

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6800825号
(P6800825)

(45) 発行日 令和2年12月16日(2020.12.16)

(24) 登録日 令和2年11月27日(2020.11.27)

(51) Int.Cl. F I
G O 6 F 16/35 (2019.01) G O 6 F 16/35
G O 6 F 16/34 (2019.01) G O 6 F 16/34

請求項の数 13 (全 26 頁)

<p>(21) 出願番号 特願2017-192750 (P2017-192750) (22) 出願日 平成29年10月2日 (2017.10.2) (65) 公開番号 特開2019-67191 (P2019-67191A) (43) 公開日 平成31年4月25日 (2019.4.25) 審査請求日 令和1年8月13日 (2019.8.13)</p>	<p>(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 (74) 代理人 110002147 特許業務法人酒井国際特許事務所 (72) 発明者 布目 光生 東京都港区芝浦一丁目1番1号 株式会社 東芝内 (72) 発明者 山崎 智弘 東京都港区芝浦一丁目1番1号 株式会社 東芝内 審査官 鹿野 博嗣</p>
---	--

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理方法およびプログラム

(57) 【特許請求の範囲】

【請求項1】

文書群に含まれる複数のキーフレーズを含むキーフレーズ群を階層クラスタリングするクラスタリング部と、

前記キーフレーズ群を複数の候補クラスタに分割する分割部と、

前記文書群を分類するための予め定められた複数の項目のうちの何れか1つ項目の選択操作を受け付ける項目選択部と、

前記複数の候補クラスタのそれぞれについて、選択された項目に対する有用性を表すスコアを算出するスコア算出部と、

前記複数の候補クラスタのうち、前記スコアが所定の順位の候補クラスタを参照クラスタとして決定する決定部と、

前記参照クラスタを複数のサブクラスタに分割するサブクラスタ生成部と、

選択された項目の下位層の予め定められた複数のサブ項目を抽出するサブ項目抽出部と

、
 前記複数のサブ項目のそれぞれ毎且つ前記複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、展開画像の提示を制御する展開画像制御部と、

を備える情報処理装置。

【請求項2】

前記決定部は、前記複数の候補クラスタのうち前記スコアが最も高い1つの候補クラスタを前記参照クラスタとして決定する

10

20

請求項 1 に記載の情報処理装置。

【請求項 3】

前記スコアは、対応する候補クラスタにおける選択された項目に分類される文書の情報量を表す

請求項 1 または 2 に記載の情報処理装置。

【請求項 4】

前記スコアは、前記複数のサブ項目のそれぞれ毎且つ対応する候補クラスタを分割した複数のサブクラスタのそれぞれ毎の文書の情報量の分散を表す

請求項 1 または 2 に記載の情報処理装置。

【請求項 5】

前記スコアは、対応する候補クラスタをユーザが過去に操作により選択した頻度または割合を表す

請求項 1 または 2 に記載の情報処理装置。

【請求項 6】

前記スコア算出部は、第 1 パラメータと、第 2 パラメータと、第 3 パラメータとを合計して前記スコアを算出し、

前記第 1 パラメータは、対応する候補クラスタにおける選択された項目に分類される文書の情報量に応じた値を表し、

前記第 2 パラメータは、前記複数のサブ項目のそれぞれ毎且つ対応する候補クラスタを分割した複数のサブクラスタのそれぞれ毎の文書の情報量の分散に応じた値を表し、

前記第 3 パラメータは、対応する候補クラスタをユーザが過去に操作により選択した頻度または割合に応じた値を表す

請求項 1 または 2 に記載の情報処理装置。

【請求項 7】

前記複数の項目のそれぞれ毎且つ前記複数の候補クラスタのそれぞれ毎の文書の情報量を表すように、初期画像の提示を制御する初期画像制御部をさらに備える

請求項 1 から 6 の何れか 1 項に記載の情報処理装置。

【請求項 8】

前記複数のサブ項目のうち何れか 1 つのサブ項目が選択された場合、前記スコア算出部、前記決定部、前記サブ項目抽出部および前記展開画像制御部は、前記複数のサブ項目を前記複数の項目と置き換え、前記複数のサブクラスタを前記複数の候補クラスタと置き換えて、再度処理を実行して展開画像を提示する

請求項 1 から 7 の何れか 1 項に記載の情報処理装置。

【請求項 9】

前記展開画像制御部は、前記展開画像とともに、前記複数の候補クラスタのうちの何れか 1 つを選択させるためのメニュー画像を提示し、

前記メニュー画像により何れか 1 つの新たな候補クラスタが選択された場合、

前記サブクラスタ生成部は、前記新たな候補クラスタを前記参照クラスタとして、新たな複数のサブクラスタに分割し、

前記展開画像制御部は、前記複数のサブ項目のそれぞれ毎且つ新たな複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、新たな展開画像の提示を制御する

請求項 1 から 8 の何れか 1 項に記載の情報処理装置。

【請求項 10】

文書抽出部と、選択クラスタリング部と、選択分割部とをさらに備え、

前記項目選択部は、前記文書群を分類するための予め定められた複数の第 1 観点項目のうち何れか 1 つの第 1 観点項目、および、前記文書群を分類するための予め定められた複数の第 2 観点項目のうち何れか 1 つの第 2 観点項目の選択操作を受け付け、

前記文書抽出部は、前記文書群から、選択された第 1 観点項目および選択された第 2 観点項目の両者に分類される複数の文書を含む選択文書群を抽出し、

前記選択クラスタリング部は、前記選択文書群に含まれる複数のキーフレーズを含む選

10

20

30

40

50

択キーフレーズ群を階層クラスタリングし、

前記選択分割部は、前記選択キーフレーズ群を複数の候補クラスタに分割し、

前記スコア算出部は、前記複数の候補クラスタのそれぞれについて、選択された第1観点項目および第2観点項目に対する有用性を表す前記スコアを算出し、

前記決定部は、前記複数の候補クラスタのうち、前記スコアが所定の順位の2つの候補クラスタを第1参照クラスタおよび第2参照クラスタとして決定し、

前記サブクラスタ生成部は、前記第1参照クラスタを複数の第1サブクラスタに分割し、前記第2参照クラスタを複数の第2サブクラスタに分割し、

前記展開画像制御部は、前記複数の第1サブクラスタのそれぞれ毎且つ前記複数の第2サブクラスタのそれぞれ毎の文書の情報量を表すように、クラスタ展開画像の提示を制御する

10

請求項1に記載の情報処理装置。

【請求項11】

前記複数の第1観点項目のそれぞれ毎且つ前記複数の第2観点項目のそれぞれ毎の文書の情報量を表すように、項目選択画像の提示を制御する初期画像制御部をさらに備える

請求項10に記載の情報処理装置。

【請求項12】

情報処理装置により実行される情報処理方法であって、

クラスタリング部が、文書群に含まれる複数のキーフレーズを含むキーフレーズ群を階層クラスタリングし、

20

分割部が、前記キーフレーズ群を複数の候補クラスタに分割し、

項目選択部が、前記文書群を分類するための予め定められた複数の項目のうちの何れか1つ項目の選択操作を受け付け、

スコア算出部が、前記複数の候補クラスタのそれぞれについて、選択された項目に対する有用性を表すスコアを算出し、

決定部が、前記複数の候補クラスタのうち、前記スコアが所定の順位の候補クラスタを参照クラスタとして決定し、

サブクラスタ生成部が、前記参照クラスタを複数のサブクラスタに分割し、

抽出部が、選択された項目の下位層の予め定められた複数のサブ項目を抽出し、

展開画像制御部が、前記複数のサブ項目のそれぞれ毎且つ前記複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、展開画像の提示を制御する

30

情報処理方法。

【請求項13】

情報処理装置で実行されるプログラムであって、

前記情報処理装置を、

文書群に含まれる複数のキーフレーズを含むキーフレーズ群を階層クラスタリングするクラスタリング部と、

前記キーフレーズ群を複数の候補クラスタに分割する分割部と、

前記文書群を分類するための予め定められた複数の項目のうちの何れか1つ項目の選択操作を受け付ける項目選択部と、

40

前記複数の候補クラスタのそれぞれについて、選択された項目に対する有用性を表すスコアを算出するスコア算出部と、

前記複数の候補クラスタのうち、前記スコアが所定の順位の候補クラスタを参照クラスタとして決定する決定部と、

前記参照クラスタを複数のサブクラスタに分割するサブクラスタ生成部と、

選択された項目の下位層の予め定められた複数のサブ項目を抽出するサブ項目抽出部と

、前記複数のサブ項目のそれぞれ毎且つ前記複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、展開画像の提示を制御する展開画像制御部と

して機能させるプログラム。

50

【発明の詳細な説明】**【技術分野】****【0001】**

本発明の実施形態は、情報処理装置、情報処理方法およびプログラムに関する。

【背景技術】**【0002】**

業務において用いた文書をデータベースに記録し、その文書を他の業務において再利用させる文書管理システムが知られている。また、文書管理システムにより管理している文書群から、データマイニングおよびテキストマイニング等により知識を抽出し、抽出した知識を業務の分析および改善に役立てることも行われている。

10

【0003】

また、文書管理システムにより管理している文書群から目的の文書を検索する手段として、キーワード（単語および単語列等）検索およびファセット検索が知られている。ファセット検索は、文書を分類するための複数の項目および階層構造を予め定義しておき、ユーザに上位層から下位層に向かい順次に項目を選択させて、文書を絞り込む方法である。

【0004】

文書管理システムにより管理している文書群の特徴をユーザに参照させる手段も、様々な方法が提案されている。例えば、文書群の特徴をユーザに参照させる手段として、OLAP（Online Analytical Processing）機能が知られている。OLAP機能は、文書群の全体の特徴を俯瞰的に参照させたり、全体を表す情報から詳細を表す情報にドリルダウンさせながら文書群の特徴を参照させたりする方法である。また、文書群の特徴をユーザに参照させる手段として、ヒートマップも知られている。ヒートマップは、異なる2つの観点から分類した情報のそれぞれの特徴を2つの軸を有するマップ上に表す方法である。

20

【0005】

ところで、ファセット検索では、予め項目の構造を定義しておかなければならない。例えば、しかし、項目の構造の設計、および、対応するデータベースの設計には、非常に多くのコストがかかる。また、文書管理システムの運用が進んだ段階で、新しい観点で文書を検索および文書群の特徴を参照したいというニーズが生じても、既に定義した項目の階層構造およびデータベースの構造を変えることは難しかった。

30

【0006】

一方で、クラスタリングにより分類の項目を自動生成するといった手法も知られている。この方法であれば、事前に、項目の構造を設計する必要はない。しかし、クラスタリングにより分類の項目を自動生成する方法では、実際に利用できる項目には大きな限定がかかる。例えば、クラスタリングにより分類の項目を自動生成する方法では、数量表現、色および形等の離散属性、および、ソースコードのパッケージ名等の、階層情報および構造が文書内に記述されている項目にしか分類ができなかった。

【先行技術文献】**【特許文献】**

40

【0007】

【特許文献1】特開2017-068534号公報

【発明の概要】**【発明が解決しようとする課題】****【0008】**

発明が解決しようとする課題は、文書群を小さいコストで適切に分類した情報を提示することにある。

【課題を解決するための手段】**【0009】**

実施形態に係る情報処理装置は、クラスタリング部と、分割部と、項目選択部と、スコ

50

ア算出部と、決定部と、サブクラスタ生成部と、サブ項目抽出部と、展開画像制御部と、を備える。前記クラスタリング部は、文書群に含まれる複数のキーフレーズを含むキーフレーズ群を階層クラスタリングする。前記分割部は、前記キーフレーズ群を複数の候補クラスタに分割する。前記項目選択部は、前記文書群を分類するための予め定められた複数の項目のうちの何れか1つ項目の選択操作を受け付ける。前記スコア算出部は、前記複数の候補クラスタのそれぞれについて、選択された項目に対する有用性を表すスコアを算出する。前記決定部は、前記複数の候補クラスタのうち、前記スコアが所定の順位の候補クラスタを参照クラスタとして決定する。前記サブクラスタ生成部は、前記参照クラスタを複数のサブクラスタに分割する。前記サブ項目抽出部は、選択された項目の下位層の予め定められた複数のサブ項目を抽出する。前記展開画像制御部は、前記複数のサブ項目のそれぞれ毎且つ前記複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、展開画像の提示を制御する。

10

【図面の簡単な説明】

【0010】

【図1】第1実施形態に係る文書管理システムの構成図。

【図2】第1実施形態に係る事前処理部の構成図。

【図3】第1実施形態に係る第1画像制御部および第2画像制御部の構成図。

【図4】情報処理装置の処理の流れを示すフローチャート。

【図5】事前処理の詳細な処理の流れを示すフローチャート。

【図6】キーフレーズ群の構造を表すデンドログラムの一例を示す図。

20

【図7】デンドログラムおよび候補クラスタの一例を示す図。

【図8】複数の項目の階層構造を示す図。

【図9】初期画像の表示処理の詳細な処理の流れを示すフローチャート。

【図10】初期画像の一例を示す図。

【図11】展開画像の表示処理の詳細な処理の流れを示すフローチャート。

【図12】初期画像および展開画像の一例を示す図。

【図13】初期画像、展開画像および新たな展開画像の一例を示す図。

【図14】メニュー画像が追加された展開画像を示す図。

【図15】複数の候補クラスタのそれぞれのスコアを示す図。

【図16】第2実施形態に係る第1画像制御部および第2画像制御部の構成図。

30

【図17】複数の第1観点項目および複数の第2観点項目の一例を示す図。

【図18】項目選択画像の一例を示す図。

【図19】選択キーフレーズ群の構造を表すデンドログラムの一例を示す図。

【図20】デンドログラムおよびクラスタ展開画像の一例を示す図。

【図21】列選択および行選択をした場合の項目選択画像を示す図。

【発明を実施するための形態】

【0011】

以下、図面を参照しながら本実施形態に係る文書管理システム10について説明する。なお、以下の実施形態では、同一の参照符号を付した部分は略同一の構成および動作をするので、相違点を除き重複する説明を適宜省略する。

40

【0012】

(第1実施形態)

図1は、第1実施形態に係る文書管理システム10の構成を示す図である。文書管理システム10は、業務等で作成された複数の文書を含む文書群を管理する。また、文書管理システム10は、ユーザの操作に応じて、文書群を分類し、分類された複数の文書毎の情報量を表示する。

【0013】

文書は、コンピュータにより情報内容を検索することが可能であれば、どのようなデータであってもよい。例えば、文書は、テキストを含むデータであってもよいし、プログラムコードを含むデータであってもよい。文書のファイル形式は、文書管理システム10に

50

より取り扱いが可能であれば、どのようなものであってもよい。

【0014】

また、複数の文書の情報量は、複数の文書の数であってもよいし、複数の文書に含まれる文字の数であってもよいし、複数の文書の合計のデータ量であってもよい。

【0015】

また、文書管理システム10は、1つの文書から1または複数のキーフレーズを抽出する。キーフレーズは、その文書に含まれる情報およびその文書に関連する情報等の、その文書の特徴を表す情報である。キーフレーズは、例えば、1つの単語であってもよいし、複数の単語が並んだ単語列であってもよい。また、キーフレーズは、プログラムコード中のコード列であってもよい。

10

【0016】

文書管理システム10は、表示装置12と、入力装置14と、記憶装置16と、情報処理装置20とを備える。

【0017】

表示装置12は、画像を表示することにより、画像をユーザに提示する。表示装置12は、情報処理装置20により生成された画像を受け取り、受け取った画像を表示する。表示装置12は、例えば、液晶表示器等の表示デバイスである。

【0018】

入力装置14は、ユーザからの指示および操作を受け付ける。入力装置14は、例えば、マウスまたはトラックボール等のポインティングデバイス、あるいはキーボード等の入力デバイスである。

20

【0019】

記憶装置16は、情報処理装置20からデータを受け取り、受け取ったデータを記憶する。また、記憶装置16は、記憶しているデータが情報処理装置20により読み出される。記憶装置16は、フラッシュメモリ等の半導体メモリ素子、ハードディスク、光ディスク等である。記憶装置16は、ネットワークを介して情報処理装置20と接続可能なサーバ装置であってもよい。

【0020】

情報処理装置20は、例えば、専用または汎用コンピュータである。情報処理装置20は、PC、あるいは、情報を保存および管理するサーバに含まれるコンピュータであってもよい。情報処理装置20は、一台の装置により実現されてもよいし、連携して動作する複数台の装置により実現されてもよい。また、情報処理装置20は、ネットワーク上に実現される仮想的な装置(例えばクラウド)等であってもよい。

30

【0021】

情報処理装置20は、表示装置12を制御して、表示装置12に画像を表示させる。また、情報処理装置20は、入力装置14から情報を受け取って、ユーザから与えられた指示内容および操作内容を識別する。また、情報処理装置20は、データを記憶装置16に書き込み、記憶装置16に記憶されたデータを読み出す。

【0022】

情報処理装置20は、通信部22と、記憶回路24と、処理回路30とを有する。表示装置12、入力装置14、記憶装置16、通信部22、記憶回路24および処理回路30は、バスを介して接続される。

40

【0023】

通信部22は、有線または無線で接続された外部装置と情報の入出力を行うインターフェースである。通信部22は、ネットワークに接続して通信を行ってもよい。

【0024】

記憶回路24は、RAM(Random Access Memory)およびROM(Read Only Memory)である。記憶回路24は、起動用プログラムを読み出すスタートプログラムが記憶されている。また、記憶回路24は、処理回路30の作業領域として機能する。

50

【 0 0 2 5 】

処理回路 30 は、1 または複数のプロセッサを含む。処理回路 30 は、情報処理を実行し、プログラムを読み出して記憶回路 24 に展開して実行し、各部を制御してデータの入出力を行ったり、データの加工を行ったりする。プロセッサは、例えば、CPU (Central Processing Unit) である。プロセッサは、CPU に限らず、プログラムを実行する他の種類のデータ処理デバイスまたは専用の処理デバイスであってもよい。

【 0 0 2 6 】

このようなハードウェア構成の文書管理システム 10 は、記憶装置 16 が、文書記憶部 42、クラスタ記憶部 44 および項目記憶部 46 として機能する。また、このようなハードウェア構成の文書管理システム 10 は、処理回路 30 が、プログラムを実行することにより、事前処理部 32、第 1 画像制御部 34 および第 2 画像制御部 36 として機能する。

10

【 0 0 2 7 】

図 2 は、第 1 実施形態に係る事前処理部 32 の構成を文書記憶部 42、クラスタ記憶部 44 および項目記憶部 46 とともに示す図である。

【 0 0 2 8 】

事前処理部 32 は、文書取得部 52 と、キーフレーズ生成部 54 と、文書登録部 56 と、クラスタリング部 58 と、分割部 60 と、項目取得部 62 と、項目登録部 64 とを有する。

【 0 0 2 9 】

文書取得部 52 は、他の装置から文書を取得する。キーフレーズ生成部 54 は、文書取得部 52 により取得された文書に対して形態素解析および複合語抽出処理等を行って、取得した文書に対する 1 または複数のキーフレーズを生成する。文書登録部 56 は、文書取得部 52 により取得された文書と、キーフレーズ生成部 54 により生成された 1 または複数のキーフレーズとを対応付けて、文書記憶部 42 に記憶させる。

20

【 0 0 3 0 】

文書取得部 52、キーフレーズ生成部 54 および文書登録部 56 は、複数の文書のそれぞれ毎に、これらの処理を実行する。これにより、文書記憶部 42 は、複数の文書を含む文書群を記憶することができる。この文書群は、データベース化されており、任意のキーフレーズを指定することにより、指定されたキーフレーズに対応付けられた 1 または複数の文書を抽出することができる。

30

【 0 0 3 1 】

クラスタリング部 58 は、文書記憶部 42 から、文書群に含まれる複数のキーフレーズを含むキーフレーズ群を取得する。クラスタリング部 58 は、取得したキーフレーズ群を階層クラスタリングする。例えば、クラスタリング部 58 は、キーフレーズ群に含まれる複数のキーフレーズを複数のクラスタにクラスタリングする。さらに、クラスタリング部 58 は、それぞれのキーフレーズをベクトル化する。そして、クラスタリング部 58 は、クラスタ中心とのベクトル距離に応じて、キーフレーズの類似度を算出する。さらに、クラスタリング部 58 は、それぞれのクラスタ内で同様の処理を繰り返して、階層化した複数のクラスタを生成する。

40

【 0 0 3 2 】

クラスタリング部 58 は、階層化した複数のクラスタのそれぞれにラベルを付与してもよい。例えば、クラスタリング部 58 は、クラスタ中心に近いキーフレーズをラベルとしてもよい。クラスタリング部 58 は、階層クラスタリングしたキーフレーズ群をクラスタ記憶部 44 に記憶させる。

【 0 0 3 3 】

分割部 60 は、階層クラスタリングされたキーフレーズ群をクラスタ記憶部 44 から読み出し、読み出したキーフレーズ群を複数の候補クラスタに分割する。例えば、分割部 60 は、階層クラスタリングされたキーフレーズ群を表すデンドログラムを描き、描いたデンドログラムにおける、所定個 (例えば、4 個以上で最小) のクラスタに分割される高さ

50

を決定する。そして、分割部 6 0 は、決定した高さでデンドログラムを切断した場合に生成される複数の階層クラスタを、複数の候補クラスタとする。分割部 6 0 は、生成した複数の候補クラスタをクラスタ記憶部 4 4 に記憶させる。

【 0 0 3 4 】

項目取得部 6 2 は、文書群を分類するための予め定められた複数の項目を他の装置から取得する。項目取得部 6 2 は、ユーザにより入力された複数の項目を取得してもよい。複数の項目は、木構造により階層化されている。項目登録部 6 4 は、項目取得部 6 2 により取得された複数の項目を項目記憶部 4 6 に記憶させる。

【 0 0 3 5 】

図 3 は、第 1 実施形態に係る第 1 画像制御部 3 4 および第 2 画像制御部 3 6 の構成を文書記憶部 4 2、クラスタ記憶部 4 4 および項目記憶部 4 6 とともに示す図である。

10

【 0 0 3 6 】

第 1 画像制御部 3 4 は、開始受付部 7 2 と、第 1 算出部 7 4 と、初期画像制御部 7 6 とを有する。開始受付部 7 2 は、ユーザによる開始操作を、入力装置 1 4 から受け付ける。

【 0 0 3 7 】

第 1 算出部 7 4 は、開始受付部 7 2 が開始操作を受け付けると、複数の候補クラスタをクラスタ記憶部 4 4 から取得する。また、第 1 算出部 7 4 は、開始受付部 7 2 が開始操作を受け付けると、項目記憶部 4 6 から予め定められた複数の項目のうち最上位の複数の項目を取得する。

【 0 0 3 8 】

20

そして、第 1 算出部 7 4 は、文書記憶部 4 2 にアクセスして、予め定められた複数の項目のそれぞれ毎、且つ、複数の候補クラスタのそれぞれ毎の文書の情報量を算出する。すなわち、第 1 算出部 7 4 は、文書群を最上位の複数の項目に従って複数の初期文書群に分類する。そして、第 1 算出部 7 4 は、複数の初期文書群のそれぞれについて、複数の候補クラスタのそれぞれに分類される文書の情報量を算出する。例えば、最上位の複数の項目が 4 個、および、候補クラスタが 5 個の場合、第 1 算出部 7 4 は、 $4 \times 5 = 20$ 個の文書の情報量を算出する。

【 0 0 3 9 】

初期画像制御部 7 6 は、予め定められた複数の項目のそれぞれ毎、且つ、複数の候補クラスタのそれぞれ毎の文書の情報量を表すように、初期画像の提示を制御する。すなわち、初期画像制御部 7 6 は、文書群を最上位の複数の項目に従って分類した複数の初期文書群のそれぞれについて、複数の候補クラスタのそれぞれに分類される文書の情報量を表す初期画像を生成する。そして、初期画像制御部 7 6 は、生成した初期画像を表示装置 1 2 へ出力して、表示装置 1 2 に初期画像を表示させる。

30

【 0 0 4 0 】

第 2 画像制御部 3 6 は、項目選択部 7 8 と、スコア算出部 8 0 と、決定部 8 2 と、サブクラスタ生成部 8 4 と、サブ項目抽出部 8 6 と、第 2 算出部 8 8 と、展開画像制御部 9 0 とを有する。

【 0 0 4 1 】

項目選択部 7 8 は、ユーザによる、文書群を分類するための予め定められた複数の項目のうち何れか 1 つ項目の選択操作を、入力装置 1 4 から受け付ける。例えば、初期画像が表示された後、項目選択部 7 8 は、初期画像に情報量が表示された最上位の複数の項目のうち、何れか 1 つの項目の選択操作を受け付ける。

40

【 0 0 4 2 】

スコア算出部 8 0 は、複数の候補クラスタのそれぞれについて、選択された項目に対する有用性を表すスコアを算出する。なお、スコアについては、詳細を後述する。

【 0 0 4 3 】

決定部 8 2 は、複数の候補クラスタのうち、算出されたスコアが所定の順位の候補クラスタを、参照クラスタとして決定する。例えば、決定部 8 2 は、複数の候補クラスタのうち、有用性が最も高いスコアの候補クラスタを、参照クラスタとして決定する。

50

【 0 0 4 4 】

サブクラスタ生成部 8 4 は、参照クラスタを、複数のサブクラスタに分割する。例えば、サブクラスタ生成部 8 4 は、参照クラスタを所定個（例えば、4 個以上で最小）に分割して、複数のサブクラスタを生成する。

【 0 0 4 5 】

サブ項目抽出部 8 6 は、項目記憶部 4 6 にアクセスして、項目選択部 7 8 により選択された項目の下位層の予め定められた複数のサブ項目を抽出する。

【 0 0 4 6 】

第 2 算出部 8 8 は、サブクラスタ生成部 8 4 から、複数のサブクラスタを取得する。また、第 2 算出部 8 8 は、サブ項目抽出部 8 6 から予め定められた複数のサブ項目を取得する。

10

【 0 0 4 7 】

そして、第 2 算出部 8 8 は、文書記憶部 4 2 にアクセスして、予め定められた複数のサブ項目のそれぞれ毎、且つ、複数のサブクラスタのそれぞれ毎の文書の情報量を算出する。すなわち、第 2 算出部 8 8 は、選択された 1 つの最上位項目に分類される初期文書群を、さらに複数のサブ項目に従って複数のサブ文書群に分類する。そして、第 2 算出部 8 8 は、複数のサブ文書群のそれぞれについて、複数のサブクラスタのそれぞれに分類される文書の情報量を算出する。例えば、複数のサブ項目が 5 個、および、複数のサブクラスタが 6 個の場合、第 2 算出部 8 8 は、 $5 \times 6 = 30$ 個の文書の情報量を算出する。

【 0 0 4 8 】

20

展開画像制御部 9 0 は、予め定められた複数のサブ項目のそれぞれ毎、且つ、複数のサブクラスタのそれぞれ毎の文書の情報量を表すように、展開画像の提示を制御する。すなわち、展開画像制御部 9 0 は、複数のサブ文書群のそれぞれについて、複数のサブクラスタのそれぞれに分類される文書の情報量を表す展開画像を生成する。そして、展開画像制御部 9 0 は、生成した展開画像を表示装置 1 2 に出力して、表示装置 1 2 に展開画像を表示させる。

【 0 0 4 9 】

なお、展開画像が表示された後、項目選択部 7 8 は、展開画像に情報量が表示された複数のサブ項目のうちの、何れか 1 つのサブ項目の選択操作を受け付けてもよい。複数のサブ項目のうち何れか 1 つのサブ項目が選択された場合、スコア算出部 8 0、決定部 8 2、サブクラスタ生成部 8 4、サブ項目抽出部 8 6、第 2 算出部 8 8 および展開画像制御部 9 0 は、複数のサブ項目を複数の項目と置き換え、複数のサブクラスタを複数の候補クラスタと置き換えて、再度処理を実行して新たな展開画像の提示を制御する。

30

【 0 0 5 0 】

図 4 は、情報処理装置 2 0 の処理の流れを示すフローチャートである。まず、S 1 1 において、情報処理装置 2 0 は、事前処理を実行する。続いて、S 1 2 において、情報処理装置 2 0 は、初期画像の表示処理を実行する。そして、S 1 3 において、情報処理装置 2 0 は、展開画像の表示処理を実行する。以下、S 1 1、S 1 2 および S 1 3 の処理について詳細に説明する。

【 0 0 5 1 】

40

図 5 は、事前処理 (S 1 1) の詳細な処理の流れを示すフローチャートである。情報処理装置 2 0 は、S 1 1 の事前処理において、以下の S 2 1 から S 2 6 の処理を実行する。

【 0 0 5 2 】

S 2 1 において、情報処理装置 2 0 は、他の装置から文書を取得する。続いて、S 2 2 において、情報処理装置 2 0 は、取得した文書に対して形態素解析および複合語抽出処理等を行って、取得した文書に対する 1 または複数のキーフレーズを生成する。続いて、S 2 3 において、情報処理装置 2 0 は、取得した文書と、生成した 1 または複数のキーフレーズとを対応付けて、文書記憶部 4 2 に登録する。

【 0 0 5 3 】

情報処理装置 2 0 は、S 2 1 から S 2 3 までの処理を、複数の文書のそれぞれに対して

50

実行する。これにより、文書記憶部 4 2 は、複数の文書を含む文書群を記憶することができる。この文書群は、データベース化されている。情報処理装置 2 0 は、文書群に対して任意のキーフレーズを指定することにより、指定されたキーフレーズに対応付けられた 1 または複数の文書を文書群から抽出することができる。

【 0 0 5 4 】

続いて、S 2 4 において、情報処理装置 2 0 は、文書記憶部 4 2 から文書群に含まれる複数のキーフレーズを含むキーフレーズ群を取得する。そして、情報処理装置 2 0 は、取得したキーフレーズ群を階層クラスタリングする。情報処理装置 2 0 は、階層化した複数のクラスタのそれぞれにラベルを付与してもよい。情報処理装置 2 0 は、階層クラスタリングしたキーフレーズ群をクラスタ記憶部 4 4 に登録する。

10

【 0 0 5 5 】

続いて、S 2 5 において、情報処理装置 2 0 は、階層クラスタリングされたキーフレーズ群をクラスタ記憶部 4 4 から読み出し、読み出したキーフレーズ群を複数の候補クラスタに分割する。そして、情報処理装置 2 0 は、生成した複数の候補クラスタをクラスタ記憶部 4 4 に登録する。

【 0 0 5 6 】

続いて、S 2 6 において、情報処理装置 2 0 は、文書群を分類するための予め定められた複数の項目を他の装置から取得する。情報処理装置 2 0 は、ユーザにより入力された複数の項目を取得してもよい。そして、情報処理装置 2 0 は、取得した複数の項目および複数の項目の階層構造を項目記憶部 4 6 に記憶させる。

20

【 0 0 5 7 】

図 6 は、S 2 4 において階層クラスタリングされたキーフレーズ群の構造を表すデンドログラムの一例を示す図である。キーフレーズ群を階層クラスタリングした場合、例えば、情報処理装置 2 0 は、図 6 に示すようなデンドログラムにより表される階層構造の複数のクラスタを生成することができる。図 6 のデンドログラムは、末端ノードにキーフレーズが対応付けられる。また、高さ方向は、キーフレーズ間の類似度を表す。

【 0 0 5 8 】

図 7 は、キーフレーズ群の構造を表すデンドログラム、および、S 2 5 において生成された候補クラスタの一例を示す図である。例えば、情報処理装置 2 0 は、階層クラスタリングしたキーフレーズ群を分割して、複数の候補クラスタを生成する。例えば、情報処理装置 2 0 は、図 7 に示すようなデンドログラムに基づき、所定個（例えば、4 個以上で最小）の候補クラスタを生成する。情報処理装置 2 0 は、複数の候補クラスタのそれぞれにラベルを付加してもよい。例えば、情報処理装置 2 0 は、候補クラスタの中心位置近傍のキーフレーズを、その候補クラスタのラベルとしてもよい。

30

【 0 0 5 9 】

図 7 の例においては、情報処理装置 2 0 は、類似度が 1 . 0 0 でクラスタを切断して、4 個の候補クラスタを生成している。具体的には、情報処理装置 2 0 は、ラベルが「アクション」の候補クラスタ、ラベルが「現象」の候補クラスタ、ラベルが「用語」の候補クラスタ、および、ラベルが「その他」の候補クラスタを生成している。

【 0 0 6 0 】

図 8 は、S 2 6 において取得された複数の項目の階層構造を示す図である。例えば、情報処理装置 2 0 は、図 8 に示すような、木構造で階層化された複数の項目を取得する。複数の項目の内容および階層構造は、例えばユーザにより予め定められている。情報処理装置 2 0 は、このような複数の項目を他の装置から取得する。また、情報処理装置 2 0 は、ユーザにより入力された複数の項目を取得してもよい。

40

【 0 0 6 1 】

これらの複数の項目のそれぞれは、文書群を分類するための情報である。情報処理装置 2 0 は、記憶装置 1 6 に記憶された文書群に対して、何れかの項目を指定することにより、その項目に関連付けられた文書を取得することができる。

【 0 0 6 2 】

50

図9は、初期画像の表示処理(S12)の詳細な処理の流れを示すフローチャートである。情報処理装置20は、S12の初期画像の表示処理において、以下のS31からS34の処理を実行する。

【0063】

S31において、情報処理装置20は、ユーザによる開始操作を、入力装置14から受け付ける。続いて、S32において、情報処理装置20は、複数の候補クラスタをクラスタ記憶部44から取得する。また、情報処理装置20は、項目記憶部46から予め定められた複数の項目のうち最上位の複数の項目を取得する。

【0064】

続いて、S33において、情報処理装置20は、文書記憶部42にアクセスして、予め定められた複数の項目のそれぞれ毎、且つ、複数の候補クラスタのそれぞれ毎の文書の情報量を算出する。すなわち、情報処理装置20は、文書群を最上位の複数の項目に従って複数の初期文書群に分類する。そして、情報処理装置20は、複数の初期文書群のそれぞれについて、複数の候補クラスタのそれぞれに分類される文書の情報量を算出する。

【0065】

続いて、S34において、情報処理装置20は、複数の項目のそれぞれ毎、且つ、複数の候補クラスタのそれぞれ毎の文書の情報量に基づき、初期画像を生成する。そして、情報処理装置20は、生成した初期画像を表示装置12へと出力して、表示装置12に初期画像を表示させる。

【0066】

図10は、S34で表示される初期画像の一例を示す図である。情報処理装置20は、S34において、例えば、図10に示すような、ヒートマップ状の初期画像を生成する。

【0067】

初期画像は、一方の軸(項目軸)が項目を表し、他方の軸(クラスタ軸)が候補クラスタを表す2次元の格子状となっている。図10の例では、項目軸が縦軸、クラスタ軸が横軸となっている。そして、初期画像は、複数の格子の内部のそれぞれの輝度または濃度が、対応する項目且つ対応する候補クラスタにより分類された文書の情報量を表す。

【0068】

例えば、図10の例の初期画像では、項目軸が、「機器」、「建屋」および「部品」の3個の項目を表す。また、この初期画像では、クラスタ軸が、「アクション」、「現象」、「用語」および「その他」の4個の候補クラスタを表す。そして、この初期画像では、「機器」および「アクション」の両者に対応する格子の内部の輝度または濃度が、項目が「機器」且つ候補クラスタが「アクション」に分類される文書の情報量を表す。他の格子の内部の輝度または濃度も同様である。

【0069】

なお、初期画像は、図10に示すようなヒートマップ状の画像に限られない。例えば、初期画像は、輝度または濃度に代えて、情報量を色によって表してもよい。また、初期画像は、情報量を数値または文字で表してもよい。また、初期画像は、オブジェクトまたはアイコンの種別の違いにより表してもよいし、ラベル文字の大小、フォント種別、描画線または線の傾き等の図形の変化により表してもよい。また、初期画像は、2次元のマップ状に限らず、情報量を短文のコメントまたはリストで表した文字情報であってもよいし、立体形状等の3次元以上の形状で情報量を表した画像であってもよい。

【0070】

また、表示装置12に表示された初期画像は、ユーザが、入力装置14を用いて項目軸に表示された複数の項目のうちの何れか1つの項目を選択することが可能である。選択操作は、マウス等のポインティングデバイスにより行われてもよいし、音声等で行われてもよい。

【0071】

図11は、展開画像の表示処理(S13)の詳細な処理の流れを示すフローチャートである。情報処理装置20は、S13の展開画像の表示処理において、以下のS41からS

10

20

30

40

50

48の処理を実行する。

【0072】

S41において、情報処理装置20は、ユーザによる、初期画像に表示された複数の項目のうち何れか1つ項目の選択操作を、入力装置14から受け付ける。

【0073】

続いて、S42において、情報処理装置20は、複数の候補クラスタのそれぞれについて、S41で選択された項目に対する有用性を表すスコアを算出する。なお、スコアについては、図15を参照して詳細を後述する。

【0074】

続いて、S43において、情報処理装置20は、複数の候補クラスタのうち、算出されたスコアが所定の順位の候補クラスタを、参照クラスタとして決定する。例えば、決定部82は、複数の候補クラスタのうち、有用性が最も高いスコアの候補クラスタを、参照クラスタとして決定する。

10

【0075】

続いて、S44において、情報処理装置20は、参照クラスタを、複数のサブクラスタに分割する。例えば、情報処理装置20は、参照クラスタを所定個（例えば、4個以上で最小）に分割して、複数のサブクラスタを生成する。また、この場合、情報処理装置20は、複数のサブクラスタのそれぞれにラベルを付加してもよい。ラベルは、そのサブクラスタの中心近傍のサブフレーズ等であってもよい。

【0076】

20

続いて、S45において、情報処理装置20は、項目記憶部46から、選択された項目の下位層の予め定められた複数のサブ項目を抽出する。

【0077】

続いて、S46において、情報処理装置20は、文書記憶部42にアクセスして、複数のサブ項目のそれぞれ毎、且つ、複数のサブクラスタのそれぞれ毎の文書の情報量を算出する。すなわち、情報処理装置20は、選択された項目の初期文書群をさらに複数のサブ項目に従って分割して複数のサブ文書群を生成する。そして、情報処理装置20は、複数のサブ文書群のそれぞれについて、複数のサブクラスタのそれぞれに分類される文書の情報量を算出する。

【0078】

30

続いて、S47において、情報処理装置20は、複数のサブ項目のそれぞれ毎、且つ、複数のサブクラスタのそれぞれ毎の文書の情報量に基づき、展開画像を生成する。そして、情報処理装置20は、生成した展開画像を表示装置12へと出力して、表示装置12に展開画像を表示させる。

【0079】

続いて、S48において、情報処理装置20は、展開画像を表示した後に、複数のサブ項目のうち、何れか1つのサブ項目の選択操作を受け付けたか否かを判断する。何れか1つのサブ項目の選択操作を受け付けた場合（S48のYes）、情報処理装置20は、処理をS42に戻す。そして、情報処理装置20は、複数のサブ項目を複数の項目と置き換え、複数のサブクラスタを複数の候補クラスタと置き換えて、再度、S42～S47の

40

【0080】

展開画像を表示した後、何れのサブ項目も選択されずに終了操作がされた場合（S48のNo）、情報処理装置20は、本フローの処理を終了する。

【0081】

図12は、初期画像および展開画像の一例を示す図である。情報処理装置20は、初期画像を表示している状態において、何れか1つの項目が選択された場合、図12の右側に示すようなヒートマップ状の展開画像を生成して表示装置12に表示させる。

【0082】

展開画像は、一方の軸（項目軸）がサブ項目を表し、他方の軸（クラスタ軸）がサブク

50

ラスタを表す2次元の格子状となっている。図12の例では、項目軸が縦軸、クラスタ軸が横軸となっている。そして、展開画像は、初期画像と同様に、複数の格子の内部のそれぞれの輝度または濃度が、対応するサブ項目且つ対応するサブクラスタにより分類された文書の情報量を表す。

【0083】

例えば、図12の例では、初期画像の項目軸に表示された複数の項目のうち、「機器」の項目がユーザにより選択された。「機器」の項目が選択されたことに応じて、情報処理装置20は、「機器」の下位層の「タービン」、「ノズル」、「ポンプ」、「配管」および「ロータ」の5個のサブ項目を抽出した。

【0084】

また、図12の例では、情報処理装置20は、「機器」の項目が選択されたことに応じて、「アクション」、「現象」、「用語」および「その他」の4個の候補クラスタのそれぞれについて、選択された項目である「機器」に対する有用性を表すスコアを算出した。そして、情報処理装置20は、スコアが最も高い「アクション」の候補クラスタを、参照クラスタとして決定した。さらに、情報処理装置20は、「アクション」の参照クラスタを分割して、「検査」、「溶接」、「拡大」、「加工」および「位置」の5個のサブクラスタを生成した。

【0085】

そして、情報処理装置20は、図12に示すような展開画像を生成した。図12の例の展開画像では、項目軸が、「機器」の下位層のサブ項目である、「タービン」、「ノズル」、「ポンプ」、「配管」および「ロータ」の5個のサブ項目を表す。また、この展開画像では、クラスタ軸が、「検査」、「溶接」、「拡大」、「加工」および「位置」の5個のサブクラスタを表す。

【0086】

なお、展開画像も、初期画像と同様に、ヒートマップ状の画像に限られない。また、表示装置12に表示された展開画像は、ユーザが、入力装置14を用いて項目軸に表示された複数のサブ項目のうちの何れか1つのサブ項目を選択することが可能である。

【0087】

図13は、初期画像、展開画像および新たな展開画像の一例を示す図である。情報処理装置20は、展開画像を表示している状態において、何れか1つのサブ項目が選択された場合、新たな展開画像を生成して表示装置12に表示させる。この場合、情報処理装置20は、複数のサブ項目を複数の項目と置き換え、複数のサブクラスタを複数の候補クラスタと置き換えて、再度、展開画像を生成する処理を実行し、新たな展開画像を生成する。

【0088】

図14は、メニュー画像が追加された展開画像を示す図である。情報処理装置20の展開画像制御部90は、展開画像とともに、複数の候補クラスタのうちの何れか1つを選択させるためのメニュー画像の提示を制御してもよい。メニュー画像は、例えば、プルダウンメニュー92のような、ユーザに操作により複数の候補クラスタのうちの何れか1つを選択させるためのユーザインターフェイス画像である。

【0089】

情報処理装置20の展開画像制御部90は、メニュー画像により何れか1つの新たな候補クラスタが選択された場合、新たな候補クラスタを参照クラスタとして、再度、展開画像を提示する。具体的には、情報処理装置20のサブクラスタ生成部84は、新たな候補クラスタを参照クラスタとして、新たな複数のサブクラスタに分割する。

【0090】

情報処理装置20の第2算出部88は、サブクラスタ生成部84から、新たな複数のサブクラスタを取得する。第2算出部88は、文書記憶部42にアクセスして、予め定められた複数のサブ項目のそれぞれ毎、且つ、新たな複数のサブクラスタのそれぞれ毎の文書の情報量を算出する。そして、展開画像制御部90は、複数のサブ項目のそれぞれ毎且つ新たな複数のサブクラスタのそれぞれ毎の文書の情報量を表す新たな展開画像を生成し、

10

20

30

40

50

表示装置 12 に表示させる。

【0091】

図 15 は、複数の候補クラスタのそれぞれのスコアを示す図である。情報処理装置 20 は、初期画像において、何れかの項目が選択された場合、クラスタ軸に表示する複数のサブクラスタを自動的に選択する。この場合、情報処理装置 20 は、選択された項目に対して、最も有用なサブクラスタがクラスタ軸に表示されるように、選択された項目に対する有用性を表すスコアを、複数の候補クラスタのそれぞれに対して算出する。

【0092】

n 番目 (n は 1 以上の整数) の候補クラスタを “C n” とした場合、情報処理装置 20 は、下記の式 (1) を演算して、選択された項目に対する n 番目の候補クラスタのスコア (V (C n)) を算出する。

$$V (C n) = \dots (1)$$

【0093】

は、スコアの第 1 パラメータである。は、対応する候補クラスタにおける選択された項目に分類される文書の情報量を表す。は、このような文書の情報量に係数等乗じた値であってもよい。

【0094】

例えば、項目として「機器」が選択され、「アクション」の候補クラスタのスコアを算出する場合、は、アクションの候補クラスタに含まれる複数の文書の情報量のうちの、「機器」に分類される複数の文書の情報量に応じた値を表す。例えば、スコアは、対応する候補クラスタにおける選択された項目に分類される文書の情報量が大きい程、大きくなる。

【0095】

は、スコアの第 2 パラメータである。は、複数のサブ項目のそれぞれ毎、且つ、対応する候補クラスタを分割した複数のサブクラスタのそれぞれ毎の文書の情報量の分散を表す。は、このような分散に係数等乗じた値であってもよい。

【0096】

例えば、項目として「機器」が選択され、「アクション」の候補クラスタのスコアを算出する場合、選択された項目の下位層の複数のサブ項目は、「タービン」、「ノズル」、「ポンプ」、「配管」および「ロータ」となる。また、対応する「アクション」の候補クラスタを分割した複数のサブクラスタは、「検査」、「溶接」、「拡大」、「加工」および「位置」となる。情報処理装置 20 は、「タービン」、「ノズル」、「ポンプ」、「配管」および「ロータ」の 5 個のサブ項目のそれぞれと、「検査」、「溶接」、「拡大」、「加工」および「位置」の 5 個のサブクラスタのそれぞれとの組み合わせ毎の文書の情報量を算出する。は、これらの組み合わせ毎の情報量の分散を表す。例えば、スコアは、このような分散が大きい程、大きくなる。

【0097】

は、スコアの第 3 パラメータである。は、対応する候補クラスタをユーザが過去に操作により選択した頻度または割合を表す。は、このような頻度または割合に係数等乗じた値であってもよい。

【0098】

例えば、項目として「機器」が選択され、「アクション」の候補クラスタのスコアを算出する場合、は、過去に、項目として「機器」が選択された後に、例えば、図 14 に示すメニュー画面等により「アクション」の候補クラスタが選択された頻度または割合を表す。また、は、時間的に近い選択操作が、時間的に遠い選択操作よりも大きな影響を与えるように、頻度または割合に重みが加えられてもよい。例えば、スコアは、対応する候補クラスタを過去にユーザが操作により選択した頻度または割合が多い程、大きくなる。

【0099】

また、スコアは、およびのうち、何れか 1 つのパラメータに基づく値、または、何れか 2 つのパラメータに基づく値であってもよい。また、スコアは、および

10

20

30

40

50

に代えて、選択された項目に対する候補クラスタの有用性を表す他のパラメータにより表されてもよい。

【0100】

(第1実施形態の効果)

以上のように、第1実施形態に係る情報処理装置20は、文書群を、検索目的およびユーザの関心点に応じた2つの軸により分類して情報量を提供する。この場合において、情報処理装置20は、一方の軸(項目軸)については予め定めた項目で分類するが、他方の軸(クラスタ軸)についてはクラスタリングにより分類する。そして、情報処理装置20は、ユーザにより予め定められた複数の項目のうち何れか1つの項目が選択された場合、選択された項目に対して有用な複数のサブクラスタを自動的に抽出して、展開画像を生成する。

10

【0101】

例えば、情報処理装置20は、選択された項目に分類される文書の情報量が多くなるように、複数のサブクラスタを自動的に抽出する。また、例えば、情報処理装置20は、文書の情報量の分散が大きくなるように複数のサブクラスタを自動的に抽出する。また、例えば、情報処理装置20は、過去にユーザが高い頻度で選択した複数のサブクラスタを自動的に抽出する。

【0102】

このように情報処理装置20は、他方の軸(クラスタ軸)をクラスタリングにより分類するので、小さいコストで文書群を分類することができる。さらに、情報処理装置20は、他方の軸(クラスタ軸)を、選択した項目に対して有用な複数のサブクラスタに分類するので、文書群を適切に分類することができる。以上のように、情報処理装置20は、文書群を適切に分類した情報を小さいコストで提示することができる。

20

【0103】

(第2実施形態)

つぎに、第2実施形態に係る文書管理システム10について説明をする。第2実施形態に係る文書管理システム10は、第1実施形態に係る文書管理システム10と略同一の構成および機能を有する。第2実施形態の説明では、第1実施形態で説明したユニットと、略同一の機能および構成のユニットには同一の符号を付けて、相違点を除き詳細な説明を省略する。

30

【0104】

図16は、第2実施形態に係る第1画像制御部34および第2画像制御部36の構成を文書記憶部42および項目記憶部46とともに示す図である。

【0105】

第2実施形態において、項目記憶部46は、文書群を第1観点により分類するための予め定められた複数の第1観点項目を記憶する。さらに、項目記憶部46は、文書群を第1観点とは異なる第2観点により分類するための複数の第2観点項目を記憶する。

【0106】

第2実施形態に係る第1画像制御部34は、開始受付部72と、第1算出部74と、初期画像制御部76とを有する。

40

【0107】

第1算出部74は、開始受付部72が開始操作を受け付けると、項目記憶部46から、複数の第1観点項目および複数の第2観点項目を取得する。そして、第1算出部74は、文書記憶部42にアクセスして、予め定められた複数の第1観点項目のそれぞれ毎、且つ、予め定められた複数の第2観点項目のそれぞれ毎の文書の情報量を算出する。すなわち、第1算出部74は、文書群を複数の第1観点項目に従って複数の第1観点文書群に分類する。そして、第1算出部74は、複数の第1観点文書群のそれぞれについて、複数の第2観点項目のそれぞれに分類される文書の情報量を算出する。

【0108】

初期画像制御部76は、予め定められた複数の第1観点項目のそれぞれ毎、且つ、予め

50

定められた複数の第2観点項目のそれぞれ毎の文書の情報量を表すように、項目選択画像の提示を制御する。すなわち、初期画像制御部76は、文書群を複数の第1観点項目に従って分類した複数の第1観点文書群のそれぞれについて、複数の第2観点項目のそれぞれに分類される文書の情報量を表す項目選択画像を生成する。そして、初期画像制御部76は、生成した項目選択画像を表示装置12に出力して、表示装置12に項目選択画像を表示させる。

【0109】

第2実施形態に係る第2画像制御部36は、項目選択部78と、文書抽出部94と、選択クラスタリング部96と、選択分割部98と、スコア算出部80と、決定部82と、サブクラスタ生成部84と、第2算出部88と、展開画像制御部90とを有する。

10

【0110】

項目選択部78は、ユーザによる、複数の第1観点項目のうちの何れか1つの第1観点項目、および、複数の第2観点項目のうちの何れか1つの第2観点項目の選択操作を、入力装置14から受け付ける。

【0111】

文書抽出部94は、文書記憶部42にアクセスして、文書群から、選択された第1観点項目および選択された第2観点項目の両者に分類される複数の文書を含む選択文書群を抽出する。選択クラスタリング部96は、文書抽出部94により抽出された選択文書群に含まれる複数のキーフレーズを含む選択キーフレーズ群を取得する。そして、選択クラスタリング部96は、取得した選択キーフレーズ群を階層クラスタリングする。選択分割部98は、階層クラスタリングされた選択キーフレーズ群を複数の候補クラスタに分割する。

20

【0112】

スコア算出部80は、選択分割部98により分割された複数の候補クラスタのそれぞれについて、選択された第1観点項目および第2観点項目に対する有用性を表すスコアを算出する。決定部82は、複数の候補クラスタのうち、算出されたスコアが所定の順位の上の2つの候補クラスタを、第1参照クラスタおよび第2参照クラスタとして決定する。例えば、決定部82は、複数の候補クラスタのうち、有用性が最も高いスコアの候補クラスタを第1参照クラスタとして決定し、有用性が2番目に高いスコアの候補クラスタを第2参照クラスタとして決定する。

【0113】

サブクラスタ生成部84は、第1参照クラスタを、複数の第1サブクラスタに分割する。例えば、サブクラスタ生成部84は、第1参照クラスタを所定個（例えば、4個以上で最小）に分割して、複数の第1サブクラスタを生成する。また、サブクラスタ生成部84は、第2参照クラスタを、複数の第2サブクラスタに分割する。例えば、サブクラスタ生成部84は、第2参照クラスタを所定個（例えば、4個以上で最小）に分割して、複数の第2サブクラスタを生成する。

30

【0114】

第2算出部88は、サブクラスタ生成部84から、複数の第1サブクラスタおよび複数の第2サブクラスタを取得する。そして、第2算出部88は、複数の第1サブクラスタのそれぞれ毎且つ複数の第2サブクラスタのそれぞれ毎の文書の情報量を算出する。すなわち、第2算出部88は、複数の第1サブクラスタに従って複数の第1サブ文書群に分類する。そして、第2算出部88は、複数の第1サブ文書群のそれぞれについて、複数の第2サブクラスタのそれぞれに分類される文書の情報量を算出する。

40

【0115】

展開画像制御部90は、複数の第1サブクラスタのそれぞれ毎且つ複数の第2サブクラスタのそれぞれ毎の文書の情報量を表すように、クラスタ展開画像の提示を制御する。すなわち、展開画像制御部90は、複数の第1サブ文書群のそれぞれについて、複数の第2サブクラスタのそれぞれに分類される文書の情報量を表すクラスタ展開画像を生成する。そして、展開画像制御部90は、生成したクラスタ展開画像を表示装置12に出力して、表示装置12にクラスタ展開画像を表示させる。

50

【 0 1 1 6 】

図 1 7 は、項目記憶部 4 6 に記憶された複数の第 1 観点項目および複数の第 2 観点項目の一例を示す図である。項目記憶部 4 6 は、文書群を異なる観点で分類するための複数の第 1 観点項目および複数の第 2 観点項目を記憶する。複数の第 1 観点項目および複数の第 2 観点項目は、例えばユーザにより予め定められている。情報処理装置 2 0 は、このような複数の第 1 観点項目および複数の第 2 観点項目を他の装置から取得する。また、情報処理装置 2 0 は、ユーザにより入力された複数の第 1 観点項目および複数の第 2 観点項目を取得してもよい。

【 0 1 1 7 】

複数の第 1 観点項目と複数の第 2 観点項目とは、文書群を異なる観点で分類するための情報である。従って、情報処理装置 2 0 は、文書群を何れかの第 1 観点項目で絞り込んだ後、さらに何れかの第 2 観点項目で絞り込むことができる。

【 0 1 1 8 】

図 1 8 は、項目選択画像の一例を示す図である。情報処理装置 2 0 は、例えば、図 1 8 に示すような、ヒートマップ状の項目選択画像を生成する。

【 0 1 1 9 】

項目選択画像は、一方の軸（第 1 観点軸）が第 1 観点項目を表し、他方の軸（第 2 観点軸）が第 2 観点項目を表す 2 次元の格子状となっている。図 1 8 の例では、第 1 観点軸が縦軸、第 2 観点軸が横軸となっている。そして、項目選択画像は、複数の格子の内部のそれぞれの輝度または濃度が、対応する第 1 観点項目且つ対応する第 2 観点項目により分類された文書の情報量を表す。

【 0 1 2 0 】

例えば、図 1 8 の例の項目選択画像では、第 1 観点軸が、「機器」、「建屋」および「部品」の 3 個の第 1 観点項目を表す。また、この項目選択画像では、第 2 観点軸が、「品質部門」、「設計部門」および「製造部門」の 3 個の第 2 観点項目を表す。そして、この項目選択画像では、「機器」および「品質部門」の両者に対応する格子の内部の輝度または濃度が、第 1 観点項目が「機器」且つ第 2 観点項目が「品質部門」に分類される文書の情報量を表す。他の格子の内部の輝度または濃度も同様である。

【 0 1 2 1 】

なお、項目選択画像は、第 1 実施形態で説明した初期画像と同様に、図 1 8 に示すようなヒートマップ状の画像に限られない。また、表示装置 1 2 に表示された項目選択画像は、ユーザが、入力装置 1 4 を用いて、第 1 観点軸に表示された複数の第 1 観点項目のうちの何れか 1 つの第 1 観点項目と、第 2 観点軸に表示された複数の第 2 観点項目のうちの何れか 1 つの第 2 観点項目とを同時に選択することができる。

【 0 1 2 2 】

図 1 9 は、項目選択画像および選択キーフレーズ群の構造を表すデンドログラムの一例を示す図である。項目選択画像が表示された後、情報処理装置 2 0 は、ユーザによる、複数の第 1 観点項目のうちの何れか 1 つの第 1 観点項目、および、複数の第 2 観点項目のうちの何れか 1 つの第 2 観点項目の選択操作を、入力装置 1 4 から受け付ける。例えば、2 次元の複数の格子が表示されている項目選択画像における何れか 1 つの格子（タイル）が選択された場合、情報処理装置 2 0 は、そのタイルに対応する第 1 観点項目および第 2 観点項目の選択操作を受け付ける。

【 0 1 2 3 】

第 1 観点項目および第 2 観点項目が選択された場合、情報処理装置 2 0 は、文書記憶部 4 2 にアクセスして、文書群から、選択された第 1 観点項目および選択された第 2 観点項目の両者に分類される複数の文書を含む選択文書群を抽出する。さらに、情報処理装置 2 0 は、文書抽出部 9 4 により抽出された選択文書群に含まれる複数のキーフレーズを含む選択キーフレーズ群を取得する。そして、情報処理装置 2 0 は、取得した選択キーフレーズ群を階層クラスタリングする。例えば、情報処理装置 2 0 は、選択キーフレーズ群に対して、図 1 9 に示すデンドログラムに表されるような階層クラスタリングを行う。

10

20

30

40

50

【 0 1 2 4 】

図 2 0 は、選択キーフレーズ群の構造を表すデンドログラム、および、クラスタ展開画像の一例を示す図である。例えば、情報処理装置 2 0 は、図 2 0 に示すようなデンドログラムにより表される階層構造のクラスタを分割して、複数の候補クラスタを生成する。例えば、情報処理装置 2 0 は、デンドログラムに基づき、階層クラスタリングされた選択キーフレーズ群を分割して、所定個（例えば、4 個以上で最小）の候補クラスタを生成する。

【 0 1 2 5 】

情報処理装置 2 0 は、複数の候補クラスタのそれぞれにラベルを付加してもよい。例えば、情報処理装置 2 0 は、候補クラスタの中心位置近傍のキーフレーズを、その候補クラスタのラベルとしてもよい。

10

【 0 1 2 6 】

図 2 0 の例においては、情報処理装置 2 0 は、3 個の候補クラスタを生成している。具体的には、情報処理装置 2 0 は、ラベルが「A」の候補クラスタ、ラベルが「B」の候補クラスタ、ラベルが「C」の候補クラスタ、および、ラベルが「D」の候補クラスタを生成している。

【 0 1 2 7 】

続いて、情報処理装置 2 0 は、分割した複数の候補クラスタのそれぞれについて、選択された第 1 観点項目および第 2 観点項目に対する有用性を表すスコアを算出する。続いて、情報処理装置 2 0 は、複数の候補クラスタのうち、有用性が最も高いスコアの候補クラスタを第 1 参照クラスタとして決定し、有用性が 2 番目に高いスコアの候補クラスタを第 2 参照クラスタとして決定する。例えば、図 2 0 の例では、ラベルが「B」の候補クラスタが第 1 参照クラスタとして決定され、ラベルが「A」の候補クラスタが第 2 参照クラスタとして決定されている。

20

【 0 1 2 8 】

続いて、情報処理装置 2 0 は、第 1 参照クラスタを複数の第 1 サブクラスタに分割し、第 2 参照クラスタを複数の第 2 サブクラスタに分割する。そして、情報処理装置 2 0 は、複数の第 1 サブクラスタのそれぞれ毎且つ複数の第 2 サブクラスタのそれぞれ毎の文書の情報量を表すクラスタ展開画像を生成する。

【 0 1 2 9 】

クラスタ展開画像は、一方の軸が第 1 サブクラスタを表し、他方の軸が第 2 サブクラスタを表す 2 次元の格子状となっている。図 2 0 の例では、縦軸が第 1 サブクラスタを表し、横軸が第 2 サブクラスタを表す。そして、クラスタ展開画像は、初期画像と同様に、複数の格子の内部のそれぞれの輝度または濃度が、対応する第 1 サブクラスタ且つ対応する第 2 サブクラスタにより分類された文書の情報量を表す。

30

【 0 1 3 0 】

縦軸が、「コンプレッサ」、「フランジボルト」、「スパッタ」および「ステータ」の 4 個の第 1 サブクラスタを表す。また、このクラスタ展開画像では、横軸が、「破損」、「遅延」、「損傷」、「溶解」および「溶接」の 5 個の第 2 サブクラスタを表す。

【 0 1 3 1 】

なお、クラスタ展開画像も、初期画像と同様に、ヒートマップ状の画像に限られない。また、クラスタ展開画像は、ユーザが、入力装置 1 4 を用いて縦軸に表示された複数の第 1 サブクラスタのうちの何れか 1 つの第 1 サブクラスタを選択することが可能である。この場合、情報処理装置 2 0 は、選択された第 1 サブクラスタのさらに下位層の複数のクラスタを縦軸に表示する。また、同様に、クラスタ展開画像は、ユーザが、入力装置 1 4 を用いて横軸に表示された複数の第 2 サブクラスタのうちの何れか 1 つの第 2 サブクラスタを選択することが可能である。この場合、情報処理装置 2 0 は、選択された第 2 サブクラスタのさらに下位層の複数のクラスタを横軸に表示する。

40

【 0 1 3 2 】

図 2 1 は、列選択および行選択をした場合の項目選択画像を示す図である。項目記憶部

50

46に記憶された複数の第1観点項目および複数の第2観点項目は、木構造に階層構造化されていてもよい。

【0133】

この場合、情報処理装置20は、ユーザによる、第1観点軸の複数の第1観点項目のうち何れか1つの第1観点項目の選択操作を、受け付けることができる。情報処理装置20は、ユーザによる、第2観点軸の複数の第2観点項目のうち何れか1つの第2観点項目の選択操作を、受け付けることができる。

【0134】

例えば、ユーザにより項目選択画像の何れかの列が選択された場合、情報処理装置20は、選択された列に対応する第2観点項目が選択されたと判断する。この場合、情報処理装置20は、項目選択画像における横軸の項目を、選択された第2観点項目の下位層に展開される複数の項目に置き換える。

10

【0135】

また、例えば、ユーザにより項目選択画像の何れかの行が選択された場合、情報処理装置20は、選択された行に対応する第1観点項目が選択されたと判断する。この場合、情報処理装置20は、項目選択画像における縦軸の項目を、選択された第1観点項目の下位層に展開される複数の項目に置き換える。

【0136】

(第2実施形態の効果)

以上のように、第2実施形態に係る情報処理装置20は、文書群を、検索目的およびユーザの関心点に応じた2つの軸により分類して情報量を提供する。この場合において、情報処理装置20は、文書群を、予め定められた複数の第1観点項目および予め定められた複数の第2観点項目で分類して、項目選択画像を表示する。続いて、情報処理装置20は、文書群から、ユーザにより選択された第1観点項目および選択された第2観点項目の両者に分類される複数の文書を含む選択文書群を抽出する。続いて、情報処理装置20は、選択文書群をクラスタリングして複数の候補クラスタを生成する。続いて、情報処理装置20は、選択された第1観点項目および第2観点項目に対して有用な2つの候補クラスタを自動的に決定する。そして、情報処理装置20は、決定した2つの候補クラスタの一方を分類した複数の第1サブクラスタを一方の軸とし、他方を分類した複数の第2サブクラスタを他方の軸とした、クラスタ展開画像を生成する。

20

30

【0137】

このように情報処理装置20は、2つの軸をクラスタリングにより分類するので、小さいコストで文書群を分類することができる。さらに、情報処理装置20は、2つの軸を、抽出した選択文書群に対して有用な複数のサブクラスタに分類するので、文書群を適切に分類することができる。以上のように、情報処理装置20は、文書群を適切に分類した情報を小さいコストで提示することができる。

【0138】

(プログラム)

情報処理装置20で実行されるプログラムは、インストール可能な形式または実行可能な形式のファイルでCD-ROM、フレキシブルディスク(FD)、CD-R、DVD等のコンピュータで読み取り可能な記録媒体に記録されて提供される。また、情報処理装置20で実行されるプログラムを、インターネット等のネットワークに接続されたコンピュータ上に格納し、ネットワーク経由でダウンロードさせることにより提供するように構成してもよい。また、情報処理装置20で実行されるプログラムをインターネット等のネットワーク経由で提供または配布するように構成してもよい。また、プログラムを、ROM等に予め組み込んで提供するように構成してもよい。

40

【0139】

情報処理装置20で実行されるプログラムは、事前処理モジュール(文書取得モジュール、キーフレーズ生成モジュール、文書登録モジュール、クラスタリングモジュール、分割モジュール、項目取得モジュールおよび項目登録モジュール)と、第1画像制御モジュ

50

ール（開始受付モジュール、第1算出モジュールおよび初期画像制御モジュール）と、第2画像制御モジュール（項目選択モジュール、スコア算出モジュール、決定モジュール、サブクラスタ生成モジュール、サブ項目抽出モジュール、第2算出モジュールおよび展開画像制御モジュール）とを有する。情報処理装置20は、プロセッサ（処理回路30）が記憶媒体（記憶装置16等）からプログラムを読み出して実行することにより各モジュールが主記憶装置（記憶回路24）上にロードされる。これにより、プロセッサ（処理回路30）は、事前処理部32（文書取得部52、キーフレーズ生成部54、文書登録部56、クラスタリング部58、分割部60、項目取得部62および項目登録部64）、第1画像制御部34（開始受付部72、第1算出部74および初期画像制御部76）、第2画像制御部36（項目選択部78、スコア算出部80、決定部82、サブクラスタ生成部84、サブ項目抽出部86、第2算出部88および展開画像制御部90）として機能する。なお、これらの一部または全部がプロセッサ以外のハードウェアにより実現されてもよい。

10

【0140】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、請求の範囲に記載された発明とその均等の範囲に含まれる。

【符号の説明】

【0141】

20

10 文書管理システム

12 表示装置

14 入力装置

16 記憶装置

20 情報処理装置

22 通信部

24 記憶回路

30 処理回路

32 事前処理部

34 第1画像制御部

30

36 第2画像制御部

42 文書記憶部

44 クラスタ記憶部

46 項目記憶部

52 文書取得部

54 キーフレーズ生成部

56 文書登録部

58 クラスタリング部

60 分割部

62 項目取得部

40

64 項目登録部

72 開始受付部

74 第1算出部

76 初期画像制御部

78 項目選択部

80 スコア算出部

82 決定部

84 サブクラスタ生成部

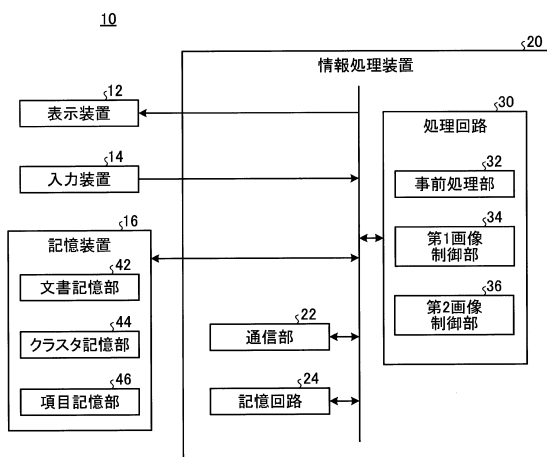
86 サブ項目抽出部

88 第2算出部

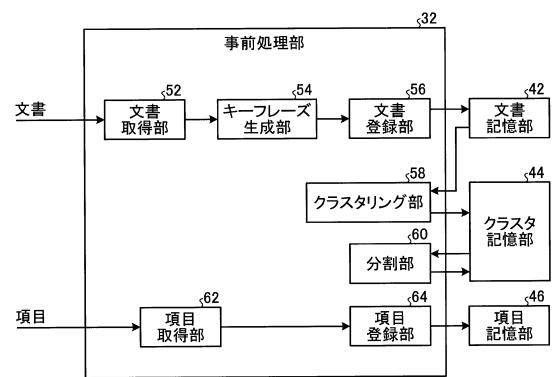
50

- 9 0 展開画像制御部
- 9 4 文書抽出部
- 9 6 選択クラスタリング部
- 9 8 選択分割部

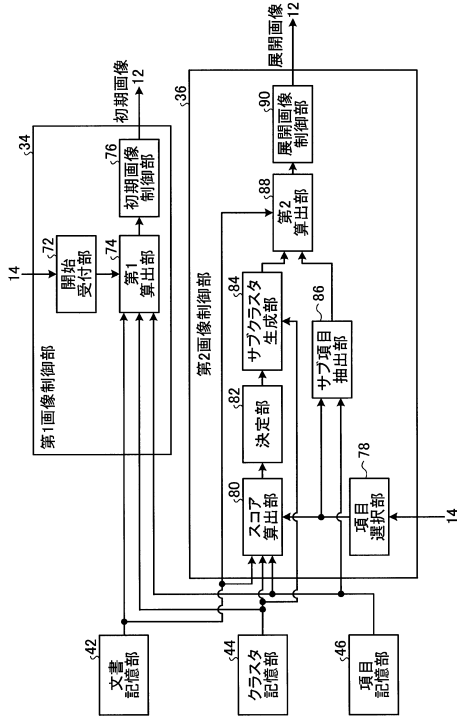
【図 1】



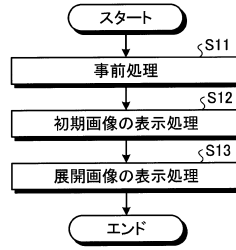
【図 2】



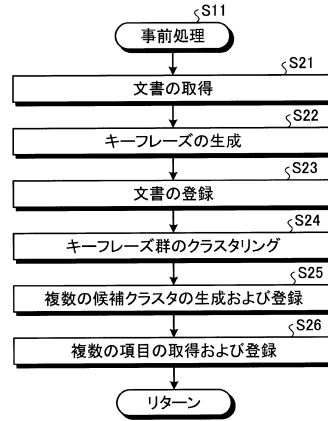
【図3】



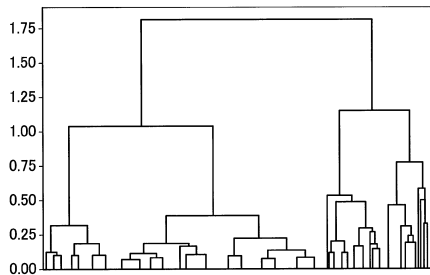
【図4】



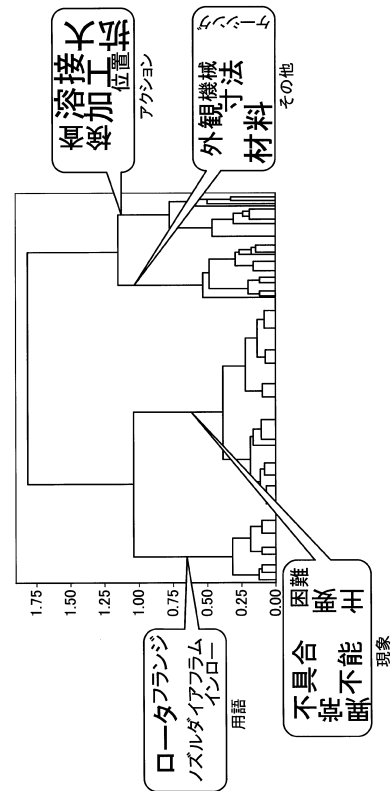
【図5】



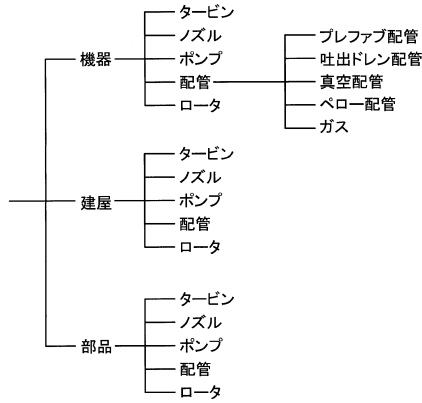
【図6】



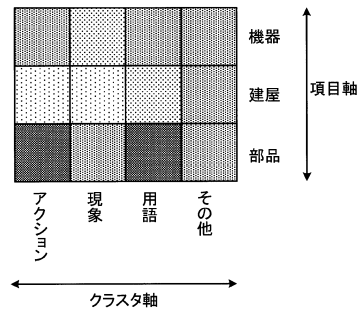
【図7】



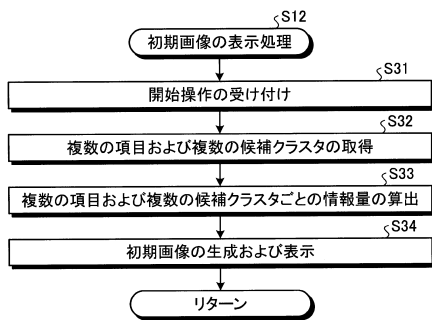
【図8】



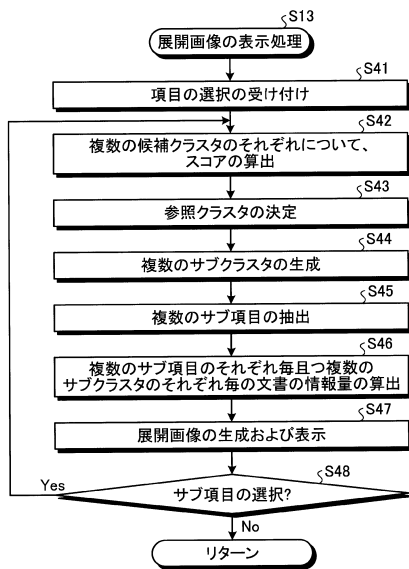
【図10】



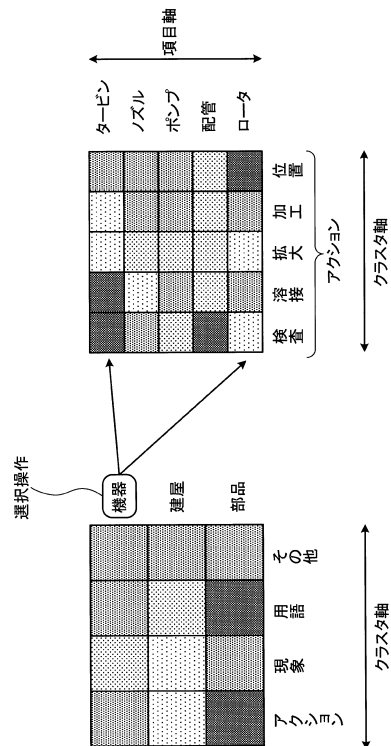
【図9】



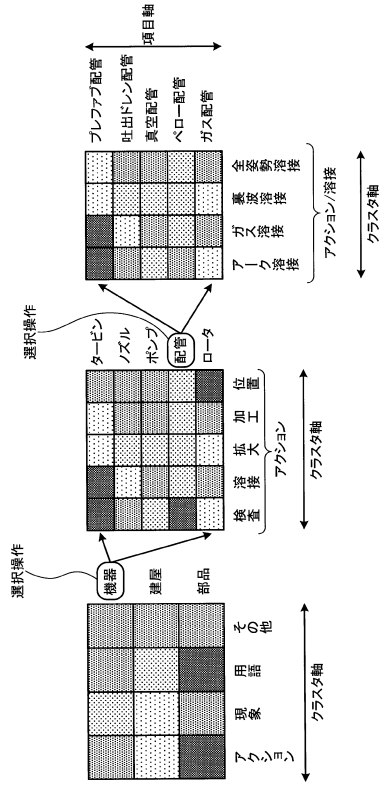
【図11】



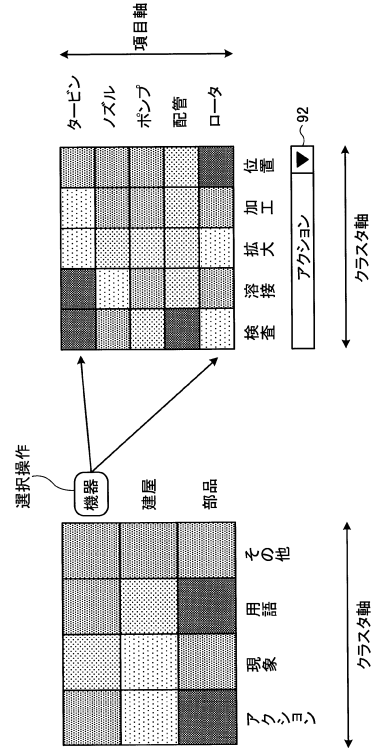
【図12】



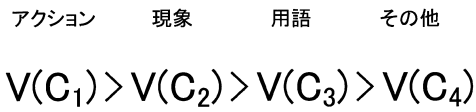
【図13】



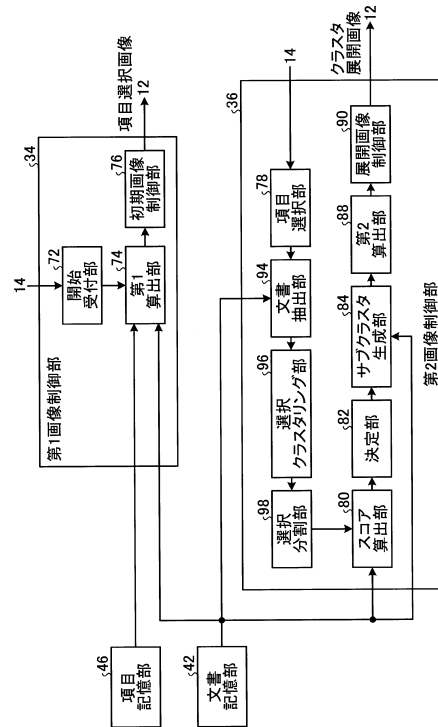
【図14】



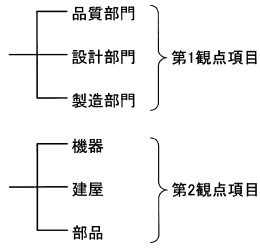
【図15】



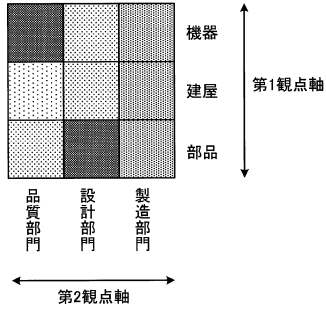
【図16】



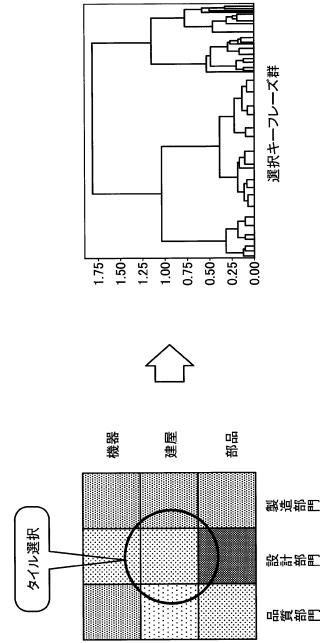
【図17】



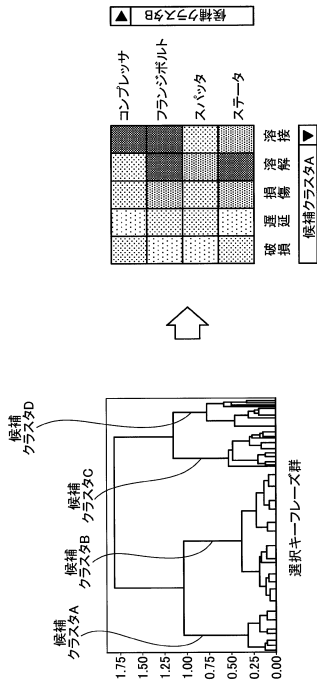
【図18】



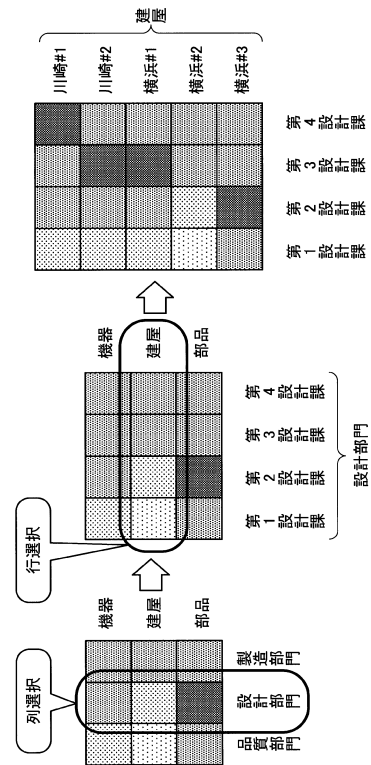
【図19】



【図20】



【図21】



フロントページの続き

(56)参考文献 特開2011-128705(JP,A)
国際公開第2007/069663(WO,A1)

(58)調査した分野(Int.Cl., DB名)
G06F 16/35
G06F 16/34