

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第4845420号
(P4845420)

(45) 発行日 平成23年12月28日 (2011.12.28)

(24) 登録日 平成23年10月21日 (2011.10.21)

(51) Int.Cl.

G O 6 F 17/30 (2006.01)

F I

G O 6 F 17/30 3 5 O C

G O 6 F 17/30 3 7 O Z

請求項の数 8 外国語出願 (全 16 頁)

(21) 出願番号	特願2005-134488 (P2005-134488)	(73) 特許権者	500046438
(22) 出願日	平成17年5月2日 (2005.5.2)		マイクロソフト コーポレーション
(65) 公開番号	特開2005-322244 (P2005-322244A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成17年11月17日 (2005.11.17)		2-6399 レッドモンド ワン マイ
審査請求日	平成20年4月30日 (2008.4.30)		クロソフト ウェイ
(31) 優先権主張番号	10/837, 540	(74) 代理人	100077481
(32) 優先日	平成16年4月30日 (2004.4.30)		弁理士 谷 義一
(33) 優先権主張国	米国 (US)	(74) 代理人	100088915
			弁理士 阿部 和夫
		(72) 発明者	ベンユー チャン
			アメリカ合衆国 98052 ワシントン
			州 レッドモンド ワン マイクロソフト
			ウェイ マイクロソフト コーポレーシ
			ョン内

最終頁に続く

(54) 【発明の名称】 情報の多様性および豊富さを向上させるよう検索結果のドキュメントを順位付ける方法およびシステム

(57) 【特許請求の範囲】

【請求項 1】

検索結果のドキュメントを順位付けるコンピューティングデバイスであって、
コンピュータ実行可能な命令を有するメモリと、
前記メモリに格納されたコンピュータ実行可能な命令を実行する中央処理装置と
を備え、前記コンピュータ実行可能な命令は、
問い合わせをユーザから受信するステップと、
前記受信された問い合わせに対して検索結果としてのドキュメントを識別するステップと、
前記識別されたドキュメントのキーワードを識別するステップと、
前記検索結果の識別されたドキュメントの各組に対して、組のうちのドキュメントが
有する、前記組のうちのもう一方のドキュメントに対する類似度を計算するステップであ
って、前記類似度は、次の式

【数 1】

$$aff(d_i, d_j) = \frac{\vec{d_i} \cdot \vec{d_j}}{\|\vec{d_i}\|}$$

に基づいて計算され、ただし、 $aff(d_i, d_j)$ はドキュメント d_i のドキュメント d_j に対する類似度であり、ドキュメント d_i は前記組のうちのドキュメントを表して、
ドキュメント d_j は前記組のうちのもう一方のドキュメントを表して、

【数 2】

\vec{d}_i

はドキュメント d_i のベクトルを表して、

【数 3】

\vec{d}_j

はドキュメント d_j のベクトルを表して、および

【数 4】

$\|\vec{d}_i\|$

10

はドキュメント d_i のベクトルの長さを表して、各ベクトルは、対応するドキュメントの識別された少なくとも 1 つのキーワードに対する少なくとも 1 つのエントリを有する、ステップと、

前記検索結果の識別された各ドキュメントに対して、類似度の行列から得られる正規化された類似度の行列の要素に基づいて、前記識別されたドキュメントに対する情報の豊富さを計算するステップであって、前記情報の豊富さは、次の式

【数 5】

$$InfoRich(d_i) = \sum_{all\ j \neq i} InfoRich(d_j) \cdot \tilde{M}_{ji}$$

20

に基づいて計算され、ただし、 $InfoRich(d_i)$ はドキュメント d_i の情報の豊富さであり、ドキュメント d_i は前記検索結果の識別されたドキュメントを表して、ドキュメント d_j は前記検索結果における別のドキュメントを表して、

【数 6】

\tilde{M}_{ji}

は前記正規化された類似度の行列の要素であり、前記類似度の行列は前記検索結果における識別されたドキュメントおよび別のドキュメントの計算された類似度を含む、ステップと、

30

前記検索結果のドキュメントについて計算された情報の豊富さに基づいて順位付けされた検索結果のドキュメントの表示を表示装置に出力するステップと

を含むことを特徴とするコンピューティングデバイス。

【請求項 2】

前記ドキュメントを識別するステップは、

検索エンジンサービスへ前記問い合わせを送信するステップと、

検索結果としてのドキュメントを受信するステップと

を含むことを特徴とする請求項 1 に記載のコンピューティングデバイス。

【請求項 3】

前記ドキュメントは、ウェブページであることを特徴とする請求項 1 に記載のコンピューティングデバイス。

40

【請求項 4】

コンピュータシステムにおいて、ドキュメントの集まり内のドキュメントの情報の豊富さを中央処理装置およびメモリにより計算する方法であって、

前記集まり内の各ドキュメントが有する、前記集まり内の別のドキュメントに対する類似度を前記中央処理装置が識別するステップであって、前記類似度は、前記集まり内のドキュメントであるドキュメントの組のうちの 1 つのドキュメント、および前記集まり内の別のドキュメントである前記組のうちのもう一方のドキュメントを含むドキュメントの各組に対して識別されて、1 つのドキュメントの情報が別のドキュメントの情報にどの程度含まれているかを示す、ステップと、

50

前記集まり内の各ドキュメントに対して、類似度の行列から得られる正規化された類似度の行列の要素に基づいて前記情報の豊富さを前記中央処理装置が決定するステップと、

前記集まり内のドキュメントの前記決定された情報の豊富さを前記メモリに前記中央処理装置が格納するステップと、

前記格納された情報の豊富さに基づいて前記ドキュメントを前記中央処理装置が順位付けするステップであって、各ドキュメントに対して前記決定された情報の豊富さは、次の式

【数 7】

$$InfoRich(d_i) = \sum_{all\ j \neq i} InfoRich(d_j) \cdot \tilde{M}_{ji}$$

10

として定義され、ただし、 $InfoRich(d_i)$ はドキュメント d_i の情報の豊富さであり、ドキュメント d_i は前記集まり内のドキュメントを表して、ドキュメント d_j は前記集まり内の別のドキュメントを表して、

【数 8】

$$\tilde{M}_{ji}$$

は前記正規化された類似度の行列の要素であり、前記類似度の行列は、前記集まり内のドキュメントの、前記集まり内の別のドキュメントに対する識別された類似度を含み、前記集まり内のドキュメントはウェブページである、ステップと

20

を備えたことを特徴とする方法。

【請求項 5】

ウェブページの集まり内のドキュメントの情報の豊富さを計算するコンピューティングデバイスであって、

コンピュータ実行可能な命令を有するメモリと、

前記メモリに格納されたコンピュータ実行可能な命令を実行する中央処理装置と

を備え、前記コンピュータ実行可能な命令は、

前記集まり内の各ウェブページが有する、前記集まり内の別のウェブページに対する類似度を識別するステップであって、前記類似度は、前記ウェブページの集まり内のウェブページであるウェブページの組のうちの 1 つのウェブページ、および前記ウェブページの集まり内の別のウェブページである前記組のうちのもう一方のウェブページを含むウェブページの各組に対して識別されて、1 つのウェブページの情報が別のウェブページの情報にどの程度含まれているかを示す、ステップと、

30

前記集まり内の各ウェブページに対して、類似度の行列から得られる正規化された類似度の行列の要素に基づいて前記情報の豊富さを決定するステップと、

前記集まり内のウェブページの前記決定された情報の豊富さを前記メモリに格納するステップと、

前記格納された情報の豊富さに基づいて前記ウェブページを順位付けするステップであって、各ウェブページに対して前記決定された情報の豊富さは、次の式

【数 9】

$$InfoRich(d_i) = \sum_{all\ j \neq i} InfoRich(d_j) \cdot \tilde{M}_{ji}$$

40

として定義され、ただし、 $InfoRich(d_i)$ はウェブページ d_i の情報の豊富さであり、ウェブページ d_i は前記集まり内のウェブページを表して、ウェブページ d_j は前記集まり内の別のウェブページを表して、

【数 10】

$$\tilde{M}_{ji}$$

は前記正規化された類似度の行列の要素であり、前記類似度の行列は、前記集まり内のウ

50

ウェブページの、前記集まり内の別のウェブページに対する識別された類似度を含む、ステップと

を含むことを特徴とするコンピューティングデバイス。

【請求項6】

前記集まり内の別のウェブページに対する類似度を識別するステップは、類似度の図式を生成するステップであって、前記類似度の図式は、ドキュメントをノードとして表現し、類似度の値をノードの間の有向のエッジの重みとして表現する、ステップを含むことを特徴とする請求項5に記載のコンピューティングデバイス。

【請求項7】

前記ウェブページを順位付けするステップは、
検索エンジンサービスへ問い合わせを送信するステップと、
検索結果としてのウェブページを受信するステップと
を含むことを特徴とする請求項5に記載のコンピューティングデバイス。

10

【請求項8】

前記ウェブページを順位付けするステップは、前記問い合わせに対する前記検索結果の各ドキュメントの適合度にさらに基づくことを特徴とする請求項7に記載のコンピューティングデバイス。

【発明の詳細な説明】

【技術分野】

【0001】

20

本発明は、一般には、検索エンジンサービスに送信された検索要求により識別される検索結果のドキュメントを順位付ける技術に関する。

【背景技術】

【0002】

GoogleやOvertureなど多数の検索エンジンサービスにより、インターネットを介してアクセス可能な情報検索が提供されている。これらの検索エンジンは、ユーザに、ユーザが関心を持つであろうウェブページなどの表示ページを検索することを可能にする。ユーザが、検索語を含む検索要求を送信した後、検索エンジンサービスは、これらの検索語に関連する可能性があるウェブページを識別する。関連するウェブページを速やかに識別するために、検索エンジンサービスは、キーワードのウェブページへのマッピングを保持する。このマッピングは、各ウェブページのキーワードを抽出するために、ウェブ（すなわちワールドワイドウェブ）を「クロール（crawling）」することによって作成される。ウェブをクロールするために、検索エンジンサービスは、ルートウェブページ（root web page）のリストを使用して、ルートウェブページを介してアクセス可能なすべてのウェブページの識別を行うことができる。様々な周知の情報検索技術、例えば、ヘッドラインの単語、ウェブページのメタデータに与えられている単語、およびハイライトされている単語などを識別することなどを用いて、あらゆる特定のウェブページのキーワードを抽出する。検索エンジンサービスは、各ウェブページが検索要求にどの程度関連しているかを示す適合度の得点を、各組合せの緊密性、およびウェブページの人気度（例えば、GoogleのPageRank）などに基づいて計算することができる。次いで、検索エンジンサービスは、ユーザに、この適合度に基づいた順序で、各ウェブページへのリンクを表示する。より一般的には、検索エンジンは、あらゆるドキュメントの集まりにおける情報検索を可能にする。例えば、ドキュメントの集まりには、すべての米国特許、すべての連邦の法廷意見、会社のすべての保管文書などが含まれる。

30

40

【0003】

ウェブベースの検索エンジンサービスによって提供される検索結果のうちで、最も高い順位のウェブページを、人気のある同一の話題にすべて方向付けることがある。例えば、ユーザが「スピルバーク（Spiegelberg）」という検索語で検索要求を送信する場合、検索結果のうちで最も高い順位のウェブページは、おそらくスティーブン・スピルバ

50

ーグ (Steven Spielberg) に関係するはずである。しかし、ユーザが、スティーブン・スピルバーグに関心が無く、代わりに、同姓の数学の教授のホームページの所在位置を突き止めることに関心がある場合、ウェブページの順位は、ユーザの役には立たない。教授のホームページは、検索結果に含まれているかも知れないが、ユーザは、教授のホームページへのリンクの所在位置を突き止めるために、検索結果のウェブページへのリンクからなるページを何ページも綿密に見る必要がある場合がある。一般に、所望のドキュメントを検索結果の最初のページ上で識別しない場合、ユーザが、所望のドキュメントの所在位置を突き止めることは困難になることがある。さらに、ユーザは、関心のあるドキュメントを探すために、検索結果を複数のページにわたり見る必要がある場合、挫折することがある。

10

【0004】

【特許文献1】米国特許出願第____号明細書、"Method and System for Calculating Importance of a Block Within a Display Page"、____日出願

【発明の開示】

【発明が解決しようとする課題】

【0005】

最も高い順位のドキュメントのうちで、より多様な話題を提供する、ドキュメントを順位付けるための技術を有することが望ましく、さらに、このような最も高い順位の各ドキュメントは、その話題に関係する情報内容を非常に豊富に有することが望ましい。

【課題を解決するための手段】

20

【0006】

システムは、情報の豊富さおよび話題の多様性に基づいて検索結果のドキュメントを順位付ける。順位付けのシステムにより、検索結果のドキュメントをそれらの関連性に基づいてグループ化し、ドキュメントを同様な話題に関連付けることを意味する。順位付けのシステムにより、最も高い順位のドキュメントには、各トピックに及ぶ少なくとも1つのドキュメントを含み得るように、ドキュメントを順位付ける。次いで、順位付けのシステムにより、最も高い順位のドキュメントのうちの1つとして、グループ内で最も豊富な情報を有するドキュメントを、各グループから選択する。

【発明を実施するための最良の形態】

【0007】

30

情報の豊富さおよび話題の多様性に基づいて検索結果のドキュメントを順位付ける方法およびシステムを提供する。一実施形態では、順位付けのシステムは、検索結果に含まれる各ドキュメントの情報の豊富さを決定する。情報の豊富さは、ドキュメントがその話題に関係する情報をどの程度含むのかを示す尺度である。情報の豊富さが高いドキュメント（例えばウェブページ）は、同一の話題に関係し、情報の豊富さがより低いドキュメントの情報を含む情報を含むことがありそうであろう。順位付けのシステムは、検索結果のドキュメントをそれらの関連性に基づいてグループ化し、ドキュメントは同様な話題に関連付けられることを意味する。順位付けのシステムにより、最も高い順位のドキュメントには、各話題に及ぶ少なくとも1つのドキュメントを含むことができ、すなわち、各グループから1つのドキュメントを含むことができ得るように、ドキュメントを順位付ける。順位付けのシステムは、グループ内でドキュメントに属する情報の豊富さが最も高いドキュメントを、各グループから選択する。ドキュメントを順位付けの順序でユーザに提示する場合、ユーザは、おそらく、検索結果の最初のページで、人気のある単一の話題だけよりもむしろ、様々な話題に及ぶドキュメントを見つけるであろう。例えば、検索要求が、検索語の「スピルバーグ」を含む場合、検索結果の最初のページ上のある1つのドキュメントは、スティーブン・スピルバーグに関係しているとともに、最初のページ上の別のドキュメントは、スピルバーグ教授に関係していることが可能である。このようにして、検索結果の最初のページで多様性のある話題に及ぶドキュメントをユーザに提示する可能性は高くなり、また、関心のある話題が検索要求に関係する最も人気のある話題ではない場合でも、ユーザが挫折する可能性は低くなる。さらに、順位付けのシステムは、情報の豊富

40

50

さがより低いドキュメントよりも、情報の豊富さがより高いドキュメントを順位付けるので、ユーザは、検索結果の最初のページで提示されるドキュメントのうちから所望の情報を見つける可能性が高くなる。

【 0 0 0 8 】

一実施形態では、順位付けのシステムは、類似度の図式に基づいて、検索結果のドキュメントの情報の豊富さを計算する。類似度は、1つのドキュメントの情報が、別のドキュメントの情報に、どの程度含まれているかを示す尺度である。例えば、スピルバーグの映画のうちの1つを表面的に説明したドキュメントは、スピルバーグの映画のすべてについて詳しく説明したドキュメントに対して、高い類似度を有することがある。反対に、スピルバーグの映画のすべてについて詳しく説明したドキュメントは、スピルバーグの映画のうちの1つを表面的に説明したドキュメントに対して、比較的低い類似度を有する可能性がある。大きく異なる話題に関係するドキュメントは、互いに類似度を有していないはずである。他のすべてのドキュメントに対する各ドキュメントの類似度の集まりにより、類似度の図式を表現する。ドキュメントに対して高い類似度を有する他の多数のドキュメントがあるドキュメントは、高い情報の豊富さを有することがありそうであろう。理由は、そのドキュメントの情報は他の多数のドキュメントの情報を含むからである。さらに、高い類似度を有するこれらの他のドキュメントが、ドキュメント自体に比較的高い情報の豊富さをも有する場合、ドキュメントの情報の豊富さは、さらに高い。

【 0 0 0 9 】

一実施形態では、順位付けのシステムは、類似度の図式をやはり使用して、検索結果で高い順位のドキュメントの多様性を得る助けとなる。順位付けのシステムは、従来技術の順位付けの技術（例えば適合度）、情報の豊富さの技術、または他のいくつかの順位付けの技術に基づいたドキュメントの初期の順位の得点を有することができる。初めに、順位付けのシステムは、最も高い初期の順位の得点を有するドキュメントを、最も高い最終順位の得点を有するドキュメントとして選択する。次いで、順位付けのシステムは、選択したドキュメントに対して高い類似度を有する各ドキュメントの順位の得点を減少させる。それらのドキュメントの内容は、選択したドキュメントにおそらく含まれている冗長な情報であるので、順位付けのシステムは、順位の得点を減少させる。次いで、順位付けのシステムは、次に最も高い順位の得点を有するドキュメントを残りのドキュメントから選択する。順位付けのシステムは、新しく選択したドキュメントに対して高い類似度を有する各ドキュメントの順位の得点を減少させる。順位付けのシステムは、所望の個数のドキュメントが最終順位の得点を得る、すべてのドキュメントが最終順位の得点を得る、または、他の何らかの終了条件が満たされるまで、このプロセスを繰り返す。一実施形態では、多様性とは、ドキュメントの集まりにおける異なる話題の数を表現し、集まりにおけるドキュメントの情報の豊富さは、集まり全体に対するドキュメントの情報を提供する度合いを意味している。

【 0 0 1 0 】

情報の豊富さと多様性との組合せではなくて、情報の豊富さのみまたは多様性のみに基づいて、検索結果のドキュメントを順位付けることができることは当事業者には理解されよう。検索エンジンサービスは、例えば、同様の話題に関係するドキュメントのグループを識別し、そのグループ内の各ドキュメントの情報の豊富さを決定することにより、情報の豊富さのみを用いることがある。次いで、検索エンジンサービスは、決定した情報の豊富さをドキュメントの順位付けの計算に入れ、グループで最も高い情報の豊富さを有するドキュメントを、グループ内の他のドキュメントよりも高く順位付けることがありそうにすることがある。検索エンジンサービスは、例えば、同様の話題に関係するドキュメントのグループを識別し、各グループからの少なくとも1つのドキュメントを、その情報の豊富さに関わらず、検索結果として高く順位付け得ることによって、多様性のみを用いることがある。例えば、検索エンジンサービスは、検索結果の最初のページに、各グループから、グループのうちで最も高い適合度を有するドキュメントを表示するように選択することができる。

【 0 0 1 1 】

類似度の図式では、ドキュメントをノードとして表現し、類似度の値を、ノードの間の有向のエッジの重みとして表現する。順位付けのシステムでは、各ドキュメントを、ドキュメントの集まり内の他のすべてのドキュメントにマッピングする正方行列によって、類似度の図式を表現する。順位付けのシステムは、行列の要素の値を、対応するドキュメントの類似度に設定する。Mを行列とする場合、 M_{ij} は、ドキュメントjに対するドキュメントiの類似度を表現する。順位付けのシステムは、各ドキュメントをベクトルとして表現することにより、ドキュメントの類似度を計算する。ベクトルは、ドキュメントの情報内容を表現する。例えば、各ベクトルは、ドキュメントの最も重要な25個のキーワードを含むことができる。順位付けのシステムは、次の式に従って、類似度を計算することができる。

10

【 0 0 1 2 】

【数1】

$$\text{aff}(d_i, d_j) = \frac{\vec{d}_i \cdot \vec{d}_j}{\|\vec{d}_i\|} \quad (1)$$

【 0 0 1 3 】

ただし、 $\text{aff}(d_i, d_j)$ は、ドキュメント d_j に対するドキュメント d_i の類似度であり、 \vec{d}_i は、ドキュメント d_i のベクトルを表現し、 \vec{d}_j は、ドキュメント d_j のベクトルを表現し、

20

【 0 0 1 4 】

【数2】

$$\|\vec{d}_i\|$$

【 0 0 1 5 】

は、ベクトル \vec{d}_i の長さ表現する。式1では、類似度を、 \vec{d}_j から \vec{d}_i への射影の長さに設定する。類似度は、異なる多くの方法で定義することができることは、当事業者には理解されよう。例えば、別のドキュメントに対するひとつのドキュメントの類似度を、その他のドキュメントのキーワードにおいて、そのひとつのドキュメントのキーワードが占める割合に基づいて定義することができる。他のドキュメントに対するひとつのドキュメントの類似度は、集合論の用語では、2つのドキュメントの積集合のキーワード数を、その他のドキュメント内のキーワード数で割ったものとして表現することができる。行列Mの各要素は、類似度の図式において、ひとつのドキュメントのノードからその他のドキュメントのノードへの有向のエッジを表現する。一実施形態では、順位付けのシステムは、類似度の閾値（例えば、2）未満の類似度の値を零に設定する。概念上、類似度が低い場合、類似度の図式において、ひとつのドキュメントのノードからその他のドキュメントのノードへの有向のエッジが存在しないことを意味する。類似度の行列は、次の式で表現する。

30

【 0 0 1 6 】

【数3】

$$M_{ij} = \begin{cases} \text{aff}(d_i, d_j), & \text{add}(d_i, d_j) \geq \text{aff}_t \text{ の場合} \\ 0, & \text{add}(d_i, d_j) < \text{aff}_t \text{ の場合} \end{cases} \quad (2)$$

40

【 0 0 1 7 】

ただし、 M_{ij} は、行列の要素であり、 aff_t は、類似度の閾値である。ノードの間に多数のエッジを有するノードのグループは、単一的话题を表現する可能性がある。理由は、グループ内のドキュメントの多くが、互いに閾値の類似度よりも大きい類似度を有するからである。対照的に、ノードの間にリンクを有していないノードは、異なる話題に方向付けられたドキュメントを表現する。

【 0 0 1 8 】

50

順位付けのシステムは、類似度の図式に対して、エッジ (e d g e) 解析のアルゴリズムを適用することによって、各ドキュメントの情報の豊富さを計算する。順位付けのシステムは、各行の値の合計が 1 になるよう類似度の行列を正規化する。正規化された類似度の行列は、次の式で表現する。

【 0 0 1 9 】

【数 4】

$$\tilde{M}_{ij} = \begin{cases} M_{i,j} / \sum_{j=1}^n M_{ij}, & \sum_{j=1}^n M_{ij} \neq 0 \text{ の場合} \\ 0 & , \sum_{j=1}^n M_{ij} = 0 \text{ の場合} \end{cases} \quad (3)$$

10

【 0 0 2 0 】

ただし、

【 0 0 2 1 】

【数 5】

\tilde{M}_{ij}

【 0 0 2 2 】

は、正規化された類似度の行列の要素である。順位付けのシステムは、次の式に従って情報の豊富さを計算する。

20

【 0 0 2 3 】

【数 6】

$$\text{InfoRich}(d_i) = \sum_{\text{all } j \neq i} \text{InfoRich}(d_j) \cdot \tilde{M}_{ji} \quad (4)$$

【 0 0 2 4 】

ただし、 $\text{InfoRich}(d_i)$ は、ドキュメント d_i の情報の豊富さである。したがって、情報の豊富さは、帰納的に定義される。式 4 は、次の式により行列の形式で表現することができる。

【 0 0 2 5 】

30

【数 7】

$$\lambda = \tilde{M}^T \lambda \quad (5)$$

【 0 0 2 6 】

ただし、 $\lambda = [\text{InfoRich}(d_i)]_{n \times 1}$ は、正規化された類似度の行列

【 0 0 2 7 】

【数 8】

\tilde{M}^T

【 0 0 2 8 】

40

の固有ベクトルである。正規化された類似度の行列

【 0 0 2 9 】

【数 9】

\tilde{M}

【 0 0 3 0 】

は、通常、疎行列であるので、場合によっては、すべてが零である行が、行列の中に現れることがあり、いくつかのドキュメントは、いくつかのドキュメントに対して重要な類似度を有する他のドキュメントがないことを意味する。順位付けのシステムは、意味のある固有ベクトルを計算するために、ドキュメントの人気度に基づくドキュメントの順位とす

50

ることができる、ダンピング因子 (d u m p i n g f a c t o r) (例 え ば 、 8 5) を使用する。ダンピング因子を用いた情報の豊富さは、次の式で表現する。

【 0 0 3 1 】

【 数 1 0 】

$$\text{InfoRich}(d_i) = c \cdot \sum_{\text{all } j \neq i} \text{InfoRich}(d_j) \cdot \tilde{M}_{ji} + \frac{(1-c)}{n} \quad (6)$$

【 0 0 3 2 】

ただし、 c は、ダンピング因子であり、 n は、集まりの中のドキュメントの数である。
式 6 は、次のように行列の形式で表現することができる。

【 0 0 3 3 】

【 数 1 1 】

$$\hat{\lambda} = c \tilde{M}^T \hat{\lambda} + \frac{(1-c)}{n} \bar{e} \quad (7)$$

【 0 0 3 4 】

ただし、

【 0 0 3 5 】

【 数 1 2 】

\bar{e}

【 0 0 3 6 】

は、すべての成分が 1 に等しい単位ベクトルである。情報の豊富さの計算は、情報のフローおよびシンク (s i n k) モデルから類推することができる。このモデルでは、情報が、各反復でノードの間を流れる。ドキュメント d_i には、次によって表現される類似度を有するドキュメントの集合 $A(d_i)$ がある。

【 0 0 3 7 】

$$A(d_i) = \{ d_j \mid j \neq i, \text{aff}(d_i, d_j) > \text{aff}_t \} \quad (8)$$

各反復において、情報を、以下の規則のうち 1 つに従って流すことができる。

1 . 確率 c (すなわち、ダンピング因子) で、情報は、 $A(d_i)$ のうちの 1 つのドキュメントに流れ込み、ドキュメント d_j に流れ込む確率は、 $\text{aff}(d_i, d_j)$ に比例する。

2 . $1 - c$ の確率で、情報は、集まり内のあらゆるドキュメントに無作為に流れ込む。

【 0 0 3 8 】

上述のプロセスからマルコフ連鎖を帰納することができ、状態はドキュメントによって与えられ、推移 (すなわち、フロー) 行列は、

【 0 0 3 9 】

【 数 1 3 】

$$c \tilde{M}^T + \frac{(1-c)}{n} U \quad (a)$$

【 0 0 4 0 】

によって与えられる。ただし、

【 0 0 4 1 】

【 数 1 4 】

$$U = \begin{bmatrix} 1 \\ n \end{bmatrix}_{n \times n}$$

【 0 0 4 2 】

である。各状態の定常確率分布は、推移行列の主固有ベクトルによって与えられる。

10

20

30

40

50

【 0 0 4 3 】

一実施形態では、順位付けのシステムは、同一の話題に方向付けられる複数のドキュメントが、他の話題に方向付けられるドキュメントを排除して、すべてが高く順位付けられないよう、情報の豊富さと類似性のペナルティを組合せることによって、類似度の順位を計算する。類似性のペナルティの使用により、最も高い順位の複数のドキュメントの間における話題の多様性が増大する。順位付けのシステムは、切望の反復アルゴリズムを使用して、類似性のペナルティを計算することができ、ドキュメントの初期の類似度の順位を、ドキュメントの情報の豊富さに設定する。各反復では、アルゴリズムは、その次に最も高い類似度の順位を有するドキュメントを選択し、同一の話題に方向付けられるドキュメントの類似度の順位を、類似性のペナルティによって減少させる。したがって、ドキュメントを選択した後、同一の話題に方向付けられる他のすべてのドキュメントが減少された類似度の順位を有することにより、最も高い順位のドキュメントが様々な話題を表す機会を向上させるであろう。順位付けのシステムは、次の式に従って、ドキュメントの類似度の順位を減少させることができる。

【 0 0 4 4 】

【数 1 5】

$$AR_j = AR_j - \tilde{M}_j \cdot \text{InfoRich}(d_i) \quad (10)$$

【 0 0 4 5 】

ただし、 AR_j はドキュメント j の類似度の順位を表現し、 i は選択されたドキュメントである。類似性のペナルティは、類似度の行列に基づくので、ドキュメントが選択したドキュメントに類似しているほど、ドキュメントの類似性のペナルティは大きくなる。

【 0 0 4 6 】

一実施形態では、順位付けのシステムは、全体の順位を生成するために、類似度の順位と、テキストベースの順位（例えば、従来技術の適合度）とを組合せることができる。順位付けを、得点または順位に基づいて組合せることができる。組合せた得点に関して、ドキュメントの最終得点を表現する全体的な得点を得るために、テキストベースの得点を類似度の順位と組合せる。組合せた得点は、テキストベースの得点と類似度の順位との線形結合に基づくことができる。得点は異なるオーダの大きさを有するので、順位付けのシステムは得点を正規化する。組合せた得点は、次の式で表現することができる。

【 0 0 4 7 】

【数 1 6】

$$\text{Score}(q, d_i) = \alpha \cdot \frac{\text{Sim}(q, d_i)}{\text{Sim}_{\Theta}(q)} + \beta \cdot \frac{\log AR_{\Theta}}{\log AR_i}, \forall d_i \in \Theta \quad (11)$$

【 0 0 4 8 】

ただし、 $\alpha + \beta = 1$ であり、 $\text{Sim}(q, d_i)$ は検索要求 q に対する検索結果を表現し、 $\text{Sim}(q, d_i)$ は検索要求 q に対するドキュメント d_i の類似性を表現し、

【 0 0 4 9 】

【数 1 7】

$$\text{Sim}_{\Theta}(q) = \text{Max}_{\forall d_i \in \Theta} \text{Sim}(q, d_i) \quad (12)$$

【 0 0 5 0 】

であり、

【 0 0 5 1 】

【数 1 8】

$$AR_{\Theta} = \text{Max}_{\forall d_i \in \Theta} AR_i \quad (13)$$

【 0 0 5 2 】

である。

【 0 0 5 3 】

組合せた順位に関して、ドキュメントの最終順位を得るために、テキストベースの順位を、類似度の順位と組合せる。組合せた順位は、テキストベースの順位と類似度の順位との線形結合に基づることができる。組合せた順位は、次の式で表現することができる。

【 0 0 5 4 】

【 数 1 9 】

$$\text{Score}(q, d_i) = \alpha \cdot \text{Rank}_{\text{Sim}(q, d_i)} + \beta \cdot \text{Rank}_{\text{AR}_i}, \forall d_i \in \Theta \quad (14)$$

【 0 0 5 5 】

ただし、Score は、検索要求 q に対するドキュメント d_i の最終順位を表現し、

10

【 0 0 5 6 】

【 数 2 0 】

$\text{Rank}_{\text{Sim}(q, d_i)}$

【 0 0 5 7 】

は、テキストベースの順位を表現し、

【 0 0 5 8 】

【 数 2 1 】

$\text{Rank}_{\text{AR}_i}$

20

【 0 0 5 9 】

は、類似度の順位を表現する。結合のアルゴリズムの中の α および β の両方は、調整可能なパラメータである。 $\alpha = 1$ かつ $\beta = 0$ の場合、再順位付けは実行されず、検索結果は、テキストベースの検索に基づいて順位付けされる。 $\beta > 0$ の場合、再順位付けを行う際に、類似度の順位付けに、より重みかけられる。 $\alpha = 1$ かつ $\beta = 0$ の場合、類似度の順位付けだけに基いて、再順位付けが行われる。

【 0 0 6 0 】

図 1 は、一実施形態における類似度の図式を例示する図である。類似度の図式 1 0 0 には、ノード 1 1 1 ~ 1 1 5、ノード 1 2 1 ~ 1 2 4、およびノード 1 3 1 を含み、各々はドキュメントを表現する。ノードの間の有向のエッジは、別のノードに対するひとつのノードの類似度を示す。例えば、ノード 1 1 1 は、ノード 1 1 5 に対する類似度を有するが、ノード 1 1 5 は、ノード 1 1 1 に対する類似度を有していない（または閾値のレベルを下回る類似度を有する）。この例では、ノードグループ 1 1 0 は、同一の話題に方向付けられたノード 1 1 1 ~ 1 1 5 を含む。理由は、このノードグループのノードの間には多数のエッジが存在するからである。同様に、ノードグループ 1 2 0 は、同一の話題に方向付けられたノード 1 2 1 ~ 1 2 4 を含む。ノードグループ 1 3 0 には、1 つのノードしかないのは、このノードが、他のどのノードに対しても類似度を有しておらず、このノードに対する類似度を有するノードもないからである。ノード 1 1 5 は、おそらくノードグループ 1 1 0 のすべてのノードのうちで最も高い情報の豊富さを有し、ノード 1 2 4 は、おそらくノードグループ 1 2 0 のすべてのノードのうちで最も高い情報の豊富さを有する。理由は、各ノードは、ノードに対する類似度を有するノードの数が最も多いからである。

30

40

【 0 0 6 1 】

図 2 は、一実施形態における順位付けのシステムのコンポーネントを例示するブロック図である。順位付けのシステム 2 0 0 は、データストア 2 0 1 ~ 2 0 4、およびコンポーネント 2 1 1 ~ 2 1 6 を含む。ドキュメントストア 2 0 1 は、ドキュメントの集まりを収容するが、このストアは、インターネットを介して利用可能なすべてのウェブページを表現する場合もある。類似度の図式を生成するコンポーネント 2 1 1 は、ドキュメントストアのドキュメントに基づいて類似度の図式を生成する。類似度の図式を生成するコンポーネントは、類似度の図式ストア 2 0 2 内に類似度を格納する。情報の豊富さを計算するコンポーネント 2 1 2 は、類似度の図式ストアからの類似度の図式を入力し、各ドキュメン

50

トの情報の豊富さの得点を計算する。このコンポーネントは、計算した情報の豊富さの得点を、情報の豊富さストア203に格納する。一実施形態では、類似度の図式を生成するコンポーネント、および情報の豊富さを計算するコンポーネントは、検索を実施するのに先立って、類似度の図式および情報の豊富さの得点を生成するためにオフラインで実行することができる。検索を実施するコンポーネント213は、ユーザからの検索要求を受信し、ドキュメントストアのドキュメントから検索結果を識別する。検索を実施するコンポーネントは、検索結果を、検索要求に対する検索結果の各ドキュメントの適合度の表示と共に検索結果ストア204に格納する。類似性のペナルティを計算するコンポーネント214は、類似性ペナルティを計算して、検索結果ストア、類似度の図式ストア、および情報の豊富さストアの情報に基づく類似度の順位に適用する。類似度の順位を計算するコンポーネント215は、検索結果の中にある各ドキュメントの類似度の順位を生成する。類似度の順位を計算するコンポーネントは、ドキュメントの情報の豊富さ、類似度の図式ストア、および検索結果を計算に入れる。最終得点を計算するコンポーネント216は、類似度の順位と適合度の得点とを組合せて、最終得点を計算する。

10

【0062】

順位付けのシステムが実装されるコンピューティングデバイスには、中央処理装置、メモリ、入力装置（例えば、キーボード、ポインティングデバイス）、出力装置（例えば、表示装置）、および記憶装置（例えば、ディスク装置）が含まれる。メモリおよび記憶装置は、順位付けのシステムを実装する命令を含むことができるコンピュータ可読媒体である。さらに、データ構造およびメッセージ構造は、通信リンク上の信号などのデータ伝送媒体を介して格納または伝送することができる。様々な通信リンク、例えば、インターネット、LAN、WAN、または、ポイントツーポイントのダイヤルアップ接続などを使用することができる。

20

【0063】

順位付けのシステムは、様々な動作環境で実装することができる。使用するのに適切となり得る、周知の様々なコンピューティングシステム、環境、および構成には、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドまたはラップトップデバイス、マルチプロセッサシステム、マイクロプロセッサベースのシステム、プログラム可能家庭用電化製品、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、および、あらゆる上述のシステムまたは装置を含む分散コンピューティング環境などが含まれる。

30

【0064】

1つまたは複数のコンピュータまたは他の装置によって実行される、プログラムモジュールなどコンピュータ実行可能命令の一般的コンテキストにおいて、順位付けのシステムを説明することができる。一般に、プログラムモジュールには、特定のタスクを実行または特定の抽象データ型などを実装する、ルーチン、プログラム、オブジェクト、コンポーネント、およびデータ構造などが含まれる。一般に、様々な実施形態において、プログラムモジュールの諸機能を、所望に応じて、結合または分散することができる。

【0065】

図3は、一実施形態における順位付けのシステムの全体的な処理を例示する流れ図である。順位付けのシステムには、検索結果を表現することができるドキュメントの集まりが提供される。ブロック301では、コンポーネントが、ドキュメントの集まりについての類似度の図式を構成する。コンポーネントは、オフラインでドキュメントの集成の中にあるすべてのドキュメント（例えば、すべてのウェブページ）に及ぶ、またはリアルタイムで集まりのドキュメントのみに及ぶ類似度の図式を構成することができる。ブロック302では、コンポーネントは、集まりの各ドキュメントの情報の豊富さを計算する。ブロック303では、コンポーネントは、集まりのドキュメントを順位付けし、次いで完了となる。

40

【0066】

図4は、一実施形態における類似度の図式を構成するコンポーネントの処理を例示する流れ図である。コンポーネントは、ドキュメントの集まりが伝わり、これらのドキュメン

50

トについての類似度の図式を構成する。ブロック401～403では、コンポーネントは、ドキュメント集まりの中にある各ドキュメントについてのドキュメントベクトルを生成するループを実行する。ブロック401では、コンポーネントは、集まり内の次のドキュメントを選択する。決定ブロック402では、集まり内のすべてのドキュメントを既に選択した場合、ブロック404に進み、そうでなければ、ブロック403に進む。ブロック403では、コンポーネントは、選択されたドキュメントに対するドキュメントベクトルを生成し、次いで、集まり内の次のドキュメントを選択するためにブロック401へのループを実行する。ブロック404～408では、コンポーネントは、集まり内のドキュメントの各組についての類似度を計算する。ブロック404では、コンポーネントは、第1のドキュメントから開始し、集まりの中にある次のドキュメントを選択する。決定ブロック405では、すべてのドキュメントが既に選択された場合、コンポーネントは、類似度の図式を返し、そうでない場合、ブロック406に進む。ブロック406～408では、コンポーネントは、集まりの各ドキュメントを選ぶループを実行する。ブロック406では、コンポーネントは、第1のドキュメントから開始し、集まりの中にある次のドキュメントを選択する。決定ブロック407では、集まり内のすべてのドキュメントが既に選ばれた場合、コンポーネントは、集まり内の次のドキュメントを選択するためにブロック404へのループを実行し、そうでない場合、ブロック408に進む。ブロック408では、コンポーネントは、式1に従って(ブロック406で)選んだドキュメントに対する(ブロック404で)選択されたドキュメントの類似度を計算し、次いで、集まり内の次のドキュメントを選ぶブロック406へのループを実行する。

【0067】

図5は、一実施形態におけるドキュメントを順位付けるコンポーネントの処理を例示する流れ図である。コンポーネントには、類似度の図式が生成され、各ドキュメントの情報の豊富さが計算されたドキュメントの集まりを伝える。ブロック501～503では、コンポーネントは、集まり内の各ドキュメントの類似度の順位を、その情報の豊富さに初期設定するループを実行する。ブロック501では、コンポーネントは、集まり内の次のドキュメントを選択する。決定ブロック502では、すべてのドキュメントを既に選択した場合、ブロック504に進み、そうでない場合、ブロック503に進む。ブロック503では、コンポーネントは、選択したドキュメントの類似度の順位を、選択したドキュメントの情報の豊富さに設定し、集まり内の次のドキュメントを選択するブロック501へのループを実行する。ブロック504～508では、コンポーネントは、ドキュメントの組を識別し、類似性のペナルティによって類似度の順位を調整するループを実行する。ブロック504では、コンポーネントは、次に最も高い類似度の順位を有するドキュメントを選択する。決定ブロック505では、終了条件に到達した場合、コンポーネントは、順位付けしたドキュメントを返し、そうでない場合、ブロック506に進む。ブロック506～508では、コンポーネントは、ドキュメントを選び、類似性のペナルティにより類似度を調整するループを実行する。ブロック506では、コンポーネントは、選ばれたドキュメントから選択されたドキュメントへの類似度についての類似度の図式において、零以外の値で指示されるように、選択されたドキュメントに対する類似度を有する、次のドキュメントを選ぶ。決定ブロック507では、このようなドキュメントがすべて既に選ばれた場合、コンポーネントは、次に最も高い類似度の順位を有するドキュメントを選択するブロック504へのループを実行する。ブロック508では、コンポーネントは、式10に従って、類似性のペナルティによって選ばれたドキュメントに対する類似度の順位を調整する。次いで、コンポーネントは、選択されたドキュメントに対する類似度を有する次のドキュメントを選ぶブロック506へのループを実行する。

【0068】

本明細書では、例示の目的のために、順位付けのシステムの特定の実施形態について説明したが、本発明の精神および範囲から逸脱することなく様々な変形形態を構成することができることは当事業者には理解されよう。一実施形態では、順位付けのシステムは、ドキュメントごとではなくブロックごとを基礎にして類似度および情報の豊富さを計算する

ことができる。ブロックは、単一の話題に一般的に係るウェブページの情報を表現する。ウェブページの順位付けは、そのウェブページに対するブロックの重要度に部分的に基づくことができる。ブロックの重要度に関しては文献に記載されている（例えば、特許文献 1 参照。この文献を参照により本明細書に援用する。）。したがって、本発明は添付の特許請求の範囲を除いて限定されるものではない。

【図面の簡単な説明】

【 0 0 6 9 】

【図 1】一実施形態における類似度の図式を例示する図である。

【図 2】一実施形態における順位付けのシステムのコンポーネントを例示するブロック図である。

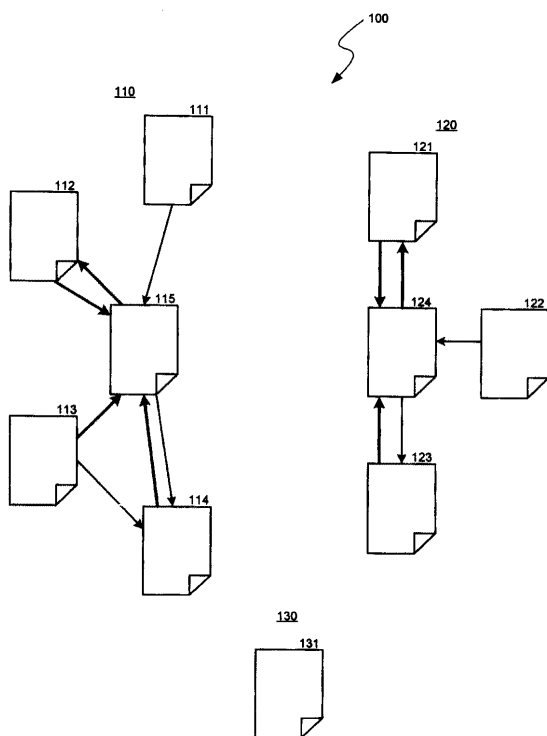
【図 3】一実施形態における順位付けのシステムの全体的な処理を例示する流れ図である。

【図 4】一実施形態における類似度の図式を構成するコンポーネントの処理を例示する流れ図である。

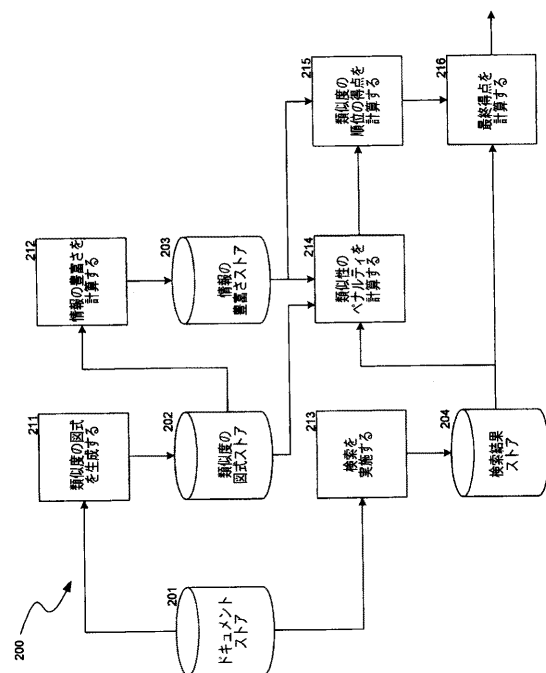
【図 5】一実施形態におけるドキュメントを順位付けるコンポーネントの処理を例示する流れ図である。

10

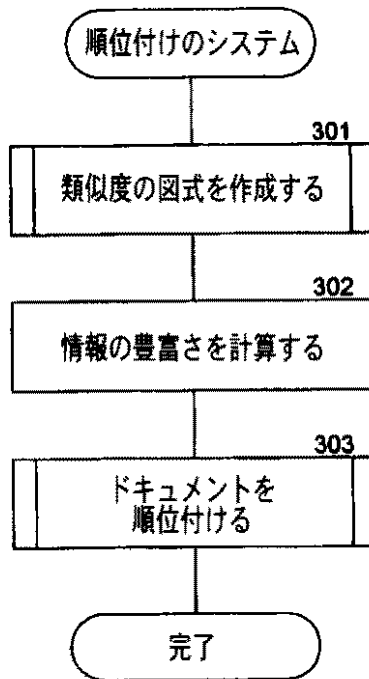
【圖 1】



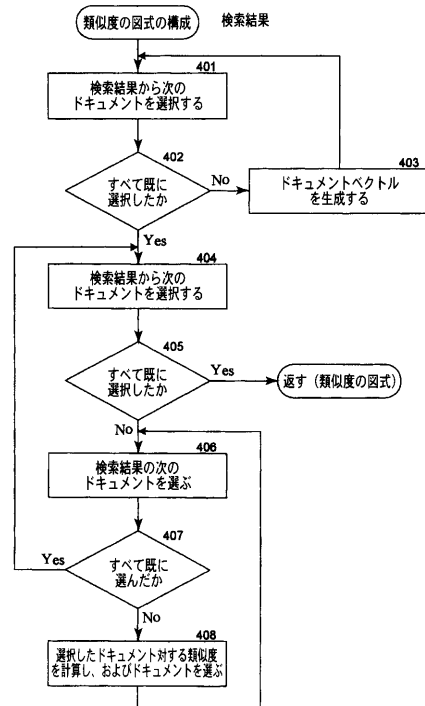
【圖 2】



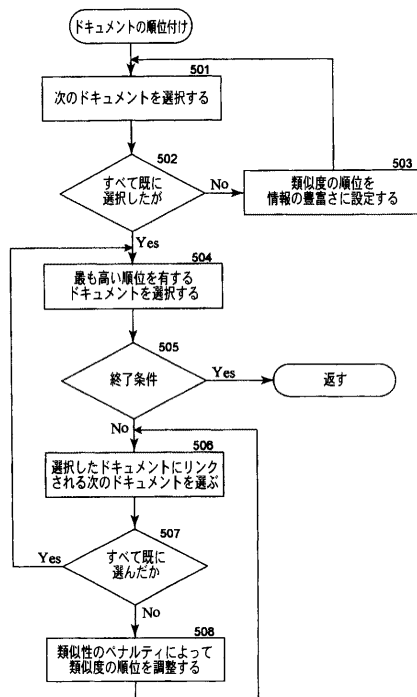
【図 3】



【図 4】



【図 5】



フロントページの続き

- (72)発明者 ホア - ジュン チュン
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 ウェイ - イン マ
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 チョン チェン
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内

審査官 打出 義尚

- (56)参考文献 特開2004-318528(JP, A)
Carbonell J, The use of MMR, diversity-based reranking for reordering documents and pr
oducing summaries, Proceedings of the 21st annual international ACM SIGIR conference o
n Research and development in information retrieval, 米国, ACM, 1998年, pp.
335 - 336, URL, <http://doi.acm.org/10.1145/290941.291025>
Goldstein J, Multi-document summarization by sentence extraction, NAACL-ANLP 2000 Work
shop on Automatic summarization - Volume 4, 米国, Association for Computational Lingui
stics, 2000年, pp. 40 - 48, URL, <http://dx.doi.org/10.3115/1117575.1117580>
Inderjeet Mani, 自動要約, 日本, 共立出版株式会社, 2003年 6月20日, 初版, pp.
181 - 183

- (58)調査した分野(Int.Cl., DB名)
G06F 17/30