

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4781466号
(P4781466)

(45) 発行日 平成23年9月28日(2011.9.28)

(24) 登録日 平成23年7月15日(2011.7.15)

(51) Int.Cl. F I
G06F 17/30 (2006.01) G O 6 F 17/30 3 5 O C
 G O 6 F 17/30 1 7 O A

請求項の数 7 (全 17 頁)

<p>(21) 出願番号 特願2009-508735 (P2009-508735) (86) (22) 出願日 平成19年3月16日 (2007.3.16) (86) 国際出願番号 PCT/JP2007/055352 (87) 国際公開番号 W02008/126184 (87) 国際公開日 平成20年10月23日 (2008.10.23) 審査請求日 平成21年5月15日 (2009.5.15)</p>	<p>(73) 特許権者 000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番 1号 (74) 代理人 100101856 弁理士 赤澤 日出夫 (72) 発明者 斎藤 郁生 神奈川県川崎市中原区上小田中4丁目1番 1号 富士通株式会社内 審査官 池田 聡史</p>
--	---

最終頁に続く

(54) 【発明の名称】 文書重要度算出プログラム

(57) 【特許請求の範囲】

【請求項1】

N個の文書それぞれの重要度の算出をコンピュータに実行させる文書重要度算出プログラムであって、

全ての要素が正の実数であるN次の正方行列D、正の実数e、前記N個の文書それぞれの重要度を要素とする縦ベクトルuがペロン・フロベニウスの定理により $e^{-1} u = D^{-1} u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定ステップと、

縦ベクトルuの近似値として求める縦ベクトルvの初期化を行い、 $(D^{-1} v) / |D^{-1} v|$ を新たな縦ベクトルvとする計算を行い、縦ベクトルvが所定の条件を満たすまで該計算を反復し、縦ベクトルvの要素をそれぞれ対応する文書の重要度とする反復ステップと

をコンピュータに実行させる文書重要度算出プログラム。

【請求項2】

請求項1に記載の文書重要度算出プログラムにおいて、

前記文書のパラメータは、該文書へのリンク関係、該文書の作成日時、該文書の参照者による該文書の評価、該文書の参照回数の少なくともいずれかを含むことを特徴とする文書重要度算出プログラム。

【請求項3】

請求項1または請求項2に記載の文書重要度算出プログラムにおいて、

前記行列決定ステップは、前記N個の文書のうちm番目の文書の前記パラメータ、または前記N個の文書のうちm番目の文書及びn番目の文書の前記パラメータに基づいて、Dのm行n列目の要素を算出することを特徴とする文書重要度算出プログラム。

【請求項4】

請求項1乃至請求項3のいずれかに記載の文書重要度算出プログラムにおいて、

前記所定の条件は、前記計算前の縦ベクトル v と前記計算後の縦ベクトル v との差のノルムが所定の閾値を下回る場合であることを特徴とする文書重要度算出プログラム。

【請求項5】

請求項1乃至請求項4のいずれかに記載の文書重要度算出プログラムにおいて、

更に行列決定ステップの前に、前記文書毎の前記パラメータを取得するパラメータ取得ステップをコンピュータに実行させることを特徴とする文書重要度算出プログラム。

10

【請求項6】

N個の文書それぞれの重要度の算出を行う文書重要度算出装置であって、

全ての要素が正の実数であるN次の正方行列D、正の実数 e 、前記N個の文書それぞれの重要度を要素とする縦ベクトル u がペロン・フロベニウスの定理により $e \cdot u = D \cdot u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定部と、

縦ベクトル u の近似値として求める縦ベクトル v の初期化を行い、 $(D \cdot v) / |D \cdot v|$ を新たな縦ベクトル v とする計算を行い、縦ベクトル v が所定の条件を満たすまで該計算を反復し、縦ベクトル v の要素をそれぞれ対応する文書の重要度とする反復部と

20

を備える文書重要度算出装置。

【請求項7】

N個の文書それぞれについての重要度の算出を行う文書重要度算出方法であって、

全ての要素が正の実数であるN次の正方行列D、正の実数 e 、前記N個の文書それぞれの重要度を要素とする縦ベクトル u がペロン・フロベニウスの定理により $e \cdot u = D \cdot u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定ステップと、

縦ベクトル u の近似値として求める縦ベクトル v の初期化を行い、 $(D \cdot v) / |D \cdot v|$ を新たな縦ベクトル v とする計算を行い、縦ベクトル v が所定の条件を満たすまで該計算を反復し、縦ベクトル v の要素をそれぞれ対応する文書の重要度とする反復ステップと

30

を実行する文書重要度算出方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数の文書についてそれぞれの重要度を算出する文書重要度算出プログラム、文書重要度算出装置、文書重要度算出方法に関するものである。

【背景技術】

【0002】

従来、Webページ等のリンクを有する文書の検索を行う検索システムにおいて、検索結果の表示順位を決定するための技術がある(例えば、特許文献1、特許文献2、非特許文献1、非特許文献2参照)。特許文献1の技術は、検索対象文書のリンク関係だけから各文書の重要度を決定し、重要度の高い順に検索結果を表示する。

40

【特許文献1】米国特許第6285999号明細書

【特許文献2】米国特許出願公開第2005/0071741号明細書

【非特許文献1】"The PageRank Citation Ranking: Bringing Order to the Web": Lawrence Page, Sergey Brin, Rajeev Motwani, Terry Winograd (1998年1月29日)

【非特許文献2】"Efficient Computation of PageRank": Taher H. Haveliwala (1999年10月18日)

【発明の開示】

50

【発明が解決しようとする課題】

【0003】

しかしながら、特許文献1の技術は、重要度算出のための制約が大きく、検索対象文書のリンク関係だけから検索対象文書の重要度を算出することしかできなかった。

【0004】

また、従来技術を基にして、複数の種類のパラメータの影響を考慮した重要度を算出するための幾つかの方法がある。特許文献2の技術においては、リンク関係に基づいて重要度を算出し、その重要度を他のパラメータに基づいて修正する方法が示されているが、リンク関係により計算された情報が他のパラメータに基づく修正により失われてしまう。また、特許文献2の技術において、多くの種類のパラメータの影響を考慮する具体的な方法は示されていない。

10

【0005】

また、パラメータ毎に重要度を算出し、得られた複数の重要度に重み付けを与える方法等が考えられるが、それぞれのパラメータに対する重要度を算出するために、時間のかかる繰り返し計算をパラメータの数だけ行うことになるため、パラメータ数が増えると計算量が膨大になってしまう。

【0006】

本発明は上述した問題点を解決するためになされたものであり、文書に関するパラメータの影響を考慮した文書の重要度を効率的に算出する文書重要度算出プログラム、文書重要度算出装置、文書重要度算出方法を提供することを目的とする。

20

【課題を解決するための手段】

【0007】

上述した課題を解決するため、本発明は、N個の文書それぞれの重要度の算出をコンピュータに実行させる文書重要度算出プログラムであって、全ての要素が正の実数であるN次の正方行列D、正の実数e、前記N個の文書それぞれの重要度を要素とする縦ベクトルuがペロン・フロベニウスの定理により $e \cdot u = D \cdot u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定ステップと、縦ベクトルuの近似値として求める縦ベクトルvの初期化を行い、 $(D \cdot v) / |D \cdot v|$ を新たな縦ベクトルvとする計算を行い、縦ベクトルvが所定の条件を満たすまで該計算を反復し、縦ベクトルvの要素をそれぞれ対応する文書の重要度とする反復ステップとをコンピュータに実行させる。

30

【0008】

また、本発明は、N個の文書それぞれの重要度の算出を行う文書重要度算出装置であって、全ての要素が正の実数であるN次の正方行列D、正の実数e、前記N個の文書それぞれの重要度を要素とする縦ベクトルuがペロン・フロベニウスの定理により $e \cdot u = D \cdot u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定部と、縦ベクトルuの近似値として求める縦ベクトルvの初期化を行い、 $(D \cdot v) / |D \cdot v|$ を新たな縦ベクトルvとする計算を行い、縦ベクトルvが所定の条件を満たすまで該計算を反復し、縦ベクトルvの要素をそれぞれ対応する文書の重要度とする反復部とを備える。

40

【0009】

また、本発明は、N個の文書それぞれについての重要度の算出を行う文書重要度算出方法であって、全ての要素が正の実数であるN次の正方行列D、正の実数e、前記N個の文書それぞれの重要度を要素とする縦ベクトルuがペロン・フロベニウスの定理により $e \cdot u = D \cdot u$ を満たすことを用いるために、前記文書に関する少なくとも1種類のパラメータに基づいて行列Dの要素を決定する行列決定ステップと、縦ベクトルuの近似値として求める縦ベクトルvの初期化を行い、 $(D \cdot v) / |D \cdot v|$ を新たな縦ベクトルvとする計算を行い、縦ベクトルvが所定の条件を満たすまで該計算を反復し、縦ベクトルvの要素をそれぞれ対応する文書の重要度とする反復ステップとを実行する。

【図面の簡単な説明】

50

【 0 0 1 0 】

【図 1】本実施の形態に係る検索システムの構成の一例を示すブロック図である。

【図 2】本実施の形態に係る文書リストファイルの一例を示す表である。

【図 3】本実施の形態に係る作成日時リストファイルの一例を示す表である。

【図 4】本実施の形態に係るリンク関係ファイルの一例を示す表である。

【図 5】本実施の形態に係る評価リストファイルの一例を示す表である。

【図 6】本実施の形態に係る参照回数リストファイルの一例を示す表である。

【図 7】本実施の形態に係る重要度算出部の動作の一例を示すフローチャートである。

【図 8】本実施の形態に係る重要度リストファイルの一例を示す表である。

【図 9】従来の反復計算処理の一例を示すフローチャートである。

10

【図 10】本実施の形態に係る反復計算処理の一例を示すフローチャートである。

【図 11】本実施の形態に係る行列決定処理により決定された行列 D の一例を示す表である。

【図 12】本実施の形態に係る反復計算処理により計算された計算結果の一例を示す表である。

【図 13】本実施の形態に係る集計データの一例を示す表である。

【図 14】本実施の形態に係るリンク関係の一例を示す表である。

【図 15】本実施の形態に係るリンク関係データの一例を示す表である。

【図 16】本実施の形態に係る第 1 の実施例において行列決定処理により決定された重みデータを示す表である。

20

【図 17】本実施の形態に係る第 1 の実施例において行列決定処理により決定された行列 D を示す表である。

【図 18】本実施の形態に係る第 1 の実施例において反復計算処理により計算された計算結果を示す表である。

【図 19】本実施の形態に係る第 2 の実施例において行列決定処理により決定された重みデータを示す表である。

【図 20】本実施の形態に係る第 2 の実施例において行列決定処理により決定された行列 D を示す表である。

【図 21】本実施の形態に係る第 2 の実施例において反復計算処理により計算された計算結果を示す表である。

30

【図 22】本実施の形態に係る第 3 の実施例において行列決定処理により決定された重みデータを示す表である。

【図 23】本実施の形態に係る第 3 の実施例において行列決定処理により決定された行列 D を示す表である。

【図 24】本実施の形態に係る第 3 の実施例において反復計算処理により計算された計算結果を示す表である。

【図 25】本実施の形態に係る第 4 の実施例において行列決定処理により決定された $R_{and}()$ の値を示す表である。

【図 26】本実施の形態に係る第 4 の実施例において行列決定処理により決定された行列 D を示す表である。

40

【図 27】本実施の形態に係る第 4 の実施例において反復計算処理により計算された計算結果を示す表である。

【図 28】本実施の形態に係る第 5 の実施例において行列決定処理により決定された行列 D を示す表である。

【図 29】本実施の形態に係る第 5 の実施例において反復計算処理により計算された計算結果を示す表である。

【発明を実施するための最良の形態】

【 0 0 1 1 】

以下、本発明の実施の形態について図面を参照しつつ説明する。

【 0 0 1 2 】

50

本実施の形態においては、本発明の文書重要度算出装置を適用した検索システムについて説明する

【0013】

まず、本実施の形態に係る検索システムの構成について説明する。

【0014】

図1は、本実施の形態に係る検索システムの構成の一例を示すブロック図である。この検索システムは、端末1、検索サーバ2、文書管理部3、検索インデックス管理部4、重要度算出部5を備える。端末1、検索サーバ2、文書管理部3、検索インデックス管理部4、重要度算出部5は、それぞれコンピュータで実現される。

【0015】

文書管理部3は、検索サーバ2による検索の対象となる文書である検索対象文書を管理する。検索インデックス管理部4は、検索対象文書の検索インデックスを登録すると共に、重要度算出部5から送信される検索対象文書の重要度リストファイルを検索インデックスに反映する。端末1は、参照者の操作に従って、検索の指示、文書の参照の指示を検索サーバ2へ送信する。また、端末1は、参照者による検索対象文書の評価を検索サーバ2へ送信する。

【0016】

検索サーバ2は、端末1からの検索の指示に従って、検索インデックス管理部4の検索インデックスを用いて検索対象文書を検索し、検索結果を端末1へ送信する。また、検索サーバ2は、端末1からの文書の参照の指示に従って、指示された文書を文書管理部3から取得し、端末1へ送信すると共に、文書の参照の履歴を参照履歴情報として重要度算出部5へ送信する。また、検索サーバ2は、端末1から送信された検索対象文書の評価を重要度算出部5へ送信する。

【0017】

重要度算出部5は、CPU及びメモリを有し、CPUがメモリ上で実行するプログラムとして、文書リスト作成部12、作成日時抽出部13、リンク解析部14、評価受付部15、参照回数算出部16、行列決定部17、反復計算部18を備える。なお、パラメータ取得部は、実施の形態における文書リスト作成部12、作成日時抽出部13、リンク解析部14、評価受付部15、参照回数算出部16に対応する。

【0018】

文書リスト作成部12は、文書管理部3に管理された検索対象文書のリストである文書リストファイルを作成する。図2は、本実施の形態に係る文書リストファイルの一例を示す表である。作成日時抽出部13は、文書管理部3に管理された検索対象文書の作成日時のリストである作成日時リストファイルを作成する。図3は、本実施の形態に係る作成日時リストファイルの一例を示す表である。リンク解析部14は、文書管理部3に管理された検索対象文書のリンクを解析し、解析結果をリンク関係ファイルとして生成する。図4は、本実施の形態に係るリンク関係ファイルの一例を示す表である。評価受付部15は、検索サーバ2から送信された検索対象文書の参照者による評価を受信し、評価リストファイルとして生成する。図5は、本実施の形態に係る評価リストファイルの一例を示す表である。参照回数算出部16は、参照履歴情報から検索対象文書の参照回数を算出し、参照回数リストファイルとして生成する。図6は、本実施の形態に係る参照回数リストファイルの一例を示す表である。

【0019】

次に、重要度算出部5の動作について説明する。

【0020】

図7は、本実施の形態に係る重要度算出部の動作の一例を示すフローチャートである。まず、文書リスト作成部12、作成日時抽出部13、リンク解析部14、評価受付部15、参照回数算出部16は、各文書のパラメータを取得し、文書リストファイル、作成日時リストファイル、リンク関係ファイル、評価リストファイル、参照回数リストを生成するパラメータ取得処理を行う(S11、パラメータ取得ステップ)。次に、行列決定部17

10

20

30

40

50

は、これらの各ファイルに基づいて、検索対象文書の重要度の算出のための行列 D を決定する行列決定処理を行う (S 1 2、行列決定ステップ)。次に、反復計算部 1 8 は、行列 D を用いて検索対象文書の重要度を算出する反復計算処理を行い、算出結果を重要度リストファイルとして検索インデックス管理部 4 へ送信し (S 1 3、反復計算ステップ)、このフローは終了する。図 8 は、本実施の形態に係る重要度リストファイルの一例を示す表である。

【 0 0 2 1 】

次に、反復計算処理の原理について説明する。

【 0 0 2 2 】

まず、検索対象文書の総数を N とし、文書管理部 3 により検索対象文書に 1 から N までの ID が振られたとする。また、 ID が n である検索対象文書を P_n とする。また、正の実数を u_n ($n = 1, 2, \dots, N$) とするとき、 u_n を要素とする縦ベクトルを u とする。 d_i_j ($i, j = 1, 2, \dots, N$) を任意の正の実数とし、第 (i, j) 成分が d_i_j である行列を D とする。

10

【 0 0 2 3 】

線形代数のペロン・フロベニウスの定理より、任意の D に対して、 D は正の固有値 (e $u = D u$ を満たす 0 でないベクトル u が存在するような数 e のこと。このときの u を e に対する D の固有ベクトルという) を持ち、そのうちの最大のものを D のフロベニウス根という。フロベニウス根 e に対する D の固有ベクトル u は正のベクトルであり、そのうち $|u| = 1$ となるものはただ 1 つである。ここでのノルムは無量大ノルム ($|u| = \max\{|u_n|; n = 1, 2, \dots, N\}$) とする。

20

【 0 0 2 4 】

検索対象文書 P_n ($n = 1, 2, \dots, N$) に与えられた様々なパラメータ (例: リンク関係、作成日時、参照者による評価、文書の参照回数など) に応じて、システムの設計者は自由に D を定義することができる。このように定義した D に対して、 D のフロベニウス根 e と、 $|u| = 1$ を満たすような e に対する正の固有ベクトル u はただ 1 組存在するので、このときの u_n を P_n の重要度と定義する。

【 0 0 2 5 】

ここで、 P_n から L_n 個のリンクがあると仮定し、次式のように d_m_n を定義すれば、 $u = D u$ 、 $|u| = 1$ を満たす u は、特許文献 1 の技術と同様な値を算出することができる。ここでのノルムは、1 ノルム ($|u| = u_1 + u_2 + \dots + u_N$) とし、 P_n から出るリンクの数を L_n とする。また、定数 $c = 0.15$ とする。

30

【 0 0 2 6 】

$d_m_n = c / N$
(P_n から P_m へのリンクがないとき)
 $d_m_n = c / N + (1 - c) / L_n$
(P_n から P_m へのリンクがあるとき)

【 0 0 2 7 】

次に、特許文献 1 等の技術を用いた従来の反復計算処理と、本実施の形態に係る反復計算部 1 8 による反復計算処理とを比較する。

40

【 0 0 2 8 】

まず、従来の反復計算処理について説明する。図 9 は、従来の反復計算処理の一例を示すフローチャートである。まず、従来の反復計算処理は、繰り返し回数 $k = 1$ とし、長さ N の縦ベクトル v_k の要素を全て 1 としして初期化を行う (S 2 1)。次に、従来の反復計算処理は、 k を 1 増加させる (S 2 2)。次に、従来の反復計算処理は、 $v_k = D v_{(k-1)}$ を計算する (S 2 3)。

【 0 0 2 9 】

次に、従来の反復計算処理は、今回得られた v_k と前回得られた $v_{(k-1)}$ から差のノルム $|v_k - v_{(k-1)}|$ を計算し、 $|v_k - v_{(k-1)}| < s$ であるか否かの判断を行う (S 3 1)。ここで、しきい値 s は、十分小さい正の整数であり、

50

例えば $s = 0.00001$ である。 $|v_k - v_{(k-1)}| < s$ でない場合 (S31, No)、処理 S22 へ戻り、処理 S22 以降を繰り返す。 $|v_k - v_{(k-1)}| < s$ である場合 (S31, Yes)、従来の反復計算処理は、 v_k を重要度ベクトル u とし (S32)、このフローは終了する。

【0030】

但し、従来の反復計算処理に用いる D の各列の要素はリンク関係に基づく確率であり、 D の各列の要素の和が 1 でなければならないという制約がある。この制約を満たさない場合、 v_k は発散し、重要度を算出することができなくなる。

【0031】

次に、本実施の形態に係る反復計算処理について説明する。図 10 は、本実施の形態に係る反復計算処理の一例を示すフローチャートである。この図において、図 9 と同一符号は、図 9 に示された対象と同一処理を示しており、ここでの説明を省略する。従来の反復計算処理と比較すると、本実施の形態の反復計算処理は、処理 S23 の代わりに処理 S24, S25 を実行する。

10

【0032】

まず、反復計算部 18 は、従来の反復計算処理と同様の処理 S21, S22 を実行する。次に、反復計算部 18 は、 $e_{(k-1)} = |D v_{(k-1)}|$ を計算し (S24)、 $v_k = (1 / e_{(k-1)}) D v_{(k-1)}$ を計算する (S25)。

【0033】

次に、反復計算部 18 は、従来の反復計算処理と同様の処理 S31 を実行する。 $|v_k - v_{(k-1)}| < s$ でない場合 (S31, No)、反復計算部 18 は、処理 S22 へ戻り、処理 S22 以降を繰り返す。 $|v_k - v_{(k-1)}| < s$ である場合 (S31, Yes)、反復計算部 18 は、従来の反復計算処理と同様の処理 S32 を実行し、このフローは終了する。

20

【0034】

本実施の形態に係る反復計算処理によれば、 k が大きくなると、 e_k は e に、 v_k は u に収束していく。また、処理 S31 により、 v_k の変化が十分小さくなったところで繰り返し計算を終了する。つまり、反復計算部 18 は、 $e u = D u$ を満たす u の近似値である v_k を算出する。

【0035】

本実施の形態に係る反復計算処理は、処理 S24, S25 により、 D に従来の制約を与えなくても v_k を収束させることができる。

30

【0036】

次に、本実施の形態に係る反復計算処理の具体例について説明する。

【0037】

図 11 は、本実施の形態に係る行列決定処理により決定された行列 D の一例を示す表である。ここでは、 N を 6 とし、この表のように D を決定したとする。この D は、従来の D のような制約はなく、各要素が正の実数であれば良い。図 12 は、本実施の形態に係る反復計算処理により計算された計算結果の一例を示す表である。この表は、処理 S22 により設定される繰り返し回数 k 、処理 S24 により算出される $D v_{(k-1)}$ 、 $e_{(k-1)}$ 、処理 S25 により算出される v_k 、処理 S31 により算出される $v_k - v_{(k-1)}$ 、 $v_{(k-1)} - v_k$ 、 $|v_k - v_{(k-1)}|$ の計算結果を示し、それぞれの計算結果を繰り返し回数 k 毎に示す。この例においては、繰り返し回数 $k = 10$ の時点で、 $|v_k - v_{(k-1)}|$ がしきい値を下回ったため、反復計算部 18 は、繰り返し計算を終了し、 v_{10} を D に対する重要度ベクトル u とする。

40

【0038】

次に、行列決定処理の概要について説明する。

【0039】

本実施の形態においては、上述の反復計算処理により、行列 D は自由に定義することができるため、検索システムの設計者は、検索対象文書の新しさや参照者からの評価などを

50

Dに反映させることができる。ここでは、検索対象文書の作成日時、評価、参照回数、リンク関係に基づいてDを決定する例を用いて、行列決定処理の概要を説明する。

【0040】

P_mの作成日時が、行列決定処理時のp日前とする。ここで、行列決定部17は、現在の日時と作成日時リストファイルからpを算出する。検索対象文書の作成日時が計算時のp日前のときの重みをx_m(正の実数)とする。x_mは、pが小さいほど大きい値になるように設定される。

【0041】

また、検索対象文書P_mの評価をy_m(正の実数)と置く。ここで、行列決定部17は、y_mは評価リストファイルから算出する。y_mは、評価が高いほど大きい値になるように設定される。

10

【0042】

検索対象文書P_mの参照回数をz_m(正の実数)と置く。ここで、行列決定部17は、参照回数リストファイルからz_mを算出する。z_mは、参照回数が多いほど大きい値になるように設定される。

【0043】

次に、行列決定部17は、リンク関係ファイルからP_nからP_mへのリンクの有無を判断し、リンクの有無によってDの要素d_{m_n}を決定する。また、ここでは、d_{m_n}の算出にx_m、y_m、z_mの積を用いる場合を示す。d_{m_n}は、次式で定義される。

20

【0044】

$$d_{m_n} = c \cdot x_m \cdot y_m \cdot z_m$$

(m, n = 1, 2, ... N)
(P_nからP_mへのリンクがないとき)

$$d_{m_n} = c \cdot x_m \cdot y_m \cdot z_m + x_m \cdot x_n \cdot y_m \cdot y_n \cdot z_m \cdot z_n$$

(m, n = 1, 2, ... N)
(P_nからP_mへのリンクがあるとき)

【0045】

ここでは、x_m、y_m、z_mがd_{m_n}へ与える影響は、x_m、y_m、z_m、x_n、y_n、z_nの積で表したが、和や他の演算によって表しても良い。

30

【0046】

このように、システムの設計者は、複数のパラメータ(ここでは、作成日、評価、参照回数、リンク関係)の影響を考慮したDを自由に決定することができる。

【0047】

次に、行列決定処理に用いるデータの具体例について説明する。

【0048】

まず、行列決定部17は、データの集計を行い、集計したデータを集計データとする。図13は、本実施の形態に係る集計データの一例を示す表である。ここでは、検索対象文書数Nを6とし、データの集計日を2007年2月7日とする。行列決定部17は、作成日時リストファイルから作成日を取得し、データの集計日(現在)と作成日の日付の差pを算出し、評価リストファイルから所定の基準を満たす有用な評価の数である有用評価数qを取得し、参照回数リストファイルから参照回数rを取得する。この表は、P_n(n = 1, 2, ... 6)のそれぞれについて、作成日、日付の差p、有用評価数q、参照回数rを集計した結果を示す。

40

【0049】

図14は、本実施の形態に係るリンク関係の一例を示す表である。ここでは、P_n(n = 1, 2, ... 6)の間に、この図に示されたリンク関係が存在するとする。行列決定部17は、このリンク関係を示したリンク関係ファイルからリンク関係を集計し、リンク関

50

係データとする。図15は、本実施の形態に係るリンク関係データの一例を示す表である。この表において、“from”の欄に示された文書は、リンク元であり、“to”の欄に示された文書は、リンク先を示す。また、“ ”が付いた欄に対応するリンク元とリンク先の間リンクが存在することを示す。また、 P_n から出るリンクの数 L_n がリンク元毎に集計されている。

【0050】

次に、重要度算出部5の動作の実施例について説明する。以下の実施例においては、上述した集計データ及びリンク関係データの値の例を用いる。

【0051】

まず、重要度算出部5による動作の第1の実施例として、作成日を重要視した重要度を算出する場合について説明する。

【0052】

まず、行列決定処理において、行列決定部17は、集計データから、検索対象文書 P_m についてのパラメータ毎の重みデータとして、作成日の重みを表す重みデータ x_m 、有用評価数の重みを表す重みデータ y_m 、参照回数の重みを表す重みデータ z_m を算出する。 x_m 、 y_m 、 z_m は、それぞれ次式で定義される。

【0053】

$$x_m = \exp(-p/100) * 1.5$$

$$y_m = \log(q+2) * 0.2$$

$$z_m = \log(r+2) * 0.03$$

【0054】

この式により x_m が高く設定される。図16は、本実施の形態に係る第1の実施例において行列決定処理により決定された重みデータを示す表である。次に、行列決定部17は、リンク関係データから P_n から P_m へのリンクの有無を判断し、リンクの有無と重みデータに従ってDの要素 d_{mn} を決定する。ここでは、 d_{mn} の算出に x_m 、 y_m 、 z_m の和を用いる場合を示す。 d_{mn} は、次式で定義される。

【0055】

$$d_{mn} = c/N + x_m + y_m + z_m$$

$$(m, n = 1, 2, \dots, N)$$

(P_n から P_m へのリンクがないとき)

$$d_{mn}$$

$$= c/N + x_m + y_m + z_m$$

$$+ (1-c) * (x_m + y_m + z_m) / L_n$$

$$(m, n = 1, 2, \dots, N)$$

(P_n から P_m へのリンクがあるとき)

$$c = 0.15$$

【0056】

次に、行列決定部17は、上述したリンク関係データの値と重みデータの値に対して上式を適用することによりDを決定する。図17は、本実施の形態に係る第1の実施例において行列決定処理により決定された行列Dを示す表である。この図において、行に付された番号はmを示し、列に付された番号はnを示し、mとnで指定される要素は d_{mn} の値を示す。

【0057】

次に、反復計算部18は、上述した反復計算処理のフローを実行することにより、重要度ベクトルuを算出する。図18は、本実施の形態に係る第1の実施例において反復計算処理により計算された計算結果を示す表である。 $k=10$ で $|v_k - v_{(k-1)}|$ がしきい値を下回ったため、反復計算部18は、繰り返し計算を終了し、 v_{10} をDに対する重要度ベクトルuとする。また、反復計算部18は、得られたuの要素である重要度 u_m ($m=1, 2, \dots, 6$)の大きさの順位を、 P_m ($m=1, 2, \dots, 6$)の表示順位(ランク)としても良い。上述した重み付けにより、作成日から集計日までの期間pが

10

20

30

40

50

小さいP__1、P__2の表示順位が高くなることが分かる。上述した第1の実施例によれば、作成日、有用評価数、参照回数、リンク関係の影響を考慮しつつ、作成日の影響を強めた重要度及び表示順位を得ることができる。

【0058】

次に、重要度算出部5による動作の第2の実施例として、有用評価数を重要視した重要度を算出する場合について説明する。

【0059】

まず、行列決定処理において、重みデータx__m、y__m、z__mは、それぞれ次式で定義される。

【0060】

$$\begin{aligned}x_m &= \exp(-p/100) * 0.2 \\ y_m &= \log(q+2) * 1.2 \\ z_m &= \log(r+2) * 0.03\end{aligned}$$

【0061】

この式を、第1の実施例と比較すると、最後の係数の大きさが異なり、y__mが高く設定される。図19は、本実施の形態に係る第2の実施例において行列決定処理により決定された重みデータを示す表である。次に、行列決定部17は、この重みデータを用い、第1の実施例と同様にしてDを決定する。図20は、本実施の形態に係る第2の実施例において行列決定処理により決定された行列Dを示す表である。

【0062】

次に、反復計算部18は、上述した反復計算処理のフローを実行することにより、重要度ベクトルuを算出する。図21は、本実施の形態に係る第2の実施例において反復計算処理により計算された計算結果を示す表である。k=10で|v__k - v__(k-1)|がしきい値を下回ったため、反復計算部18は、繰り返し計算を終了し、v__10をDに対する重要度ベクトルuとする。上述した重み付けにより、有用評価数qが大きいP__4、P__6の表示順位が上位になることが分かる。上述した第2の実施例によれば、作成日、有用評価数、参照回数、リンク関係の影響を考慮しつつ、有用評価数の影響を強めた重要度及び表示順位を得ることができる。

【0063】

次に、重要度算出部5による動作の第3の実施例として、参照回数を重要視した重要度を算出する場合について説明する。

【0064】

まず、行列決定処理において、重みデータx__m、y__m、z__mは、それぞれ次式で定義される。

【0065】

$$\begin{aligned}x_m &= \exp(-p/100) * 0.2 \\ y_m &= \log(q+2) * 0.2 \\ z_m &= \log(r+2) * 0.2\end{aligned}$$

【0066】

この式を、第1の実施例と比較すると、最後の係数の大きさが異なり、z__mが高く設定される。図22は、本実施の形態に係る第3の実施例において行列決定処理により決定された重みデータを示す表である。次に、行列決定部17は、この重みデータを用い、第1の実施例と同様にしてDを決定する。図23は、本実施の形態に係る第3の実施例において行列決定処理により決定された行列Dを示す表である。

【0067】

次に、反復計算部18は、上述した反復計算処理のフローを実行することにより、重要度ベクトルuを算出する。図24は、本実施の形態に係る第3の実施例において反復計算処理により計算された計算結果を示す表である。k=11で|v__k - v__(k-1)|がしきい値を下回ったため、反復計算部18は、繰り返し計算を終了し、v__11をDに対する重要度ベクトルuとする。上述した重み付けにより、参照回数rが大きいP__4、

10

20

30

40

50

P__3の表示順位が上位になることが分かる。上述した第3の実施例によれば、作成日、有用評価数、参照回数、リンク関係の影響を考慮しつつ、参照回数の影響を強めた重要度及び表示順位を得ることができる。

【0068】

次に、重要度算出部5による動作の第4の実施例として、重みデータをランダムに設定して重要度を算出する場合について説明する。

【0069】

まず、行列決定部17は、d__m__nの重みデータをそれぞれランダムな数として設定する。次に、行列決定部17は、リンク関係ファイルからP__nからP__mへのリンクの有無を判断し、リンクの有無と重みデータに従ってDの要素d__m__nを決定する。d__m__nは、次式で定義される。

【0070】

$$d_{m,n} = c / N + \text{Rand}() \\ (m, n = 1, 2, \dots, N) \\ (P_{n} \text{から} P_{m} \text{へのリンクがないとき}) \\ d_{m,n} \\ = c / N + (1 - c) / L_n + \text{Rand}() \\ (m, n = 1, 2, \dots, N) \\ (P_{n} \text{から} P_{m} \text{へのリンクがあるとき}) \\ c = 0.15$$

【0071】

ここで、Rand()は、0以上1未満のランダムな実数を生成する関数である。図25は、本実施の形態に係る第4の実施例において行列決定処理により決定されたRand()の値を示す表である。図26は、本実施の形態に係る第4の実施例において行列決定処理により決定された行列Dを示す表である。

【0072】

次に、反復計算部18は、上述した反復計算処理のフローを実行することにより、重要度ベクトルuを算出する。図27は、本実施の形態に係る第4の実施例において反復計算処理により計算された計算結果を示す表である。k=8で|v__k - v__(k-1)|がしきい値を下回ったため、反復計算部18は、繰り返し計算を終了し、v__8をDに対する重要度ベクトルuとする。上述した第4の実施例によれば、リンク関係の影響を考慮しつつ、乱数の影響を与えた重要度及び表示順位を得ることができる。

【0073】

次に、重要度算出部5による動作の第5の実施例として、重みデータを用いず、リンク関係のみから重要度を算出する場合について説明する。

【0074】

まず、行列決定部17は、リンク関係ファイルからP__nからP__mへのリンクの有無を判断し、リンクの有無に従ってDの要素d__m__nを決定する。d__m__nは、次式で定義される。

【0075】

$$d_{m,n} = c / N \\ (m, n = 1, 2, \dots, N) \\ (P_{n} \text{から} P_{m} \text{へのリンクがないとき}) \\ d_{m,n} \\ = c / N + (1 - c) / L_n \\ (m, n = 1, 2, \dots, N) \\ (P_{n} \text{から} P_{m} \text{へのリンクがあるとき}) \\ c = 0.15$$

【0076】

図28は、本実施の形態に係る第5の実施例において行列決定処理により決定された行

10

20

30

40

50

列Dを示す表である。このDは、リンク関係のみを反映した行列である。

【0077】

次に、反復計算部18は、上述した反復計算処理のフローを実行することにより、重要度ベクトルuを算出する。図29は、本実施の形態に係る第5の実施例において反復計算処理により計算された計算結果を示す表である。k=14で $|v_k - v_{(k-1)}|$ がしきい値を下回ったため、反復計算部18は、繰り返し計算を終了し、 v_{14} をDに対する重要度ベクトルuとする。上述した第5の実施例によれば、重み付けを用いないことにより、特許文献1の技術と同様、リンク関係のみからも重要度及び表示順位を得ることができる。

【0078】

本実施の形態によれば、検索システムの設計者は、上述した反復計算処理のアルゴリズムを用いることにより、従来より少ない制約の中で複数のパラメータの重み付けを定義し、行列Dを決定することができる。また、パラメータの数に関わらず、1回の行列演算で重要度を算出することができるため、全てのパラメータの情報を失わず、且つ高速に重要度を算出することができる。

【0079】

また、本実施の形態に係る文書重要度算出装置は、情報処理装置に容易に適用することができ、情報処理装置の性能をより高めることができる。ここで、情報処理装置には、例えばサーバ、PC(Personal Computer)等が含まれ得る。

【0080】

更に、文書重要度算出装置を構成するコンピュータにおいて上述した各ステップを実行させるプログラムを、文書重要度算出プログラムとして提供することができる。上述したプログラムは、コンピュータにより読取り可能な記録媒体に記憶させることによって、文書重要度算出装置を構成するコンピュータに実行させることが可能となる。ここで、上記コンピュータにより読取り可能な記録媒体としては、ROMやRAM等のコンピュータに内部実装される内部記憶装置、CD-ROMやフレキシブルディスク、DVDディスク、光磁気ディスク、ICカード等の可搬型記憶媒体や、コンピュータプログラムを保持するデータベース、或いは、他のコンピュータ並びにそのデータベースや、更に回線上の伝送媒体をも含むものである。

【産業上の利用可能性】

【0081】

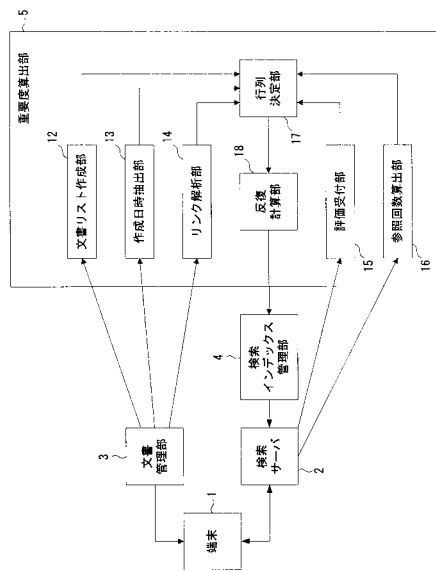
本発明によれば、文書に関するパラメータの影響を考慮した文書の重要度を効率的に算出することができる。

10

20

30

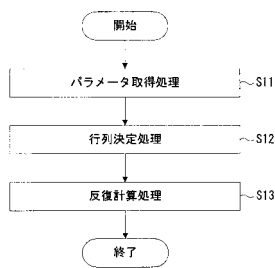
【図1】



【図2】

番号	URL
1	http://server1.domain1.com/document1.html
2	http://server2.domain2.com/document2.html
3	http://server3.domain3.com/document3.html
...	...

【図7】



【図8】

URL	重要度
http://server1.domain1.com/document1.html	1.000000
http://server2.domain2.com/document2.html	0.980003
http://server3.domain3.com/document3.html	0.971333
...	...

【図3】

URL	作成日時
http://server1.domain1.com/document1.html	2005/9/21
http://server2.domain2.com/document2.html	2005/9/22
http://server3.domain3.com/document3.html	2005/9/23
...	...

【図4】

リンク元URL	リンク先URL
http://server1.domain1.com/document1.html	http://server2.domain2.com/document2.html, http://server3.domain3.com/document3.html
http://server2.domain2.com/document2.html	http://server3.domain3.com/document3.html
http://server3.domain3.com/document3.html	http://server1.domain1.com/document1.html
...	...

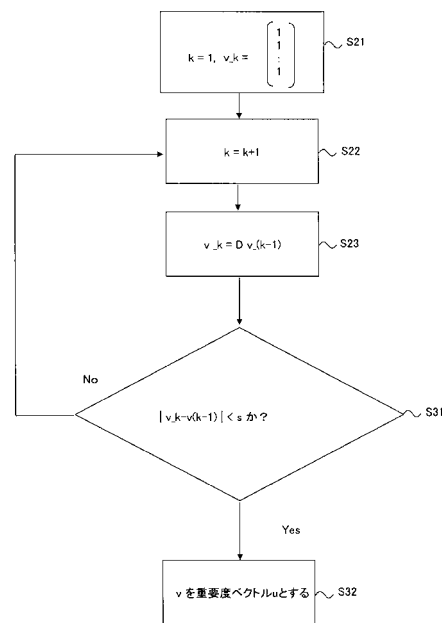
【図5】

URL	評価
http://server1.domain1.com/document1.html	5
http://server2.domain2.com/document2.html	2
http://server3.domain3.com/document3.html	3
...	...

【図6】

URL	参照回数
http://server1.domain1.com/document1.html	141
http://server2.domain2.com/document2.html	21
http://server3.domain3.com/document3.html	33
...	...

【図9】



フロントページの続き

- (56)参考文献 米国特許第6285999(US, B1)
米国特許出願公開第2005/0071741(US, A1)
米国特許出願公開第2003/0204502(US, A1)
特開2001-84256(JP, A)
特開2002-41524(JP, A)

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