

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7670951号
(P7670951)

(45)発行日 令和7年5月1日(2025.5.1)

(24)登録日 令和7年4月22日(2025.4.22)

(51)国際特許分類

F I

G 0 6 F 16/33 (2025.01)

G 0 6 F 16/33

請求項の数 11 (全21頁)

(21)出願番号	特願2020-178643(P2020-178643)	(73)特許権者	390002761
(22)出願日	令和2年10月26日(2020.10.26)		キャノンマーケティングジャパン株式会
(65)公開番号	特開2022-69790(P2022-69790A)		社
(43)公開日	令和4年5月12日(2022.5.12)		東京都港区港南2丁目16番6号
審査請求日	令和5年10月25日(2023.10.25)	(73)特許権者	592135203
			キャノンITソリューションズ株式会
			社
			東京都港区港南2丁目16番6号
		(74)代理人	100189751
			弁理士 木村 友輔
		(72)発明者	松田 雄介
			東京都港区港南2丁目16番6号 キャ
			ノンITソリューションズ株式会社内
		(72)発明者	福田 直之
			東京都港区港南2丁目16番6号 キャ
			ノンITソリューションズ株式会社内
			最終頁に続く

(54)【発明の名称】 情報処理装置、情報処理方法、プログラム

(57)【特許請求の範囲】

【請求項1】

文書のカテゴリ毎に、文書のフィールドと、当該フィールドから抽出する文字列と、当該フィールドの重みと、を定義した定義情報を受け付ける定義情報受付手段と、複数の検索対象の文書のそれぞれについて、当該検索対象の文書に係るカテゴリに対応する前記定義情報に基づき、当該検索対象の文書から当該検索対象の文書のフィールド毎に前記文字列を抽出する抽出手段と、

ユーザから前記検索対象の文書を検索するための検索語を受け付ける検索語受付手段と、複数の検索対象の文書のそれぞれについて、前記検索語受付手段により受け付けた検索語と、前記抽出手段により当該検索対象の文書から抽出された文字列と、前記定義情報により定義された当該文字列に係るフィールドの重みと、に基づき、当該検索対象の文書のフィールド毎に算出された検索スコアに基づく当該検索対象の文書の検索スコアを算出するスコア算出手段と、

を備えることを特徴とする情報処理システム。

【請求項2】

前記スコア算出手段は、前記フィールド毎に算出された検索スコアの合計値を当該検索対象の文書の検索スコアとして算出することを特徴とする請求項1に記載の情報処理システム。

【請求項3】

前記スコア算出手段は、前記フィールド毎の検索スコアを、当該フィールドにおける前

記検索語と一致する前記抽出された文字列の数と、当該フィールドの重みとに基づき算出することを特徴とする請求項 1 または 2 に記載の情報処理システム。

【請求項 4】

前記スコア算出手段は、前記フィールド毎の検索スコアを、当該フィールドにおける前記検索語と一致する前記抽出された文字列の数に、当該フィールドの重みを掛け合わせることで算出することを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の情報処理システム。

【請求項 5】

前記定義情報は、前記フィールドから抽出する文字列の抽出方法として、前記フィールド毎に、形態素解析、キーワードマッチ、正規表現パターンのいずれかの方法が定義されていることを特徴とし、

10

前記抽出手段は、前記定義情報に定義された抽出方法に基づき、当該検索対象の文書のフィールド毎に前記文字列を抽出することを特徴とする請求項 1 乃至 4 のいずれか 1 項に記載の情報処理システム。

【請求項 6】

前記スコア算出手段により算出された検索スコアに基づき、前記検索語による検索結果を表示するよう制御する表示制御手段をさらに備えることを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の情報処理システム。

【請求項 7】

前記検索語により検索された文書の閲覧実績に基づき、当該検索された文書に係るカテゴリに対応する定義情報での前記フィールドの重みを更新する更新手段をさらに備えることを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の情報処理システム。

20

【請求項 8】

前記更新手段は、閲覧された文書に係るカテゴリに対応する定義情報での前記フィールドの重みを高くすることを特徴とする請求項 7 に記載の情報処理システム。

【請求項 9】

前記更新手段は、閲覧されなかった文書に係るカテゴリに対応する定義情報での前記フィールドの重みを低くすることを特徴とする請求項 7 または 8 に記載の情報処理システム。

【請求項 10】

情報処理システムの定義情報受付手段が、文書のカテゴリ毎に、フィールドと、当該フィールドから抽出する文字列と、当該フィールドの重みと、を定義した定義情報を受け付ける定義情報受付工程と、

30

前記情報処理システムの抽出手段が、複数の検索対象の文書のそれぞれについて、当該検索対象の文書に係るカテゴリに対応する前記定義情報に基づき、当該検索対象の文書から当該検索対象の文書のフィールド毎に前記文字列を抽出する抽出工程と、

前記情報処理システムの検索語受付手段が、ユーザから前記検索対象の文書を検索するための検索語を受け付ける検索語受付工程と、

前記情報処理システムのスコア算出手段が、複数の検索対象の文書のそれぞれについて、前記検索語受付工程により受け付けた検索語と、前記抽出工程により抽出された文字列と、前記定義情報により定義された当該文字列に係るフィールドの重みと、に基づき、当該検索対象の文書のフィールド毎に算出された検索スコアに基づく当該検索対象の文書の検索スコアを算出するスコア算出工程と、

40

を備えることを特徴とする情報処理方法。

【請求項 11】

コンピュータを請求項 1 乃至 9 のいずれか 1 項に記載の各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】

【0001】

本発明は、情報処理装置、情報処理方法、プログラムに関する。

50

【背景技術】

【0002】

従来の単語頻度のみによる全文検索では単語の重要度や意味というものが考慮されない。そのため、出現頻度は低い重要な単語ではヒットしても検索上位に現れなかったり、字面は同じだがニュアンスが異なる単語にヒットした文書が検索結果に現れたりするという問題があった。

【0003】

特許文献1には、文書データのフィールド情報を検索スコアの計算に用いて、ユーザの検索意図に近い検索結果を得るための技術について開示されている。

【先行技術文献】

【特許文献】

【0004】

【文献】特開2005-063468号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

特許文献1には、文書データのフィールド情報を検索スコアの計算に用いて、ユーザの検索意図に近い検索結果を得るための技術が記載されている。

【0006】

しかし、フィールドごとのスコア反映割合を検索のたびにユーザが入力する必要があり、フィールド数が多くなった場合に煩雑である。また、フィールドが事前に文書のメタデータとして用意されていない部分についてはフィールドに格納されない問題がある。さらに文書のカテゴリについての概念がないため、各カテゴリに応じたフィールド情報の抽出やスコア計算を行うことができないという課題がある。

【0007】

そのため、文書データにカテゴリ情報を付与し、カテゴリごとにフィールド抽出情報を定義することが望まれる。

【0008】

そこで、本発明は、文書データのカテゴリとフィールド情報を検索スコアの計算に用いて、検索ユーザの意図に近い検索結果を得られるよう検索精度の向上を行うことを目的とする。

【課題を解決するための手段】

【0009】

本発明の情報処理システムは、文書のカテゴリ毎に文書に含まれる各フィールドに設定される重みを記憶する記憶手段と、前記フィールドに設定された重みと、当該フィールドと検索語との関係とに基づき、当該文書のスコアを算出する算出手段と、を備えることを特徴とする。

【0010】

本発明の情報処理方法は、文書のカテゴリ毎に文書に含まれる各フィールドに設定される重みを記憶する記憶ステップと、前記フィールドに設定された重みと、当該フィールドと検索語との関係とに基づき、当該文書のスコアを算出する算出ステップと、を備えることを特徴とする。

【0011】

本発明のプログラムは、コンピュータを、文書のカテゴリ毎に文書に含まれる各フィールドに設定される重みを記憶する記憶手段と、前記フィールドに設定された重みと、当該フィールドと検索語との関係とに基づき、当該文書のスコアを算出する算出手段として機能させることを特徴とする。

【発明の効果】

【0012】

本発明によれば、文書データのカテゴリとフィールド情報を検索スコアの計算に用いて

10

20

30

40

50

、検索ユーザの意図に近い検索結果を得られるよう検索精度の向上を行うことが可能となる。

【図面の簡単な説明】

【 0 0 1 3 】

【図 1】情報処理システムのシステム構成を示す図である。

【図 2】情報処理装置のハードウェア構成を示す図である。

【図 3】本実施例での処理の流れを示すフローチャートである。

【図 4】本発明の実施形態における、文書登録処理の一例を示すフローチャートである。

【図 5】本発明の実施形態における、フィールド抽出処理の一例を示すフローチャートである。

10

【図 6】本発明の実施形態における、キーワードによるフィールド抽出処理の一例を示すフローチャートである。

【図 7】本発明の実施形態における、パターンによるフィールド抽出処理の一例を示すフローチャートである。

【図 8】本発明の実施形態における、形態素解析によるフィールド抽出処理の一例を示すフローチャートである。

【図 9】本発明の実施形態における、検索処理の一例を示すフローチャートである。

【図 10】本発明の実施形態における、検索セッション統計情報の更新処理の一例を示すフローチャートである。

【図 11】本発明の実施形態における、フィールド重みの更新処理の一例を示すフローチャートである。

20

【図 12】本発明の実施形態における、抽出定義一覧画面の一例を示す図である。

【図 13】本発明の実施形態における、抽出定義詳細画面の一例を示す図である。

【図 14】本発明の実施形態における、フィールド名とキーワードの距離についての説明の図である。

【図 15】本発明の実施形態における、フィールド重み更新処理の一例を示す図である。

【図 16】本発明の実施形態における、検索セッション統計情報のテーブルの一例を示す図である。

【図 17】本発明の実施形態における、フィールドスコアの計算の一例を示す図である。

【発明を実施するための形態】

30

【 0 0 1 4 】

以下、図面を参照して、本発明の実施形態を詳細に説明する。なお、以下に説明する実施形態は、本発明を具体的に実施した場合の一例を示すもので、特許請求の範囲に記載した構成の一例である。

【 0 0 1 5 】

図 1 は、本発明の実施形態における文書検索システムのシステム構成の一例を示す図である。本発明における実施形態における文書検索システム 2000 は、情報処理装置 100 と、文書 DB 107、クライアント PC 108 で構成される。情報処理装置 100 は、文書登録処理部 101、文書検索処理部 102、形態素解析辞書 103、登録文書インデックス 104、抽出定義 DB 105、検索セッション統計情報 106 から構成され、外部の文書 DB 107 や文書検索システムに文書を登録する際に使用するクライアント PC とネットワークを介して通信可能に接続されている。

40

【 0 0 1 6 】

文書登録処理部 101 では、ユーザから受け付けた文書に係る処理を実行する機能部である。具体的には、テキスト抽出処理やカテゴリ付与やフィールドの抽出処理を行い、検索インデックスを作成し、登録文書インデックス 104 に格納するなどの処理を行う。

【 0 0 1 7 】

文書検索処理部 102 では、ユーザから受け付けた検索語を用いて、インデックス済みの文書を検索する機能部である。ユーザから検索語を受け付けると、インデックス済みの文書から本文スコアとフィールドスコアを計算して、それぞれを合算して検索結果に反映

50

させる処理を行う。

【 0 0 1 8 】

形態素解析辞書 1 0 3 は、形態素解析を行う際に使用される辞書である。

【 0 0 1 9 】

登録文書インデックス 1 0 4 は、登録対象となる文書から抽出した本文及び各フィールドに対する検索インデックスを格納する D B である。本 D B を用いて、検索処理部 1 0 2 による処理が行われる。

【 0 0 2 0 】

抽出定義 D B 1 0 5 は、カテゴリ毎に定義づけられる抽出定義を記憶しておく D B である。本抽出定義 D B に記憶される当該カテゴリの抽出定義として設定された抽出方式により、フィールドの抽出を行う。抽出方式は、キーワードによる抽出を行うか、パターンによる抽出を行うか、形態素解析による抽出などがある。

10

【 0 0 2 1 】

検索セッション統計情報 1 0 6 は、ユーザの検索セッション統計情報を更新する D B である。ユーザの検索セッション統計情報の更新を行い、抽出定義のフィールド重みを更新する際に利用する。

【 0 0 2 2 】

文書 D B 1 0 7 は、文書が記憶されている D B である。クラウドサービスなどの外部 D B も含まれる。

【 0 0 2 3 】

20

クライアント P C 1 0 8 は、ユーザから文書登録を受付ける際に使用される。

【 0 0 2 4 】

図 2 は、本発明の実施形態における情報処理装置のハードウェア構成の一例を示すブロック図である。

【 0 0 2 5 】

図 2 に示すように、情報処理装置は、システムバス 2 0 0 を介して C P U (C e n t r a l P r o c e s s i n g U n i t) 2 0 1、ROM (R e a d O n l y M e m o r y) 2 0 2、RAM (R a n d o m A c c e s s M e m o r y) 2 0 3、記憶装置 2 0 4、入力コントローラ 2 0 5、音声入力コントローラ 2 0 6、ビデオコントローラ 2 0 7、メモリコントローラ 2 0 8、および通信 I / F コントローラ 2 0 9 が接続される。

30

【 0 0 2 6 】

C P U 2 0 1 は、システムバス 2 0 0 に接続される各デバイスやコントローラを統括的に制御する。

【 0 0 2 7 】

ROM 2 0 2 あるいは外部メモリ 2 1 3 は、C P U 2 0 1 が実行する制御プログラムである BIOS (B a s i c I n p u t / O u t p u t S y s t e m) や OS (O p e r a t i n g S y s t e m) や、本情報処理方法を実現するためのコンピュータ読み取り実行可能なプログラムおよび必要な各種データ (データテーブルを含む) を保持している。

【 0 0 2 8 】

RAM 2 0 3 は、C P U 2 0 1 の主メモリ、ワークエリア等として機能する。C P U 2 0 1 は、処理の実行に際して必要なプログラム等を ROM 2 0 2 あるいは外部メモリ 2 1 3 から RAM 2 0 3 にロードし、ロードしたプログラムを実行することで各種動作を実現する。

40

【 0 0 2 9 】

入力コントローラ 2 0 5 は、キーボード 2 1 0 や不図示のマウス等のポインティングデバイス等の入力装置からの入力を制御する。入力装置がタッチパネルの場合、ユーザがタッチパネルに表示されたアイコンやカーソルやボタンに合わせて押下 (指等でタッチ) することにより、各種の指示を行うことができることとする。

【 0 0 3 0 】

また、タッチパネルは、マルチタッチスクリーンなどの、複数の指でタッチされた位置

50

を検出することが可能なタッチパネルであってもよい。

【 0 0 3 1 】

ビデオコントローラ 2 0 7 は、ディスプレイ 2 1 2 などの外部出力装置への表示を制御する。ディスプレイは本体と一体になったノート型パソコンのディスプレイも含まれるものとする。なお、外部出力装置はディスプレイに限ったものはなく、例えばプロジェクタであってもよい。また、前述のタッチ操作を受け付け可能な装置については、入力装置も提供する。

【 0 0 3 2 】

なおビデオコントローラ 2 0 7 は、表示制御を行うためのビデオメモリ (V R A M) を制御することが可能で、ビデオメモリ領域として R A M 2 0 3 の一部を利用することもできるし、別途専用のビデオメモリを設けることも可能である。

10

【 0 0 3 3 】

メモリコントローラ 2 0 8 は、外部メモリ 2 1 3 へのアクセスを制御する。外部メモリとしては、ブートプログラム、各種アプリケーション、フォントデータ、ユーザファイル、編集ファイル、および各種データ等を記憶する外部記憶装置 (ハードディスク)、フレキシブルディスク (F D)、或いは P C M C I A カードスロットにアダプタを介して接続されるコンパクトフラッシュ (登録商標) メモリ等を利用可能である。

【 0 0 3 4 】

通信 I / F コントローラ 2 0 9 は、ネットワークを介して外部機器と接続・通信するものであり、ネットワークでの通信制御処理を実行する。例えば、 T C P / I P を用いた通信や I S D N などの電話回線、および携帯電話の 3 G 回線を用いた通信が可能である。

20

【 0 0 3 5 】

尚、 C P U 2 0 1 は、例えば R A M 2 0 3 内の表示情報用領域へアウトラインフォントの展開 (ラスタライズ) 処理を実行することにより、ディスプレイ 2 1 2 上での表示を可能としている。また、 C P U 2 0 1 は、ディスプレイ 2 1 2 上の不図示のマウスカーソル等でのユーザ指示を可能とする。

【 0 0 3 6 】

本発明を実現するための後述する各種プログラムは、外部メモリ 2 1 3 に記憶されており、必要に応じて R A M 2 0 3 にロードされることにより C P U 2 0 1 によって実行されるものである。さらに上記プログラムの実行時に用いられる定義ファイル及び各種情報テーブル等も外部メモリ 2 1 3 に格納されており、これらについての詳細な説明も後述する。

30

【 0 0 3 7 】

次に図 3 を用いて、本願発明における処理の流れについて説明する。

【 0 0 3 8 】

ステップ S 3 0 1 では、事前設定として、カテゴリ毎のフィールド抽出定義情報 (フィールド重みを含む) とデフォルトの抽出定義 (カテゴリが設定されていないファイルやフィールド重みセットで指定しなかったフィールドに使う抽出定義) の設定を受け付ける。フィールド抽出定義情報とは、抽出定義詳細画面 1 3 0 0 に示すように、カテゴリ毎に、フィールド名と当該フィールドを抽出する方法と抽出定義が対応付けられた情報である。例えば図 1 3 に示す抽出定義情報によれば、「工事概要」というカテゴリの文書については、「事務所」や「病院」といったキーワードにより抽出されるフィールドを「建物用途」というフィールドとして抽出することが可能となる。

40

【 0 0 3 9 】

設定された抽出定義情報は、抽出定義 D B 1 0 5 に保存される。

【 0 0 4 0 】

ステップ S 3 0 2 では、ユーザから受け付けた文書 (検索対象文書) に対して、文書登録処理を実行する。文書登録処理では、検索対象文書の本文抽出やカテゴリの付与、検索対象文書のフィールド抽出、本文及びフィールドに対する検索インデックスの構築などが行われる。文書登録処理の詳細については、図 4 を用いて後述する。

【 0 0 4 1 】

50

ステップ S 3 0 3 では、ユーザから受け付けた検索語に基づき、文書検索処理を実行する。文書検索処理では、ステップ S 3 0 2 で構築した検索インデックスを用いた検索処理が行われる。文書検索処理の詳細については、図 9 を用いて後述する。

【 0 0 4 2 】

次に図 4 ~ 図 8 のフローチャートを用いて、本発明の実施形態における文書登録処理部が実行する文書登録処理について説明する。

【 0 0 4 3 】

図 4 のフローチャートは、文書登録処理部 1 0 1 において文書を登録する処理を示すフローチャートである。

【 0 0 4 4 】

ステップ S 4 0 1 では、登録対象となる文書全てに対して処理が終了したかどうかを判定する。処理が終了していれば (S 4 0 1 の Y e s) 該フローチャートの処理を終了し、処理の終了していない文書が残っていれば (S 4 0 1 の N o) ステップ S 4 0 2 に進む。

【 0 0 4 5 】

ステップ S 4 0 2 では該文書に対してテキスト抽出処理を行う。該テキスト抽出処理は一般に開示されている技術により実現されるものであり、どのような技術・方法を用いても構わない。

【 0 0 4 6 】

ステップ S 4 0 3 では該文書に対するカテゴリ付与を行う。カテゴリとは、その文書がいかなるタイプの文書であるかを分類するために付与され、本実施例であれば工事概要、注文書、議事録などがカテゴリの分類例である。ここでのカテゴリ付与は計算機によって自動で行ってもよいし、ユーザによって手動で行っても構わない。

【 0 0 4 7 】

ステップ S 4 0 4 ではフィールド抽出処理を行う。フィールド抽出処理については、図 5 を使い後述する。

【 0 0 4 8 】

ステップ S 4 0 5 ではステップ S 4 0 2 で抽出したテキスト及びステップ S 4 0 4 で抽出した各フィールドに対する検索インデックスの作成を行い登録文書インデックス 1 0 4 に格納する。検索インデックスとは、図 9 で示す文書検索処理の処理時に使用する検索インデックスである。

【 0 0 4 9 】

図 5 のフローチャートは、文書からフィールドを抽出する処理を示すフローチャートである。

【 0 0 5 0 】

ステップ S 5 0 1 では、ステップ S 4 0 3 で付与された該文書のカテゴリを取得する。

【 0 0 5 1 】

ステップ S 5 0 2 では、抽出定義 D B 1 0 5 からステップ S 3 0 1 で設定された該カテゴリの抽出定義情報を取得する。ステップ S 5 0 1 でカテゴリが取得できなかった場合はデフォルトの抽出定義を取得する。

【 0 0 5 2 】

ステップ S 5 0 3 では、該抽出定義に定義された全てのフィールドに対して抽出処理が終了したかどうかを判断する。終了していれば (S 5 0 3 の Y e s) 該フローチャートを終了し、そうでなければ (S 5 0 3 の N o) 処理をステップ S 5 0 4 に進める。

【 0 0 5 3 】

ステップ S 5 0 4 では、該抽出定義情報に設定された処理対象のフィールドの抽出方式に応じて、処理を分岐する。例えば、図 1 3 の例では、「住所」のフィールドについては形態素解析により抽出することを意味している。抽出方式が「キーワード」であればステップ S 5 0 5 に、「パターン」であればステップ S 5 0 6 に、「形態素解析」であればステップ S 5 0 7 に処理を進める。

【 0 0 5 4 】

10

20

30

40

50

ステップS505では、キーワードによる抽出処理を行う。キーワードによる抽出処理の詳細は、図6のフローチャートを用いて後述する。

【0055】

ステップS506では、パターンによる抽出処理を行う。パターンによる抽出処理の詳細は、図7のフローチャートを用いて後述する。

【0056】

ステップS507では、形態素解析による抽出処理を行う。形態素解析による抽出処理の詳細は、図8のフローチャートを用いて後述する

ステップS508では、抽出されたフィールドを該文書のフィールドとして記録しておく。このとき、該抽出定義情報のフィールド名と関連付けて記録する。

【0057】

図6のフローチャートは、文書からキーワード方式でフィールドを抽出する処理を示すフローチャートである。

【0058】

ステップS601では、該抽出定義の全てのキーワードを処理したかどうかを判断する。全て処理していれば(S601のYes)該フローチャートを終了し、そうでなければ(S601のNo)処理をステップS602に進める。

【0059】

ステップS602では、該抽出定義から未処理のキーワードを取得する。

【0060】

ステップS603では、該文書に対するキーワードマッチを実行する。このキーワードマッチにはどのような手法を用いても構わない。

【0061】

ステップS604では、該文書に該キーワードが存在するかどうかを判断する。存在しない場合(S604のNo)処理をステップS601に進め、存在する場合(S604のYes)処理をステップS605に進める。

【0062】

ステップS605では、存在キーワードの近くにフィールド名が存在するかどうかを判定する。キーワードによる抽出処理に関する抽出定義は、図13に示すように、フィールド名とキーワードとが対応付けて登録されたものである。検出されたキーワードの近くに、当該キーワードに対応付けられたフィールド名が存在する場合(ステップS605: YES)は処理をステップS606に進め、存在しない場合(ステップS605: NO)は処理をステップS601に戻す。

【0063】

ここで、キーワード(Value)とフィールド名(Key)の距離について、図14を用いて具体的に説明する。

【0064】

図13のようにフィールド名「建物用途」には事務所、病院、飲食店、駐車場、ホテルの5つのキーワードが対応付けられているため、図14に示す文書において抽出されるキーワードは、V1「病院」、V2「事務所」、V3「駐車場」となる。このうちV1とV2はキーであるK1「建物用途」と同じ行にあり、距離的に近いと言える。一方でV3はK1と5行離れており、距離的には遠く、「建物用途」とは異なる文脈で使用されていると考えられる。したがって、キーワード抽出の際にはこのキーワードとフィールド名の距離を考慮し、遠いものを抽出対象としないようにすることで誤抽出を防ぐことができる。

【0065】

図6の説明に戻る。

【0066】

ステップS606では、ステップS605で抽出されたキーワードを該文書のフィールドの値として記録する。

【0067】

10

20

30

40

50

図 7 は、正規表現パターンによるフィールドの抽出処理を示すフローチャートである。

【 0 0 6 8 】

ステップ S 7 0 1 では、抽出定義情報に設定された正規表現パターンを取得する。正規表現パターンの例としては、図 1 3 のフィールド名 1 3 0 2 の関連法令であれば、抽出定義 1 3 0 4 の「 . + (法 | 条例) 」となる。これは、「法」または「条例」が後方一致する文字列を検出するための正規表現であり、この条件によれば例えば、「建築基準法」「騒音対策条例」などが抽出可能となる。

【 0 0 6 9 】

ステップ S 7 0 2 では、該文書に対してステップ S 7 0 1 で取得した正規表現のパターンマッチを行う。

10

【 0 0 7 0 】

ステップ S 7 0 3 では、ステップ S 7 0 2 でマッチした部分全てについて処理が行われたかどうかを判断する。全てのパターンで処理が行われた場合 (S 7 0 3 の Y e s) 該フローチャートの処理を終了し、そうでない場合 (S 7 0 3 の N o) ステップ S 7 0 4 へ処理を進める。

【 0 0 7 1 】

ステップ S 7 0 4 では、マッチした部分を該文書のフィールドの値として記録する。また、グループや名前付き前方参照といった正規表現の機能を用いてマッチした部分の一部をフィールドの値として使ってもよい。

【 0 0 7 2 】

20

図 8 は、形態素解析で得られた品詞によるフィールドの抽出処理を示すフローチャートである。

【 0 0 7 3 】

ステップ S 8 0 1 では、抽出定義を取得する。例えば、本実施例であれば、図 1 3 のフィールド名 1 3 0 2 の住所の抽出定義 1 3 0 4 に定められる品詞の並びを取得する。この場合であれば、抽出定義は [名詞 - 固有名詞 - 地域] の並びで定められている。つまり、これは、名詞の中の固有名詞の中の地域カテゴリに属する単語の並びを抽出することを意味し、「東京都港区港南」といった文字列が抽出される。

【 0 0 7 4 】

ステップ S 8 0 2 では、該文書に形態素解析を実行する。

30

【 0 0 7 5 】

ステップ S 8 0 3 では、ステップ S 8 0 1 で取得した抽出定義に合致する品詞の並びがあるかどうかを判断する。品詞の並びがない場合 (S 8 0 3 の N o) 該フローチャートの処理を終了し、そうでない場合 (S 8 0 3 の Y e s) 処理をステップ S 8 0 4 に進める。

【 0 0 7 6 】

ステップ S 8 0 4 では、マッチした部分を該文書のフィールドの値として記録する。

【 0 0 7 7 】

続けて、図 9、図 1 7 を用いて、本発明の実施形態における文書検索処理部が実行する処理について説明する。

【 0 0 7 8 】

40

図 9 は、検索処理部 1 0 2 において、ユーザからの検索語を入力として受けとり、インデックス済みの文書を検索する処理を示すフローチャートである。

【 0 0 7 9 】

ステップ S 9 0 1 では、ユーザからの検索語を取得する。

【 0 0 8 0 】

ステップ S 9 0 2 では、インデックス済みの全文書に対して文書スコアが未計算の文書が存在するかどうかを判断する。文書スコアが未計算の文書が存在する場合 (S 9 0 2 の Y e s) 処理をステップ S 9 0 3 に進め、そうでない場合 (S 9 0 2 の N o) 処理をステップ S 9 0 8 に進める。

【 0 0 8 1 】

50

ステップS903では、文書スコア未計算の文書を取得する。

【0082】

ステップS904では、該文書の本文に対する検索スコアを計算する。検索スコアとは、検索語との関連度合いを数値で表した値である。本文に対する検索スコアを、本文スコアと呼ぶ。なお、本実施例においては、本文スコアは公知の検索スコア算出方法により算出される値とする。

【0083】

ステップS905では、フィールドスコアが未計算のフィールドが存在するかどうかを判断する。存在する場合(S905のYes)処理をステップS806に進め、そうでない場合(S905のNo)処理をステップS907に進める。

【0084】

ステップS906では、フィールドスコア未計算のフィールドを取得し、該フィールドに対する検索スコアを計算する。このスコアをフィールドスコアと呼ぶ。

【0085】

フィールドスコアの計算の方法の一例を、図17を用いて説明する。ユーザから「AAA株式会社 大阪」という検索語を受け付けた場合について説明する。。

【0086】

図17Aは、大阪府警担当者議事録というタイトルの文書を示した図で、当該文書をフィールド毎に分け、各フィールドの値と重みに対応付けられている。図17Bは、プロジェクト概要というタイトルの文書を示した図で、図17Aと同様に、フィールド毎に値と重みとが対応付けてある。なお、重みは、当該文書のカテゴリによって定まる値である。なお、図17において各フィールドの値として示している内容は、説明の為に抽出定義に合致しない文字列も含めて示しているが、ステップS506、ステップS604、ステップS704で説明した通り、各フィールドの値として登録されるのは、抽出定義に合致した文字列である。

【0087】

まず、検索語の出現回数をフィールド毎にカウントする。

【0088】

図17Aの文書であれば、タイトルフィールド1803には「大阪」は1回出現、人名フィールド1804には「大阪」は0回出現、本文フィールド1705には「大阪」は3回出現している。そして、各フィールドでの検索語の出現回数をフィールド毎に設定されている重みとを掛けてフィールドスコアを求める。

【0089】

タイトルフィールド1803に設定されている重みは1806に示すように2で大阪は1回出現なので、 $1 \times 2 = 2$ となる。同様に、人名フィールド1804は $0 \times 5 = 0$ 、本文フィールド1805は $3 \times 1 = 3$ となる。これらの合計値($2 + 0 + 3 = 5$)が「大阪府警担当者議事録」という文書のフィールドスコアとして算出される。

【0090】

同様に図17Bの、プロジェクト概要．PDFのフィールドスコアを計算すると、会社名フィールドで、検索語AAA株式会社が1回出現しているので $1 \times 5 = 5$ 、住所フィールドで大阪が1回出現しているので $1 \times 5 = 5$ 、本文フィールドでAAA株式会社と大阪がそれぞれ1回ずつ出現しているので $2 \times 1 = 2$ となる。これらの合計値($5 + 5 + 2 = 12$)がプロジェクト概要．PDFのフィールドスコアとして算出される。

【0091】

ステップS907では、ステップS904で算出した該文書の本文スコアと、ステップS906で算出した該文書のフィールドスコアを合算する。この値を文書スコアと呼ぶ。

【0092】

なお、本実施例においては、本文スコアとフィールドスコアとを合算したスコアを文書スコアとしたが、各フィールドの重みを考慮したスコアであるフィールドスコアのみを用いても良い。

10

20

30

40

50

【 0 0 9 3 】

ステップ S 9 0 8 では、文書スコアの降順で検索結果をユーザに示す。なお、本実施例では検索語との関連性が強い文書の文書スコアが高くなる計算方法を用いたため、降順で検索結果をユーザに示したが、検索語との関連性が強い文書の文書スコアが小さくなる算出方法を用いる場合は、昇順により表示する。すなわち、検索語との関連性が強い文書が検索結果の上位に表示されるようソートして表示する。

【 0 0 9 4 】

以上のように、抽出定義情報で「人名」や「会社名」や「住所」など、当該カテゴリの文書の特徴付けるフィールドに対して大きな重みを設定し、設定されたフィールド毎の重みを考慮して検索スコアを算出することで、検索語が同じ数だけ含まれる文書であっても、よりユーザ（検索者）の意図に合った（ユーザが探し求めている）文書を上位に表示することが可能となる。

10

【 0 0 9 5 】

ステップ S 9 0 9 では、検索セッション統計情報更新処理を行う。図 1 0 のフローチャートを用いて後述する。

【 0 0 9 6 】

ステップ S 9 1 0 では、フィールド重み更新処理を行う。図 1 1 のフローチャートを用いて後述する。

【 0 0 9 7 】

図 1 0 は、ユーザの検索セッションでの統計情報を更新する処理を示すフローチャートである。なお、検索セッションとはユーザが検索結果を取得して、該検索結果を破棄するまでの期間のことを言う。

20

【 0 0 9 8 】

ステップ S 1 0 0 1 では、検索セッション統計情報テーブル図 1 5 の 1 4 0 0 の初期化を行う。該検索結果に含まれる全ての文書情報について、文書 ID、カテゴリを設定しセッション閲覧数を 0 に設定する。

【 0 0 9 9 】

ステップ S 1 0 0 2 では、検索セッションが終了しているかどうかを判断する。終了している場合（S 1 0 0 2 の Yes）該フローチャートの処理を終了し、そうでない場合（S 1 0 0 2 の No）ステップ S 1 0 0 3 に処理を進める。

30

【 0 1 0 0 】

ステップ S 1 0 0 3 では、ユーザが検索結果の文書を選択したかどうかを判断する。選択していない場合（S 1 0 0 3 の No）処理をステップ S 1 0 0 2 に進め、そうでない場合（S 9 0 3 の Yes）は処理をステップ S 1 0 0 4 に進める。

【 0 1 0 1 】

ステップ S 1 0 0 4 では、ユーザが選択した文書の情報を取得する。

【 0 1 0 2 】

ステップ S 1 0 0 5 では、検索セッション統計情報テーブルの該文書のエントリを更新する。この場合、該テーブルのセッション閲覧数に 1 を加える。

【 0 1 0 3 】

40

図 1 1 は、検索セッション統計情報を利用して抽出定義のフィールド重みを更新する処理を示すフローチャートである。検索の情報に応じてフィールド重みを更新していくことで、より検索精度が向上していくことが見込まれる。

【 0 1 0 4 】

ステップ S 1 1 0 1 では、検索セッション統計情報テーブル図 1 5 の 1 4 0 0 を取得する。なお、ここで取得するのは検索セッションの終了した検索セッション統計情報テーブルのみである。

【 0 1 0 5 】

ステップ S 1 1 0 2 では、ヒット文書のカテゴリごとに閲覧数を集計する。ヒット文書とは、検索処理部により検索された文書である。検索の結果ヒットした文書をユーザが閲

50

覧したかを集計することで、次の検索精度を上げるために利用される。

【0106】

ステップS1103では、ステップS1102で集計したカテゴリの中に未処理のカテゴリがあるかどうかを判断する。未処理のカテゴリがある場合(S1103のYes)処理をステップS1104に進め、そうでない場合(S1103のNo)処理をステップS1109に進める。

【0107】

ステップS1104では、未処理のカテゴリの抽出定義を取得する。

【0108】

ステップS1105では、検索語に含まれる未処理のフィールド情報(当該カテゴリの抽出定義として設定されたフィールドのうち、検索語として用いられたワードが該当するフィールドであって、未処理のフィールド)があるかどうかを判断する。未処理のフィールド情報がある場合(S1105のYes)処理をステップS1106へ進め、そうでない場合(S1105のNo)処理をステップS1103に進める。

10

【0109】

ステップS1106では、該カテゴリのセッション閲覧数が0より大きいかどうかを判断する。0より大きい場合(S1106のYes)処理をステップS1107に進め、そうでない場合(S1106のNo)処理をステップS1108に進める。

【0110】

ステップS1107では、該カテゴリのセッション閲覧数が0より大きく、該フィールドが検索に貢献できたと考え、該フィールドのフィールド重みを(セッション閲覧数) \times 0.01だけ加算する。この計算式はあくまでも一例であり、その他の計算方法を用いても構わない。

20

【0111】

ステップS1108では、該カテゴリのセッション閲覧数が0であり、該フィールドが検索に貢献していないと考え、該フィールドのフィールド重みを0.01だけ減算する。この計算式はあくまでも一例であり、その他の計算方法を用いても構わない。

【0112】

ステップS1109では、不要となった該検索セッション統計情報テーブルを破棄する。

【0113】

ここで、図11を用いて、フィールド重み更新処理の一例を説明する。まず、検索語に「住所」と「建物用途」を含む検索語が使われたとし、検索セッション終了時の検索セッション統計情報テーブルが図16の1400であったとする。また、カテゴリ「工事概要」「注文書」「議事録」の抽出定義がそれぞれ、1500、1600、1700であったとする。またフィールド重み更新式はセッション閲覧数が0より大きい場合は(セッション閲覧数) \times 0.01を加算、0の場合は0.01の減算とする。この場合、テーブル1300より検索結果のカテゴリごとのセッション閲覧数は、工事概要が2、注文書と議事録が0となる。工事概要のフィールド重みの更新は、フィールド「住所」(図16の1501)と「建物用途」(図16の1502)が両方とも定義されていることから、 $2 \times 0.01 = 0.02$ が加算され、更新後のフィールド重みはそれぞれ3.02と2.02となる。注文書のフィールド重みの更新は「住所」(図16の1601)のみが定義されていることから、0.01の減算となり、更新後のフィールド重みは0.09となる。議事録のフィールド重みの更新は「住所」「建物用途」ともに定義されていないため行われない。

30

40

【0114】

図12は、現在定義されている抽出定義の確認と、追加、削除を行う画面である。抽出定義一覧画面1200は抽出定義追加ボタン1201、一括削除ボタン1202、チェックボックス1203、編集ボタン1204、個別削除ボタン1205からなる。

【0115】

抽出定義追加ボタン1201は、押下することで抽出定義詳細画面(図13)に遷移し

50

、新規に抽出定義を作成するためのものである。

【 0 1 1 6 】

一括削除ボタン 1 2 0 2 は、押下することでチェックボックス 1 2 0 3 が有効になっている全ての抽出定義を一括削除するものである。

【 0 1 1 7 】

チェックボックス 1 2 0 3 は、有効にすることで一括削除ボタン 1 2 0 2 を用いて一括削除を行えるようにするためのものである。

【 0 1 1 8 】

編集ボタン 1 2 0 4 は、押下することで抽出定義詳細画面（図 1 3 ）に遷移し、選択した抽出定義を編集するためのものである。

【 0 1 1 9 】

個別削除ボタン 1 2 0 5 は、押下することで選択した抽出定義を削除するためのものである。

【 0 1 2 0 】

図 1 3 は、抽出定義の詳細の追加、確認、編集を行う画面である。抽出定義詳細画面 1 3 0 0 は、カテゴリ名テキストボックス 1 3 0 1、フィールド名テキストボックス 1 3 0 2、抽出方式プルダウンリスト 1 3 0 3、抽出定義テキストボックス 1 3 0 4、フィールド重みテキストボックス 1 3 0 5、フィールド削除ボタン 1 3 0 6、抽出定義フィールド追加ボタン 1 3 0 7 からなる。

【 0 1 2 1 】

なお、抽出定義一覧画面 1 1 0 0 の抽出定義追加ボタン 1 1 0 1 を押下して本画面に遷移した場合は、カテゴリ名テキストボックス 1 3 0 1 は空欄で、フィールド名テキストボックス 1 3 0 2、抽出方式プルダウンリスト 1 3 0 3、抽出定義テキストボックス 1 3 0 4、フィールド重みテキストボックス 1 3 0 5 は初期状態では表示されていない。また、抽出定義一覧画面 1 2 0 0 の編集ボタンから本画面に遷移した場合、該抽出定義の内容がカテゴリ名テキストボックス 1 3 0 1、フィールド名テキストボックス 1 3 0 2、抽出方式プルダウンリスト 1 3 0 3、抽出定義テキストボックス 1 3 0 4、フィールド重みテキストボックス 1 3 0 5 に表示される。

【 0 1 2 2 】

カテゴリ名テキストボックス 1 3 0 1 は、この抽出定義につける名称を設定するためのものである。

【 0 1 2 3 】

フィールド名テキストボックス 1 3 0 2 は、フィールドの名称を設定するためのものである。

【 0 1 2 4 】

抽出方式プルダウンリスト 1 3 0 3 は、抽出方式を選択するためのものである。ここでは「キーワード」「パターン」「形態素解析」から選択する。

【 0 1 2 5 】

抽出定義テキストボックス 1 3 0 4 は、抽出の定義を設定するためのものである。抽出方式が「キーワード」の場合は抽出するキーワードのリスト、「パターン」の場合は正規表現パターン、「形態素解析」の場合は抽出したい形態素の並びを設定する。

【 0 1 2 6 】

フィールド重みテキストボックス 1 3 0 5 は、フィールド重みを設定するためのものである。

【 0 1 2 7 】

フィールド削除ボタン 1 3 0 6 は、押下することで該フィールドの抽出定義を削除するためのものである。

【 0 1 2 8 】

抽出定義フィールド追加ボタン 1 3 0 7 は、押下することで空欄のフィールド名テキストボックス 1 3 0 2、抽出方式プルダウンリスト 1 3 0 3、抽出定義テキストボックス 1

10

20

30

40

50

３０４、フィールド重みテキストボックス１３０５、フィールド削除ボタン１３０６が最下行に追加され新しいフィールドの定義ができるようになる。

【０１２９】

このようにして、カテゴリごとに抽出定義を設定することにより検索精度の向上が見込まれる。例えば、登録文書内の建築設計書の工事概要と注文書を比較した場合、工事概要の住所（建設場所）の情報は地形や適用される自治体の条例が異なるなど非常に重要な項目であるが、注文書の住所は特に重要な情報でないため、工事概要ではフィールド重みを高め（例えば３）に、注文書では低め（例えば０．１）に設定することで、同じフィールドでのカテゴリごとと重要度の違いを表現できる。このように設定することで、住所で検索を行った場合、検索スコアが高めになる工事概要が検索結果上位に、検索スコアが低めになる注文書は検索下位に表示されることが見込み、検索ユーザの意図に沿った検索結果となりやすい。

10

【０１３０】

仮にカテゴリごとにフィールド重みを設定しなかった場合、住所で検索した場合、どのカテゴリの文書でも住所を重視するよう設定した場合、工事概要と注文書の両方が検索結果に混在することになり利便性が低下すると考えられる。

【０１３１】

図１５の検索セッション統計情報テーブル１４００は、検索セッションの統計情報を保持するためのテーブルであり、文書ＩＤ１４０１、カテゴリ１４０２、セッション閲覧数１４０３の項目からなる

20

文書ＩＤ１４０１には、検索でヒットした文書を特定するための項目であり、ヒットした文書のＩＤが登録される。

【０１３２】

カテゴリ１４０２には、該文書のカテゴリが登録される。

【０１３３】

セッション閲覧数１４０３には、ユーザが検索セッション中に該文書を閲覧した回数を記録する。

【０１３４】

以上、本実施形態について示したが、本発明は、例えば、システム、装置、方法、プログラムもしくは記録媒体等としての実施態様をとることが可能である。具体的には、複数の機器から構成されるシステムに適用しても良いし、また、一つの機器からなる装置に適用しても良い。

30

【０１３５】

また、本発明におけるプログラムは、図３～図１１に示すフローチャートの処理方法をコンピュータが実行可能なプログラムであり、本発明の記憶媒体は図３～図１１の処理方法をコンピュータが実行可能なプログラムが記憶されている。

【０１３６】

以上のように、前述した実施形態の機能を実現するプログラムを記録した記録媒体を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ（またはＣＰＵやＭＰＵ）が記録媒体に格納されたプログラムを読み出し、実行することによっても本発明の目的が達成されることは言うまでもない。

40

【０１３７】

この場合、記録媒体から読み出されたプログラム自体が本発明の新規な機能を実現することになり、そのプログラムを記録した記録媒体は本発明を構成することになる。

【０１３８】

プログラムを供給するための記録媒体としては、例えば、フレキシブルディスク、ハードディスク、光ディスク、光磁気ディスク、ＣＤ－ＲＯＭ、ＣＤ－Ｒ、ＤＶＤ－ＲＯＭ、磁気テープ、不揮発性のメモリカード、ＲＯＭ、ＥＥＰＲＯＭ、シリコンディスク等を用いることができる。

【０１３９】

50

また、コンピュータが読み出したプログラムを実行することにより、前述した実施形態の機能が実現されるだけでなく、その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）等が実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

【0140】

さらに、記録媒体から読み出されたプログラムが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPU等が実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

10

【0141】

また、本発明は、複数の機器から構成されるシステムに適用しても、ひとつの機器から成る装置に適用しても良い。また、本発明は、システムあるいは装置にプログラムを供給することによって達成される場合にも適用できることは言うまでもない。この場合、本発明を達成するためのプログラムを格納した記録媒体を該システムあるいは装置に読み出すことによって、そのシステムあるいは装置が、本発明の効果を享受することが可能となる。

【0142】

さらに、本発明を達成するためのプログラムをネットワーク上のサーバ、データベース等から通信プログラムによりダウンロードして読み出すことによって、そのシステムあるいは装置が、本発明の効果を享受することが可能となる。なお、上述した各実施形態およびその変形例を組み合わせた構成も全て本発明に含まれるものである。

20

【符号の説明】

【0143】

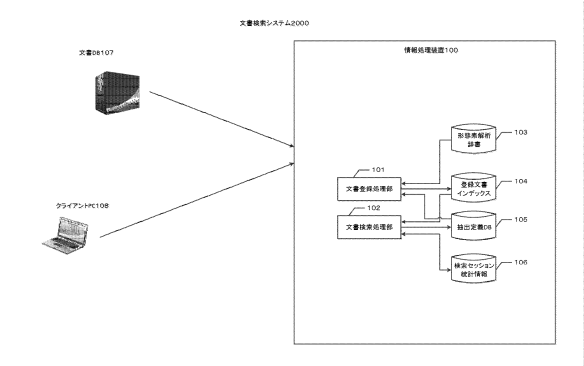
- 2000 文書検索システム
- 100 情報処理装置
- 101 文書登録処理部
- 102 文書検索処理部
- 103 形態素解析辞書
- 104 登録文書インデックス
- 105 抽出定義DB
- 106 検索セッション統計情報
- 107 文書DB
- 108 クライアントPC

30

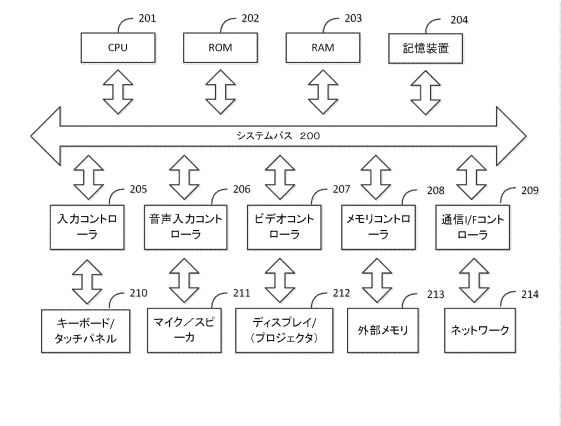
40

50

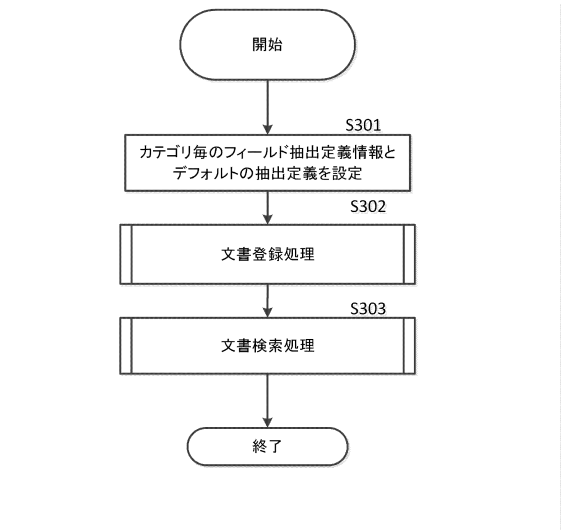
【図面】
【図 1】



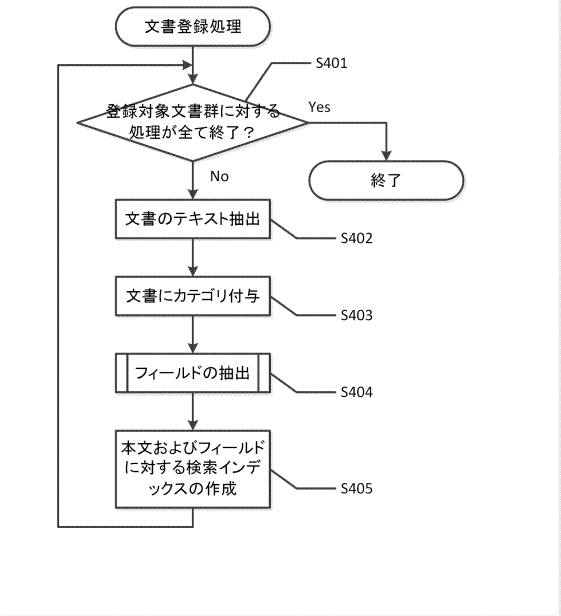
【図 2】



【図 3】



【図 4】



10

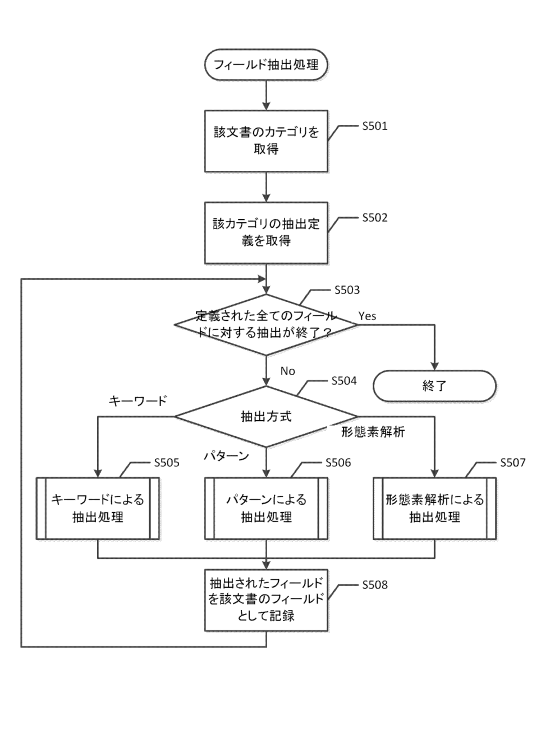
20

30

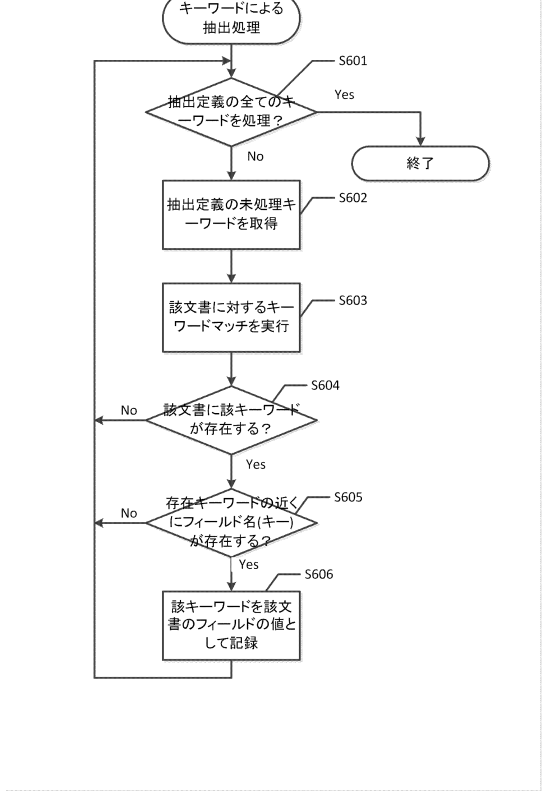
40

50

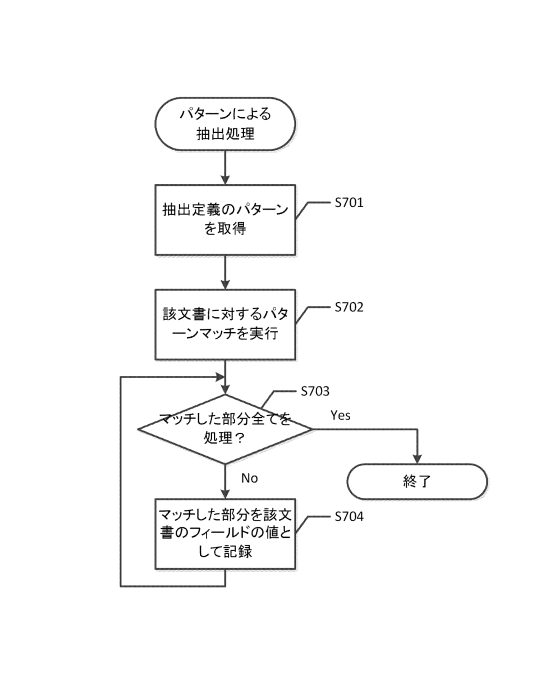
【図 5】



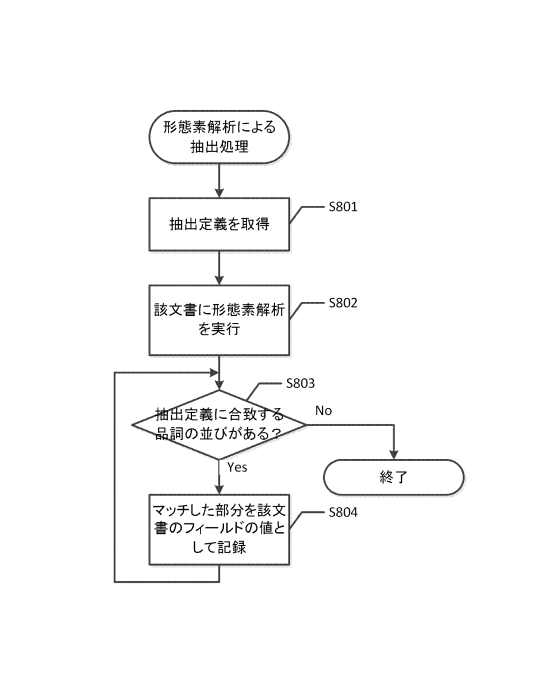
【図 6】



【図 7】



【図 8】



10

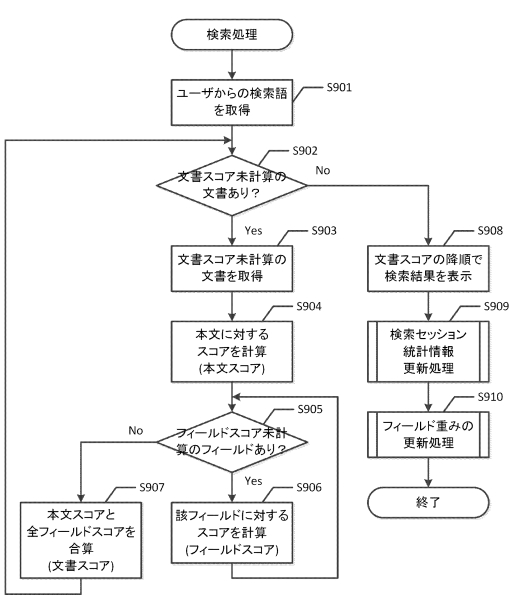
20

30

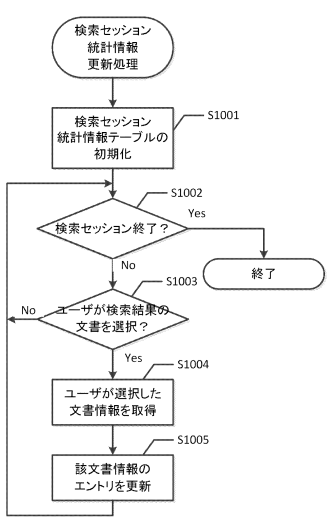
40

50

【図 9】



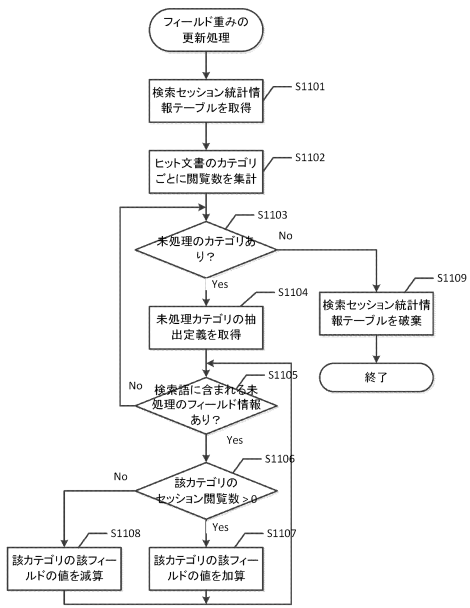
【図 10】



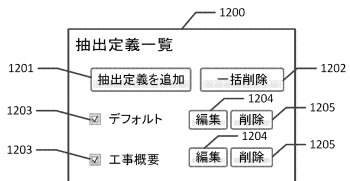
10

20

【図 11】



【図 12】



30

40

50

【図 1 3】

抽出定義詳細 - 工事概要

カテゴリ名

フィールド名

抽出方式

抽出定義

フィールド重み

住所

形態素解析

品詞:名詞-固有名詞-地域:の並び

3

削除

建物用途

キーワード

事務所,病院,飲食店,駐車場,ホテル

2

削除

関連法令

パターン

・(法|条例)

3

削除

抽出定義フィールドを追加

1300

1301

1302

1303

1304

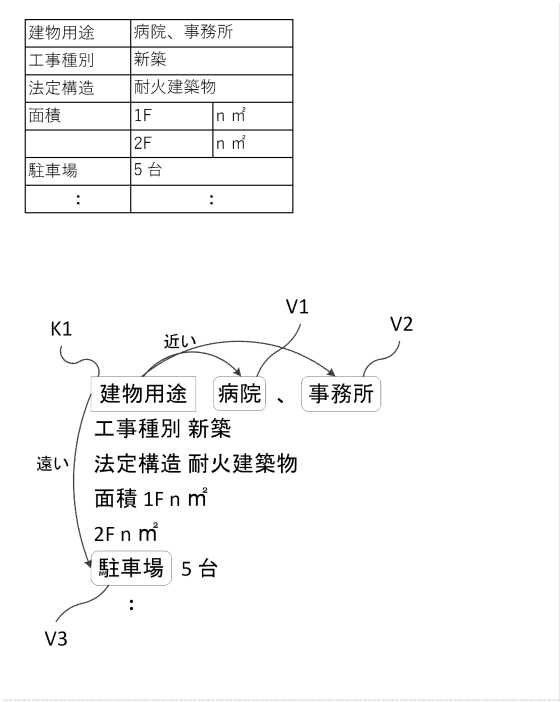
1305

1306

1307

1308

【図 1 4】



10

20

【図 1 5】

文書ID	カテゴリ	セッション閲覧数
AAA	工事概要	1
BBB	注文書	0
CCC	議事録	0
DDD	工事概要	1
:	:	:

1400

1401

1402

1403

1404

【図 1 6】

抽出定義詳細 - 注文書

検索語に「住所」と「建物用途」を含む検索セッション統計情報

文書ID	カテゴリ	セッション閲覧数
AAA	工事概要	1
BBB	注文書	0
CCC	議事録	0
DDD	工事概要	1

抽出定義詳細 - 工事概要

フィールド名

抽出方式

抽出定義

フィールド重み

住所

形態素解析

品詞:名詞-固有名詞-地域:の並び

3

削除

建物用途

キーワード

事務所,病院,飲食店,駐車場,ホテル

2

削除

関連法令

パターン

・(法|条例)

3

削除

抽出定義詳細 - 注文書

フィールド名

抽出方式

抽出定義

フィールド重み

住所

形態素解析

品詞:名詞-固有名詞-地域:の並び

0.1

削除

会社名

パターン

・株式会社|株式会社・

2

削除

抽出定義詳細 - 議事録

フィールド名

抽出方式

抽出定義

フィールド重み

人名

形態素解析

品詞:名詞-人名:の並び

3

削除

会社名

パターン

・株式会社|株式会社・

2

削除

カテゴリ「工事概要」のセッション閲覧数合計: 2
→ カテゴリ「工事概要」のフィールド重み
・住所: $3.00 + 0.02 = 3.02$
・建物用途: $2.00 + 0.02 = 2.02$
カテゴリ「注文書」のセッション閲覧数合計: 0
→ カテゴリ「注文書」のフィールド重み
・住所: $0.10 - 0.01 = 0.09$
・建物用途: 変更なし (抽出定義なし)
カテゴリ「議事録」のセッション閲覧数合計: 0
→ カテゴリ「議事録」のフィールド重み
・住所: 変更なし (抽出定義なし)
・建物用途: 変更なし (抽出定義なし)

1500

1501

1502

1600

1601

1700

30

40

50

【図 17】

1800	1801	図 17A		1802
	フィールド	値	重み	
1803	タイトル	大阪府警担当者議事録.docx	2	1806
1804	人名	佐藤 一郎 高橋 二郎	5	1807
1805	本文	大阪府警担当者議事録 日時: 20XX/XX/XX 出席者: 佐藤 一郎 (大阪府警)、 高橋 二郎 (A社) 場所: 会議室A 議題1: 大阪での○○について ...	1	1808

1806

図 17B

フィールド	値	重み
タイトル	○○プロジェクト概要.pdf	2
会社名	AAA株式会社	5
人名	鈴木 太郎	5
住所	大阪府○○市○○0-0-0	5
本文	名称: ○○プロジェクト 顧客会社名: AAA株式会社 顧客責任者: 鈴木 太郎 住所: 大阪府○○市○○0-0-0 ...	1

10

20

30

40

50

フロントページの続き

審査官 酒井 恭信

(56)参考文献 米国特許出願公開第 2 0 1 9 / 0 3 2 5 2 1 2 (U S , A 1)

中国特許出願公開第 1 1 0 3 9 0 0 9 4 (C N , A)

特開平 1 0 - 0 4 9 5 4 9 (J P , A)

米国特許第 0 6 1 5 4 7 3 7 (U S , A)

(58)調査した分野 (Int.Cl. , D B 名)

G 0 6 F 1 6 / 0 0 - 1 6 / 9 5 8