

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4912384号
(P4912384)

(45) 発行日 平成24年4月11日(2012.4.11)

(24) 登録日 平成24年1月27日(2012.1.27)

(51) Int. Cl. F 1
G 0 6 F 17/30 (2006.01)
 G 0 6 F 17/30 3 2 0 D
 G 0 6 F 17/30 3 8 0 E
 G 0 6 F 17/30 3 4 0 B

請求項の数 3 (全 13 頁)

(21) 出願番号	特願2008-297847 (P2008-297847)	(73) 特許権者	000004226
(22) 出願日	平成20年11月21日(2008.11.21)		日本電信電話株式会社
(65) 公開番号	特開2010-123036 (P2010-123036A)		東京都千代田区大手町二丁目3番1号
(43) 公開日	平成22年6月3日(2010.6.3)	(74) 代理人	100086232
審査請求日	平成22年2月4日(2010.2.4)		弁理士 小林 博通
		(74) 代理人	100104938
			弁理士 鶴澤 英久
		(74) 代理人	100140361
			弁理士 山口 幸二
		(74) 代理人	100096459
			弁理士 橋本 剛
		(72) 発明者	村田 眞哉
			東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 文書検索装置、文書検索方法、および文書検索プログラム

(57) 【特許請求の範囲】

【請求項1】

ユーザ端末から検索指示されたクエリを含む電子文書を検索するときに検索エンジンの検索ログを利用する文書検索装置であって、

前記検索ログに含まれたクエリに応じた検索結果のタイトルおよび概要文から拡張語を生成し、該拡張語を前記クエリの拡張情報として保存するクエリ情報保存手段と、

前記検索ログからクリックされた検索結果の検索時のクエリを判別し、該クエリに関連する拡張語を前記クエリ情報保存手段から求め、検索結果がクリックされたときに投入されたクエリの拡張語群を検索結果の拡張情報として保存する検索結果情報保存手段と、

ユーザ端末から検索指示されたクエリについて、クエリの拡張語を前記クエリ情報保存手段から取得し、ユーザ端末に送る照合処理手段と、

前記検索指示されたクエリを前記拡張語で拡張した拡張クエリの検索で得られた検索結果に対して、前記検索結果情報保存手段に保存された拡張語群を付与して検索結果を拡張し、拡張された検索結果を前記拡張クエリとの類似度により並び替え、この並び替えた結果をリスト化した最終検索結果を前記ユーザ端末に送る検索結果処理手段と、

を備えることを特徴とする文書検索装置。

【請求項2】

ユーザ端末から検索指示されたクエリを含む電子文書を検索するときに検索エンジンの検索ログを利用する文書検索方法であって、

前記検索ログに含まれたクエリに応じた検索結果のタイトルおよび概要文から拡張語を

10

20

生成し、該拡張語をクエリ情報保存手段に前記クエリの拡張情報として保存するクエリ情報保存ステップと、

検索結果情報保存手段が、前記検索ログからクリックされた検索結果の検索時のクエリを判別し、該クエリに関連する拡張語を前記クエリ情報保存手段から求め、検索結果がクリックされたときに投入されたクエリの拡張語群を検索結果の拡張情報として保存する検索結果情報保存ステップと、

照合処理手段が、ユーザ端末から検索指示されたクエリについてクエリの拡張語を前記クエリ情報保存手段から取得し、ユーザ端末に送る照合処理ステップと、

検索結果処理手段が、前記検索指示されたクエリを前記拡張語で拡張した拡張クエリの検索で得られた検索結果に対して、前記検索結果情報保存手段に保存された拡張語群を付与して検索結果を拡張し、拡張された検索結果を前記拡張クエリとの類似度により並び替え、この並び替えた結果をリスト化した最終検索結果を前記ユーザ端末に送る検索結果処理ステップと、

を有することを特徴とする文書検索方法。

【請求項3】

請求項1記載の文書検索装置を構成する各手段としてコンピュータを機能させることを特徴とする文書検索プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、電子文書群中からクエリに該当する電子文書を検索する技術に関する。

【背景技術】

【0002】

インターネットなどで接続されたクライアント端末から検索語（クエリ）を受信して検索結果を返信する検索エンジンの運営サーバには、時々刻々、検索のログが保存される。このログには、投入されたクエリの情報や、検索結果に対するユーザのクリックの情報などが保存される。

【0003】

このような検索エンジンのログを利用した文書検索システムが非特許文献1に提案されている。この文書検索システムは、検索エンジンのログのうち、特にクリックに関する情報（クリックログ）を利用して検索結果の精度向上を行うものである。

【0004】

すなわち、この文書検索システムでは、クリックログを解析することにより、多くの検索結果中からクリックが集中しているサイト（アクセス集中サイト）を的確に判別している。このアクセス集中サイトのタイトルとスニペット（概要文）には有用な情報（キーワード）が含まれていると考えられ、この情報でクエリを拡張し、各クエリの情報要求に沿った高精度な検索を実現するものである。

【非特許文献1】“Improving Mobile Web - IR Using Access Concentration Sites in Search Results.” Masaya Murata, et al. Proc. of WISE 2008, pp 221 - 234, 2008.

【発明の開示】

【発明が解決しようとする課題】

【0005】

非特許文献1の文書検索システムでは、クリックログを解析し、検索結果中においてアクセスが集中しているサイトを的確に判別することで、そこから有用な情報（キーワード）を抽出する。そしてこのキーワード群を基にクエリを拡張し、検索を行う。これは、通常1語か2語であるクエリの少ない単語数を、その情報要求を的確に表現する他のキーワード群で補い、ユーザの検索を補助するものである。

【0006】

10

20

30

40

50

一方、各検索結果が満たすことのできる情報要求の表現については、それらの内容、すなわち文書の作成者が記した本文のみが利用され、ユーザがその検索結果を実際に見て、それが要求通りか否かをどのように判断したのかに関する情報は考慮されていない。このような文書を利用するユーザ側からの情報を基に各検索結果を拡張することで、検索エンジンに対してクエリを発行したユーザが持つ情報要求（クエリの情報要求）と各検索結果が満たすことのできる情報要求レベルでのマッチングが可能となる。

【0007】

そこで本発明は、このような問題に鑑み、クエリの情報要求と、検索システムに登録されている各文書が満たすことのできる情報要求とを考慮した検索の実現を解決課題としている。

【課題を解決するための手段】

【0008】

本発明は、前記課題を解決するため、検索エンジンのログを利用してクエリおよび検索結果を拡張することによりそれぞれの満たす情報要求の表現を行い、これらの間の関係性を基に検索結果のランキングを行うことで、高精度の検索を実現している。

【0009】

具体的には、請求項1記載の発明は、ユーザ端末から検索指示されたクエリを含む電子文書を検索するときに検索エンジンの検索ログを利用する文書検索装置であって、前記検索ログに含まれたクエリに応じた検索結果のタイトルおよび概要文から拡張語を生成し、該拡張語を前記クエリの拡張情報として保存するクエリ情報保存手段と、前記検索ログからクリックされた検索結果の検索時のクエリを判別し、該クエリに関連する拡張語を前記クエリ情報保存手段から求め、検索結果がクリックされたときに投入されたクエリの拡張語群を検索結果の拡張情報として保存する検索結果情報保存手段と、ユーザ端末から検索指示されたクエリについて、クエリの拡張語を前記クエリ情報保存手段から取得し、ユーザ端末に送る照合処理手段と、前記検索指示されたクエリを前記拡張語で拡張した拡張クエリの検索で得られた検索結果に対して、前記検索結果情報保存手段に保存された拡張語群を付与して検索結果を拡張し、拡張された検索結果を前記拡張クエリとの類似度により並び替え、この並び替えた結果をリスト化した最終検索結果を前記ユーザ端末に送る検索結果処理手段と、を備えることを特徴としている。

【0012】

請求項2記載の発明は、ユーザ端末から検索指示されたクエリを含む電子文書を検索するときに検索エンジンの検索ログを利用する文書検索方法であって、前記検索ログに含まれたクエリに応じた検索結果のタイトルおよび概要文から拡張語を生成し、該拡張語をクエリ情報保存手段に前記クエリの拡張情報として保存するクエリ情報保存ステップと、検索結果情報保存手段が、前記検索ログからクリックされた検索結果の検索時のクエリを判別し、該クエリに関連する拡張語を前記クエリ情報保存手段から求め、検索結果がクリックされたときに投入されたクエリの拡張語群を検索結果の拡張情報として保存する検索結果情報保存ステップと、照合処理手段が、ユーザ端末から検索指示されたクエリについてクエリの拡張語を前記クエリ情報保存手段から取得し、ユーザ端末に送る照合処理ステップと、検索結果処理手段が、前記検索指示されたクエリを前記拡張語で拡張した拡張クエリの検索で得られた検索結果に対して、前記検索結果情報保存手段に保存された拡張語群を付与して検索結果を拡張し、拡張された検索結果を前記拡張クエリとの類似度により並び替え、この並び替えた結果をリスト化した最終検索結果を前記ユーザ端末に送る検索結果処理ステップと、を有することを特徴としている。

【0015】

請求項3記載の発明は、文書検索プログラムであって、請求項1記載の文書検索装置を構成する各手段としてコンピュータを機能させることを特徴としている。

【発明の効果】

【0016】

請求項1～3記載の発明によれば、検索エンジンのログを利用してクエリおよび検索結

10

20

30

40

50

果を拡張し、これらの関係性に基づいて検索結果をランキングすることで、検索の精度が向上する。

【発明を実施するための最良の形態】

【0017】

本発明は、検索エンジンのログ（クリックログ）を利用することで、クエリ拡張によるクエリの情報要求の表現、検索結果の拡張による検索結果の満たす情報要求の表現を行い、これらの間の関係性を基に検索結果のランキングを行っている。

【0018】

すなわち、多くのユーザが有用だと判断し、アクセスが集中しているサイトのタイトルとスニペットを拡張語の取得源とみなすことにより、クエリに対する高い適合性を持った拡張語を抽出する。この拡張語を用いたクエリ拡張により、クエリの情報要求を表現する。

10

【0019】

また、クリックログを解析することで、検索結果がクリック（閲覧）されたときに投入されたクエリを判別し、この判別したクエリの拡張語で検索結果を拡張することにより、この検索結果が満たすことができる情報要求を表現する。

【0020】

そして、それぞれ拡張されたクエリと検索結果同士をキーワードベースで比較することで、クエリと検索結果が潜在的に持つ情報要求に沿った高精度な検索を可能としている。以下、図面に基づき本発明の実施形態に係る文書検索装置1を説明する。

20

【0021】

図1は、本発明の実施形態に係る文書検索装置1の構成例を示している。この文書検索装置1は、インターネット経由で複数のユーザ端末13とネットワーク接続されている。このユーザ端末13をもってユーザはクエリを送信し文書検索を行う。

【0022】

前記文書検索装置1は、主に2つの処理部、すなわちユーザから投入されたクエリに対する検索結果を取得する検索エンジン100と、前記検索エンジン100の取得した検索結果を適切なランキングに並べ替え（re-ranking）、前記ユーザ端末13へ返信する支援処理部125とで構成されている。

【0023】

30

前記検索エンジン100は、「World Wide Web（WWW）」もしくは「Mobile Web（MW）」101から各サイトのデータを随時ダウンロードし、そのインデックスをインデックスDB102に格納する。そして、前記ユーザ端末13からの検索指示に従って前記インデックスDB102を検索し、検索結果を取得する。

【0024】

前記支援処理部125は、前記検索エンジン100の検索結果を適切なランキングに並べ替えて前記ユーザ端末13へ返信する。この支援処理部125は、図1に示すように、ログDB110、解析処理部111、クエリ情報要求生成部114、検索結果情報要求生成部115、検索結果情報要求DB116、クエリ情報要求DB117、照合処理部119、ランキング処理部123として機能している。

40

【0025】

ここで、前記ログDB110には、ユーザの検索ログに含まれたクエリと、該クエリの検索結果から実際にユーザがクリックして閲覧した電子文書のURLとを対応付けたクリックログが格納されている。

【0026】

前記解析処理部111は、前記クリックログを解析して、使用頻度が上位のクエリ（以下、頻度上位クエリとする）112を求める。そして、該頻度上位クエリ112を前記検索エンジン100に送信し、各クエリに対する検索結果集合113を取得する。この取得した検索結果集合113を前記各情報要求生成部114、115に送信する。

【0027】

50

前記クエリ情報要求生成部 114 は、前記クリックログを用いて前記検索結果集合 113 の解析を行い、クエリの情報要求を生成して前記クエリ情報要求 DB 117 へ格納する。

【0028】

ここでは、クエリの情報要求は、前記各頻度上位クエリ 112 に対する拡張語の集合として求められる。すなわち、前記クエリ情報要求 DB 117 には、前記各頻度上位クエリ 112 と、該各クエリに対する拡張語の集合とが対応して格納される。

【0029】

前記検索結果情報要求生成部 115 は、前記クリックログを用いて前記検索結果集合 113 の解析を行い、検索結果の情報要求を生成して前記検索結果情報要求 DB 116 へ格納する。

10

【0030】

ここでは、前記検索結果集合 113 からユーザの閲覧した文書の URL を含む検索結果を探索し、該検索結果の検索時のクエリを判別する。そして、この判別したクエリに対する拡張語の集合を該検索結果の情報要求として求める。すなわち、前記検索結果情報要求 DB 116 には、検索結果の電子文書の URL と、該電子文書がクリックされたときに投入されたクエリの拡張語の集合とが対応して格納される。

【0031】

ここまでの前記両 DB 116、117 の生成処理は、図 2 に示すように、前記ユーザ端末 13 と未接続のオフライン状態で実施される。この生成処理後に前記ユーザ端末 13 と

20

【0032】

前記ユーザ端末 13 は、ネットワークに接続可能なブラウザなどのユーザインタフェース 130 を備えていればよい。例えば、パーソナルコンピュータ (PC) や携帯電話などが該当する。前記ユーザインタフェース 130 には、ユーザがクエリを入力するクエリ入力画面 131、および検索結果を表示する検索結果表示画面 132 が表示される。

【0033】

ここでは、ユーザは前記クエリ入力画面 131 にてクエリ 118 を投入する。投入されたクエリ 118 は前記照合処理部 119 へ送信される。前記照合処理部 119 は、前記クエリ 118 を受信すると、対応する情報要求 120 を前記クエリ情報要求 DB 117 から

30

【0034】

前記検索エンジン 100 は、受信した拡張クエリ 121 をもって前記インデックス DB 102 を検索し、検索結果 122 を取得する。そして、取得した検索結果 122 を前記ランキング処理部 123 へ送信する。

【0035】

前記ランキング処理部 123 は、前記検索結果 122 に応じた情報要求を前記検索結果情報要求 DB 116 から取得し、該情報要求を用いて前記検索結果 122 を拡張する。そして、拡張した検索結果 122 を前記拡張クエリ 121 との関係性に基づいて並べ替え、ランク付けされた最終検索結果 124 を生成する。そして、生成した最終検索結果 124 を前記ユーザ端末 13 へ返信して検索結果表示画面 132 に表示させ、ユーザに提示する。

40

【0036】

ユーザは、前記検索結果表示画面 132 に表示された最終検索結果 124 をクリックして任意の電子文書を閲覧する。クエリ 118 の投入から最終検索結果 124 のクリックまでの操作情報は、検索ログ記録部 133 で随時取得され、該取得情報は前記ログ DB 110 へ蓄積される。ここまでの処理フローを図 3 に示す。

50

【 0 0 3 7 】

前記文書検索装置 1 の各機能ブロック 1 0 0 . 1 0 2 . 1 1 0 . 1 1 1 . 1 1 4 ~ 1 1 7 . 1 1 9 . 1 2 3 の機能は、コンピュータのハードウェアとソフトウェアの協働で実現されている。また、前記文書検索装置 1 は、コンピュータの通常の構成要素、例えば図示省略の処理データなどを一時記憶する書き換え可能なメモリ (R A M) と、前記ユーザ端末 1 3 とのネットワーク接続に使用する通信デバイスと、ハードディスクドライブ装置などの保存部などを備え、前記各 D B 1 0 2 . 1 1 0 . 1 1 6 . 1 1 7 は前記ハードディスクドライブ装置上に構築されている。以下、前記文書検索装置 1 の動作例を説明する。

【 0 0 3 8 】

< 動作例 >

前記文書検索装置 1 が実行する一連の処理は、主にオフラインで行われる情報要求生成フェーズと、オンラインで行われる検索実行フェーズから構成されている。以下、両フェーズの処理について、図 4 ~ 8 に基づき説明する。

【 0 0 3 9 】

(1) 情報要求生成フェーズ

情報要求生成フェーズでは、クリックログを解析して、クエリの情報要求および該クエリの検索結果が満たす情報要求を生成する。この情報要求生成フェーズは通常、前記ユーザ端末 1 3 と接続されないオフライン状態で、ユーザからの検索要求を受け付ける前に行われる。

【 0 0 4 0 】

図 4 . 5 は、情報要求生成フェーズの処理フローを示している。まず、前記解析処理部 1 1 1 は、前記ログ D B 1 1 0 に格納されたクリックログに含まれるクエリを使用頻度順に並べ、頻度上位クエリ 1 1 2 を得る。そして、該頻度上位クエリ 1 1 2 を前記検索エンジン 1 0 0 に送信し、それぞれのクエリの上位 K 件の検索結果集合 1 1 3 を得る。ここでは、前記頻度上位クエリ 1 1 2 の任意のクエリ q に対する検索結果集合 1 1 3 を例に説明する。

【 0 0 4 1 】

前記検索結果集合 1 1 3 は、前記検索エンジン 1 0 0 から前記クエリ情報要求生成部 1 1 4 および前記検索結果情報要求生成部 1 1 5 にそれぞれ送信され、クエリの情報要求、および検索結果が満たすことができる情報要求の算出が開始される。

【 0 0 4 2 】

< クエリの情報要求 >

前記クエリ情報要求生成部 1 1 4 は、前記クエリ q に対する検索結果集合 1 1 3 の各検索結果 $s r_i$ (s e a r c h r e s u l t) ($i = 1 , \dots , k$) において、そのタイトルとスニペットを形態素解析して内容語 (キーワード) t を抽出し、このキーワード集合をベクトル $V (s r_i)$ で表現する。

【 0 0 4 3 】

このベクトル $V (s r_i)$ には、キーワード t の $t f (t) \cdot i d f (t)$ に基づく重みが含まれている。 $t f (t)$ は、その検索結果 $s r_i$ のタイトルとスニペットにおけるキーワード t の出現頻度 (T e r m F r e q u e n c y) 、 $i d f$ はあるドキュメント集合におけるキーワード t の出現頻度 (I n v e r s e D o c u m e n t F r e q u e n c y) である。

【 0 0 4 4 】

このベクトル $V (s r_i)$ を、検索結果 $s r_i$ に対するアクセス集中度合 $A C D (s r_i)$ (A c c e s s C o n c e n t r a t i o n D e g r e e) で加重平均する。その結果得られるベクトルを、クエリ q の情報要求ベクトル $V_{IN} (q)$ (I n f o r m a t i o n N e e d V e c t o r) とする。具体的には、クエリ q の情報要求ベクトル $V_{IN} (q)$ は以下の式 (1) で与えられる。この式 (1) は、プログラムなどに定義されていればよい。

【 0 0 4 5 】

10

20

30

40

50

【数1】

$$V_{IN}(q) = \frac{\sum_i (ACD(sr_i) \times V(sr_i))}{\sum_i ACD(sr_i)} \quad \dots\dots (1)$$

【0046】

このように算出されたクエリ q の情報要求ベクトル $V_{IN}(q)$ は、クエリ q の拡張語の集合を表現するベクトルであり、前記クエリ情報要求DB117へ格納される。なお、アクセス集中度合 $ACD(sr_i)$ は、例えば非特許文献1の手法により求めることができる。また、ここでは処理を簡単にするため、クリックログを解析することで得られる文書の絶対的クリック回数 $C(sr_i)$ をアクセス集中度合 $ACD(sr_i)$ に置き換えてもよい。

10

【0047】

< 検索結果の情報要求 >

前記検索結果情報要求生成部115は、前記クリックログを解析することで、前記検索結果集合113のある検索結果 sr がクリックされたときに投入されたクエリの集合 q_j ($j = 1, \dots, m$)を求める。また、該集合の各クエリから検索結果 sr に対するアクセス集中度合 $ACD(sr, q_j)$ を求める。

20

【0048】

次に、クエリ q_j をもって前記クエリ情報要求DB117を検索し、該クエリ q_j に対応する情報要求ベクトル $V_{IN}(q_j)$ を取得する。そして、この情報要求ベクトル $V_{IN}(q_j)$ を前記アクセス集中度合 $ACD(sr, q_j)$ で加重平均して得られるベクトルを、検索結果 sr が満たすことができる(できた)情報要求ベクトル $V_{IN}(sr)$ とみなす。この情報要求ベクトル $V_{IN}(sr)$ は、以下の式(2)で与えられる。この式(2)もプログラムなどに定義されていればよい。

【0049】

【数2】

$$V_{IN}(sr) = \frac{\sum_j (ACD(sr, q_j) \times V_{IN}(q_j))}{\sum_j ACD(sr, q_j)} \quad \dots\dots (2)$$

30

【0050】

このように算出された検索結果の情報要求ベクトル $V_{IN}(sr)$ は、検索結果 sr がクリックされたときに投入されたクエリの拡張語集合を表現するベクトルであり、前記検索結果情報要求DB116へ格納される。なお、処理を簡単にするため、クリックログを解析することで得られる文書の絶対的クリック回数 $C(sr, q_j)$ をアクセス集中度合 $ACD(sr, q_j)$ に置き換えてもよい。また、クリックログを解析することで得られるクエリ q_j の全投入回数 $TIN(q_j)$ (Total Input Number)で絶対的クリック回数 $C(sr, q_j)$ を正規化した値で置き換えてもよい。

40

【0051】

ここで、前記クエリ q_j をもって前記クエリ情報要求DB117を検索した際、該クエリ q_j に対応する情報要求ベクトル $V_{IN}(q_j)$ が存在しない場合もある。そのような場合、前記検索結果情報要求生成部115は、このクエリ q_j を前記解析処理部111へ送信し、前述したクエリの情報要求の生成処理と同様の手順で、該クエリ q_j に対応する情報要求ベクトル $V_{IN}(q_j)$ の算出を行えばよい。この算出した情報要求ベクトル $V_{IN}(q_j)$

50

)を前記クエリ情報要求DB117へ格納し、これを用いて前記式(2)により検索結果srの情報要求ベクトル $V_{IN}(sr)$ を算出すればよい。以上の情報要求生成フェーズにおけるデータ例を図6に示す。

【0052】

(2)検索実行フェーズ

検索実行フェーズでは、情報要求生成フェーズで生成したクエリおよび検索結果の情報要求を用いて、ユーザの投入したクエリを拡張するとともに、該拡張クエリに対する検索結果も拡張する。そして、拡張されたクエリと検索結果との間の関係性に基づき検索結果を並べ替え、最終検索結果を生成している。この検索実行フェーズは、前記ユーザ端末13と接続されたオンライン状態で行われる。

10

【0053】

図7は、検索実行フェーズの処理フローを示している。まず、前記照合処理部119は、ユーザが前記クエリ入力画面131をもって投入したクエリ118を受信する。

【0054】

前記照合処理部119は、前記クエリ情報要求DB117から前記クエリ118に対応する情報要求ベクトル $V_{IN}(q)$ を取得し、これをクエリの情報要求120として前記ユーザ端末13へ返信しクエリ入力画面131へ表示させる。このクエリの情報要求120によって、前記クエリ118が拡張される。

【0055】

すなわち、前記クエリ118と前記クエリの情報要求120との組み合わせが、拡張クエリ121として前記クエリ入力画面131から前記検索エンジン100へ送信される。前記検索エンジン100は、受信した拡張クエリ121を用いて検索を行い、取得した検索結果122を前記ランキング処理部123へ送信する。

20

【0056】

ここで、通常の実験結果は、ランク付けされた検索結果のタイトルの一部、本文の一部(スニペット)、およびURLが返されるが、ここではランク付けされた検索結果のタイトルの全文、全本文、およびURLが返される。この時点でのランク付けは、従来の全文検索アルゴリズムに沿って行われる。

【0057】

前記ランキング処理部123は、前記検索結果122の各検索結果srに対応する情報要求ベクトル $V_{IN}(sr)$ を前記検索結果情報要求DB116から取得する。そして、取得した情報要求ベクトル $V_{IN}(sr)$ で表現される拡張語群を各検索結果srのタイトルと本文に付与する。これにより、検索結果122の拡張が行われる。

30

【0058】

そして、このように拡張された検索結果122を、拡張クエリ121とのキーワードベースでの類似度を考慮に入れて並べ替える(re-ranking)。そして、この並べ替えた結果をリスト化した最終検索結果124を前記ユーザ端末13へ返信し、検索結果表示画面132へ表示させてユーザへ提示する。ユーザは、提示された前記最終検索結果124から任意の電子文書をクリックして閲覧する。

【0059】

なお、検索実行フェーズにおけるユーザのクエリの投入から最終検索結果のクリックまでの行動情報は、前記検索ログ記録部133で常に監視・取得される。この行動情報は、前記文書検索装置1へ送信され、新たなクリックログとして前記ログDB110に蓄積される。これにより、ユーザの行動情報が以降の情報要求生成フェーズに随時反映され、時々刻々と変化するユーザの情報要求を適切に把握することが可能となる。なお、この検索ログ記録部133は、前記文書検索装置1内に実装されていてもよい。ここまでの検索実行フェーズにおけるデータ例を図8に示す。

40

【0060】

<発明の効果>

以上のように、文書検索システムに本発明の前記文書検索装置1を配置し、検索結果の

50

精度評価を行った実験の結果を表1に示す。

【0061】

【表1】

	5	10	15	20	30	50
BM25 (baseline1)	42.4 **	37.5 **	37.6 **	37.6 **	36.5 **	36.1 *
RCN (baseline2)	34.9 **	34.1 **	33.7 **	31.2 **	30.7 **	30.2 **
QEC (baseline3)	46.3 *	43.7 **	43.5	42.4 *	42.5	39.4
本発明	54.1	47.3	44.8	45.2	42.9	40.2

10

【0062】

精度評価の指標は「Precision@X」と呼ばれるものを使用している。これは、クエリに対して正解であるサイトが検索結果の上位X件に多く入るほど高い数値を出す指標である。

20

【0063】

なお、本発明との比較対象の手法は、(1)BM25、(2)クリック回数の多いサイトで検索結果を並べ替える方法(Re-ranking by Click Number, RCN)、(3)クリック回数に基づくクエリ拡張法(Query Expansion method using Click number, QEC)の3種類としている。

【0064】

(1)のBM25は、クエリ-サイト間のキーワードマッチングベースのランキングとして幅広く用いられている手法である。(2)のRCNは、単純に検索結果をそのクリック回数順で並べ替える手法である。(3)のQECは、クリック回数の多い検索結果のタイトルとスニペットからキーワードを抽出し、抽出したキーワードでクエリ拡張を行い、その拡張されたクエリとサイト間の類似度を基にランキングを行う手法である。

30

【0065】

表1中の「*」、「**」は、それぞれウィルコクソンの符号付順位和検定において、本発明の手法と各比較手法との統計的有意差が5%、1%であった結果である。太文字になっている手法と数値が「Precision@X」に対する最大値である。

【0066】

表1に示すように、本発明の手法は全ての検索結果ランク@Xにおいて最大の精度を達成しており、特に検索結果の上位ランク1 X 20の領域における精度を著しく向上させている。

40

【0067】

なお、本発明は、コンピュータを前記文書検索装置1の各機能ブロック100・102・110・111・114~117・119・123として機能させる文書検索プログラムとしても提供することができる。このプログラムは、各機能ブロック100・102・110・111・114~117・119・123の全ての機能を実現させるものでもよく、あるいは一部の機能を実現させるものであってもよい。

【0068】

このプログラムは、Webサイトなどからのダウンロードによってコンピュータに提供される。また、前記プログラムは、CD-ROM、DVD-ROM、CD-R、CD-R

50

W, DVD-R, DVD-RW, MO, HDD, Blu-ray Disk (登録商標) などの記録媒体に格納してコンピュータに提供してもよい。

【図面の簡単な説明】

【0069】

【図1】本発明の実施形態に係る文書検索装置の構成図。

【図2】同 情報要求生成フェーズの概略図。

【図3】同 検索実行フェーズの概略図。

【図4】同 情報要求生成フェーズ前半の処理フロー図。

【図5】同 情報要求生成フェーズ後半の処理フロー図。

【図6】同 情報要求生成フェーズのデータ例。

10

【図7】同 検索実行フェーズの処理フロー図。

【図8】同 検索実行フェーズのデータ例。

【符号の説明】

【0070】

1 ... 文書検索装置

13 ... ユーザ端末

100 ... 検索エンジン

101 ... World Wide WebもしくはMobile Web

102 ... インデックスDB

110 ... ログDB

20

111 ... 解析処理部

112 ... 頻度上位クエリ

113 ... 検索結果集合

114 ... クエリ情報要求生成部

115 ... 検索結果情報要求生成部

116 ... 検索結果情報要求DB

117 ... クエリ情報要求DB

118 ... クエリ

119 ... 照合処理部

120 ... クエリの情報要求

30

121 ... クエリとクエリの情報要求(拡張クエリ)

122 ... 検索結果

123 ... ランキング処理部(検索結果処理手段)

124 ... 最終検索結果

125 ... 支援処理部

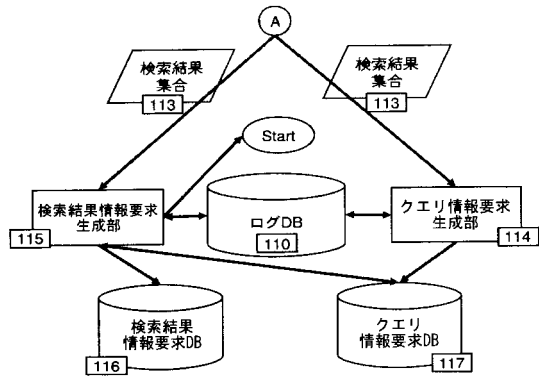
130 ... ユーザインタフェース

131 ... クエリ入力画面

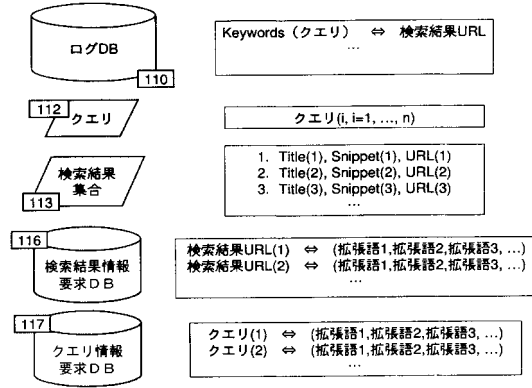
132 ... 検索結果表示画面

133 ... 検索ログ記録部

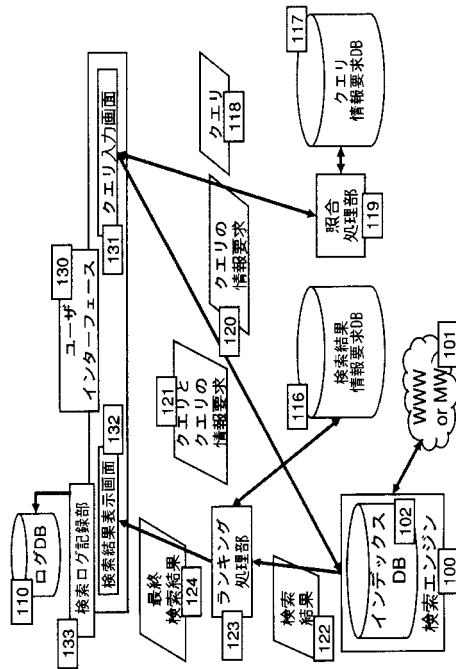
【図5】



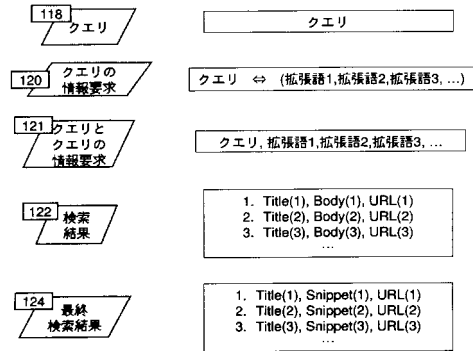
【図6】



【図7】



【図8】



フロントページの続き

- (72)発明者 戸田 浩之
東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
- (72)発明者 松浦 由美子
東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
- (72)発明者 片岡 良治
東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内

審査官 岩間 直純

- (56)参考文献 村田 眞哉, ほか, 検索結果中のアクセス集中サイトを利用したクエリ拡張法の提案, データベースとWeb情報システムに関するシンポジウム 情報処理学会シンポジウムシリーズ, 日本, 社団法人情報処理学会, 2007年11月27日, Vol. 2007, No. 3, pp. 1-7
- 村田 眞哉, ほか, 検索結果中のアクセス集中サイトを利用したクエリ拡張法の提案, 日本データベース学会 Letters, 日本, 日本データベース学会, 2008年 3月21日, Vol. 6, No. 4, pp. 45-48

- (58)調査した分野(Int.Cl., DB名)
G06F 17/30