

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5743938号
(P5743938)

(45) 発行日 平成27年7月1日(2015.7.1)

(24) 登録日 平成27年5月15日(2015.5.15)

(51) Int.Cl.	F I
G06F 17/30 (2006.01)	G06F 17/30 340Z
	G06F 17/30 170A

請求項の数 15 (全 16 頁)

(21) 出願番号	特願2012-69750 (P2012-69750)	(73) 特許権者	000005108
(22) 出願日	平成24年3月26日(2012.3.26)		株式会社日立製作所
(65) 公開番号	特開2013-200795 (P2013-200795A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成25年10月3日(2013.10.3)	(74) 代理人	100091096
審査請求日	平成26年6月4日(2014.6.4)		弁理士 平木 祐輔
		(74) 代理人	100105463
			弁理士 関谷 三男
		(74) 代理人	100102576
			弁理士 渡辺 敏章
		(72) 発明者	今一 修
			東京都国分寺市東恋ヶ窪一丁目280番地
			株式会社日立製作所 中央研究所内
		審査官	野崎 大進

最終頁に続く

(54) 【発明の名称】 連想検索システム、連想検索サーバ及びプログラム

(57) 【特許請求の範囲】

【請求項1】

検索要求文書を入力する入力手段と、検索された検索結果を表示する検索結果表示手段とを少なくとも有する検索クライアントと、

複数の文書を格納した文書データベースと、

受信した検索要求文書に対する関連度が高い文書を前記文書データベースから検索する検索手段と、与えられた文書群から特徴単語群を抽出すると共に、抽出された特徴単語群を各単語の重要度とその出現位置情報とに基づいて1つ又は複数の特徴単語群に分類する特徴単語抽出手段とを有し、複数の特徴単語群が抽出された場合、分類後の特徴単語群のそれぞれについて関連度が高い文書を前記文書データベースから検索する連想検索サーバと

を有する連想検索システム。

【請求項2】

請求項1に記載の連想検索システムにおいて、

特徴単語群の分類数は、ユーザがインタフェースを通じて任意に入力することを特徴とする連想検索システム。

【請求項3】

請求項2に記載の連想検索システムにおいて、

前記インタフェースは、特徴単語群の分類数の変更を指示するボタンを有することを特徴とする連想検索システム。

10

20

【請求項 4】

請求項 1 に記載の連想検索システムにおいて、
特徴単語群の分類数は、分類スコアと閾値の比較により自動設定される
ことを特徴とする連想検索システム。

【請求項 5】

請求項 1 に記載の連想検索システムにおいて、
前記特徴単語抽出手段は、複数回出現する単語が被覆する範囲を求め、被覆度の少ない
箇所の特徴単語群を分類する
ことを特徴とする連想検索システム。

【請求項 6】

請求項 1 に記載の連想検索システムにおいて、
前記特徴単語抽出手段は、複数回出現する単語の重心位置を求め、その重心位置を中心
に特徴単語群を分類する
ことを特徴とする連想検索システム。

【請求項 7】

請求項 1 に記載の連想検索システムにおいて、
前記検索クライアントは、検索された文書群の特徴単語を表示する特徴単語表示手段を
有する
ことを特徴とする連想検索システム。

【請求項 8】

検索クライアントから入力された検索要求文書に類似する文書を複数の文書を格納した
文書データベースから検索する連想検索サーバにおいて、
受信した検索要求文書に対する関連度が高い文書を前記文書データベースから検索する
検索手段と、
与えられた文書群から特徴単語群を抽出すると共に、抽出された特徴単語群を各単語の
重要度とその出現位置情報とに基づいて 1 つ又は複数の特徴単語群に分類する特徴単語抽
出手段と
を有し、複数の特徴単語群が抽出された場合、分類後の特徴単語群のそれぞれについて
関連度が高い文書を前記文書データベースから検索する連想検索サーバ。

【請求項 9】

請求項 8 に記載の連想検索サーバにおいて、
前記特徴単語抽出手段は、ユーザがインタフェースを通じて任意に指定した分類数に基
づいて特徴単語群を分類する
ことを特徴とする連想検索サーバ。

【請求項 10】

請求項 9 に記載の連想検索サーバにおいて、
前記インタフェースは、特徴単語群の分類数の変更を指示するボタンを有する
ことを特徴とする連想検索サーバ。

【請求項 11】

請求項 8 に記載の連想検索サーバにおいて、
前記特徴単語抽出手段は、分類スコアと閾値の比較により、特徴単語群の分類数を自動
設定する。
ことを特徴とする連想検索サーバ。

【請求項 12】

検索クライアントから入力された検索要求文書に類似する文書を複数の文書を格納した
文書データベースから検索する連想検索サーバとして機能するコンピュータに、
受信した検索要求文書に対する関連度が高い文書を前記文書データベースから検索する
第 1 の処理と、
与えられた文書群から特徴単語群を抽出すると共に、抽出された特徴単語群を各単語の
重要度とその出現位置情報とに基づいて 1 つ又は複数の特徴単語群に分類する第 2 の処理

10

20

30

40

50

と、

複数の特徴単語群が抽出された場合、分類後の特徴単語群のそれぞれについて関連度が高い文書を前記文書データベースから検索する第3の処理と
を実行させるプログラム。

【請求項13】

請求項12に記載のプログラムにおいて、
前記第2の処理は、ユーザがインタフェースを通じて任意に指定した分類数に基づいて特徴単語群を分類することを特徴とするプログラム。

【請求項14】

請求項13に記載のプログラムにおいて、
前記インタフェースは、特徴単語群の分類数の変更を指示するボタンを有することを特徴とするプログラム。

【請求項15】

請求項12に記載のプログラムにおいて、
前記第2の処理は、分類スコアと閾値の比較により、特徴単語群の分類数を自動設定することを特徴とするプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、検索要求として与えられた文書に関連する文書を検索する連想検索システムに関し、特に、与えられた文書中の特徴単語の出現位置の情報を用いる連想検索システム、連想検索サーバ及びそれらを実現するプログラムに関する。

【背景技術】

【0002】

コンピュータやインターネットの普及に伴い、文書情報の電子化が急速に進んでいる。一方、入手可能な情報の増加に伴い、それらの中から必要な情報を探し出すことが重要な課題となってきた。また、複数の文書データベース間での文書群の関連性を調べたいという要求も高まっている。例えば、興味のある新聞記事に対し、それらに関連する百科事典の項目を検索したいという要求は多い。

【0003】

現在実用化されているキーワード検索技術の場合、複数の文書データベースを切り替えて検索することは可能であるが、ある文書データベースに含まれる文書群に対し、それに関連する文書群を、同一の文書データベース、あるいは、別の文書データベースから検索すること（文書連想検索と呼ばれる検索方式）は不可能である。

【0004】

同一の文書データベースに限れば、文書間の類似度を予め計算しておくことにより、文書群を検索入力とした文書連想検索を実現することはできる。しかし、複数の文書データベース間での文書連想検索を実現しようとする、予め計算すべき文書間の関連度の組み合わせ数が、文書データベース数の増加に伴って爆発的に増加する。このため、文書間の類似度を予め計算する方法による文書連想検索の現実化は不可能である。

【0005】

これに対し、特許文献1には、利用者が指定した文書データベース中の任意の文書群に対して、その文書群に関連する文書群を任意の文書データベースから効率よく検索するための方法が開示されている。

【0006】

特許文献1に開示の方法は、文書群として入力された検索入力内の特徴的な単語群（特徴単語群）のみを使用し、高速な文書連想検索を実現する。この方法を用いれば、利用者は、複数の異なる種類の文書データベースを切り替えながら、文書群の関連性を調べるこ

10

20

30

40

50

とができ、高精度かつ効率的に文書を検索することができる。また、この方法は、検索結果として得られた文書群に出現する特徴単語群を抽出し、それらを検索結果の概観（要約）として利用者に提示することにより、利用者による検索結果の可否の判断を支援する技術も提供する。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開2000-155758号公報

【発明の概要】

【発明が解決しようとする課題】

10

【0008】

一般に、単語に基づく文書検索では、文書中に出現する単語によって文書のインデックス付けを行ない、文書検索を実現する。特許文献1の場合も同様であり、文書から特徴単語群を抽出する際には、文書に含まれる単語の統計的尺度（tf*idf法などが代表的）を用いて重要度を計算し、重要度の高い順に単語を抽出し、連想検索を実現する。

【0009】

しかし、従来の連想検索では、特徴単語を抽出する対象は文書全体である。このため、文書に複数の話題が含まれている場合には、複数の話題の特徴単語が混在した状態のまま単語が抽出される。つまり、複数の話題を総合的に判断して類似した文書が検索される。このため、利用者が望んだ結果が必ずしも得られるとは限らない。

20

【0010】

この技術課題を鑑み、本発明は、検索入力となる文書群中に含まれる話題ごとに類似する文書を検索できる連想検索システムを提供する。

【課題を解決するための手段】

【0011】

このために、本発明においては、連想検索における検索入力文書から特徴単語群を抽出する際に、各単語の重要度だけでなく、その単語の検索入力文書中での位置情報も付加して抽出処理を実行する。次に、抽出した特徴単語群を、各単語の重要度と出現位置に基づいて分類する。特徴単語群の分類数は、検索入力文書中での特徴単語の重要度と距離に応じて分類する際の分類スコアに閾値を設定して自動的に設定してもよいし、利用者がユーザインタフェース上で分類数を任意に設定してもよい。最後に、分類結果として得られた特徴単語群のそれぞれを検索入力として検索を実行する。

30

【発明の効果】

【0012】

本発明によれば、複数の話題が含まれている文書群を検索入力とする場合でも、文書群全体として類似した文書ではなく、分類された特徴単語群（文書中に含まれる話題に相当）毎に類似した文書を連想検索結果として得ることができる。これにより、利用者の希望により近い結果を提示することができる。前述した以外の課題、構成及び効果は、以下の実施の形態の説明により明らかにされる。

【図面の簡単な説明】

40

【0013】

【図1】連想検索システムの概念的な構成例を示す図。

【図2】連想検索スレーブサーバの構成例を示す図。

【図3】検索手段の構成例を示す図。

【図4】特徴単語抽出手段の構成例を示す図。

【図5】文書中における特徴単語の分布例を示す図。

【図6】特徴単語の抽出及び分類方法の例を示す図。

【図7】特徴単語の抽出及び分類方法の例を示す図。

【図8】検索クライアントにおける初期画面の例を示す図。

【図9】検索クライアントにおける検索結果の表示例を示す図。

50

【図10】検索クライアントにおける検索結果の表示例を示す図。

【図11】検索クライアントにおける検索結果の表示例を示す図。

【図12】特徴単語群分類結果の確認画面の表示例を示す図。

【図13】インデックス付けの一例を示す図。

【図14】検索クライアント、連想検索マスタサーバ、連想検索スレーブサーバ間におけるデータ及び処理の流れを示す図。

【図15】検索クライアント、連想検索マスタサーバ、連想検索スレーブサーバ間におけるデータ及び処理の流れを示す図。

【図16】検索クライアント、連想検索マスタサーバ、連想検索スレーブサーバ間におけるデータ及び処理の流れを示す図。

10

【発明を実施するための形態】

【0014】

以下、図面に基づいて、本発明の実施の形態を説明する。なお、本発明の実施の態様は、後述する形態例に限定されるものではなく、その技術思想の範囲において、種々の変形が可能である。

【0015】

図1は、形態例に係る連想検索システムの概略構成を示している。このシステムは、利用者による検索要求の入力及び検索結果の表示に使用される検索クライアント20と、文書データベースを検索する連想検索スレーブサーバ40、50、60と、連想検索クライアント20と連想検索スレーブサーバ40、50、60を仲介する連想検索マスタサーバ30と、これらを接続する通信ネットワーク10とで構成される。

20

【0016】

図1の例では、文書データベースを検索する連想検索スレーブサーバが通信ネットワーク10に3台接続されている場合を表しているが、通信ネットワーク10に接続される連想検索スレーブサーバの数は任意でよい。検索クライアント20の数も任意である。

【0017】

また、図1の例では、検索クライアント10と、連想検索マスタサーバ30と、連想検索スレーブサーバ40、50、60とを通信ネットワーク10を介して接続しているが、これらのうちの幾つかを、あるいは、全てを同一の計算機上に構成してもよい。

【0018】

図2に、連想検索スレーブサーバ40の構成例を示す。他の連想検索スレーブサーバ50、60の構成も、連想検索スレーブサーバ40と同じである。連想検索スレーブサーバ40は、メモリ装置491、演算処理装置492、インタフェース装置493、補助記憶装置494、入力装置495、出力装置496を有し、それぞれがバス490を介して相互に接続されている。

30

【0019】

メモリ装置491は、補助記憶装置494からプログラムを読み出して記憶するRAM(Random Access Memory)等の記憶装置である。メモリ装置491には、検索手段410と特徴単語抽出手段420に対応するプログラム、その実行に必要な検索インデックス430と文書データベース440に対応するファイルやデータ等が記憶される。

40

【0020】

演算処理装置492は、メモリ装置491に格納されたプログラムを実行するCPU(Central Processing Unit)等の演算処理装置である。インタフェース装置493は、外部ネットワーク等に接続するためのインタフェース装置である。補助記憶装置494は、検索手段410と特徴単語抽出手段420に対応するプログラム、検索インデックス430と文書データベース440に対応するファイルやデータ等を記憶するHDD(Hard Disk Drive)等の記憶装置である。入力装置495は、ユーザインタフェースを提供する装置(例えば、キーボード、マウス)である。出力装置496は、ユーザインタフェースを提供する出力装置(例えば、ディスプレイ装置)である。

【0021】

50

図2は、連想検索スレーブサーバの構成を示す図であるが、検索クライアント20と連想検索マスタサーバ30の構成も、補助記憶装置に記憶されるプログラムやデータの違いを除き、同様に構成される。

【0022】

図3に、連想検索スレーブサーバ40が備える検索手段410の機能ブロック構成を示す。プログラムとしての検索手段410は、単語頻度取得手段411、位置情報取得手段412、関連度計算手段413、近似性計算手段414、スコア計算手段415の各機能により構成される。これらの検索手段410を構成する各手段もプログラム処理を通じて提供される。

【0023】

連想検索スレーブサーバ40は、連想検索マスタサーバ30が備える検索要求発行手段320から送られてきた検索要求に対し、関連度の高い文書群を文書データベース440から検索し、その検索結果を関連度のスコア付きで連想検索マスタサーバ30に返す。ここでの検索は、例えば公知のキーワード検索手法により実現することができる。

【0024】

キーワード検索手法では、検索処理の効率を上げるために、文書データベースに含まれる文書を単語に分割し（日本語の文書に対しては形態素解析を実行し、英語の文書に対してはステミング処理を実行する）、どの文書にどの単語が含まれているかを示す検索インデックスを事前に作成する。後述する本実施例の検索方法のように、検索時に位置情報も用いる場合には、各単語の出現位置もインデックスに格納しておく。検索実行時には、事前に作成された検索インデックスを用いることで、検索処理を高速に実行することができる。

【0025】

図1の場合には、連想検索スレーブサーバ40、50、60が有する文書データベース440、540、640のそれぞれについて、検索インデックス430、530、630を事前に作成し、検索処理に利用する。

【0026】

検索要求と検索対象文書間の関連度の計算は、以下の手順で実行される。まず、検索手段410が、連想検索マスタサーバ30の検索要求発行手段320から送信された検索要求を受信する。検索手段410は、受信した検索要求に含まれる単語群を含む文書を検索する。単語頻度取得手段411は、検索結果として得られた文書のそれぞれについて、各文書に含まれる単語群のうち検索要求に含まれる単語群の頻度情報を取得する。次に、関連度計算手段413は検索要求とその文書の関連度を計算する。関連度の計算方法は任意でよい。例えば公知の技術である $tf \cdot idf$ 法により単語の重要度を計算し、その総和を関連度とする。単語の近接性を検索スコアに反映する場合には、位置情報取得手段412が、各文書に含まれる単語群のうち検索要求に含まれる単語群の出現位置情報を取得し、近接性計算手段414が近接スコアを計算する。近接スコアの計算方法は任意でよい。例えば、検索要求に含まれる単語群がどれくらい密集して出現しているかを計算し、その計算結果を近接スコアとする。スコア計算手段415は、関連度計算手段413と近接性計算手段414のそれぞれから得られたスコアを統合し、統合後のスコアを関連度として文書に付与する。

【0027】

図4に、連想検索スレーブサーバ40が備える特徴単語抽出手段420の機能ブロック構成を示す。プログラムとしての特徴単語抽出手段420は、単語頻度取得手段421、位置情報取得手段422、重要度計算手段423、近接性クラスタリング手段424、単語追加手段425の各機能により構成される。これらの特徴単語抽出手段420を構成する各手段もプログラム処理を通じて提供される。

【0028】

特徴単語抽出手段420は、連想検索マスタサーバ30が備える特徴単語要求手段330から送られてきた文書群に対する特徴単語を、文書データベース440から抽出する。

10

20

30

40

50

特徴単語抽出手段 4 2 0 は、特徴単語の高速抽出を実現するために、検索手段 4 1 0 と同様、検索インデックス 4 3 0 を利用する。すなわち、特徴単語抽出手段 4 2 0 は、ある文書にどの単語が含まれているかを、検索インデックス 4 3 0 を参照して調べる。

【 0 0 2 9 】

特徴単語の抽出は、以下の手順で実行される。まず、特徴単語抽出手段 4 2 0 が、連想検索マスタサーバ 3 0 の特徴単語要求手段 3 3 0 から送信された文書群を受信する。単語頻度取得手段 4 2 1 は、受信した文書群に含まれる各単語の頻度情報を取得する。取得された頻度情報に基づいて、重要度計算手段 4 2 3 は、各単語の重要度を計算する。重要度の計算方法は任意でよい。例えば公知の技術である $tf \cdot idf$ 法により単語の重要度を計算する。位置情報を用いない連想検索の場合、特徴単語抽出手段 4 2 0 は、高い重要度が付された単語から順番に特徴単語として連想検索マスタサーバ 3 0 に返す。

10

【 0 0 3 0 】

本実施の形態では、位置情報取得手段 4 2 2 が、重要度付きの各単語について出現位置情報を取得する。さらに、近接性クラスタリング手段 4 2 4 が、重要度と位置情報とに基づいて検索された単語群を分類する。さらに、単語追加手段 4 2 5 が、分類結果のそれぞれに含まれる単語群に近接する単語を追加する。特徴単語抽出手段 4 2 0 は、このようにして得られた特徴単語群の集合を連想検索マスタサーバ 3 0 に返す。単語追加手段 4 2 5 の使用は任意でよい。

【 0 0 3 1 】

次に、近接性クラスタリング手段 4 2 4 の動作を図 5、図 6、図 7 を用いて説明する。図 5 は、同じ単語が含まれる二つの文書 1、文書 2 を例示している。文書 1 では、各単語 (term1 ~ term6) が文書全体に分散して分布しているのに対し、文書 2 では、term1 ~ term3 が文書中の前半に、term4 ~ term6 が文書中の後半に集中して分布している。

20

【 0 0 3 2 】

このような場合でも、従来の連想検索では、特徴単語の出現位置を考慮していないため、文書 1 を検索入力として連想検索を実行した場合の結果と、文書 2 を検索入力として連想検索を実行した場合の結果は、同じである。しかし、特徴単語の文書中での分布が偏っている場合、複数の話題について書かれている可能性があるため、特徴単語群を個々の話題に分類することが望ましい。

【 0 0 3 3 】

図 6 は、文書 1 から特徴単語群を抽出する場合の例である。この場合、各特徴単語は、文書全体に分散して分布しているため、一つの話題について書かれていると考えられる。従って、位置情報に基づいて特徴単語群を分類しても、分類することができず、一つの特徴単語群となる。

30

【 0 0 3 4 】

図 7 は、文書 2 から特徴単語群を抽出する場合の例である。この場合、各特徴単語は、文書の前半に term1 ~ term3、文書の後半に term4 ~ term6 が集中して分布しているため、二つの話題について書かれていると考えられる。従って、位置情報に基づいて特徴単語群を分類すると、term1 ~ term3 の特徴単語群と、term4 ~ term6 の特徴単語群の二つの特徴単語群が抽出される。

40

【 0 0 3 5 】

近接性クラスタリング手段 4 2 4 による特徴単語群の分類には、例えば、単語の出現位置とその重みを用いる階層的クラスタリング手法を適用すればよい。複数回出現する単語については、予め、その重心位置を求めておく。その後、各単語の位置に基づいて、最も近接する単語をまとめあげる。その際、それぞれに単語の重みを考慮して、新しい重心を決定する。この処理を繰り返すことでクラスタリング結果を得る。

【 0 0 3 6 】

あるいは、別の手法として、複数回出現する単語が文書中のどの範囲を被覆するかを求め、文書全体における被覆度の少ない箇所で特徴単語群を分類してもよい。

【 0 0 3 7 】

50

前述の説明では、近接性クラスタリング手段424における特徴単語分類手法として、二つの手法について説明したが、位置情報に基づいて特徴単語群を分類する手法であれば任意のものを用いてもよい。

【0038】

このようにして得られた特徴単語群を用いて連想検索を実行することにより、文書中に複数の話題が含まれている場合でも、利用者の望んだ検索結果を得ることが可能となる。

【0039】

図8は、検索クライアント20が備える検索要求入力手段210により提供される画面例を表している。利用者は、検索要求入力エリア211に検索要求を入力し、検索指示ボタン212をクリックすることにより検索の実行を検索クライアント20に指示する。

10

【0040】

図9は、検索クライアント20による検索結果の表示例である。検索結果は、検索結果表示手段220により表示され、検索結果から抽出された特徴単語群が特徴単語表示手段230により表示される。特徴単語表示手段230を用いるか否かは任意である。検索結果表示手段220は文書群指定手段も兼ねている。文書選択チェックボックス221により任意個の文書を選択した状態で、連想検索指示ボタン213をクリックすると、選択した文書と関連する文書を検索することができる。特徴単語表示手段230は、単語群指定手段も兼ねている。単語選択チェックボックス231により任意個の単語を選択した状態で、連想検索指示ボタン213をクリックすると、特徴単語からの検索を実行することができる。分類数指定手段240は、文書を選択して連想検索を実行する場合に、文書中に含まれる話題を何個に分割するかを指定入力するために用いられる。分類数は、数値として直接指定してもよいし、スライダーやボタン等を用いて指定してもよい。また、分類数は、分類スコアと閾値との比較を通じて自動的に設定してもよい。分類スコアは、特徴単語の重要度のスコアと近接度合のスコアを統合したスコアとして規定する。分類数を閾値により自動設定する場合には、分類数指定手段240を画面に表示しなくてもよい。

20

【0041】

図10は、検索入力として与えられた文書に二つの話題が含まれている場合の検索結果の例である。この場合、検索結果表示手段220には、二列に分けて、それぞれの話題に関する検索結果が表示される。左列の記事1～5が話題1に対応し、右列の記事A～Eが話題2に対応する。なお、図10の場合、特徴単語表示手段230には、二つの話題の検索結果を統合して、そこから特徴単語群を抽出した結果を表示している。

30

【0042】

一方、図11は、検索入力として与えられた文書に二つの話題が含まれている点は図10と同じであるが、特徴単語表示手段230に、各話題の検索結果ごとに特徴単語群を抽出し、それぞれを二列に表示している。左列の特徴ターム1～5が話題1に対応し、右列の特徴タームA～Eが話題2に対応する。図9の場合と同様、特徴単語表示手段230を用いるか否かは任意である。

【0043】

図12は、近接性クラスタリング手段424が分類した特徴単語群を確認する画面である。利用者は、この画面を用いて、分類された特徴単語群が適切かどうかを判断し、適切であれば検索指示ボタン213をクリックする。適切でなければ、利用者は、分類数指定手段240に新たな分類数を指定し、その後、分類数変更指示ボタン241をクリックし、再度、分類された特徴単語群を確認する。なお、この画面の使用は任意である。

40

【0044】

図13は、文書データベース440、540、640に含まれる文書から検索インデックス430、530、630を作成した場合の検索インデックスの例である。文書IDの列に個々の文書を識別する識別子、その識別子に該当する文書に含まれる単語の出現位置の情報が格納されている。

【0045】

次に、実施の形態に係る連想検索システムで実行される処理の流れを、図14のシーケ

50

ンス図を用いて説明する。以下では、連想検索スレーブサーバとして連想検索スレーブサーバ40を用いる場合を説明する。

【0046】

利用者は、検索クライアント20が備える検索要求入力手段210を用い、検索要求を入力する。入力された検索要求は、検索クライアント20から連想検索マスターサーバ30に送信される(T11)。

【0047】

連想検索マスターサーバ30の検索要求解析手段310は検索要求を解析し、連想検索スレーブサーバ40に送信するための検索要求を作成する。検索要求発行手段320により、検索要求が連想検索スレーブサーバ40に送信される(T12)。

10

【0048】

連想検索スレーブサーバ40が備える検索手段410は、検索インデックス430を用いて文書データベース440を検索し、その結果を連想検索マスターサーバ30に返す(T13)。

【0049】

連想検索マスターサーバ30の特徴単語要求手段330は、得られた検索結果から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T14)。

【0050】

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスターサーバ30へ返す(T15)。

20

【0051】

最後に、検索結果と特徴単語群が連想検索マスターサーバ30から検索クライアント20に送信され(T16)、検索クライアント20の検索結果表示手段220と特徴単語表示手段230によって利用者に提示される。

【0052】

次に、図15に示すシーケンス図について説明する。このシーケンス図は、検索結果として得られた文書群から連想検索を実行する場合の処理の流れを示している。

【0053】

利用者は、検索クライアント20が備える文書群指定手段220を用いて、検索入力となる文書群を選択する。選択された文書群の識別子は連想検索マスターサーバ30に送信される(T21)。

30

【0054】

連想検索マスターサーバ30の特徴単語要求手段330は、選択された文書群から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T22)。

【0055】

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスターサーバ30へ返す(T23)。

【0056】

連想検索マスターサーバ30の検索要求発行手段320は、得られた特徴単語群を連想検索スレーブサーバに送信する(T24)。

40

【0057】

連想検索スレーブサーバ40が備える検索手段410は、検索インデックス430を用いて文書データベース440を検索し、その結果を連想検索マスターサーバ30に返す(T25)。

【0058】

連想検索マスターサーバ30の特徴単語要求手段330は、得られた検索結果から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T26)。

50

【 0 0 5 9 】

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスターサーバ30へ返す(T27)。

【 0 0 6 0 】

最後に、検索結果と特徴単語群が連想検索マスターサーバ30から検索クライアント20に送信され(T28)、検索クライアント20の検索結果表示手段220と特徴単語表示手段230によって利用者に提示される。

【 0 0 6 1 】

次に、図16に示すシーケンス図について説明する。このシーケンス図は、検索結果として得られた文書群から連想検索を実行する場合の処理の流れを示しており、かつ、得られた文書群に二つの話題が含まれている場合を示している。

10

【 0 0 6 2 】

利用者は、検索クライアント20が備える文書群指定手段220を用いて、検索入力となる文書群を選択する。選択された文書群の識別子は、検索クライアント20から連想検索マスターサーバ30に送信される(T31)。

【 0 0 6 3 】

連想検索マスターサーバ30の特徴単語要求手段330は、選択された文書群から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T32)。

【 0 0 6 4 】

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスターサーバ30へ返す(T33)。

20

【 0 0 6 5 】

連想検索マスターサーバ30の検索要求発行手段320は、得られた二つの特徴単語群のうち一つ目の話題に相当する特徴単語群を連想検索スレーブサーバ40に送信する(T341)。

【 0 0 6 6 】

連想検索スレーブサーバ40が備える検索手段410は、検索インデックス430を用いて文書データベース440を検索し、その結果を連想検索マスターサーバ30に返す(T351)。

30

【 0 0 6 7 】

連想検索マスターサーバ30の特徴単語要求手段330は、得られた検索結果から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T361)。

【 0 0 6 8 】

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスターサーバ30へ返す(T371)。

【 0 0 6 9 】

次に、連想検索マスターサーバ30の検索要求発行手段320は、得られた二つの特徴単語群のうち二つ目の話題に相当する特徴単語群を連想検索スレーブサーバ40に送信する(T342)。

40

【 0 0 7 0 】

連想検索スレーブサーバ40が備える検索手段410は、検索インデックス430を用いて文書データベース440を検索し、その結果を連想検索マスターサーバ30に返す(T352)。

【 0 0 7 1 】

連想検索マスターサーバ30の特徴単語要求手段330は、得られた検索結果から特徴単語を抽出するために、特徴単語の抽出要求を連想検索スレーブサーバ40に送信する(T362)。

【 0 0 7 2 】

50

連想検索スレーブサーバ40が備える特徴単語抽出手段420は、検索インデックス430を利用して特徴単語群を抽出し、連想検索マスタサーバ30へ返す(T372)。

【0073】

最後に、検索結果と特徴単語群が連想検索マスタサーバ30から検索クライアント20に送信され(T28)、検索クライアント20の検索結果表示手段220と特徴単語表示手段230によって利用者に提示される。

【0074】

話題が三つ以上ある場合には、T33の後の検索要求発行手段 T341 検索手段 T351 特徴単語要求手段 T361 特徴単語抽出手段 T371と同様の処理を必要な回数繰り返せばよい。

10

【0075】

図10に示したように二つの話題の検索結果全体から特徴単語を抽出する場合は、図16のシーケンス図において、T351の後の特徴単語要求手段 T361 特徴単語抽出手段 T371を省略し、T352の後の特徴単語要求手段において、二つの話題の検索結果全体の文書群を連想検索スレーブサーバ40に送信すればよい。

【0076】

なお、本発明は上述した形態例に限定されるものでなく、様々な変形例が含まれる。例えば、上述した形態例は、本発明を分かりやすく説明するために詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、ある形態例の一部を他の形態例の構成に置き換えることが可能であり、また、ある形態例の構成に他の形態例の構成を加えることも可能である。また、各形態例の構成の一部について、他の構成を追加、削除又は置換することも可能である。

20

【0077】

また、上述した各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路その他のハードウェアとして実現することも可能である。

【符号の説明】

【0078】

- 10：通信ネットワーク
- 20：検索クライアント
- 210：検索要求入力手段
- 211：検索要求入力エリア
- 212：検索指示ボタン
- 213：連想検索指示ボタン
- 220：検索結果表示手段(文書群指定手段)
- 221：文書選択チェックボックス
- 230：特徴単語表示手段(単語群指定手段)
- 231：単語選択チェックボックス
- 240：分類数指定手段
- 241：分類数変更指示ボタン
- 30：連想検索マスタサーバ
- 310：検索要求解析手段
- 320：検索要求発行手段
- 330：特徴単語要求手段
- 40：連想検索スレーブサーバ
- 410：検索手段
- 411：単語頻度取得手段
- 412：位置情報取得手段
- 413：関連度計算手段
- 414：近接性計算手段
- 415：スコア計算手段

30

40

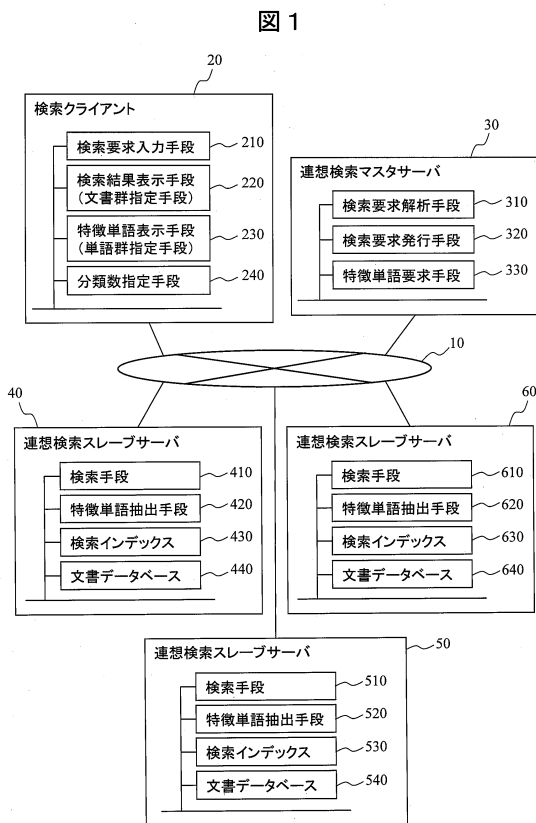
50

- 4 2 0 : 特徴単語抽出手段
- 4 2 1 : 単語頻度取得手段
- 4 2 2 : 位置情報取得手段
- 4 2 3 : 重要度計算手段
- 4 2 4 : 近接性クラスタリング手段
- 4 2 5 : 単語追加手段
- 4 3 0 : 検索インデックス
- 4 4 0 : 文書データベース
- 4 9 0 : バス
- 4 9 1 : メモリ装置
- 4 9 2 : 演算処理装置
- 4 9 3 : インタフェース装置
- 4 9 4 : 補助記憶装置
- 4 9 5 : 入力装置
- 4 9 6 : 出力装置
- 5 0 : 連想検索スレーブサーバ
- 5 1 0 : 検索手段
- 5 2 0 : 特徴単語抽出手段
- 5 3 0 : 検索インデックス
- 5 4 0 : 文書データベース
- 5 0 : 連想検索スレーブサーバ
- 5 1 0 : 検索手段
- 5 2 0 : 特徴単語抽出手段
- 5 3 0 : 検索インデックス
- 5 4 0 : 文書データベース

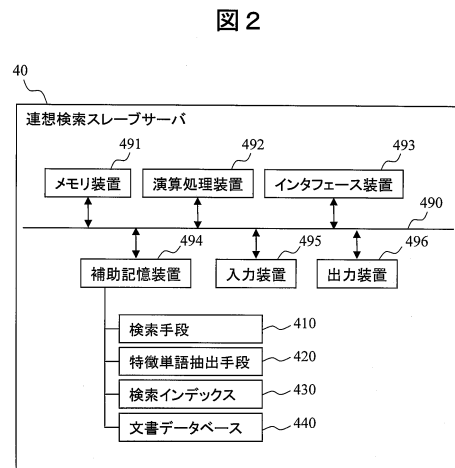
10

20

【図 1】

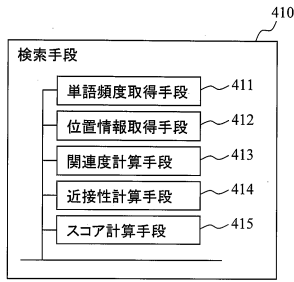


【図 2】



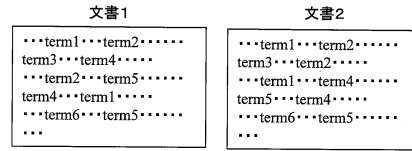
【図3】

図3



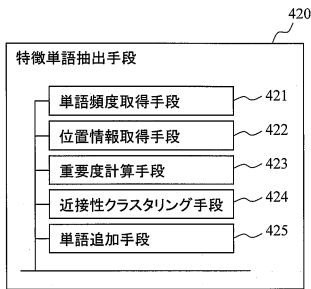
【図5】

図5



【図4】

図4



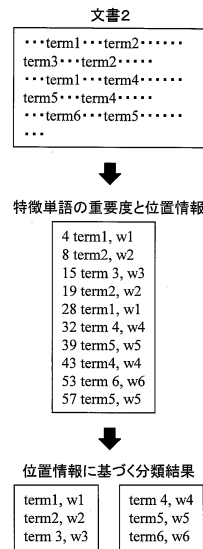
【図6】

図6



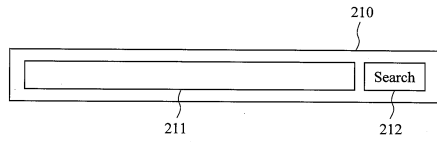
【図7】

図7



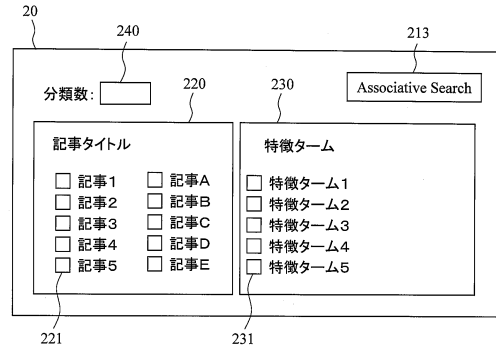
【図8】

図8



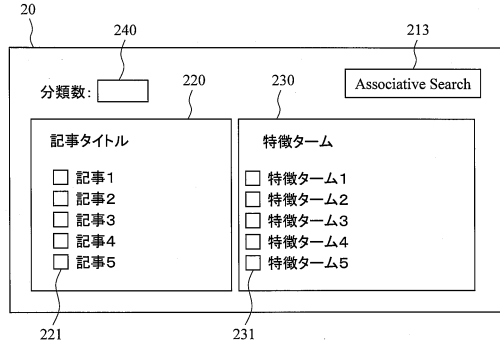
【図10】

図10



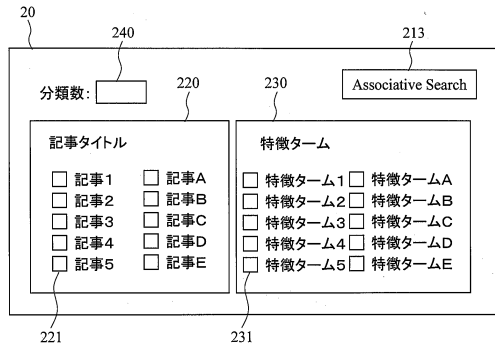
【図9】

図9



【図11】

図11



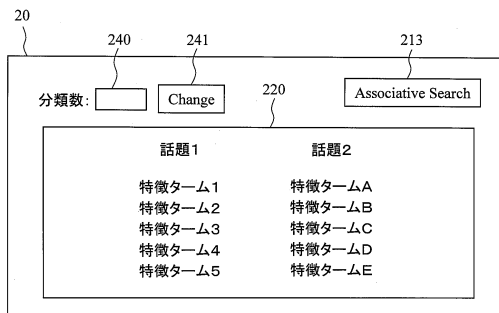
【図13】

図13

文書ID	出現位置
123	1,3,7 term ₁ 2 term ₂ 4 term ₃ ... n term _m
124	...

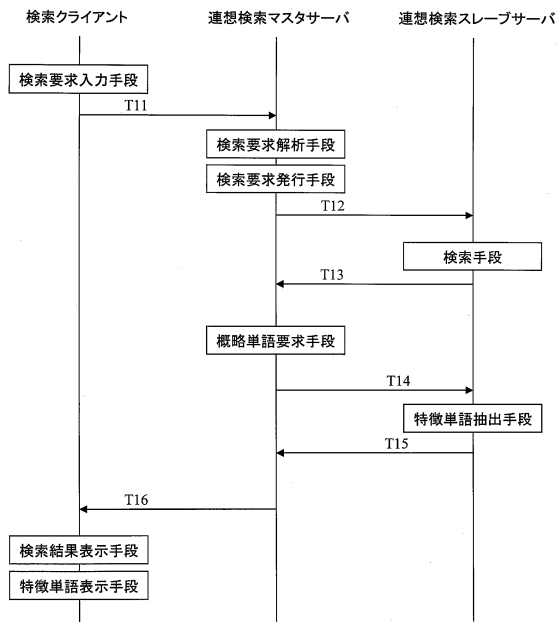
【図12】

図12



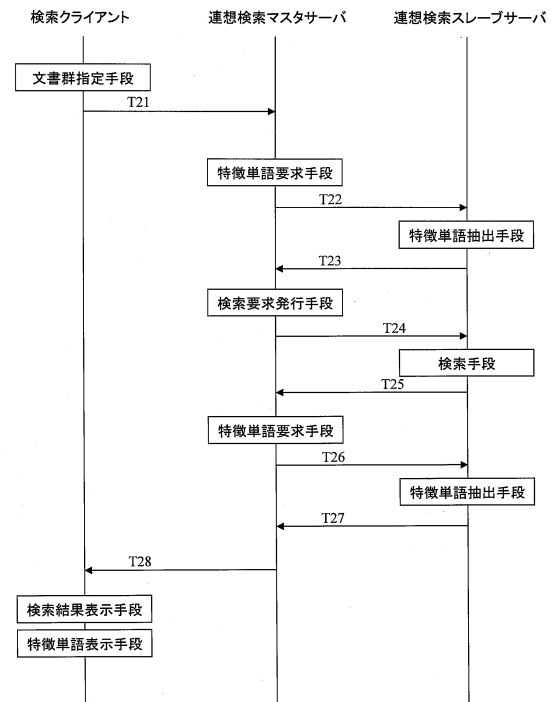
【 図 1 4 】

図 1 4



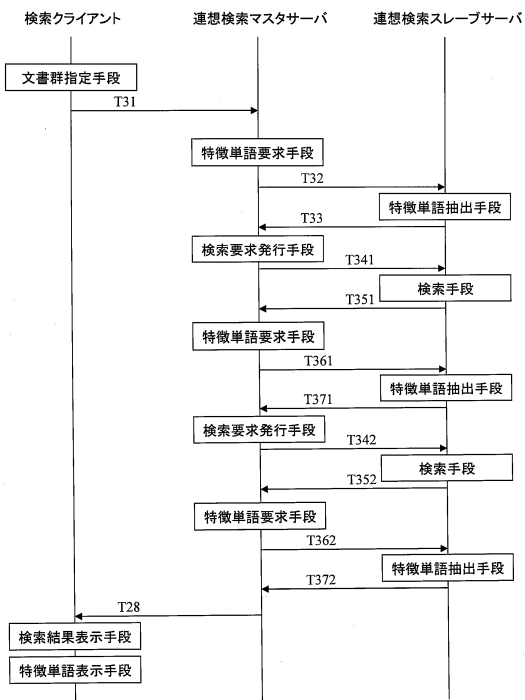
【 図 1 5 】

図 1 5



【 図 1 6 】

図 1 6



フロントページの続き

(56)参考文献 特開2000-010986(JP,A)
特開2002-269140(JP,A)
特開2007-265034(JP,A)
特許第5387870(JP,B2)
欧州特許第00730765(EP,B1)
米国特許出願公開第2009/0169110(US,A1)

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