合。 $comTerms = \{term1, term2, \dots\}$ ($term\ i$ は単語)
相違特徴語 $diffTerms(cmp)$	比較対象集合 cmp に対する相違点を表す単語の集合。 $diffTerms(cmp) = \{term1, term2, \dots\}$ ($term\ i$ は単語) 引数 cmp により、指定される比較対象の相違特徴語を返す。 また、 $diffTerms$ と表記した場合は、複数のカテゴリ cmp について、それぞれの相違特徴語 $diffTerms(cmp)$ を、要素とする集合とする。 $diffTerms = \{diffTerms(cmp1), diffTerms(cmp2), \dots\}$ ($cmp\ i$ は比較対象)
着目語集合 $tgtTerms$	複数の単語から構成される集合。 $tgtTerms = \{term1, term2, \dots\}$ ($term\ i$ は単語)

10

20

30

【0026】

なお、以下で説明する実施形態においては、複数の文書からなる文書集合を入力とし、その文書集合に含まれる文書 ID に関連した文書特徴語が、当該文書 ID の文書中で出現する文書数（文書の中で出現する特徴後の数を含む概念）を算出する処理（出現文書数算出機能）に、特に特徴があるといえる。

【0027】

従って、全文書中出现文書数算出機能と、カテゴリ文書中出现文書数算出機能と、第 1 文書数算出機能と、第 2 文書数算出機能とを例にして、出現文書数算出機能を説明していく。なぜならば、上記の各表で定義されている全文書集合、カテゴリ、比較対象集合（比較対象の集合）、比較対象（カテゴリの集合）などについて、各入力の種類は異なっても、これらはいずれも文書集合（文書 ID の集合）と換言することができるからである。

40

【0028】

図 1 は本発明の一実施形態に係る特徴語抽出装置の構成を示すブロック図であり、図 2 乃至図 4 は同装置内の各記憶部 10、20、30 を説明するための模式図である。この特徴語抽出装置は、文書記憶部 10、カテゴリ記憶部 20、特徴語記憶部 30、特徴語抽出

50

部 4 0 及びユーザ操作・提示部 5 0 を備えている。

【 0 0 2 9 】

文書記憶部 1 0 は、各部 4 0 , 5 0 から読出 / 書込可能な記憶装置であり、図 2 に示すように、文書データ 1 0 d を記憶している。文書データ 1 0 d は、各文書を識別する文書 ID としての文書 1 1 d と、内容テキスト情報 (文字列情報) 1 2 d としての文書名 1 2 d 及び / 又は本文 1 5 d とを有する複数の文書を電子化したデータであり、ここでは特許文献の例が図示されている。なお、文書データ 1 0 d は、文書 ID 1 1 d と内容テキスト情報に加え、出願日 1 3 d や出願人 1 4 d などの属性値を有していてもよい。

【 0 0 3 0 】

カテゴリ記憶部 2 0 は、各部 4 0 , 5 0 から読出 / 書込可能な記憶装置であり、図 3 に示すように、カテゴリ ID 2 1 c 毎に 1 つ以上の文書 ID からなる所属文書情報 2 2 c を関連付けて記憶している。ここで、カテゴリ ID 及び所属文書情報 2 2 c の集合をカテゴリデータ 2 0 c と呼ぶ。1 つのカテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c 内の文書 ID の集合は特徴語抽出の最小単位となる文書集合を示しており、この文書集合がカテゴリとも呼ばれる。例えば、カテゴリ ID = C 0 1 で識別されるカテゴリは、文書 ID = D 1 7、D 2 3、D 4 1 で識別される文書が所属している。このカテゴリデータ 2 0 c は、予め与えられるものである。例えば、文書クラスタリングなどの文書分類技術による分類結果をカテゴリデータ 2 0 c としてもよく、文書の作成年や作成者といった属性値によって分割される文書の集合をカテゴリデータ 2 0 c としてもよい。さらに、一つの文書がただ一つのカテゴリに所属するようなカテゴリ構造でもよく、1 つの文書が複数の 10
20
カテゴリに所属するようなカテゴリ構造でもよい。カテゴリデータ 2 0 c は、カテゴリ ID 2 1 c と所属文書情報 2 2 c 以外にも、カテゴリ名やラベルなどの属性情報を有していてもよい。

【 0 0 3 1 】

特徴語記憶部 3 0 は、各部 4 0 , 5 0 から読出 / 書込可能な記憶装置であり、図 4 に示すように、文書特徴語データ 3 0 d t 及びカテゴリ特徴語データ 3 0 c t を記憶する。

【 0 0 3 2 】

文書特徴語データ 3 0 d t は、文書記憶部 1 0 内の文書毎に、当該文書の文書 ID 3 1 d t と、当該文書の内容テキスト情報から抽出された文書特徴語 3 2 d t とを関連付けた 30
データである。この文書特徴語 3 2 d t は、文書特徴語抽出部 4 1 において、文書記憶部 1 0 に記憶された文書データの内容テキスト情報を形態素解析して得られた単語の集合から、不要語を除去して抽出された単語の集合である。不要語の除去では、名詞や未知語といった品詞で、特徴語として利用する単語の条件に合致しない単語や、" こと " や " もの " という一般性が高く特徴語として不適切な単語を排除する。反対に、文書中に 1 回しか出現しないような出現頻度が極端に少ない単語も不要語として排除してもよい。特許文献やメール文書といった特徴語抽出の対象となる文書の種類や、調査や分析といった特徴語抽出の目的などに応じて、保持する品詞の種類を変更することができる。この例では、文書特徴語データ 3 0 d t として文書特徴語 3 2 d t を単語のみで保持しているが、文書中 40
での単語の出現回数 T F を各文書特徴語 3 2 d t の当該単語に関連付けて保持してもよい。T F は、特徴語抽出において、単語の特徴語を求める際の 1 つの指標として利用することができる。

【 0 0 3 3 】

カテゴリ特徴語データ 3 0 c t は、カテゴリ記憶部 2 0 内のカテゴリ ID 2 1 c 及び文書所属情報 2 2 c と同一のカテゴリ ID 3 1 c t 及び所属文書情報 3 2 c t と、当該カテゴリ ID 3 1 c t に関連したカテゴリ特徴語 3 3 c t とを関連付けたデータである。カテゴリ特徴語 3 3 c t は、所属文書情報 3 2 c t 内の文書 ID に関連した文書特徴語 3 2 d t である各単語と、当該各単語に付加された特徴度とからなる。

【 0 0 3 4 】

特徴語抽出部 4 0 は、文書特徴語抽出部 4 1、カテゴリ特徴語抽出部 4 2、カテゴリ共 50

通特徴語抽出部 4 3 及びカテゴリ相違特徴語抽出部 4 4 を備えている。なお、カテゴリ共通特徴語抽出部 4 3 及びカテゴリ相違特徴語抽出部 4 4 は、いずれか一方があれば文書集合の分析が可能のため、いずれか一方を残し、他方を省略することも可能である。

【 0 0 3 5 】

文書特徴語抽出部 4 1 は、文書記憶部 1 0 内の文書毎に内容テキスト情報を形態素解析し、形態素解析の結果から文書特徴語を抽出し、当該抽出した文書特徴語と、当該文書特徴語に対応する文書の文書 ID とを関連付けた文書特徴語データ 3 0 d t を文書特徴語記憶部 3 0 に書き込む機能をもっている。ここで、文書特徴語の抽出は、例えば形態素解析の結果から、文書中に 1 回しか出現していないなど、特徴語抽出において不要な単語（不要語）を排除する処理により実行すればよい。

10

【 0 0 3 6 】

カテゴリ特徴語抽出部 4 2 は、以下の各機能 (f42-1) ~ (f42-5) をもっている。

(f42-1) 文書記憶部 1 0 内の文書特徴語毎に、文書記憶部 1 0 の全ての文書中で当該文書特徴語が出現する文書数を算出する全文書中出現文書数算出機能。

【 0 0 3 7 】

(f42-2) カテゴリ記憶部 2 0 内のカテゴリ ID 2 1 c 毎に、当該カテゴリ ID 2 1 c に関連付けられた文書 ID に関連した文書特徴語が当該文書 ID の文書中で出現する文書数を算出するカテゴリ文書中出現文書数算出機能。

【 0 0 3 8 】

(f42-3) 全文書中出現文書数算出機能により算出された文書数と、カテゴリ文書中出現文書数算出機能により算出された文書数とに基づいて、全文書中におけるカテゴリ ID 2 1 c に関連した文書に対する当該文書特徴語の特徴度を算出する特徴度算出機能。文書特徴語の特徴度は、カテゴリに属する文書の文書特徴語の統計情報に基づいて算出される。

20

【 0 0 3 9 】

(f42-4) この文書特徴語に当該特徴度を付加したカテゴリ特徴語 3 3 c t を作成する機能。

【 0 0 4 0 】

(f42-5) 当該作成したカテゴリ特徴語 3 3 c t と当該カテゴリ特徴語 3 3 c t に関連したカテゴリ ID 3 1 c t 及び所属文書情報 3 2 c t とを関連付けたカテゴリ特徴語データ 3 0 c t を特徴語記憶部 3 0 に書き込む機能。

30

【 0 0 4 1 】

カテゴリ共通特徴語抽出部 4 3 は、以下の各機能 (f43-1) ~ (f43-3) をもっている。

【 0 0 4 2 】

(f43-1) ユーザ操作により共通・相違特徴語提示部 5 3 が選択を受け付けた比較対象集合内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、当該各文書 ID に関連した全ての文書中で当該文書特徴語が出現する文書数を算出する第 1 文書数算出機能。

【 0 0 4 3 】

(f43-2) 選択を受け付けた比較対象集合内の比較対象毎に、当該比較対象内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語が当該各文書 ID の文書中で出現する文書数を算出する第 2 文書数算出機能。

40

【 0 0 4 4 】

(f43-3) 選択を受け付けた比較対象集合内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、第 1 文書数算出機能により算出された文書数と、第 2 文書数算出機能により算出された文書数とに基づいて、比較対象集合における当該文書特徴語の共通性を表す度合いの共通特徴度を算出し、共通特徴度が上位の文書特徴語をカテゴリ共通特徴語として共通・相違特徴語提示部 5 3 に送出するカテゴリ共通特徴語送出機能。ここで、共通特徴度は、各比較対象集合における各特徴語の共通性を表す度合いであり、比較対象集合に属する文書の文書集合の統計情報に基づい

50

て算出される。

【 0 0 4 5 】

カテゴリ相違特徴語抽出部 4 4 は、以下の各機能 (f44-1) ~ (f44-3) をもっている。

【 0 0 4 6 】

(f44-1) ユーザ操作により共通・相違特徴語提示部 5 3 が選択を受け付けた比較対象集合内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、当該各文書 ID に関連した全ての文書中で当該文書特徴語が出現する文書数を算出する第 1 文書数算出機能。

【 0 0 4 7 】

(f44-2) 選択を受け付けた比較対象集合内の比較対象毎に、当該比較対象内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語が当該各文書 ID の文書中で出現する文書数を算出する第 2 文書数算出機能。

【 0 0 4 8 】

(f44-3) 選択を受け付けた比較対象集合内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、第 1 文書数算出機能により算出された文書数と、第 2 文書数算出機能により算出された文書数とに基づいて、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として共通・相違特徴語提示部 5 3 に送出するカテゴリ相違特徴語送出機能。ここで、相違特徴度は、各比較対象集合における各特徴語の相違性を表す度合いであり、それぞれの比較対象に属する文書の文書集合の統計情報に基づいて算出される。

【 0 0 4 9 】

ユーザ操作・提示部 5 0 は、画面提示部 5 1、カテゴリ特徴語提示部 5 2、相違・共通特徴語提示部 5 3 及び関連カテゴリ提示部 5 4 を備えている。なお、関連カテゴリ提示部 5 4 は、文書集合の分析に必須ではなく、省略してもよい。

【 0 0 5 0 】

画面提示部 5 1 は、ユーザの操作に応じて、各記憶部 1 0 , 2 0 , 3 0 を参照して画面データを作成する機能と、当該画面データに基づいて画面を提示する機能とをもっている。ここで、画面データとしては、例えば、カテゴリ記憶部 2 0 内のカテゴリ ID 2 1 c 毎に、当該カテゴリ ID 2 1 c に関連付けられた文書所属情報 2 2 c 内の文書 ID の個数を含むカテゴリ個数データを各セルに提示したクロス分析画面の画面データ、選択された特徴語を分析軸にするようにカテゴリ個数データを修正して各セルに提示したクロス分析画面の画面データ、提示中のクロス分析画面の画面データに基づくグラフ表示画面の画面データ、あるいは、提示中のグラフ表示画面の画面データに基づくクロス分析画面の画面データ、などがある。表示形式は、クロス表示やグラフ表示以外にも、文書集合を平面上に楕円などで表現したマップ表示や、コンピュータのファイルシステムで使われるようなフォルダ表示でもよい。

【 0 0 5 1 】

カテゴリ特徴語提示部 5 2 は、ユーザによるカテゴリの選択を受け付け、カテゴリ特徴語記憶部 3 0 から選択されたカテゴリに対応するカテゴリ特徴語データを取得する。取得したカテゴリ特徴語データに基づき、該カテゴリにおいて特徴度が上位の特徴語をカテゴリ特徴語として、ユーザに提示する。

【 0 0 5 2 】

相違・共通特徴語提示部 5 3 は、ユーザによる比較対象集合の選択を受け付け、選択された比較対象集合を特徴語抽出部 4 0 に送出する機能と、特徴語抽出部 4 0 から受けた共通特徴語を提示する機能と、特徴語抽出部 4 0 から各々の比較対象に対するそれぞれの相違特徴語を受けると、これらの相違特徴語を各々の比較対象に対応づけて提示する機能とをもっている。

【 0 0 5 3 】

関連カテゴリ提示部 5 4 は、ユーザによる着目語集合の選択を受け付け、特徴語記憶部

10

20

30

40

50

30に記憶されたカテゴリ特徴語データに基づき、その着目語集合と各カテゴリとの関連度を算出し、関連度が大きいカテゴリを関連カテゴリとして、該当するカテゴリ個数データを強調表示する機能をもっている。ここで、関連度が大きいカテゴリとしては、関連度がしきい値以上のカテゴリとしたが、これに限らず、関連度が上位s個以内のカテゴリとしてもよい。

【0054】

次に、以上のように構成された特徴語抽出装置の動作を図5乃至図23のフローチャートや模式図を参照しながら説明する。

【0055】

(特徴語抽出部41の動作：図5)

特徴語抽出部41は、概略的には、文書記憶部10内の文書毎に内容テキスト情報を形態素解析し、当該形態素解析の結果から文書特徴語を抽出し、当該抽出した文書特徴語と、当該文書特徴語に対応する文書の文書IDとを関連付けて特徴語記憶部30に書き込む処理を実行する(S1~S4)。

【0056】

具体的には、特徴語抽出部41は、文書記憶部10から、分析対象の全ての文書データの集合である全文書集合docAllを取得する(S1)。

【0057】

次に、特徴語抽出部41は、この全文書集合docAllに含まれる文書データdoc毎に、ステップS3とステップS4の処理を繰り返す(S2)。

【0058】

すなわち、特徴語抽出部41は、文書データdoc毎に内容テキスト情報を形態素解析する(S3)。また、特徴語抽出部41は、この形態素解析の結果から、特徴語抽出の対象とする品詞以外の単語や、“こと”、“もの”などの不要語を排除して抽出した単語群を文書特徴語とする。しかる後、特徴語抽出部41は、抽出した文書特徴語と文書IDとを関連付けた文書特徴語データを特徴語記憶部30に書き込む(S4)。

【0059】

例えば、図2に示す文書データについて、本文25を分析対象の内容テキスト情報とした場合、ステップS3とステップS4の手順により、図4に示すように、文書特徴語データ30dtが特徴語記憶部30に書き込まれる。

【0060】

(カテゴリ特徴語抽出部42の動作：図6)

カテゴリ特徴語抽出部42は、概略的には、文書記憶部10内の文書特徴語毎に、文書記憶部10の全ての文書中で当該文書特徴語が出現する文書数df(t, docAll)を算出する全文書中出現文書数算出処理(S11~S13)と、カテゴリ記憶部20内のカテゴリID21c毎に、当該カテゴリID21cに関連付けられた所属文書情報22cの文書IDに関連した文書特徴語32dtが当該文書IDの文書中で出現する文書数df(t, cat)を算出するカテゴリ文書中出現文書数算出処理(S14~S18)と、全文書中出現文書数算出処理により算出された文書数df(t, docAll)と、カテゴリ文書中出現文書数算出処理により算出された文書数df(t, cat)とに基づいて、全文書中におけるカテゴリID21cに関連した文書に対する当該文書特徴語32dtの特徴度score(t, cat)を算出する特徴度算出処理(S19)と、この文書特徴語32dtに当該特徴度score(t, cat)を付加したカテゴリ特徴語33ctを作成し、当該作成したカテゴリ特徴語33ctと当該カテゴリ特徴語に関連したカテゴリID31ct(カテゴリID21cと同一値)及び1つ以上の文書IDを含む所属文書情報32ct(所属文書情報22cと同一値)とを関連付けて特徴語記憶部30に書き込む処理(S20)とを実行する。

【0061】

具体的には、カテゴリ特徴語抽出部42は、特徴語記憶部30から全ての文書docAllの文書特徴語データを取得する(S11)。

10

20

30

40

50

【0062】

次に、カテゴリ特徴語抽出部42は、ステップS11によって得られた文書特徴語データに含まれる文書特徴語t毎に、ステップS13の処理を繰り返す(S12)。

【0063】

すなわち、カテゴリ特徴語抽出部42は、文書特徴語t毎に、全ての文書docAllの文書特徴語データを参照しながら、全文書集合docAll中で当該文書特徴語tが出現する文書数df(t, docAll)を求める処理(S13)を繰り返す。

【0064】

しかる後、カテゴリ特徴語抽出部42は、カテゴリ記憶部20から全てのカテゴリデータ20cを取得する(S14)。

10

【0065】

また、カテゴリ特徴語抽出部42は、全てのカテゴリcatについて、カテゴリID21c毎に、ステップS16~S20の処理を繰り返す(S15)。

【0066】

さらに、カテゴリ特徴語抽出部42は、当該カテゴリID21cに関連付けられた所属文書情報22cの文書IDに関連した文書特徴語を特徴語記憶部30から読み出すことにより、特徴語記憶部30から、カテゴリcatに所属する複数の文書について、それぞれの文書の文書特徴語データを取得する(S16)。

【0067】

続いて、カテゴリ特徴語抽出部42は、取得した文書特徴語データに含まれる文書特徴語t毎に、ステップS18の処理を繰り返す(S17)。

20

【0068】

カテゴリ特徴語抽出部42は、文書特徴語t毎に、ステップS16で取得した文書特徴語データを参照しながら、カテゴリcatに所属する複数の文書中で、文書特徴語tが出現する文書数df(t, cat)を求める(S18)。

【0069】

カテゴリ特徴語抽出部42は、ステップS13で求めた文書数df(t, docAll)と、ステップS18で求めた文書数df(t, cat)に基づき、全文書集合docAllにおけるカテゴリcatに対する特徴語tの特徴度score(t, cat)を算出する(S19)。

30

【0070】

具体的には、特徴度score(t, cat)は、各文書数df(t, docAll)、df(t, cat)に基づいて、表3に示すように、共通パラメータa, b, c, d, nを算出した後、表4に示す如き、いずれかの統計指標として算出される。

【0071】

【表3】

[表3]

特徴度(t, cat) 算出用の 共通パラメータ	特徴度(t, cat)算出用の共通パラメータの算出式
a	$a = df(t, cat)$
b	$b = df(t, docAll) - df(t, cat)$
c	$c = cat - df(t, cat)$
d	$d = docAll - df(t, docAll) - cat + df(t, cat)$
n	$n = a + b + c + d = docAll $

40

【0072】

【表 4】

[表 4]

統計指標の例	特徴度 (t, cat) の算出式の例
対数尤度比 LLR	$LLR' = a \log(a n / ((a+b)(a+c))) + b \log(b n / ((a+b)(b+d)))$ $+ c \log(c n / ((c+d)(a+c))) + d \log(d n / ((c+d)(b+d)))$ $\text{if}((ad-bd) < 0) LLR = -LLR'$ $\text{else } LLR = LLR'$ $\text{score}(t, cat) = LLR$
ダイス係数 Dice	$Dice = 2a / ((a+b) + (a+c))$ $\text{score}(t, cat) = Dice$
イエーツ補正 χ ² 乗値 Yates'	$Yates' = \frac{n \left(\left ad - bc \right - \frac{n}{2} \right)^2}{(a+b)(c+d)(a+c)(b+d)}$ $\text{if}((ad-bd) < 0) Yates = -Yates'$ $\text{else } Yates = Yates'$ $\text{score}(t, cat) = Yates$
自己相互情報量 MI	$MI = \log(a n / ((a+b)(a+c)))$ $\text{score}(t, cat) = MI$

10

20

【 0 0 7 3 】

ここでは、対数尤度比 LLR という統計指標として、各単語 t の特徴度 score (t , cat) を算出している。但し、統計指標は、対数尤度比 LLR に限らず、例えば、ダイス係数 Dice、イエーツ補正 2乗値 Yates' 又は自己相互情報量 MI 等としてもよい。なお、各統計指標にはそれぞれ特徴があるため、各統計指標の特徴に応じて、得られる特徴語の傾向が異なる。

【 0 0 7 4 】

例えば、ダイス係数 Dice は、カテゴリ cat 内で単語 t が出現する文書数 df (t , cat) の大きい単語 t (カテゴリ cat に多く含まれる単語 (高頻度)) を高く評価する。

30

【 0 0 7 5 】

イエーツ補正 2乗値 Yates' は、全文書集合 doc All 中での出現確率に対し、カテゴリ cat 中での出現確率が高い単語を高く評価する。結果として、イエーツ補正 2乗値 Yates' は、対数尤度比 LLR やダイス係数 Dice を利用した場合よりも、比較的 low 頻度の単語が特徴語として抽出されやすい。

【 0 0 7 6 】

自己相互情報量 MI は、全文書集合 doc All 中での出現確率と、カテゴリ cat 中での出現確率とで偏りの大きい単語を高く評価する。但し、自己相互情報量 MI は、低頻度語を過大評価する傾向があるため、利用する場合、df (t , cat) が極端に小さい単語を特徴語から排除するなどの処理が必要となる。以上の各統計量の詳細については、非特許文献 1 に記載されている。

40

【 0 0 7 7 】

カテゴリ特徴語抽出部 4 2 は、ステップ S 1 9 で算出した各特徴語の特徴度 score (t , cat) を、その特徴語に付加したカテゴリ特徴語情報 3 3 c t として、カテゴリ cat のカテゴリデータ 2 0 c に付加したカテゴリ特徴語データ 3 0 c t を特徴語記憶部 3 0 に格納する (S 2 0) 。

【 0 0 7 8 】

(カテゴリ共通特徴語抽出部 4 3 の動作 : 図 7)

画面提示部 5 1 は、図 8 に示すように、カテゴリ記憶部 2 0 内のカテゴリ ID 2 1 c 毎

50

に、当該カテゴリID 21cに関連付けられた文書IDの個数を含むカテゴリ個数データを提示した画面G10を表示する。例えば、画面G10内のセルc1, c2は、図示しないカテゴリID 21毎に表示されており、各セルc1, c2内の値“75”, “50”がカテゴリ個数データに相当している。

【0079】

相違・共通特徴語提示部53は、画面提示部51によるカテゴリ個数データの提示中、ユーザの操作により、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象cmp_iからなる比較対象集合tgtSetの選択を受け付ける。例えば図9に示す場合、第1の比較対象cmp1は、実線枠f1で囲まれた5つのカテゴリ個数データ“65”, “50”, “69”, “75”, “72”の各カテゴリIDの和集合であり、第2の比較対象cmp2は、点線枠f2で囲まれた5つのカテゴリ個数データ“10”, “21”, “45”, “53”, “35”の各カテゴリIDの和集合である。

10

【0080】

カテゴリ共通特徴語抽出部43は、概略的には、相違・共通特徴語提示部53により選択を受け付けた比較対象集合tgtSet内の各カテゴリID 21cに関連付けられた所属文書情報22cの各文書IDに関連した文書特徴語毎に、当該各文書IDに関連した全ての文書(tgtDocs)中で当該文書特徴語が出現する文書数df(t, tgtDocs)を算出する第1文書数算出処理(S21~S25)と、選択を受け付けた比較対象集合tgtSet内の比較対象cmp_i毎に、当該比較対象cmp_i内の各カテゴリID 21cに関連付けられた所属文書情報22cの各文書IDに関連した文書特徴語が当該各文書IDの文書中で出現する文書数df(t, cmp)を算出する第2文書数算出処理(S26~S29)と、選択を受け付けた比較対象集合tgtSet内の各カテゴリID 21cに関連付けられた所属文書情報22cの各文書IDに関連した文書特徴語毎に、第1文書数算出処理により算出された文書数df(t, tgtDocs)と、第2文書数算出処理により算出された文書数df(t, cmp)とに基づいて、比較対象集合tgtSetにおける当該文書特徴語の共通性を表す度合いの共通特徴度com(t, tgtSet)を算出し、共通特徴度com(t, tgtSet)が上位の文書特徴語をカテゴリ共通特徴語として相違・共通特徴語提示部53に送出する処理を実行する(S30~S31)。

20

【0081】

具体的には、カテゴリ共通特徴語抽出部43は、ユーザ操作・提示部50から、ユーザによって選択された各カテゴリ個数データに対応する複数の比較対象cmp(各文書ID)からなる比較対象集合tgtSetを取得する(S21)。

30

【0082】

カテゴリ共通特徴語抽出部43は、比較対象集合tgtSetに含まれる全ての比較対象cmpの和集合をとり、比較範囲tgtDocsを求める(S22)。

【0083】

カテゴリ共通特徴語抽出部43は、比較範囲tgtDocsに含まれる全ての文書IDに関連した文書特徴語データを、特徴語記憶部30から取得する(S23)。

【0084】

カテゴリ共通特徴語抽出部43は、ステップS23で取得した文書特徴語データに含まれる全ての特徴語tについて、ステップS25を繰り返す(S24)。

40

【0085】

カテゴリ共通特徴語抽出部43は、ステップS23で取得した文書特徴語データを参照しながら、比較範囲tgtDocsに含まれる文書IDに関連した文書の中で、特徴語tが出現する文書数df(t, tgtDocs)を求める(S25)。

【0086】

カテゴリ共通特徴語抽出部43は、比較対象集合tgtSetに含まれる比較対象cmp毎に、ステップS27~S29の処理を繰り返す(S26)。

【0087】

50

カテゴリ共通特徴語抽出部 43 は、比較対象 cmp の文書 ID に関連する文書特徴語データを、特徴語記憶部 30 から取得する (S27)。

【0088】

カテゴリ共通特徴語抽出部 43 は、ステップ S27 で取得した文書特徴語データに含まれる全ての特徴語 t について、ステップ S29 の処理を繰り返す (S28)。

【0089】

カテゴリ共通特徴語抽出部 43 は、ステップ S27 で取得した文書特徴語データを参照しながら、比較対象 cmp の文書 ID に関連した文書の中で、単語 t が出現する文書数 $df(t, cmp)$ を求める (S29)。

【0090】

カテゴリ共通特徴語抽出部 43 は、ステップ S25 で算出した比較範囲 $tgtDocs$ 内で単語が出現する文書数 $df(t, tgtDocs)$ と、ステップ S29 で算出した各比較対象 cmp 内で単語が出現する文書数 $df(t, cmp)$ に基づき、比較範囲 $tgtDocs$ 内の文書に含まれる全ての単語 t について、比較対象集合 $tgtSet$ における共通特徴度 $com(t, tgtSet)$ を算出する (S30)。

【0091】

具体的には、共通特徴度 $com(t, tgtSet)$ を算出する場合、始めに、各文書数 $df(t, tgtDocs)$ 、 $df(t, cmp)$ に基づいて、表 5 に示すように、共通パラメータ a' 、 b' 、 c' 、 d' 、 n' を算出した後、表 6 に示す如き、いずれかの統計指標として評価値 $eval(t, cmp_i)$ を算出する。

【0092】

【表 5】

[表 5]

評価値 $eval(t, cmp_i)$ 算出用の 共通パラメータ	評価値 $eval(t, cmp_i)$ 算出用の共通パラメータの算出式
a'	$a' = df(t, tgtDocs)$
b'	$b' = df(t, tgtDocs) - df(t, cmp_i)$
c'	$c' = cmp - df(t, cmp_i)$
d'	$d' = tgtDocs - df(t, tgtDocs) - cmp_i + df(t, cmp_i)$
n'	$n' = a' + b' + c' + d' = tgtDocs $

【0093】

10

20

30

【表 6】

[表 6]

統計指標の例	評価値eval(t, cmp_i)の算出式の例
対数尤度比 LLR	$\begin{aligned} \text{LLR}' = & a' \log(a' n' / ((a' + b') (a' + c'))) \\ & + b' \log(b' n' / ((a' + b') (b' + d'))) \\ & + c' \log(c' n' / ((c' + d') (a' + c'))) \\ & + d' \log(d' n' / ((c' + d') (b' + d'))) \\ \text{if}((a' d' - b' d') < 0) & \text{LLR} = -\text{LLR}' \\ \text{else LLR} = & \text{LLR}' \\ \text{eval}(t, \text{cmp}_i) = & \text{LLR} \end{aligned}$
イエーツ補正 × 2乗値 Yates'	$\text{Yates}' = \frac{n' \left(a' d' - b' c' - \frac{n'}{2} \right)^2}{(a' + b')(c' + d')(a' + c')(b' + d')}$ $\begin{aligned} \text{if}((a' d' - b' d') < 0) & \text{Yates} = -\text{Yates}' \\ \text{else Yates} = & \text{Yates}' \\ \text{eval}(t, \text{cmp}_i) = & \text{Yates} \end{aligned}$
自己相互情報量 MI	$\text{MI} = \log(a' n' / ((a' + b') (a' + c')))$ $\text{eval}(t, \text{cmp}_i) = \text{MI}$

10

20

【 0 0 9 4 】

続いて、比較対象 cmp_i の評価値 $\text{eval}(t, \text{cmp}_i)$ の総和を算出し、得られた総和の値を、比較範囲 tgtDocs における単語 t の共通特徴度 $\text{com}(t, \text{tgtSet})$ とする。

【 0 0 9 5 】

この指標では、より多くの比較対象 cmp_i に特徴語として含まれ、かつそれぞれの比較対象 cmp_i で、より高い評価値 $\text{eval}(t, \text{cmp}_i)$ を持つ単語ほど、共通特徴語として高く評価される。

【 0 0 9 6 】

ここでは、例えば対数尤度比 LLR という統計指標を用いて、単語の共通特徴語 $\text{com}(t, \text{tgtSet})$ を求めている。なお、対数尤度比 LLR に代えて、前述したイエーツ 2乗値や自己相互情報量 MI などの統計指標を用いてもよい。

30

【 0 0 9 7 】

このような統計指標において、全文書集合 docAll における各特徴語 t の出現頻度 $\text{df}(t, \text{docAll})$ や、各カテゴリ cat における各単語の出現頻度 $\text{df}(t, \text{cat})$ も利用してもよい。

【 0 0 9 8 】

しかる後、カテゴリ共通特徴語抽出部 43 は、ステップ S30 で算出した各特徴語の共通特徴度 $\text{com}(t, \text{tgtSet})$ について、上位 r 個の単語を tgtSet の共通特徴語 comTerms として、ユーザ操作・提示部 50 に送出する (S31)。

40

【 0 0 9 9 】

ここで、 r とは共通特徴語、相違特徴語及びカテゴリ特徴語の提示において、提示する特徴語の個数の設定値であり、事前に設定されてもよく、特徴語抽出を行う都度設定されてもよい。また、共通特徴度 $\text{com}(t, \text{tgtSet})$ が上位 r 個以内の特徴語を共通特徴語としたが、これに限らず、共通特徴度 $\text{com}(t, \text{tgtSet})$ がしきい値以上の特徴語を共通特徴語としてもよい。

【 0 1 0 0 】

相違・共通特徴語提示部 53 は、図 10 に示すように、ステップ S31 で送出された r 個のカテゴリ共通特徴語をリスト L_{com} に提示する。

50

【 0 1 0 1 】

(カテゴリ相違特徴語抽出部 4 4 の動作 : 図 1 1)

画面提示部 5 1 は、図 8 に示したように、カテゴリ記憶部 2 0 内のカテゴリ ID 2 1 c 毎に、当該カテゴリ ID 2 1 c に関連付けられた文書 ID の個数を含むカテゴリ個数データを提示した画面 G 1 0 を表示する。

【 0 1 0 2 】

相違・共通特徴語提示部 5 3 は、画面提示部 5 1 によるカテゴリ個数データの提示中、ユーザの操作により、複数個のカテゴリ個数データの各カテゴリ ID の和集合である複数の比較対象 cmp_i からなる比較対象集合 $tgtSet$ の選択を受け付ける。

【 0 1 0 3 】

カテゴリ相違特徴語抽出部 4 4 は、概略的には、図 9 に示したように相違・共通特徴語提示部 5 3 により選択を受け付けた比較対象集合 $tgtSet$ 内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、当該各文書 ID に関連した全ての文書 ($tgtDocs$) 中で当該文書特徴語が出現する文書数 $df(t, tgtDocs)$ を算出する第 1 文書数算出処理 (S 4 1 ~ S 4 5) と、選択を受け付けた比較対象集合 $tgtSet$ 内の比較対象 cmp_i 毎に、当該比較対象 cmp_i 内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語が当該各文書 ID の文書中で出現する文書数 $df(t, cmp)$ を算出する第 2 文書数算出処理 (S 4 6 ~ S 4 9) と、選択を受け付けた比較対象集合 $tgtSet$ 内の各カテゴリ ID 2 1 c に関連付けられた所属文書情報 2 2 c の各文書 ID に関連した文書特徴語毎に、第 1 文書数算出処理により算出された文書数 $df(t, tgtDocs)$ と、第 2 文書数算出処理により算出された文書数 $df(t, cmp)$ とに基づいて、各比較対象 cmp_i における当該文書特徴語の相違性を表す度合いの相違特徴度 $diff(t, cmp)$ を算出し、相違特徴度 $diff(t, cmp)$ が上位の文書特徴語をカテゴリ相違特徴語として相違・共通特徴語提示部 5 3 に送出する処理を実行する (S 5 0 ~ S 5 1) 。

【 0 1 0 4 】

具体的には、カテゴリ相違特徴語抽出部 4 4 は、前述したステップ S 2 1 ~ S 2 9 と同様に、ステップ S 4 1 ~ S 4 9 を実行する。なお、ステップ S 4 1 ~ S 4 9 に代えて、ステップ S 2 1 ~ S 2 9 の結果をステップ S 5 0 で用いるようにカテゴリ相違特徴語抽出部 4 4 を変形してもよい。逆に、ステップ S 2 1 ~ S 2 9 に代えて、ステップ S 4 1 ~ S 4 9 の結果をステップ S 3 0 で用いるようにカテゴリ共通特徴語抽出部 4 3 を変形してもよい。

【 0 1 0 5 】

ステップ S 4 1 ~ S 4 9 の実行後、カテゴリ相違特徴語抽出部 4 4 は、ステップ S 2 5 と同様のステップ S 4 5 で算出した比較範囲 $tgtDocs$ 内で単語が出現する文書数 $df(t, tgtDocs)$ と、ステップ S 2 9 と同様のステップ S 4 9 で算出した各比較対象 cmp 内で単語が出現する文書数 $df(t, cmp)$ に基づき、比較範囲 $tgtDocs$ 内の文書に含まれる全ての特徴語 t について、各比較対象 cmp に対する相違特徴度 $diff(t, cmp)$ を算出する (S 5 0) 。

【 0 1 0 6 】

具体的には、相違特徴度 $diff(t, cmp)$ としては、各文書数 $df(t, tgtDocs)$, $df(t, cmp)$ に基づいて、表 7 に示すように、共通パラメータ a " , b " , c " , d " , n " を算出した後、表 8 に示す如き、いずれかの統計指標として相違特徴度 $diff(t, cmp)$ とする。ここでは、例えば、T 統計量を相違特徴度 $diff(t, cmp)$ とする場合について述べる。

【 0 1 0 7 】

10

20

30

40

【表 7】

[表 7]

相違特徴度 diff(t, cmp) 算出用の 共通パラメータ	相違特徴度diff(t, cmp)算出用の 共通パラメータの算出式
a''	$a'' = df(t, cmp_i)$
b''	$b'' = df(t, cmpDocs_i)$
c''	$c'' = cmp_i - df(t, cmp_i)$
d''	$d'' = cmpDocs_i - df(t, cmpDocs_i)$
n''	$n'' = a'' + b'' + c'' + d'' = tgtDocs $

10

【 0 1 0 8 】

【表 8】

[表 8]

統計指標の例	相違特徴度diff(t, cmp)の算出式の例
T統計量 tscore	$tscore = (a'' - ((a'' + b'') / (a'' + c'')) / \sqrt{a''}$ $diff(t, cmp) = tscore$
対数尤度比LLR	$LLR' = a'' \log(a'' n'' / ((a'' + b'') (a'' + c''))) + b'' \log(b'' n'' / ((a'' + b'') (b'' + d''))) + c'' \log(c'' n'' / ((c'' + d'') (a'' + c''))) + d'' \log(d'' n'' / ((c'' + d'') (b'' + d'')))$ $if((a'' d'' - b'' d'') < 0) LLR = -LLR'$ $else LLR = LLR'$ $diff(t, cmp) = LLR$
イエーツ補正 χ^2 乗値 Yates'	$Yates' = \frac{n'' \left(a'' d'' - b'' c'' - \frac{n''}{2} \right)^2}{(a'' + b'')(c'' + d'')(a'' + c'')(b'' + d'')}$ $if((a'' d'' - b'' d'') < 0) Yates = -Yates'$ $else Yates = Yates'$ $diff(t, cmp) = Yates$
自己相互情報量 MI	$MI = \log(a'' n'' / ((a'' + b'') (a'' + c'')))$ $diff(t, cmp) = MI$

20

30

40

【 0 1 0 9 】

相違特徴度 $diff(t, cmp)$ は、T 統計量を利用して、単語 t について、比較対象 cmp_i と、比較範囲 $tgtDocs$ から比較対象 cmp_i を除いた差集合 $cmpDocs_i$ との間の出現頻度の平均の差に基づき、有意性を求める指標である。これにより、比較対象 cmp_i において、比較対象以外の比較範囲 ($cmpDocs_i$) に比べ、有意に出現頻度の多い単語を相違特徴語として抽出することができる。なお、T 統計量 (T スコアともいう) については、例えば非特許文献 3 に記載されている。また、T 統計量に代えて、特徴度の算出の説明で述べたような対数尤度比 LLR や χ^2 乗値、自己相互情報量 MI などの統計指標を使ってもよい。

50

【0110】

このような統計指標において、全文書集合 $docAll$ における各特徴語 t の出現頻度 $df(t, docAll)$ や、各カテゴリ cat における各特徴語の出現頻度 $df(t, cat)$ も利用してもよい。

【0111】

カテゴリ相違特徴語抽出部 44 は、比較対象集合 $tgtSet$ に含まれるそれぞれの対象集合 cmp について、ステップ S50 で算出した各特徴語の相違特徴度 $diff(t, cmp)$ が上位 r 個の特徴語を、相違特徴語 $diffTerms(cmp)$ として、ユーザ操作・提示部 50 に送出する (S51)。ここで、上位個数 r は前述した設定値である。また、相違特徴度 $diffTerms(cmp)$ が上位 r 個以内の特徴語を相違特徴語としたが、これに限らず、相違特徴度 $diffTerms(cmp)$ がしきい値以上の特徴語を相違特徴語としてもよい。

10

【0112】

相違・共通特徴語提示部 53 は、図 12 に示すように、ステップ S51 で送出された各 r 個のカテゴリ相違特徴語をリスト $Ldif1, Ldif2$ に提示する。

【0113】

(関連カテゴリ提示部 54 の動作：図 13)

関連カテゴリ提示部 54 は、概略的には、例えば相違・共通特徴語提示部 53 によるカテゴリ相違特徴語の提示中、ユーザの操作により、複数のカテゴリ相違特徴語からなる着目語集合 $tgtTerms$ の選択を受けると、当該選択を受け付けた着目語集合 $tgtTerms$ と、カテゴリ特徴語記憶部 30 内のカテゴリ特徴語との関連度 $rel(cat, tgtTerms)$ を算出し、関連度 $rel(cat, tgtTerms)$ の高いカテゴリ特徴語に関連付けられたカテゴリ ID に関連したカテゴリ個数データを強調表示する (S61~S66)。

20

【0114】

具体的には、関連カテゴリ提示部 54 は、ユーザによって選択された複数の単語から構成される着目語集合 $tgtTerms$ を取得する (S61)。なお、着目語集合に含まれる単語としては、前述したカテゴリ相違特徴語に限らず、カテゴリ特徴語やカテゴリ共通特徴語などが適宜、選択可能となっている。

【0115】

関連カテゴリ提示部 54 は、特徴語記憶部 30 から全てのカテゴリ特徴語 $33ct$ を取得する (S62)。

30

【0116】

関連カテゴリ提示部 54 は、全てのカテゴリデータ cat について、ステップ S64 とステップ S65 の処理を繰り返す (S63)。

【0117】

関連カテゴリ提示部 54 は、カテゴリデータ cat のカテゴリ特徴語 $33ct$ に含まれる特徴語を特徴度でソートし、特徴語ランキング $termRnk$ を求める (S64)。

【0118】

関連カテゴリ提示部 54 は、着目語集合 $tgtTerms$ と、特徴語ランキング $termRnk$ に基づいて、カテゴリ cat と着目語集合 $tgtTerms$ との関連度 $rel(cat, tgtTerms)$ を求める (S65)。

40

【0119】

関連度 $rel(cat, tgtTerms)$ としては、平均精度と呼ばれる統計指標を利用することができる。この統計指標は、特徴語ランキング $termRnk$ において、着目語集合 $tgtTerms$ に含まれる単語が、より上位に多く出現する程、高い値をとる指標である。平均精度の詳細については、非特許文献 2 に記載されている。関連度 $rel(cat, tgtTerms)$ としては、平均精度以外にも、カテゴリ cat のカテゴリ特徴語において、着目語集合 $tgtTerms$ に存在する単語 t の特徴度 $score(t, cat)$ を足し合わせた値としてもよい。

50

【 0 1 2 0 】

関連カテゴリ提示部 5 4 は、ステップ S 6 5 により算出された各カテゴリの関連度 $rel(cat, tgtTerms)$ に基づき、当該関連度 $rel(cat, tgtTerms)$ がしきい値 s 以上のカテゴリ特徴語に関連付けられたカテゴリ ID を、着目語集合 $tgtTerms$ の関連カテゴリ $relCats$ として、関連カテゴリ $relCats$ に含まれるカテゴリ ID に関連したカテゴリ個数データのセルを強調表示する (S 6 6)。

【 0 1 2 1 】

なお、関連カテゴリとしては、関連度 $rel(cat, tgtTerms)$ がしきい値以上のカテゴリに限らず、関連度 $rel(cat, tgtTerms)$ が上位 t 個以内のカテゴリとしてもよい。しきい値 s や上位個数 t は、前述した上位個数 r と同様に、予め設定されていてもよく、関連カテゴリの提示を行う都度設定されてもよい。

10

【 0 1 2 2 】

(ユーザ操作・提示部 5 0 の動作：図 1 4)

次に、以上のような特徴語抽出部 4 0 や関連カテゴリ提示部 5 4 等の処理をユーザ操作に応じて用いるユーザ操作・提示部 5 0 の動作について説明する。なお、文書特徴語抽出部 4 1 及びカテゴリ特徴語抽出部 4 2 の動作 (ステップ S 1 ~ S 4 , S 1 1 ~ S 2 0) は予め完了している状態であるとする。

【 0 1 2 3 】

ユーザ操作・提示部 5 0 は、概略的には、カテゴリ ID 毎にカテゴリ個数データをセル表示し、ユーザによる選択操作に応じて、カテゴリ特徴語、カテゴリ共通特徴語及びカテゴリ相違特徴語を提示し、また、関連カテゴリを強調して提示する処理を実行する (S 1 0 0 ~ S 1 3 1)。

20

【 0 1 2 4 】

具体的には、ユーザ操作・提示部 5 0 においては、画面提示部 5 1 が、特徴語記憶部 3 0 に記憶された全てのカテゴリデータについて、それぞれのカテゴリを 1 つのセルとして表示する (S 1 0 0)。

【 0 1 2 5 】

この表示例としては、図 8 の画面 G 1 0 に示すようなクロス表示が挙げられる。この例では、文書データは図 2 に示すような特許文献とし、カテゴリとしては、特許文献の出願人 1 4 d の属性値と、出願日 1 3 d の上位 4 桁である出願年の属性値との 2 つの属性値で予め分類された文書集合を想定する。ユーザは特許文献から競合他社の技術動向を調査する作業中であるものとする。画面 G 1 0 のクロス表示において、一つのセルが 1 つのカテゴリに相当する。例えば、セル c 1 は、F 社が 2 0 0 4 年に申請した特許文献を含むカテゴリに相当する。なお、表示形式は、クロス表示に限らず、グラフ表示、マップ表示又はフォルダ表示といった任意の表示形式が使用可能となっている。

30

【 0 1 2 6 】

ステップ S 1 1 0 ~ S 1 1 2 は、ユーザによるカテゴリの選択を受け付け、該カテゴリにおけるカテゴリ特徴語を提示するカテゴリ特徴語提示部 5 2 の処理を示している。

【 0 1 2 7 】

すなわち、カテゴリ特徴語提示部 5 2 は、ステップ S 1 0 0 によりセルとして表示されたカテゴリに対して、ユーザがカテゴリ cat を選択した場合、ステップ S 1 1 1 と S 1 1 2 の処理を行う (S 1 1 0)。

40

【 0 1 2 8 】

カテゴリ特徴語提示部 5 2 は、ユーザが選択したカテゴリ cat のカテゴリ ID に関連するカテゴリ特徴語データを、特徴語記憶部 3 0 から取得する (S 1 1 1)。

【 0 1 2 9 】

カテゴリ特徴語提示部 5 2 は、取得したカテゴリ特徴語データに含まれる特徴度に基づき、特徴度 $score(t, cat)$ が上位 r 個の特徴語をカテゴリ特徴語として、ユーザに提示する (S 1 1 2)。

【 0 1 3 0 】

50

例えば、図15に示すように、ユーザがセル(カテゴリ)c2をマウスのクリックなどにより選択した場合、該カテゴリに対するカテゴリ特徴語をリストL2に表示する。これにより、ユーザは、選択したセルc2に含まれる文献の内容の特徴を把握することができる。すなわち、選択したセルc2に対応するF社の2005年の出願特許におけるカテゴリ特徴語のリストL2に“検索”や“Web”という技術用語が有意に出現していることにより、ユーザは、F社の2005年における注力技術としては、検索やWebなどがあることを把握できる。ユーザは、他に選択したセルc3があれば、同様にリストL3から、出願年及び企業名の分析軸におけるカテゴリ特徴語を把握することができる。

【0131】

ステップS120～S126は、ユーザによる比較対象集合tgtSetの選択を受け付け、選択された比較対象集合tgtSetを特徴語抽出部40に送り、特徴語抽出部40によって抽出される共通特徴語と相違特徴語を受け取り、ユーザに提示する相違・共通特徴語提示部53の処理を示している。

10

【0132】

すなわち、相違・共通特徴語提示部53は、ステップS100によって表示されたカテゴリに対して、比較対象集合tgtSetとして複数の比較対象を選択した場合、ステップS121～S126の処理を行う(S120)。

【0133】

相違・共通特徴語提示部53は、比較対象集合tgtSetを特徴語抽出部40に送る(S121)。特徴語抽出部40では、カテゴリ共通特徴語抽出部43が、前述したステップS21～S31の処理を実行し、得られた共通特徴語comTermsを相違・特徴語提示部53に送出する。

20

【0134】

相違・共通特徴語提示部53は、特徴語抽出部40から共通特徴語comTermsを受け取り、ユーザに提示する(S122)。

【0135】

相違・共通特徴語提示部53は、比較対象集合tgtSetを特徴語抽出部40に送る(S123)。特徴語抽出部40では、カテゴリ相違特徴語抽出部44が、前述したステップS41～S51の処理を実行し、得られた相違特徴語diffTermsを相違・特徴語提示部53に送出する。

30

【0136】

相違・共通特徴語提示部53は、特徴語抽出部40から相違特徴語diffTermsを取得する(S124)。

【0137】

相違・共通特徴語提示部53は、比較対象集合tgtSetに含まれる全ての比較対象cmpについて、ステップS126の処理を繰り返す(S125)。

【0138】

相違・共通特徴語提示部53は、比較対象cmpに対する相違特徴語diffTerms(cmp)をユーザに提示する(S126)。

【0139】

40

ステップS120～S126における表示例は、図12に示す通りである。ユーザは、例えば、ユーザがA社とB社の技術動向を比較したい場合、画面G10においてA社に関するカテゴリを示す複数のセルを実線枠f1のように選択することにより1つの比較対象を選択する。

【0140】

また、もう一つの比較対象として、B社に関するカテゴリを示す複数のセルを点線枠f2のように選択した場合、実線枠f1と点線枠f2で示される2つの比較対象から構成される比較範囲における共通特徴語リストLcomを表示する。

【0141】

このように、両社の出願特許における共通特徴語リストLcomに有意に出現する技術

50

用語として、“分類”や“クラスタリング”が提示され、これらの技術がA社とB社で共通する技術分野であることを把握できる。

【0142】

また、実線枠f1で示される比較対象に対する相違特徴語リストLdif1を表示し、点線枠f2で示される比較対象に対する相違特徴語リストLdif2を表示する。このような相違特徴語リストLdif1、Ldif2の表示により、A社とB社の独自性を表す技術を把握することができる。

【0143】

また、複数の比較対象に対し、共通特徴語と相違特徴語を表示することにより、単に文書集合に対する特徴語を提示するよりも、比較対象間の特徴を、より明確にユーザに示すことができる。

【0144】

ステップS130、S131は、ステップS112や、S122、S126による各特徴語の提示中に、これら各特徴語から選択された着目語からなる着目語集合tgtTermsを受け付け、その着目語集合と各カテゴリとの関連度を算出し、関連度の高いカテゴリを関連カテゴリとしてユーザに提示する関連カテゴリ提示部54の処理を示している。

【0145】

関連カテゴリ提示部54は、カテゴリ特徴語提示部52や、相違・共通特徴語提示部53によって提示された、カテゴリ特徴語または共通特徴語または相違特徴語から、ユーザが着目語を選択した場合、ステップS131の処理を行う(S130)。ここで、ユーザは複数の単語を着目語として選択できるものとし、選択された複数の着目語を着目語集合tgtTermsとする。また、本実施形態では、提示された特徴語から着目語を選択する場合について説明したが、これに限らず、Webの検索のようにユーザが任意のキーワードを着目語として入力してもよい。

【0146】

関連カテゴリ提示部54は、前述したステップS61～S66の処理を実行することにより、着目語集合tgtTermsと各カテゴリとの関連度を算出して関連度の高いカテゴリを関連カテゴリとしてユーザに提示する(S131)。

【0147】

例えば、図16に示すように、ユーザは、共通特徴語のリストLcomや相違特徴語のリストLdif1、Ldif2の中から、着目したい単語Tcom1とTdif2を選択する。ここでは、ユーザは“分類”と“XML”との単語が気になった場合、単語“分類”を示す共通語Tcom1と、単語“XML”を示す相違特徴語Tdif2とを着目語として選択する。関連カテゴリ提示部54は、ユーザの着目語の選択を受けて、着目語との関連度が高い関連カテゴリのセルc4を、背景色を変更する等して、強調表示する。

【0148】

これによって、ユーザは着目する技術について、調査すべき範囲の糸口をつかむことができる。図16に示した例では、ユーザが着目した“分類”と“XML”について、企業の観点から見ると、C社もこれらの技術に関連していることがわかる。さらに、出願年の観点からみると2006～2008年の間で、これらの技術に関連する特許が有意に出現していることがわかる。これによって、ユーザは着目している技術について、詳細に調査すべき範囲を明確化でき、効率的に先行技術調査を行うことができる。

【0149】

また、着目語の選択は、共通特徴語や相違特徴語だけでなく、カテゴリ特徴語からも選択することができる。例えば図17に示すように、着目語の選択に加え、セルc2におけるカテゴリ特徴語リストL2内のカテゴリ特徴語からも着目語を選択した場合には、この選択に応じて、関連カテゴリの表示が変化する。

【0150】

画面提示部51は、ユーザがシステムの終了を選択した場合、処理を終了し、それ以外はステップS110に処理を戻す(S140)。

10

20

30

40

50

【 0 1 5 1 】

例えば、ステップ S 1 1 0 に処理を戻し、調査を継続する場合の例について説明する。図 1 8 は比較対象の絞り込みと共通特徴語及び相違特徴語の表示例を表す図である。ユーザは、共通特徴語や相違特徴語、カテゴリ特徴語の提示や、着目語指定に対する関連カテゴリの提示を受けて、比較対象の縮小（絞り込み）や拡大といった変更を行うことができる。

【 0 1 5 2 】

例えば、図 1 6 に示す如き、特徴語や関連語カテゴリの提示中に、ユーザは、図 1 8 に示すように、比較範囲の各棒 f_1 , f_2 を出願年について 2 0 0 6 ~ 2 0 0 8 年に絞り込み、新たな点線棒 f_3 により C 社を比較企業に選択する。これら各棒 $f_1 \sim f_3$ に基づく比較対象集合 $t g t S e t$ に基づいて、特徴語抽出装置は、提示する共通特徴語や相違特徴語を変化させる。これにより、ユーザはそれまで思いつかなかったが着目すべき技術用語を発見する手がかりとなる。

10

【 0 1 5 3 】

図 1 9 は、ユーザによる着目語の変更（追加 / 削除）と、関連カテゴリの表示例を表す図である。ユーザは、共通特徴語や相違特徴語やカテゴリ特徴語を見ながら、着目語を追加したり、削除したりすることができる。

【 0 1 5 4 】

例えば、図 1 8 による特徴語の変化や、C 社の相違特徴語を受けて、新たに単語“マイニング”を示す相違特徴語 $T d i f 3$ を着目語に追加する。これを受けて、特徴語抽出装置は、提示する関連カテゴリを変化させる。

20

【 0 1 5 5 】

これにより、ユーザは、着目語を切り替えながら関連するカテゴリを概観することで、それまで気づいていなかった着目語とカテゴリの関連を発見することができる。先行技術調査であれば、着目している技術を扱っている意外な企業や、ある企業はユーザが認識しているよりも早い年代から着目している技術に関する特許を出願しているといったことを発見する手がかりとなる。

【 0 1 5 6 】

また、図 1 5 ~ 図 1 9 を用いて述べたように、特徴語抽出装置による「特徴語の提示」と「関連カテゴリ提示」、ユーザによる「比較対象の指定」と「着目語の選択」、というプロセスを繰り返すことで、分析する対象や特徴語を明確化するとともに、それまでユーザが意識していなかったキーワードや、分析対象を発見することができる。特許調査においては、新たに着目すべき技術や、注意すべき競合他社を発見する糸口となる。また、比較対象や着目語を利用することで、適切な比較範囲に対する、適切な単語による、先行技術調査を実現することができる。

30

【 0 1 5 7 】

例えば図 2 0 に示す如き、「時系列 × 企業」の分析軸によるクロス分析の画面 G 1 0 は、図 2 1 及び図 2 2 に示すように、ある企業に対する「時系列 × 特徴語」の分析軸によるクロス分析の画面 G 2 0、ある特徴語に対する「時系列 × 企業」の分析軸によるクロス分析の画面 G 3 0 などのように、任意の分析軸の画面に適用して適切な比較範囲と適切な単語による分析・調査を実現することができる。

40

【 0 1 5 8 】

また例えば、ある特徴語に対する「時系列 × 企業」のクロス分析の画面 G 3 0 は、図 2 3 に示す如き、ある特徴語に対する「時系列 × 企業」のグラフ表示の画面 G 3 1 に表示形式を変更することができる。なお、表示形式を変更できることは、他のクロス分析の画面 G 1 0 , G 2 0 でも同様である。

【 0 1 5 9 】

上述したように本実施形態によれば、カテゴリ ID に関連付けられた文書 ID の個数を含むカテゴリ個数データの提示中に、選択を受け付けたカテゴリ個数データのカテゴリ ID に関連したカテゴリ特徴語のうち、特徴度が上位のカテゴリ特徴語における文書特徴語

50

をカテゴリ特徴語として提示する。

【0160】

また、カテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象からなる比較対象集合の選択を受け付けると、選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、各比較対象における当該文書特徴語の相違性を表す度合いの相違特徴度を算出し、相違特徴度が上位の文書特徴語をカテゴリ相違特徴語として提示する。

【0161】

このように、選択したカテゴリ個数データに応じて、カテゴリ特徴語やカテゴリ相違特徴語といった特徴語を提示する構成により、事前に定義された属性に限らずに分析軸の候補として特徴語を提示できると共に、着目すべき比較範囲の絞り込みや、着目すべき特徴語の参照を支援できる。

10

【0162】

補足すると、カテゴリ特徴語を提示することにより、ユーザが各カテゴリに対するカテゴリ特徴語を確認して、文書集合の全体像や、個々のカテゴリの内容を効率よく把握できる。

【0163】

また、複数の比較対象間の相違特徴語を提示する構成により、ユーザは着目している任意の範囲における比較対象の相違点を把握することができる。さらに、比較対象集合を絞り込めば、各比較対象における相違点をより詳細に把握できる。一方、比較対象集合を拡大すれば、マクロな視点で相違点を把握でき、全体的な内容の理解を深めることができる。以上により、ユーザは各特徴語を参照しながら、文書集合について内容の理解を進め、分析すべき範囲や、着目すべきキーワードを明確化することができる。

20

【0164】

また、本実施形態によれば、カテゴリ個数データの提示中、複数個のカテゴリ個数データの各カテゴリIDの和集合である複数の比較対象からなる比較対象集合の選択を受け付けると、選択を受け付けた比較対象集合内の各カテゴリIDに関連付けられた各文書IDに関連した文書特徴語毎に、比較対象集合における当該文書特徴語の共通性を表す度合いの共通特徴度を算出し、共通特徴度が上位の文書特徴語をカテゴリ共通特徴語として提示する構成により、ユーザは自身が任意の着目している範囲における文書の共通点を把握でき、文書集合に対する理解がさらに深められ、分析すべき範囲やキーワードをより明確に捉えることができる。

30

【0165】

さらに、本実施形態によれば、例えば、カテゴリ相違特徴語の提示中、複数のカテゴリ相違特徴語からなる着目語集合の選択を受け付けると、選択を受け付けた着目語集合と、カテゴリ特徴語記憶部30内のカテゴリ特徴語との関連度を算出し、関連度の高いカテゴリ特徴語に関連付けられたカテゴリIDに関連したカテゴリ個数データを強調表示する構成により、ユーザは、着目語の選択に対して提示されるカテゴリを概観することで、自身が着目しているキーワードに関連しているカテゴリを把握でき、それまで気づいていなかった分析対象を発見できる。従って、ユーザは、分析したい事項について、適切な分析対象を把握でき、より精度の高い分析が可能となる。

40

【0166】

このように、ユーザは、相違特徴語・共通特徴語と関連カテゴリを確認しながら、比較対象の選択（絞り込みや拡大）と着目語の選択を繰り返すことで、分析する範囲や着目する特徴語を明確化することができる。これによって、複数の文書集合に対して、ユーザは漏れなく、無駄なく、目的にあった、内容把握や比較調査を効率的に行うことができる。

【0167】

なお、上記実施形態に記載した手法は、コンピュータに実行させることのできるプログラムとして、磁気ディスク（フロッピー（登録商標）ディスク、ハードディスクなど）、光ディスク（CD-ROM、DVDなど）、光磁気ディスク（MO）、半導体メモリなど

50

の記憶媒体に格納して頒布することもできる。

【0168】

また、この記憶媒体としては、プログラムを記憶でき、かつコンピュータが読み取り可能な記憶媒体であれば、その記憶形式は何れの形態であってもよい。

【0169】

また、記憶媒体からコンピュータにインストールされたプログラムの指示に基づきコンピュータ上で稼働しているOS（オペレーティングシステム）や、データベース管理ソフト、ネットワークソフト等のMW（ミドルウェア）等が上記実施形態を実現するための各処理の一部を実行してもよい。

【0170】

さらに、本発明における記憶媒体は、コンピュータと独立した媒体に限らず、LANやインターネット等により伝送されたプログラムをダウンロードして記憶または一時記憶した記憶媒体も含まれる。

【0171】

また、記憶媒体は1つに限らず、複数の媒体から上記実施形態における処理が実行される場合も本発明における記憶媒体に含まれ、媒体構成は何れの構成であってもよい。

【0172】

尚、本発明におけるコンピュータは、記憶媒体に記憶されたプログラムに基づき、上記実施形態における各処理を実行するものであって、パソコン等の1つからなる装置、複数の装置がネットワーク接続されたシステム等の何れの構成であってもよい。

【0173】

また、本発明におけるコンピュータとは、パソコンに限らず、情報処理機器に含まれる演算処理装置、マイコン等も含み、プログラムによって本発明の機能を実現することが可能な機器、装置を総称している。

【0174】

なお、本願発明は、上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態に開示されている複数の構成要素の適宜な組合せにより種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。更に、異なる実施形態に亘る構成要素を適宜組合せてもよい。

【符号の説明】

【0175】

10...文書記憶部、20...カテゴリ記憶部、30...特徴語記憶部、40...特徴語抽出部、41...文書特徴語抽出部、42...カテゴリ特徴語抽出部、43...カテゴリ共通特徴語抽出部、44...カテゴリ相違特徴語抽出部、50...ユーザ操作・提示部、51...画面提示部、52...カテゴリ特徴語提示部、53...相違・共通特徴語提示部、54...関連カテゴリ提示部。

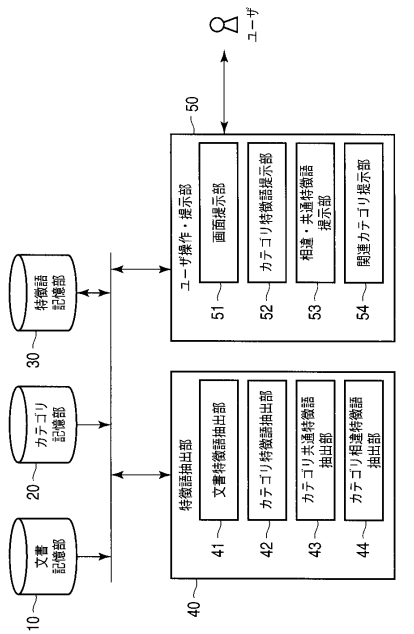
10

20

30

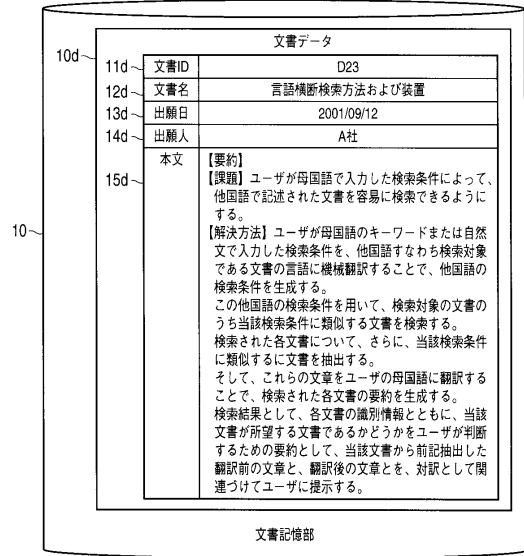
【図 1】

図 1



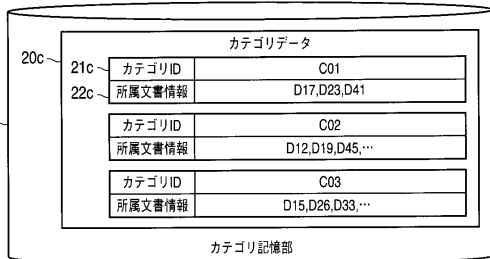
【図 2】

図 2



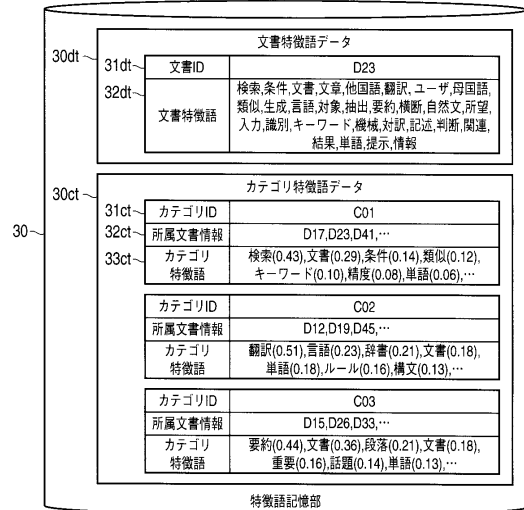
【図 3】

図 3



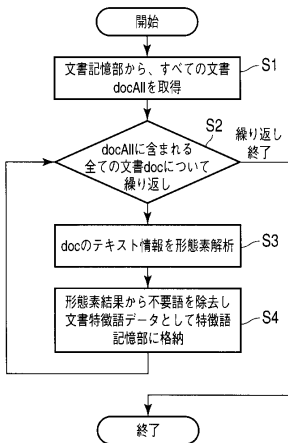
【図 4】

図 4



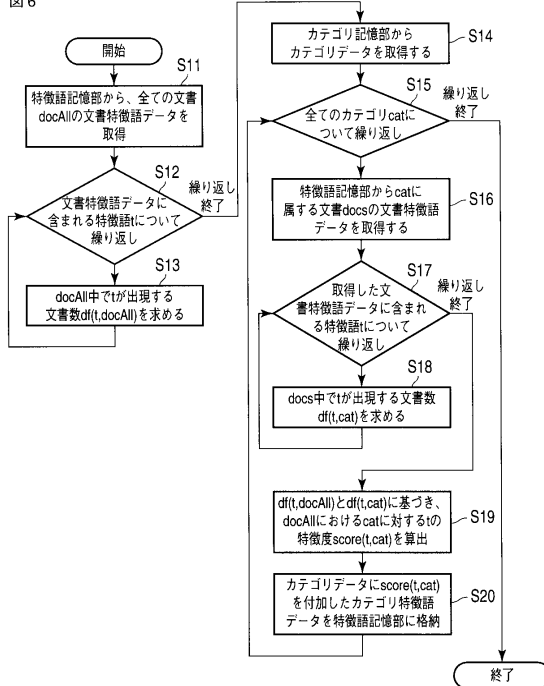
【図5】

図5



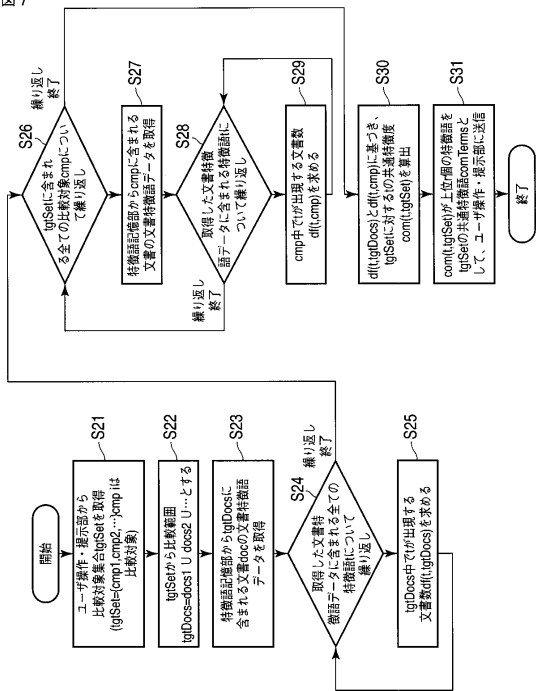
【図6】

図6



【図7】

図7



【図8】

図8

		出願年				
		2004/	2005/	2006/	2007/	2008/
企業	A社	65	50	69	75	72
	B社	10	21	45	53	35
	C社	25	10	24	36	20
	D社	16	6	12	5	10
	E社	52	50	32	54	43
	F社	75	63	80	69	78

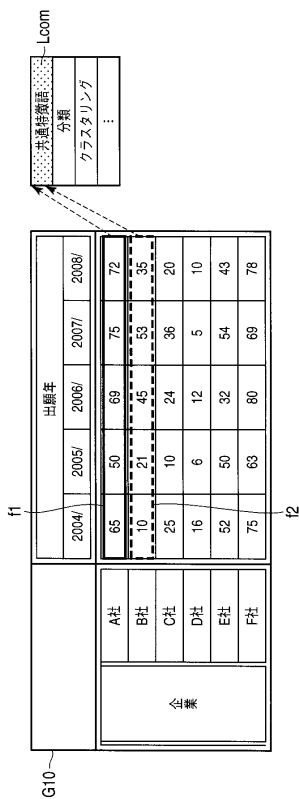
【図9】

図9

		出願年				
		2004/	2005/	2006/	2007/	2008/
企業	A社	65	50	69	75	72
	B社	10	21	45	53	35
	C社	25	10	24	36	20
	D社	16	6	12	5	10
	E社	52	50	32	54	43
	F社	75	63	80	69	78

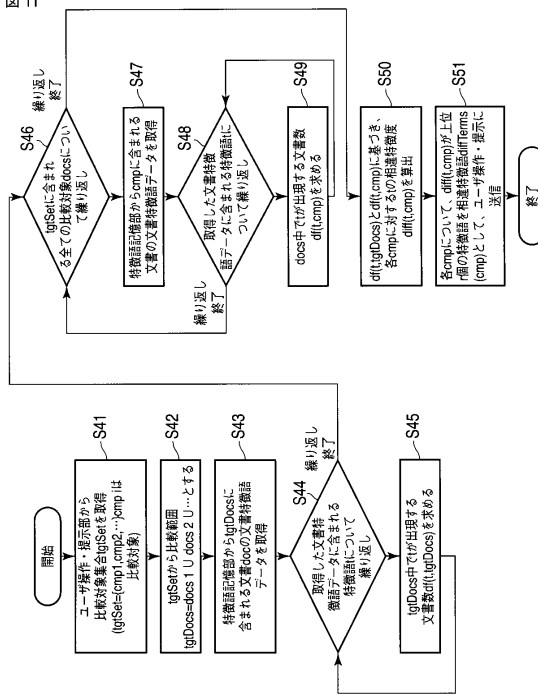
【図 10】

図 10



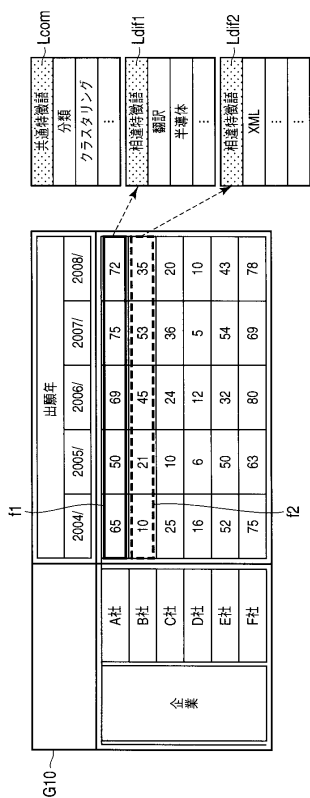
【図 11】

図 11



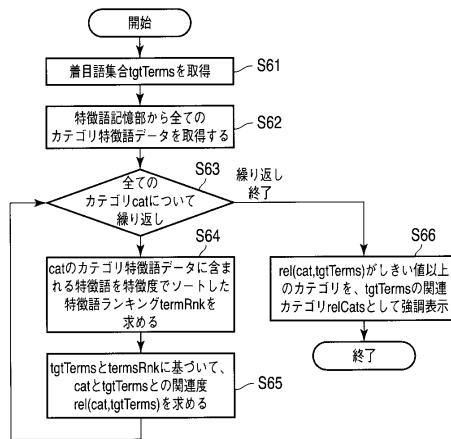
【図 12】

図 12

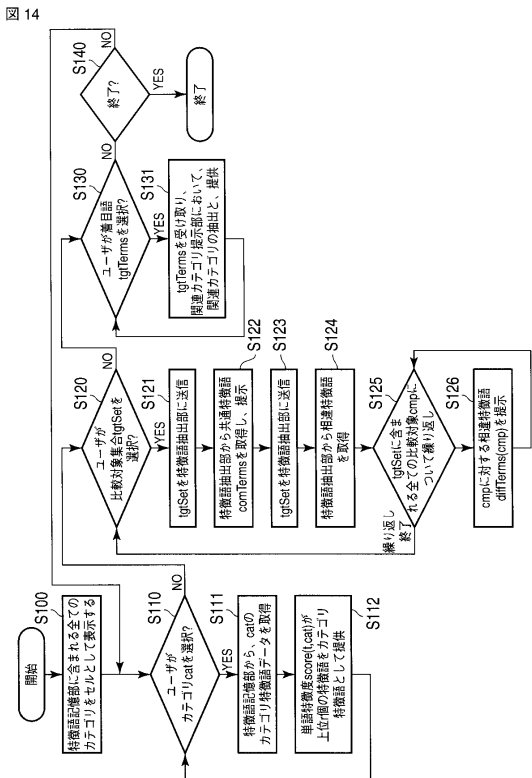


【図 13】

図 13



【 図 14 】



【 図 15 】

企業	A社	65	50	69	75	72
	B社	10	21	45	53	35
	C社	25	10	24	36	20
	D社	16	6	12	5	10
	E社	52	50	32	54	43
	F社	75	63	80	69	78
	出願年	2004/	2005/	2006/	2007/	2008/

【 図 16 】

図 16

企業	A社	65	50	69	75	72
	B社	10	21	45	53	35
	C社	25	10	24	36	20
	D社	16	6	12	5	10
	E社	52	50	32	54	43
	F社	75	63	80	69	78
	出願年	2004/	2005/	2006/	2007/	2008/

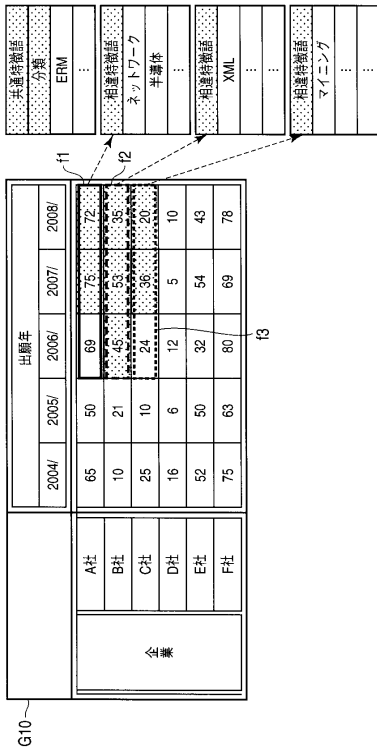
【 図 17 】

図 17

企業	A社	65	50	69	75	72
	B社	10	21	45	53	35
	C社	25	10	24	36	20
	D社	16	6	12	5	10
	E社	52	50	32	54	43
	F社	75	63	80	69	78
	出願年	2004/	2005/	2006/	2007/	2008/

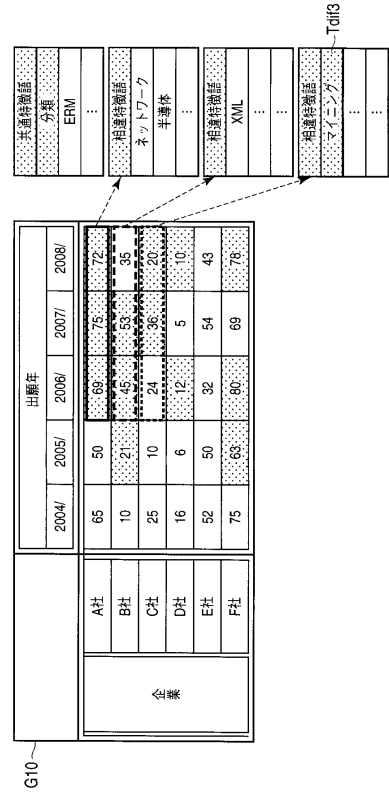
【図 18】

図 18



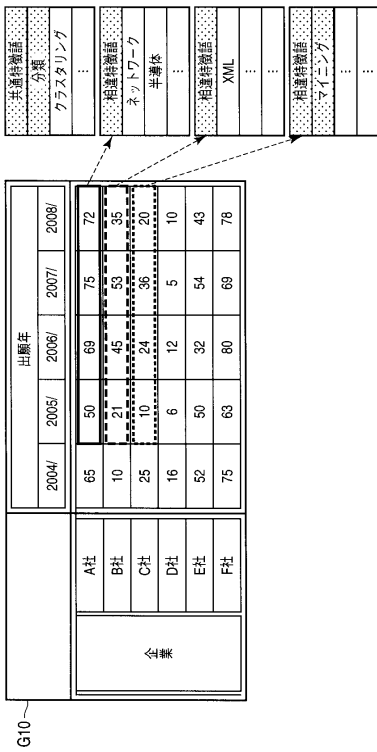
【図 19】

図 19



【図 20】

図 20



【図 21】

図 21

G20

A社		出願年			
		2005	2006	2007	2008
特許分類	分類	10	25	50	31
	クラスターリング	21	45	53	35
	XML	10	24	26	20
	マイニング	123	101	201	93

【図 22】

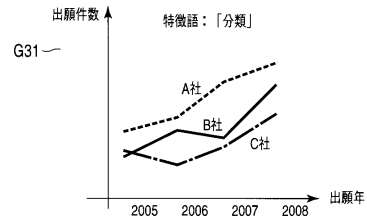
図 22

G30

「分類」		出願年			
		2005	2006	2007	2008
業企	A社	10	13	21	26
	B社	4	10	9	23
	C社	5	4	7	12

【図 23】

図 23



フロントページの続き

- (74)代理人 100075672
弁理士 峰 隆司
- (74)代理人 100095441
弁理士 白根 俊郎
- (74)代理人 100084618
弁理士 村松 貞男
- (74)代理人 100103034
弁理士 野河 信久
- (74)代理人 100119976
弁理士 幸長 保次郎
- (74)代理人 100153051
弁理士 河野 直樹
- (74)代理人 100140176
弁理士 砂川 克
- (74)代理人 100101812
弁理士 勝村 紘
- (74)代理人 100124394
弁理士 佐藤 立志
- (74)代理人 100112807
弁理士 岡田 貴志
- (74)代理人 100111073
弁理士 堀内 美保子
- (74)代理人 100134290
弁理士 竹内 将訓
- (74)代理人 100127144
弁理士 市原 卓三
- (74)代理人 100141933
弁理士 山下 元
- (72)発明者 岩崎 秀樹
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内
- (72)発明者 後藤 和之
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内
- (72)発明者 松本 茂
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内
- (72)発明者 平 博司
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内
- (72)発明者 宮部 泰成
東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

審査官 齊藤 貴孝

- (56)参考文献 特開2001-092825(JP,A)
特開2003-345810(JP,A)
国際公開第2006/115260(WO,A1)
特開2009-288999(JP,A)
特開2009-294938(JP,A)
特開2009-294939(JP,A)
特表2008-524712(JP,A)

特開2006-215675(JP,A)
特開2003-345811(JP,A)
国際公開第2007/069663(WO,A1)
特開2007-004233(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

G06Q 10/00