

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4067603号
(P4067603)

(45) 発行日 平成20年3月26日(2008.3.26)

(24) 登録日 平成20年1月18日(2008.1.18)

(51) Int.Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 210D
 G06F 17/30 350C
 G06F 17/30 170A

請求項の数 3 (全 15 頁)

<p>(21) 出願番号 特願平9-217131 (22) 出願日 平成9年7月27日(1997.7.27) (65) 公開番号 特開平11-45247 (43) 公開日 平成11年2月16日(1999.2.16) 審査請求日 平成16年7月2日(2004.7.2)</p> <p>前置審査</p>	<p>(73) 特許権者 390024350 株式会社ジャストシステム 徳島県徳島市川内町平石若松108番地4</p> <p>(74) 代理人 100096655 弁理士 川井 隆</p> <p>(74) 代理人 100091225 弁理士 仲野 均</p> <p>(72) 発明者 野村 直之 徳島県徳島市沖浜東3丁目46番地 株 式会社ジャストシステム内</p> <p>審査官 辻本 泰隆</p>
---	--

最終頁に続く

(54) 【発明の名称】 文書分類装置、文書分類プログラムが記憶された記憶媒体、及び文書分類方法

(57) 【特許請求の範囲】

【請求項1】

予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得手段と、

前記対象文書を取得する対象文書取得手段と、

前記対象文書取得手段で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得手段と、

前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定手段と、

各分類に対する評価値を担当者ごとに格納する評価関数データベースと、

対象文書を手動で分類した担当者の情報を取得する分類担当取得手段と、を具備し、

前記自動分類結果取得手段は、

前記対象文書の特徴づける対象文書ベクトルを取得する文書ベクトル取得手段と、

前記各分類を特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得手段と、

前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出手段と

を有し、前記類似度算出手段によって得られた各分類に対する類似度を分類結果とし、

前記分類決定手段は、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記類似度算出手段によって得られた各分類

10

20

に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類装置。

【請求項 2】

各分類に対する評価値を担当者ごとに格納する評価関数データベースを備えたコンピュータに、

予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得機能と、

前記対象文書を取得する対象文書取得機能と、

前記対象文書取得機能で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得機能と、

前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定機能と、

対象文書を手動で分類した担当者の情報を取得する分類担当取得機能と、を実現させるためのコンピュータ読取り可能な文書分類プログラムが記憶された記憶媒体であって、

前記自動分類結果取得機能は、

前記対象文書を特徴づける対象文書ベクトルを取得する文書ベクトル取得機能と、

前記各分類を特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得機能と、

前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出機能と

を有し、前記類似度算出機能によって得られた各分類に対する類似度を分類結果とし、

前記分類決定機能は、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記類似度算出機能によって得られた各分類に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類プログラムが記憶された記憶媒体。

【請求項 3】

各分類に対する評価値を担当者ごとに格納する評価関数データベース、手動分類結果取得手段、対象文書取得手段、自動分類結果取得手段、分類決定手段、分類担当取得手段、文書ベクトル取得手段、典型文書ベクトル取得手段、類似度算出手段を有する文書分類装置において用いられる文書分類方法であって、

前記手動分類結果取得手段が、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する第 1 ステップと、

前記対象文書取得手段が、前記対象文書を取得する第 2 ステップと、

前記自動分類結果取得手段が、前記第 2 ステップで取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る第 3 ステップと、

前記分類決定手段が、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する第 4 ステップと、

前記分類担当取得手段が、対象文書を手動で分類した担当者の情報を取得する第 5 ステップと、を有し、

前記第 3 ステップは、

前記文書ベクトル取得手段が、前記対象文書を特徴づける対象文書ベクトルを取得する第 6 ステップと、

前記典型文書ベクトル取得手段が、前記各分類を特徴づける典型文書の典型文書ベクトルを取得する第 7 ステップと、

前記類似度算出手段が、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る第 8 ステップと

を有し、前記第 8 ステップによって得られた各分類に対する類似度を分類結果とし、

10

20

30

40

50

前記第4ステップは、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記第8ステップによって得られた各分類に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、文書分類装置、文書分類プログラムが記憶された記憶媒体、及び文書分類方法に係り、詳細には、取得した対象文書に対する分類精度の向上に関する。

10

【0002】

【従来の技術】

文書をファイルしたり、電子的に配信したり、記憶媒体に記憶させたりする場合、その対象文書を予め決められたカテゴリに分類する場合がある。

このように対象文書の分類を行う場合、従来では分類者担当ものがその対象文書を読んだ後に手動分類をしたり、コンピュータシステムを使用して文書内容を解析することで対象文書を自動的に分類したりしている。

【0003】

【発明が解決しようとする課題】

しかし、従来の人手による文書の手動分類では、必ずしも正確に分類付けがされない場合があった。

20

一方、コンピュータシステムによる判断は高速に大量の文書を分類することが可能であるが、この分類も必ずしも正確であるとは限らなかった。

また、従来の手動分類と自動分類とでは、分類形態が全く異なるため両者を融合したシームレスな使い勝手が実現しないかった。

【0004】

本発明は、このような従来技術の課題を解決するために成されたもので、手動分類と自動分類の両分類結果を使用して、対象文書に対してより精度の高い分類を行うことが可能な文書分類装置を提供することを第1の目的とする。

また、本発明は、手動分類と自動分類の両分類結果を使用して、対象文書に対してより精度の高い分類を行うことが可能な文書分類プログラムが記録された記憶媒体を提供することを第1の目的とする。

30

また、本発明は、手動分類と自動分類の両分類結果を使用して、対象文書に対してより精度の高い分類を行うことが可能な文書分類方法を提供することを第3の目的とする。

【0005】

【課題を解決するための手段】

前記第1の目的を達成するために、請求項1に記載した発明では、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得手段と、前記対象文書を取得する対象文書取得手段と、前記対象文書取得手段で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得手段と、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定手段と、各分類に対する評価値を担当者ごとに格納する評価関数データベースと、対象文書を手動で分類した担当者の情報を取得する分類担当取得手段と、を具備し、前記自動分類結果取得手段は、前記対象文書の特徴づける対象文書ベクトルを取得する文書ベクトル取得手段と、前記各分類を特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得手段と、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出手段とを有し、前記類似度算出手段によって得られた各分類に対する類似度を分類結果とし、前記分類決定手段は、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記類似度算出手段によって得られた各分

40

50

類に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類装置を提供する。

前記第2の目的を達成するために、請求項2に記載した発明では、各分類に対する評価値を担当者ごとに格納する評価関数データベースを備えたコンピュータに、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得機能と、前記対象文書を取得する対象文書取得機能と、前記対象文書取得機能で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得機能と、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定機能と、対象文書を手動で分類した担当者の情報を取得する分類担当取得機能と、を実現させるためのコンピュータ読取り可能な文書分類プログラムが記憶された記憶媒体であって、前記自動分類結果取得機能は、前記対象文書を特徴づける対象文書ベクトルを取得する文書ベクトル取得機能と、前記各分類を特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得機能と、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出機能とを有し、前記類似度算出機能によって得られた各分類に対する類似度を分類結果とし、前記分類決定機能は、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記類似度算出機能によって得られた各分類に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類プログラムが記憶された記憶媒体を提供する。

前記第3の目的を達成するために、請求項3に記載した発明では、各分類に対する評価値を担当者ごとに格納する評価関数データベース、手動分類結果取得手段、対象文書取得手段、自動分類結果取得手段、分類決定手段、分類担当取得手段、文書ベクトル取得手段、典型文書ベクトル取得手段、類似度算出手段を有する文書分類装置において用いられる文書分類方法であって、前記手動分類結果取得手段が、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する第1ステップと、前記対象文書取得手段が、前記対象文書を取得する第2ステップと、前記自動分類結果取得手段が、前記第2ステップで取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る第3ステップと、前記分類決定手段が、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する第4ステップと、前記分類担当取得手段が、対象文書を手動で分類した担当者の情報を取得する第5ステップと、を有し、前記第3ステップは、前記文書ベクトル取得手段が、前記対象文書を特徴づける対象文書ベクトルを取得する第6ステップと、前記典型文書ベクトル取得手段が、前記各分類を特徴づける典型文書の典型文書ベクトルを取得する第7ステップと、前記類似度算出手段が、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る第8ステップとを有し、前記第8ステップによって得られた各分類に対する類似度を分類結果とし、前記第4ステップは、前記手動分類結果と前記取得した分類担当者に対応する評価値とに基づいて手動分類の点数を算出し、また、前記第8ステップによって得られた各分類に対する類似度に基づいて自動分類の点数を算出し、そして、前記算出された手動分類の点数と自動分類の点数との合計値を分類ごとに算出し、この算出結果に基づいて、前記対象文書に対する分類を最終決定することを特徴とする文書分類方法を提供する。

【0006】

【発明の実施の形態】

以下、本発明の文書分類装置、文書分類プログラムが記憶された記憶媒体、及び文書分類方法における好適な実施の形態について、図1から図7を参照して説明する。

(1) 実施形態の概要

本実施形態による文書分類処理では、過去に行った分類に対する正解率から求めた重み付

10

20

30

40

50

け等による評価関数を各分類担当者毎にデータベース化しておくと共に、各分類毎にその分類を特徴づける典型文書を予め用意しておく。

そして、分類担当者（人手）による対象文書の分類結果と評価関数とから、各分類に対する手動分類の点数化を行う。また、対象文書と典型文書との類似度を算出し、この類似度を用いて各分類に対する自動分類の点数化を行う。この両点数を各分類毎に合計した値が最も高い分類を最終分類結果とする。

このように、手動分類と自動分類とを融合化することで、より正確な分類結果を得ることができる。

【 0 0 0 7 】

(2) 実施の形態の詳細

本実施形態の文書分類装置は、パーソナルコンピュータやワードプロセッサ等を含むコンピュータシステムで構成するだけでなく、LAN（ローカル・エリア・ネットワーク）のサーバ、コンピュータ（パソコン）通信のホスト、インターネット上に接続されたコンピュータシステム等によって構成することも可能である。また、ネットワーク上の各機器に機能分散させ、ネットワーク全体で文書分類装置を構成することも可能である。

【 0 0 0 8 】

図 1 は、文書分類装置の構成を表したブロック図である。

文書分類装置は、図 1 に示すようにシステム全体を制御するための制御部 1 1 を備えている。この制御部 1 1 には、データバス等のバスライン 2 1 を介して、入力装置としてのキーボード 1 2 やマウス 1 3、表示装置 1 4、印刷装置 1 5、記憶装置 1 6、記憶媒体駆動装置 1 7、通信制御装置 1 8、入出力 I / F 1 9、及び文字認識装置 2 0 が接続されている。

制御部 1 1 は、CPU 1 1 1、ROM 1 1 2、RAM 1 1 3 を備えている。

ROM 1 1 2 は、CPU 1 1 1 が各種制御や演算を行うための各種プログラムやデータが予め格納されたリードオンリーメモリである。

【 0 0 0 9 】

RAM 1 1 3 は、CPU 1 1 1 にワーキングメモリとして使用されるランダムアクセスメモリである。この RAM 1 1 3 には、本実施形態による文書分類処理を行うためのエリアとして、自動分類と手動分類の分類結果を点数化して正規化等の処理を行う分類処理表が格納される分類処理表格納エリア 1 1 3 1、分類の対象となる対象文書が格納される対象文書格納エリア 1 1 3 2、抽出したキーワードの重要度等を要素値として対象文書の特徴づける対象文書ベクトルが格納される対象文書ベクトル格納エリア、典型文書の特徴づける典型文書ベクトルが格納される典型文書ベクトル格納エリア 1 1 3 4、対象文書と各典型文書との類似度が格納される類似度格納エリア 1 1 3 5、...、その他の各種エリアが確保されるようになっている。

【 0 0 1 0 】

キーボード 1 2 は、自装置内で対象文書を作成する場合の対象文書取得手段や群類担当者による分類結果を入力する場合の手動分類結果取得手段の一部を構成し、かな文字を入力するためのかなキーやテンキー、各種機能を実行するための機能キー、カーソルキー、等の各種キーが配置されている。

マウス 1 3 は、ポインティングデバイスであり、表示装置 1 4 に表示されたキーやアイコン等を左クリックすることで対応する機能の指定を行う入力装置である。

表示装置 1 4 は、例えば CRT や液晶ディスプレイ等が使用される。この表示装置には、キーボード 1 2 やマウス 1 3 による入力結果が表示されたり、最終分類結果が表示されたりするようになっている。

印刷装置 1 5 は、表示装置 1 4 に表示された文書や、記憶装置 1 6 の文書格納部 1 6 4 に格納された文書等の印刷を行うためのものである。この印刷装置としては、レーザプリンタ、ドットプリンタ、インクジェットプリンタ、ページプリンタ、感熱式プリンタ、熱転写式プリンタ、等の各種印刷装置が使用される。

【 0 0 1 1 】

10

20

30

40

50

記憶装置 16 は、読み書き可能な記憶媒体と、その記憶媒体に対してプログラムやデータ等の各種情報を読み書きするための駆動装置で構成されている。この記憶装置 16 に使用される記憶媒体としては、主としてハードディスクが使用されるが、後述の記憶媒体駆動装置 17 で使用される各種記憶媒体のうちの読み書き可能な記憶媒体を使用するようにしてもよい。

記憶装置 16 は、仮名漢字変換辞書 161、プログラム格納部 162、データ格納部 163、文書データベース 164、評価関数データベース 165、文書ベクトルデータベース 166、図示しないその他の格納部（例えば、この記憶装置 16 内に格納されているプログラムやデータ等をバックアップするための格納部）等を有している。

プログラム格納部 162 には、本実施形態における文書分類処理プログラム、文書ベクトル作成処理プログラム等の各種プログラムの他、仮名漢字変換辞書 161 を使用して入力された仮名文字列を漢字混り文に変換する仮名漢字変換プログラム等の各種プログラムが格納されている。

データ格納部 163 には、ユーザに関するデータ等の、システムが必要とする各種データが格納されている。

【0012】

文書データベース 164 には、各の分類を特徴づける典型文書や、典型文書以外の通常の文書等が格納されている。この文書データベース 164 に格納される各文書の形式は特に限定されるものではなく、テキスト形式の文書、HTML (Hyper Text Markup Language) 形式の文書、JIS 形式の文書等の各種形式の文書の格納が可能である。

この典型文書により特徴づけられる分類としては、技術動向報告、主張報告、新プロジェクト等の社内用の分類や、政治、経済、健康等の一般的な分類、図書館等弟子用される一般図書や科学技術文献に関する分類、その他各種分類が使用目的によって適宜選択可能になっている。

【0013】

図 2 は、評価関数データベース 165 の内容を概念的に表したものである。

この図 2 に示すように、評価関数は各分類担当者花子、太郎、四郎、... 毎に、各分類甲、乙、丙、... に対する、「重み」が評価値として格納されている。

「重み」は各分類に対する分類担当者の正解率（または誤り率）等に基づいて決定される。この「重み」は、各担当者が対象文書に対する分類を決定する毎に、最終分類結果と比較して、変更される。

この図 2 に示すように、分類担当者花子さんは、分類甲に対しての正解率が低く、分類丙に対する正解率が高いことが理解される。

【0014】

図 3 は、文書ベクトルデータベース 166 の内容を概念的に表したものである。

この図 3 に示されるように、文書 A_{jk} の中から自動抽出されたキーワード x に対して求められた重要度 f(x) が文書ベクトルの要素値 f(x) として格納されている。この文書ベクトルは各文書 jk (j = 1 ~、k = 1 ~) 毎に格納され、文書データベース 164 に格納されている各文書と対応づけられている。

各文書ベクトルの次元は採用するキーワード x (重要語句) の数であるが、2 文書間の類似度を両文書ベクトルから求める場合には、両文書のキーワードの和集合の数が両文書ベクトルの次元となる。この場合、一方の文書ベクトルにのみ含まれるキーワードに対する他方の文書ベクトルの要素値は、"0" に定義される。

【0015】

例えば図 3 おいて、文書 B のキーワードは「重要、重要語、重要度、...」、文書 C のキーワードは「重要、...、政治、...」であり、両文書の文書ベクトルは次の通りである。

文書 B の文書ベクトル = (1, 18, 19, ...)

文書 C の文書ベクトル = (18, ..., 21, ...)

これに対して文書 B と文書 C との類似度を算出する場合には、両文書のキーワードを「重要、重要語、重要度、...、政治、...」とし、両文書の文書ベクトルはつぎの通り定義され

10

20

30

40

50

る。

文書Aの文書ベクトル = (1 , 1 8 , 1 9 , ... , 0 , ...)、

文書Cの文書ベクトル = (1 8 , 0 , 0 , ... , 2 1 , ...)

【 0 0 1 6 】

記憶媒体駆動装置17(図1)は、CPU111が外部の記憶媒体からコンピュータプログラムや文書を含むデータ等を読み込むための駆動装置である。記憶媒体に記憶されているコンピュータプログラム等には、本実施形態の文書分類装置により実行される文書分類処理等の各種処理プログラム、および、そこで使用される辞書、データ等も含まれる。

ここで、記憶媒体とは、コンピュータプログラムやデータ等が記憶される記憶媒体をいい、具体的には、フロッピーディスク、ハードディスク、磁気テープ等の磁気記憶媒体、メモリチップやICカード等の半導体記憶媒体、CD-ROMやMO、PD(相変化書換型光ディスク)等の光学的に情報が読み取られる記憶媒体、紙カードや紙テープ等の用紙(および、用紙に相当する機能を持った媒体)を用いた記憶媒体、その他各種方法でコンピュータプログラム等が記憶される記憶媒体が含まれる。

本実施形態の文書分類装置において使用される記憶媒体としては、主として、CD-ROMやフロッピーディスク等の記憶媒体が使用される。

記憶媒体駆動装置17は、これらの各種記憶媒体からコンピュータプログラムを読み込む他に、フロッピーディスクのような書き込み可能な記憶媒体に対してRAM113や記憶装置16に格納されているデータ等を書き込むことが可能である。

【 0 0 1 7 】

本実施形態の文書分類装置では、制御部11のCPU111が、記憶媒体駆動装置17にセットされた外部の記憶媒体からコンピュータプログラムを読み込んで、記憶装置16の各部に格納(インストール)する。そして、本実施形態による文書分類処理等の各種処理を実行する場合、記憶装置16から該当プログラムをRAM113に読み込み、実行するようになっている。

但し、記憶装置16からではなく、記憶媒体駆動装置17により外部の記憶媒体から直接RAM113にプログラムを読み込んで実行することも可能である。また、文書分類装置によっては、本実施形態の文書分類処理プログラム等を予めROM112に記憶させておき、これをCPU111が実行するようにしてもよい。

さらに、本実施形態の文書分類処理プログラム等の各種プログラムやデータを、通信制御装置18を介して他の記憶媒体からダウンロードし、実行するようにしてもよい。

【 0 0 1 8 】

通信制御装置18は、文書分類装置と他のパーソナルコンピュータやワードプロセッサ等の各種電子機器との間をネットワーク接続するための制御装置である。

通信制御装置18は、これら各種電子機器が有している対象文書と同一の言語の文書、入力された他言語の文書、および同一言語や他言語の文書のデータベースを検索対象としてアクセスすることが可能になっている。対象となる文書には、テキスト形式やHTML形式等の各種形式の文書その他、ビットマップデータ等の各種データも含まれる。

入出力I/F19は、音声や音楽等の出力を行うスピーカ等の各種機器を接続するためのインターフェースである。

文字認識装置20は、用紙等に記載された文字をテキスト形式やHTML等の各種形式で認識する装置であり、イメージスキャナや文字認識プログラム等で構成されている。

【 0 0 1 9 】

本実施形態では、キーボード12の入力操作により作成した文書(RAM113の所定格納エリアに格納)の他、外部で作成して所定の記憶媒体に格納した文書で記憶媒体駆動装置17から読み込んだ文書、予め文書データベースに格納されている文書、通信制御装置18からダウンロードした文書、及び文字認識装置20で文字認識した文書、等の各種文書を検索の元になる対象文書として取得する(文書取得手段)ことが可能である。

【 0 0 2 0 】

以上のように構成された本実施形態の文書分類装置による文書分類処理の動作について、

10

20

30

40

50

図4を使用して説明する。

図4は文書分類処理のメイン動作を表したフローチャートである。

CPU111は、まず分類を希望する対象文書Tを取得しRAM113の対象文書格納エリア1132に格納する(ステップ11)。

【0021】

そして、CPU111は、分類担当者と、その分類担当者によって分類された手動分類結果を取得し、RAM113の分類処理表格納エリア1131の分類処理表に格納する(ステップ12)。

図6は、RAM1131の作業領域としてエリアが確保されている自動分類表の内容を概念的に表したものである。

分類担当者花子が対象文書を読んで決定した分類が分類甲であった場合、図6に示すように、花子の分類結果として花子a欄61における、分類甲の点数が1点で他の分類が0点となる。

【0022】

次にCPU111は、取得した対象文書Tの文書ベクトルBtが既に作成されていて文書ベクトルデータベース166中に格納されているか否かを確認し(ステップ14)、格納されていれば(; Y)、その文書ベクトルBtを読み込んでRAM113の対象文書ベクトル格納エリア1133に格納する(ステップ15)。

対象文書の文書ベクトルBtが文書ベクトルデータベース166に格納されていない場合(ステップ14 ; N)、CPU111は、対象文書に対する文書ベクトルBtを作成する(ステップ16)。

【0023】

図5は、文書ベクトル作成処理の動作を表したフローチャートである。

CPU111は、形態素解析を行うことで対象文書Tから自立語を抽出する(ステップ131)と共に、名詞句、複合名詞句等を含めた候補語(句)を対象文書Tから抽出しRAM113の所定作業領域に格納する(ステップ132)。

そして抽出した候補語(句)の対象文書Tでの出現頻度、評価関数から、各候補語(句)重要度f(x)を決定する(ステップ133)。ここで、評価関数としては、例えば、所定の重要語が予め指定されている場合にはその重要語に対する重み付け、単語、名詞句、複合名詞句等の候補語(句)の種類による重み付け等が使用される。

さらにCPU111は、決定した重要度f(x)の値から対象文書Tのキーワードa, b, ...を決定する(ステップ134)。そして、各キーワードの重要度f(x)を要素として、文書ベクトルB=(f(a), f(b), ...)をRAM113の対象文書ベクトル格納エリア1133に格納して(ステップ135)、図4の文書分類処理ルーチンにリターンする。

【0024】

次にCPU111は、対象文書Tと分類甲、乙、丙、...の各典型文書との類似度Sを算出する(ステップ17)。

すなわち、CPU111は、図7に示すように、対象文書の文書ベクトルBtと典型文書の文書ベクトルBjkとを比較し、両者ベクトルの角度に依存するコサインにより両文書間の類似度Sを算出する。

一般に、文書Axの文書ベクトルBxと文書Ayの文書ベクトルByとの間の角度をとし、両文書ベクトルの内積をBx・Byとし、両文書ベクトルの大きさをそれぞれ|Bx|、|By|とした場合、両文書ベクトルの類似度Sは次の数式1により求まる。

【0025】

【数1】

類似度 $S = \text{COS}(\quad) = (B_x \cdot B_y) / (|B_x| \times |B_y|)$

【0026】

この類似度Sの値は-1 S 1の値をとり、1に近いほど2つの文書ベクトルが互いに平行に近く、2つの文書Axと文書Ayは互いに類似していると考えられる。

10

20

30

40

50

【0027】

次に、CPU111は、各分類の典型文書に対して算出した類似度Sの合計値が1になるように正規化し、正規化後の類似度を自動分類の点数として分類処理表エリア1131の自動b欄62(図6)に格納する(ステップ18)。

【0028】

そして、CPU111は、手動分類と自動分類による点数に対して評価関数の処理を行う(ステップ19)。

すなわち、分類担当花子の評価関数のうち、分類甲に対象文書を分類した場合の評価関数(重み $w = 0.5$)を評価関数データベース165から読み出し、図6の分類処理表における、花子a欄61の各分類の点数に、乗じて花子c欄63に格納する。また、自動b欄62における各分類の点数に($1 - w = 0.5$)を乗じて、自動d欄64に格納する。

10

【0029】

さらにCPU111は、評価関数処理を行った後の手動分類の各点数(花子c欄63)と、に評価関数処理後の自動分類の各点数(自動d欄64)との合計値($c + d$)を各分類毎に求め、合計値が最大となる分類を対象文書Tに対する分類として最終決定する(ステップ20;分類決定手段)。

CPU111は、最終決定した分類により、分類目的に応じて対象文書を処理し(ステップ21)、処理を終了する。対象文書の処理の例としては、分類目的が配信であればその分類に属するユーザに対象文書を配信する。

【0030】

以上説明したように本実施形態によれば、各分類担当者による手動分類の結果から各分類に対する手動分類の点数化を行い、各分類を特徴づける典型文書の文書ベクトルと対象文書の対象文書ベクトルとの類似度から各分類に対する自動分類の点数化を行うことで、手動分類と自動分類とを融合させることができ、より正確な分類結果を得ることができる。

20

【0031】

以上、本実施形態の構成および他言語文書検索の処理について説明したが、本発明では、これらの各形態に限定されるものではなく、請求項に記載された発明の範囲内で種々の変形をすることが可能である。

例えば、典型文書は、必ずしも予め選ばれている必要がなく、文書データベース164に格納されてる通常の文書を典型文書として使用してもよい。

また、文書データベース163に格納されている文書の中から、クラスタリング処理により自動抽出した文書を典型文書として使用するようにしてもよい。

30

【0032】

説明した実施形態では、典型文書とその典型文書ベクトルとがそれぞれ文書データベース164、文書ベクトルデータベース166に格納されていることを前提に説明したが、必ずしも両者が存在する必要はない。

すなわち、典型文書に対する典型文書ベクトルが存在すれば(文書ベクトルデータベース166に格納されていれば)、対象文書Tとの類似度Sを算出することができるので、典型文書自体は必ずしも必要ではない。

40

逆に、各分類毎にその分類を特徴づける典型文書が存在すれば(文書データベース164に格納されていれば)、図5に示した文書ベクトル作成処理により、典型文書ベクトルを作成することができるので、典型文書ベクトル自体は必ずしも必要ではない。

【0033】

また、説明した実施形態では、1分類に対する典型文書の数については特に限定しなかったが、典型文書は必ずしも1分類に1典型文書である必要はなく、1分類に複数の典型文書を用意するようにしてもよい。この場合、各分類に対する対象文書の類似度としては合計値または平均値(正規化処理を行うのでどちらを使用することも可能である。)を使用する。このように1分類複数典型文書とすることで、各をよりの確に特徴づけることができ、自動分類側の精度を上げることができる。

50

【 0 0 3 4 】

また、最終分類結果と分類対象者による分類結果が異なる場合には、評価関数の重み付けを変えることで、学習を行うようにしても良い。で文書分類装置を構成することも可能である。

また、自動分類による分類結果（例えば、ステップ 18 による正規化後の類似度の値）に対して、手動分類の場合と同様に、重み付け（自動分類に対する評価関数）を規定するようにしてもよい。そして、この場合の重み付けに対しても、学習により変更するようにしてもよい。

【 0 0 3 5 】

さらに、説明した実施形態では、対象文書の言語については特に言及しなかったが、本発明では日本語に限定されるものではなく、あらゆる言語の対象文書に適用することが可能である。この場合、対象文書の言語用の形態素解析アルゴリズム等を使用するといった、本発明の構成には影響のない部分を変更するだけでよい。

但し、典型文書の言語は対象文書の言語と同一である必要がある。

【 0 0 3 6 】

以上の実施形態において説明した、各装置、各部、各動作、各処理等に対しては、それらを含む上位概念としての各手段（～手段）により、実施形態を構成することが可能である。

例えば、「CPU 111 は、図 7 に示すように、対象文書の文書ベクトル B_t と典型文書の文書ベクトル B_{jk} とを比較し、両者ベクトルの角度に依存するコサインにより両文書間の類似度 S を算出する。」との記載に対して「類似度算出手段」を構成するようにしてもよい。

同様に、その他各種動作に対して「～（動作）手段」等の上位概念で実施形態を構成するようにしてもよい。

例えば、以下のように構成するようにしてもよい。

（ 1 ） 図 8 に示すように、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得手段と、前記対象文書を取得する対象文書取得手段と、前記対象文書取得手段で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得手段と、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定手段と、を文書分類装置に具備させる。

（ 2 ） 図 9 に示すように、上記（ 1 ）に記載した文書分類装置において、前記自動分類結果取得手段は、前記対象文書の特徴づける対象文書ベクトルを取得する文書ベクトル取得手段と、前記各分類の特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得手段と、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出手段とを有し、前記類似度算出手段によって得られた各分類に対する類似度を分類結果とする。

（ 3 ） 図 10 に示すように、予め決められた複数分類のセットの範囲内で、人手によって対象文書を分類した手動分類結果を取得する手動分類結果取得機能と、前記対象文書を取得する対象文書取得機能と、前記対象文書取得機能で取得された対象文書を、前記複数分類のセットの範囲内で、自動的に分類して自動分類結果を得る自動分類結果取得機能と、前記手動分類結果と前記自動分類結果とから前記対象文書に対する分類を最終決定する分類決定機能と、をコンピュータに実現させるためのコンピュータ読取り可能な文書分類プログラムを記憶媒体に記憶させる。

（ 4 ） 図 11 に示すように、前記自動分類結果取得機能は、前記対象文書の特徴づける対象文書ベクトルを取得する文書ベクトル取得機能と、前記各分類の特徴づける典型文書の典型文書ベクトルを取得する典型文書ベクトル取得機能と、前記対象文書ベクトルと前記各典型文書ベクトルとの類似度を算出して各分類に対する類似度を得る類似度算出機能とを有し、前記類似度算出機能によって得られた各分類に対する類似度を分類結果とする。

（ 5 ） 図 12 に示すように、予め決められた複数分類のセットの範囲内で対象文書を自動

10

20

30

40

50

的に分類し、この自動分類結果と、前記複数分類のセットの範囲内で、人手によって前記対象文書を分類した手動分類結果とから前記対象文書に対する分類を最終決定する。

【 0 0 3 7 】

【発明の効果】

本発明によれば、同一の複数分類のセットの範囲内で、手動分類と自動分類を行うと共に、両分類結果を使用して対象文書に対する最終分類を決定するようにしたので、手動分類と自動分類の両分類結果を使用して、対象文書に足してより精度の高い分類を行うことができる。

【図面の簡単な説明】

【図 1】本発明の 1 実施形態における文書分類装置の構成を表したブロック図である。 10

【図 2】同上、実施形態における評価関数データベースの内容を概念的に表した説明図である。

【図 3】同上、実施形態における文書ベクトルデータベースの内容を概念的に表した説明図である。

【図 4】同上、実施形態における文書分類処理のメイン動作を表したフローチャートである。

【図 5】同上、実施形態の文書分類処理における文書ベクトル作成処理の動作を表したフローチャートである。

【図 6】同上、実施形態において分類の最終決定までの分類処理表での処理を表した説明図である。 20

【図 7】同上、実施形態において対象文書と典型文書との類似関係を文書ベクトルを用いて表した説明図である。

【図 8】請求項 1 に記載した発明のクレーム対応図である。

【図 9】請求項 2 に記載した発明のクレーム対応図である。

【図 10】請求項 3 に記載した発明のクレーム対応図である。

【図 11】請求項 4 に記載した発明のクレーム対応図である。

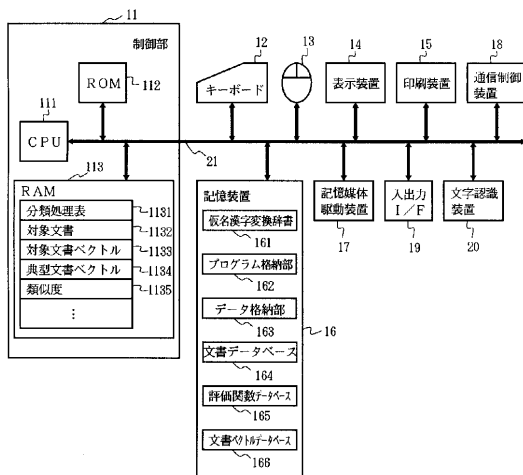
【図 12】請求項 5 に記載した発明のクレーム対応図である。

【符号の説明】

1 1	制御部	
1 1 2	ROM	30
1 1 3	RAM	
1 1 3 1	分類処理表	
1 1 3 2	対象文書格納エリア	
1 1 3 3	対象文書ベクトル格納エリア	
1 1 3 4	典型文書ベクトル格納エリア	
1 1 3 5	類似度格納エリア	
1 2	キーボード	
1 3	マウス	
1 4	表示装置	
1 5	印刷装置	40
1 6	記憶装置	
1 6 1	仮名漢字変換辞書	
1 6 2	プログラム格納部	
1 6 3	データ格納部	
1 6 4	文書データベース	
1 6 5	評価関数データベース	
1 6 6	文書ベクトルデータベース	
1 7	記憶媒体駆動装置	
1 8	通信制御装置	
1 9	入出力 I / F	50

2 0 文字認識装置

【図 1】



【図 2】

評価関数データベース ; 165

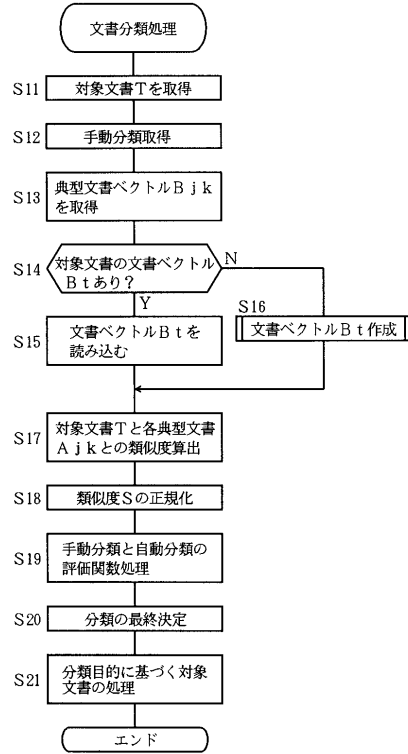
分類	花子	太郎	四郎	...
甲	0.5	0.8	0.9	...
乙	0.9	0.7	0.6	...
丙	0.9	0.9	0.8	...
.	.	.	.	
.	.	.	.	
.	.	.	.	

【図3】

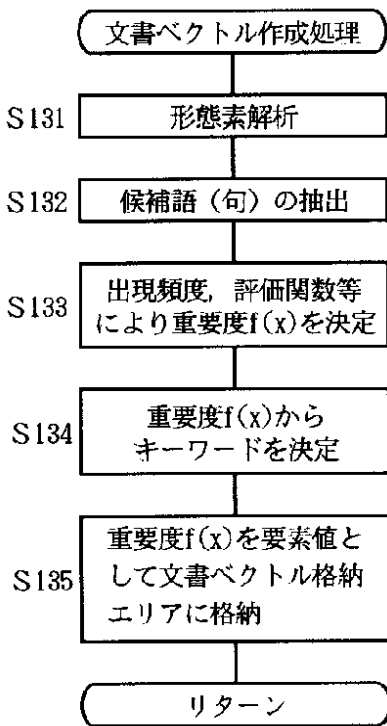
文書ベクトルデータベース; 166

文書	キーワードの要素値 f(x)					
	重要	重要語	重要度	政治
A	2	20	21
B	1	18	19
C	18	21
.....

【図4】



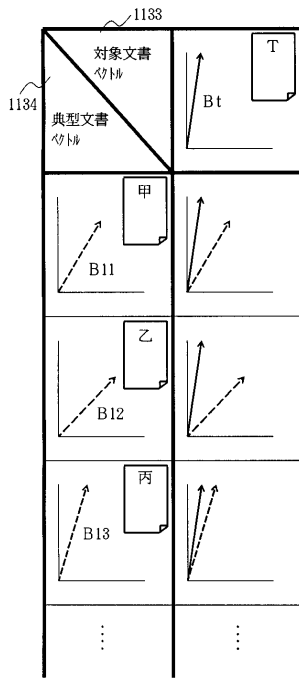
【図5】



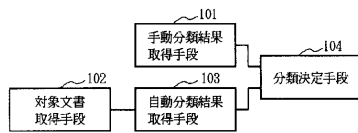
【図6】

分類	⁶¹ 花子 a	⁶² 自動 b	⁶³ 花子 c a × w	⁶⁴ 自動 d b × (1-w)	⁶⁵ 最終評価 c+d
甲	1	0.2	0.5	0.1	0.6
乙	0	0.1	0	0.05	0.05
丙	0	0.7	0	0.85	0.85
⋮	⋮	⋮	⋮	⋮	⋮

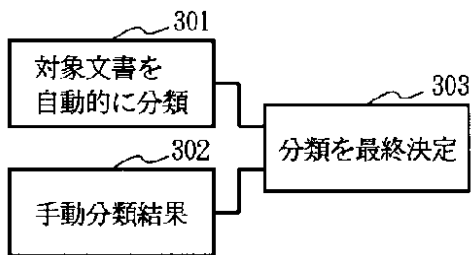
【図7】



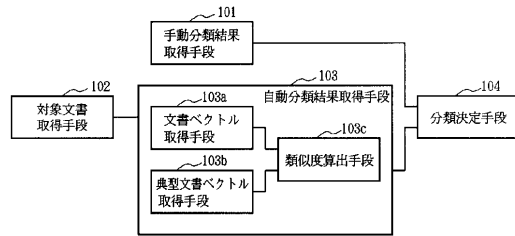
【図8】



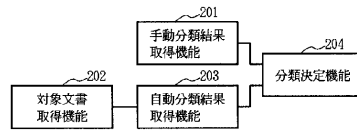
【図12】



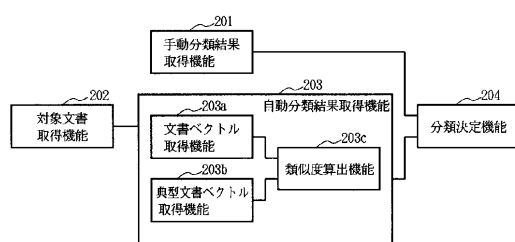
【図9】



【図10】



【図11】



フロントページの続き

(56)参考文献 特開平09-022414(JP,A)

特開平7-282078(JP,A)

特開平9-26963(JP,A)

田中 栄治, 情報探索支援システムの構築(1), 電子情報通信学会技術研究報告, 日本, 社団法人 電子情報通信学会, 1996年12月14日, 第96巻 第431号, 55~62

(58)調査した分野(Int.Cl., DB名)

G06F 17/30