

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3953295号
(P3953295)

(45) 発行日 平成19年8月8日(2007.8.8)

(24) 登録日 平成19年5月11日(2007.5.11)

(51) Int. Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 350C
 G06F 17/30 330A
 G06F 17/30 340Z
 G06F 17/30 412

請求項の数 12 (全 44 頁)

| | |
|---|---|
| <p>(21) 出願番号 特願2001-324437 (P2001-324437) (22) 出願日 平成13年10月23日(2001.10.23) (65) 公開番号 特開2003-141160 (P2003-141160A) (43) 公開日 平成15年5月16日(2003.5.16) 審査請求日 平成14年10月23日(2002.10.23)</p> | <p>(73) 特許権者 390009531 インターナショナル・ビジネス・マシー ズ・コーポレーション INTERNATIONAL BUSIN ESS MASCHINES CORPO RATION アメリカ合衆国10504 ニューヨーク 州 アーモンク ニュー オーチャード ロード (74) 代理人 100086243 弁理士 坂口 博 (74) 代理人 100091568 弁理士 市位 嘉宏</p> |
|---|---|

最終頁に続く

(54) 【発明の名称】 情報検索システム、情報検索方法、情報検索を実行させるためのプログラムおよび情報検索を実行させるためのプログラムが記録された記録媒体

(57) 【特許請求の範囲】

【請求項1】

データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索システムであって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて前記情報検索システムに保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、

前記ドキュメント - アトリビュート行列から共分散行列を生成し、ドキュメントベクトルの積和行列 (SUM(M)₁) と、ドキュメント - ベクトルの平均 (MEAN(M)₁) と、ドキュメント - ベクトルの平均の積行列 (SUM(M)₂) と、ドキュメントの全数 (M) とを保持させるための手段と、

前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するための手段と、

生成された前記共分散行列と、前記SUM(M)₁と、前記MEAN(M)₁と、前記SUM(M)₂と、前記Mと、追加された前記副行列のドキュメント数(H)とからなる前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記共分散行列を更新し、更新された前記共分散行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するための手段と、

前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うための手段と

を含む、情報検索システム。

【請求項 2】

前記すべてのドキュメント - アトリビュート行列の次元削減を実行するための手段は、追加された前記副行列を含むすべてのドキュメント - アトリビュート行列の共分散行列 C' を、下記式

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

により生成する手段を含む

請求項 1 に記載の情報検索システム。

【請求項 3】

さらに、ドキュメント - ベクトルに含まれるアトリビュートを自動的に検索し、アトリビュート・ハッシュ・テーブルを生成して前記アトリビュートを追加または削除するための手段を含む

請求項 1 または 2 に記載の情報検索システム。

【請求項 4】

データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索方法であって、前記方法は情報検索システムによって実行され、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、

前記情報検索システムが備える保持させるための手段が、前記ドキュメント - アトリビュート副行列から共分散行列を生成し、ドキュメントベクトルの積和行列 (SUM(M)₁) と、ドキュメント - ベクトルの平均 (MEAN(M)₁) と、ドキュメント - ベクトルの平均の積和行列 (SUM(M)₂) と、ドキュメントの全数 (M) とを保持させるステップと、

前記情報検索システムが備えるドキュメント - アトリビュート副行列を生成するための手段が、前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、

前記情報検索システムが備える次元削減を実行するための手段が、生成された前記共分散行列と、前記 SUM(M)₁ と、前記 MEAN(M)₁ と、前記 SUM(M)₂ と、前記 M と、追加された前記副行列のドキュメント数 (H) とからなる前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記共分散行列を更新し、更新された前記共分散行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、

前記情報検索システムが備える情報検索を行うための手段が、前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと

を含む、情報検索方法。

【請求項 5】

前記すべてのドキュメント - アトリビュート行列の次元削減を実行するステップは、前記情報検索システムが備える共分散行列 C' を生成する手段が、追加された前記副行列を含むすべてのドキュメント - アトリビュート行列の共分散行列 C' を、下記式

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

により生成するステップを含む

請求項 4 に記載の情報検索方法。

【請求項 6】

前記情報検索システムが備えるアトリビュートを追加または削除するための手段が、さらに、ドキュメント - ベクトルに含まれるアトリビュートを自動的に検索し、アトリビュ

10

20

30

40

50

ート・ハッシュ・テーブルを生成して前記アトリビュートを追加または削除するステップを含む

請求項 4 または 5 に記載の情報検索方法。

【請求項 7】

データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索方法をコンピュータに実行させるためのプログラムであって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、前記ドキュメント - アトリビュート副行列から共分散行列を生成し、ドキュメントベクトルの積和行列 ($SUM(M)_1$) と、ドキュメント - ベクトルの平均 ($MEAN(M)_1$) と、ドキュメント - ベクトルの平均の積行列 ($SUM(M)_2$) と、ドキュメントの全数 (M) とを保持させるステップと、

前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、

生成された前記共分散行列と、前記 $SUM(M)_1$ と、前記 $MEAN(M)_1$ と、前記 $SUM(M)_2$ と、前記 M と、追加された前記副行列のドキュメント数 (H) とからなる前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記共分散行列を更新し、更新された前記共分散行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、

前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと

を含む、情報検索方法を前記コンピュータに実行させるためのプログラム。

【請求項 8】

前記すべてのドキュメント - アトリビュート行列の次元削減を実行するステップにおいて、追加された前記副行列を含むすべてのドキュメント - アトリビュート行列の共分散行列 C' を、下記式

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

により生成するステップを実行させる

請求項 7 に記載のプログラム。

【請求項 9】

さらに、ドキュメント - ベクトルに含まれるアトリビュートを自動的に検索し、アトリビュート・ハッシュ・テーブルを生成して前記アトリビュートを追加または削除するステップを実行させる

請求項 7 または 8 に記載のプログラム。

【請求項 10】

データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索方法をコンピュータに実行させるためのプログラムが記録されたコンピュータ可読な記録媒体であって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、

前記ドキュメント - アトリビュート行列から共分散行列を生成し、ドキュメントベクトルの積和行列 ($SUM(M)_1$) と、ドキュメント - ベクトルの平均 ($MEAN(M)_1$) と、ドキュメント - ベクトルの平均の積行列 ($SUM(M)_2$) と、ドキュメントの全数 (M) とを保持させるステップと、

前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、

生成された前記共分散行列と、前記 $SUM(M)_1$ と、前記 $MEAN(M)_1$ と、前記 $SUM(M)_2$ と、前記 M と、追加された前記副行列のドキュメント数 (H) とからなる

10

20

30

40

50

前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記共分散行列を更新し、更新された前記共分散行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、

前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと

を含む、情報検索方法を前記コンピュータに実行させるためのプログラムが記録された記録媒体。

【請求項 1 1】

前記すべてのドキュメント - アトリビュート行列の次元削減を実行するステップにおいて、追加された前記副行列を含むすべてのドキュメント - アトリビュート行列の共分散行列 C' を、下記式

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

により生成するステップを実行させる

請求項 1 0 に記載の記録媒体。

【請求項 1 2】

さらに、ドキュメント - ベクトルに含まれるアトリビュートを自動的に検索し、アトリビュート・ハッシュ・テーブルを生成して前記アトリビュートを追加または削除するステップを実行させる

請求項 1 0 または 1 1 に記載の記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、データベースに保持されたデータの情報検索に関し、より詳細にはデータベースに対して情報が逐次的に追加され、情報が更新されるきわめて大きなデータベースのための情報検索システム、情報検索方法、情報検索を行うためのプログラムおよび該プログラムが記録されたコンピュータ可読な記録媒体に関する。

【0002】

【従来の技術】

ドキュメントのベクトル空間モデルを使用する情報検索システムは、きわめて大きな、静的なデータベースに対して比較的成功的なものである。Deerwesterらは、テキストに基づいたドキュメントの情報検索に関連する問題をより低次元のサブスペースへとマッピングするアルゴリズムである、潜在的意味解析法(Latent Semantic Indexing)を開発し、実時間における検索を可能とさせている。Deerwesterらのアルゴリズムは、ドキュメント - アトリビュート行列の特異値分解(Singular value decomposition: SVD)の効率的で、数値的に正確な計算に基づくものである。上述したDeerwesterらの手法は、精度が高く充分なものではあるが、そのために使用するサブスペースの基底ベクトルの決定は通常、逐次的なドキュメントデータの追加および削除の要求がある場合でもドキュメントデータに対する特異値分解を含む計算により実行されるのでコストが高く、例えば一晩を要するという不都合がある。

【0003】

情報検索における上述したSVDに基づく次元削減のためのアルゴリズムは、例えば、ニュース・データベース、医療データベース、消費者プロフィール・データベースといったきわめて迅速に情報が追加される、きわめて大きなデータベースから情報を検索するのに適切ではない。この理由は、主として(1)データベースは、新たに追加される情報を考慮すると共に、日付が古くなったり不要になった情報を削除するため、本明細書においてアップデートまたはダウンデートとして参照するプロセスを通じて頻繁に更新されなければならないこと、(2)データベース内のコンテンツの変更に対応するためのデータベースの更新(すなわちアップデートまたはダウンデート)のたびに、ドキュメント - アトリ

10

20

30

40

50

ビュート行列に対して再度基底ベクトルを算出する必要が生じ、このため、計算時間、労力、ハードウェア資源などに関連してコストが極めて高くなることを挙げる事ができる。

【 0 0 0 4 】

< 先行技術の開示 >

これまで上述した問題点を解決するべく、いくつかの方法が提案されている。Berry, DumaisおよびO'Brienは、“Using linear algebra for intelligent information retrieval”、1995、pp.573-595、において、およびZhaおよびSiomonは、“On updating problems in latent semantic indexing”、SIAM Journal of Scientific Computation、Vol. 21、No. 2、pp. 782-791、March 2000において、SVDに基づく次元削減につき、ベクトル空間モデルの下でドキュメントをアップデート、すなわちドキュメントを付け加える方法を提案している。しかしながら、Berryらの方法は、ベクトル空間モデルに対するアップデートの正確な計算ではなく単に近似でしかないために、信頼性において充分ではないという不都合があった。

10

【 0 0 0 5 】

Berryら(1995)およびZhaら(2000)によるDeerwesterらのアルゴリズムに対するアップデート方法も特異値分解に対する非線形な近似解であり、信頼性に欠けるという不都合がある。

【 0 0 0 6 】

一方、Witterは、“Downdating the latent semantic model for information retrieval”、M.A. Thesis、Univ. of Tennessee、Knoxville、Dec. 1997、において、Deerwesterらのアルゴリズムのためのダウンデート方法を開示している。しかしながら、Witterの方法を急速に変化するデータベースのダウンデートに対して適用すると、(1)ドキュメントが一度に削除されるので、以後に連続するダウンデートのための浮動小数点計算における誤差が蓄積し、これが重大な影響を与えること、(2)ダウンデートは、次元削減されたドキュメント-アトリビュート行列に対してのみ実行されるので、次元削減された行列は、その次元を失い元々の行列にある少量ではあるが意味的に重要ないわゆるアウトライア・ドキュメントを検出できず、ドキュメントが削除されるにつれて主要トピックスのみになってしまう可能性があることである。これについては、本願出願人らによる特許出願、特願2001-205183号にも詳細に記載されている。

20

30

【 0 0 0 7 】

上述した理由から、Witterのアルゴリズムによるダウンデートされたドキュメント-アトリビュート行列の出力は不正確であり、情報検索を含んだ多くの用途に対して適切ではない。なお、浮動小数点計算における誤差に関する検討は、最も標準的な数値解析のテキストに記載されており、このようなテキストとしては、GolubおよびVan Loan、“Matrix Computations”、第2版、John Hopkins Univ. Press、Baltimore、MD、1989およびGoldbergによる総説論文“What every computer scientist should know about floating-point arithmetic”、ACM Computing Surveys、Vol. 23、No. 1、1991年3月を挙げる事ができる。

【 0 0 0 8 】

またこれまで、急速に変化するデータベースを含み、SVDに基づいた情報検索用途に使用することができるに足る、十分な精度および効率で、次元削減されたドキュメント-アトリビュート行列のアップデートを実行するための方法は知られていない。精度および効率を必要とするアップデート用途の重要なクラスの例としては、本願出願人らによるデータベースにおける新たなトピックス/事象の検出およびその累積の追跡を挙げる事ができる。

40

【 0 0 0 9 】

一方で、SVDのプロセスについて考察すれば、一般的なドキュメント-アトリビュート行列AのSVDは、下記式により与えられる。

【 0 0 1 0 】

50

【数 1】

$$A = U \Sigma V^T$$

(上記式中、 U 、 V は、直交行列を示し、 Σ は、対角行列を示し、 T は、行列の転置を示す。以下、本発明において同様である。)

【0011】

この場合、上述した行列 A の SVD を、 $A^T A$ や $A A^T$ の固有値問題として実行する方法もあり得る。しかしながら、 $A^T A$ や $A A^T$ 生成させる場合には、小さな特異値に対して激しい精度低下を生じさせるという問題が生じ、反復してドキュメント - アトリビュート行列のアップデートに対応すると、出力精度が著しく低下してしまうという問題が生じることにもなる。

10

【0012】

【発明が解決しようとする課題】

すなわち本発明は、 SVD を使用してきわめて大きく、かつ急速に変化するデータベースから実時間での情報検索を行うための情報検索システム、情報検索方法、情報検索を実行させるためのプログラムおよび情報検索を実行させるためのプログラムが記録された記録媒体を提供することを目的とする。

【0013】

【課題を解決するための手段】

本発明は、上述した課題を解決するに際して、ドキュメント - アトリビュート行列の以前に算出した結果を一部再利用すれば、計算時間の短縮、精度向上、および効率化を達成することができる、という新奇な着想に基づいてなされたものである。

20

【0014】

本発明は、上述した以前の計算結果を使用するプロセスとして、ドキュメント - アトリビュート行列を QR 分解し、新たに付け加えられたドキュメントを含むドキュメント - アトリビュート行列を、すでに算出された QR 分解の結果のうちの直前に生成された行列 R を使用する。具体的には行列 R と付け加えられたドキュメント - アトリビュート副行列を含む行列をハウスホルダー変換し、付け加えられたドキュメント - アトリビュート副行列の結果を反映させて新たな行列 R として更新する。

30

【0015】

本発明においては、上述のようにして更新された上三角行列 R を SVD に用い、得られた特異ベクトルを使用してドキュメント - アトリビュート行列の次元削減を実行することにより、急速にデータが追加されるデータベースにおける次元削減プロセスを高精度、かつ効率的に実行させ、最新のデータベースでの情報検索を可能にするものである。

【0016】

また、本発明の別の実施の形態においては、ドキュメント - アトリビュート行列について共分散行列を生成し、追加されたドキュメントからなるドキュメント - アトリビュート行列について以前に算出された共分散行列に関連する所定の行列を使用して共分散行列を更新する。本発明においては、共分散行列の更新の際にアトリビュートを追加・削除することもできる。上述のようにして得られた共分散行列に対し SVD を実行して特異ベクトルを生成し、ドキュメント - アトリビュート行列の次元削減に使用することにより、次元削減プロセスの高精度化および効率化を達成するものである。

40

【0017】

すなわち、動的に変化するドキュメント - アトリビュート行列の本発明による特異値トリプレット(特異値と、それに対応する左および右特異ベクトル)のアップデートは、上述した先行技術において提案された結果よりも、より数値的に正確な結果を与えることができる。本発明は、特にアップデートされた行列の SVD を迅速に計算すると共に、特異値トリプレットの近似を与えるものではなく、数値的に正確な計算を実行させるものである。

50

【 0 0 1 8 】

さらに、本発明は、多数回のアップデートが行われることに対応しているので、浮動小数点誤差の蓄積を最小化させるものである。情報検索システムによる次元削減の問題に対してSVDを使用する結果の出力の品質は、一般には、算出された特異値トリプレットの精度に依存するので、本発明は、従来の技術においてこれまで開示されている方法よりもより良好な結果を与えることが可能となる。

【 0 0 1 9 】

さらに、本発明は、ドキュメント - アトリビュート行列のSVDに基づいたアップデートの正確な計算を可能とするので、本発明は、例えば新たなトピックス / 事象の検出といった、内容の著しく変化するデータベースにおけるドキュメントといったデータの検出に使用することができる。また、本発明の方法は、例えば特願2000 - 175848号、特願2001 - 157614号にも記載されているトピックス / 事象の追跡にも適用することができる。

10

【 0 0 2 0 】

すなわち、本発明によれば、データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索システムであって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて前記情報検索システムに保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、

前記ドキュメント - アトリビュート行列から所定の行列を生成して保持させるための手段と、

20

前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するための手段と、

前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記所定の行列を更新し、更新された前記所定の行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するための手段と、

前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うための手段と

を含む、情報検索システムが提供される。

【 0 0 2 1 】

本発明においては、前記ドキュメント - アトリビュート副行列をQR分解し、前記所定の行列として行列Rのみを使用する手段を含むことができる。本発明においては、前記保存された行列Rの更新を、前記ドキュメント - アトリビュート副行列に関連したハウスホルダー変換により実行させるための手段を含むことができる。

30

【 0 0 2 2 】

本発明においては、前記ドキュメント - アトリビュート副行列から前記所定の行列として共分散行列を生成し、ドキュメント・ベクトルの積和行列($SUM(M)_1$)と、ドキュメント・ベクトルの平均($MEAN(M)_1$)と、ドキュメント・ベクトルの平均の積行列($SUM(M)_2$)と、ドキュメントの全数(M)とを保持するための手段とを含むことができる。また、本発明においては、追加された前記副行列を含むすべてのドキュメント - アトリビュート行列の共分散行列C'を、下記式

40

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

(上記式中、Hは、追加された副行列のドキュメント数を示す。)

により生成する手段を含むことができる。本発明においては、さらに、ドキュメント・ベクトルに含まれるアトリビュートを自動的に検索し、アトリビュート・ハッシュ・テーブルを生成して前記アトリビュートを追加または削除するための手段を含むことができる。

【 0 0 2 3 】

本発明によれば、データベースに時間の経過と共に追加されるドキュメントを検索するた

50

めの情報検索方法であって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、
前記ドキュメント - アトリビュート副行列から所定の行列を生成して保持させるステップと、
前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、
前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記所定の行列を更新し、更新された前記所定の行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、
前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと
を含む、情報検索方法が提供される。

10

【 0 0 2 4 】

さらに本発明によれば、データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索方法を実行させるためのプログラムであって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、
前記ドキュメント - アトリビュート副行列から所定の行列を生成して保持させるステップと、
前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、
前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記所定の行列を更新し、更新された前記所定の行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、
前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと
を含む、情報検索方法を実行させるためのプログラムが提供できる。

20

【 0 0 2 5 】

また、本発明によれば、データベースに時間の経過と共に追加されるドキュメントを検索するための情報検索方法を実行させるためのプログラムが記録されたコンピュータ可読な記録媒体であって、前記ドキュメントは、ドキュメント - アトリビュート行列へと変換されて保持され、かつ前記ドキュメント - アトリビュート行列は、逐次的に追加されるドキュメント - アトリビュート副行列から構成され、
前記ドキュメント - アトリビュート副行列から所定の行列を生成して保持させるステップと、
前記データベースに所定の期間の間に追加されたドキュメントからドキュメント - アトリビュート副行列を生成するステップと、
前記ドキュメント - アトリビュート副行列に関連する情報を使用して前記所定の行列を更新し、更新された前記所定の行列を特異値分解して、データベースに保持されたすべてのドキュメント - アトリビュート行列の次元削減を実行するステップと、
前記次元削減されたドキュメント - アトリビュート行列を使用してユーザが入力したクエリーによる情報検索を行うステップと
を含む、情報検索方法を実行させるためのプログラムが記録された記録媒体が提供される。

30

40

【 0 0 2 6 】**【 発明の実施の形態 】**

以下、本発明につき図面に示した実施の形態に基づいて詳細に説明するが、本発明は後述する特定の実施の形態に限定されるものではない。

50

【0027】

図1は、本発明において、情報検索を実行するためのデータベースの概略構成を示した図である。図1に示されたデータベースにおいては、ドキュメントは、バイナリ・モデルやアトリビュート頻度・モデルといった適切な方法を使用して、ドキュメント・ベクトルへと変換されている。図1には、上述したドキュメントにより生成されたドキュメント・アトリビュート行列としてデータベースのデータ構成を示している。図1においては、行方向にドキュメント・ベクトルが並べられており、列方向には、所定のアトリビュートがドキュメントに含まれている場合には、適切な重み、所定のアトリビュートが含まれていない場合には、0とする方法により、数値要素が並べられている。

【0028】

なお、本発明においては、図1で示したバイナリ・モデルの他に、ユーザが指定する重み付け因子（ウエイト・ファクタ）を適用した、アトリビュート頻度モデル与えることもできる。以下、本発明においては、上述して得られた行列をドキュメント・アトリビュート行列Aとして参照する。アトリビュートとしては、テキスト・ドキュメント・データに用いられるキーワード・アトリビュートの他にも、タイム・スタンプ、画像、オーディオ・データなど、いかなるアトリビュートでも本発明においては使用することができる

【0029】

図1に示されるように、本発明においては、データベースには、ドキュメント/データが頻繁に追加されており、これらのドキュメントは、それぞれが含むアトリビュートに基づいて、ドキュメント・ベクトルへと変換される。これらのドキュメント・ベクトルは、例えば、日単位、週単位、月単位、或いは、データベースの管理者が設定した単位に区切られて、それぞれの副行列A₁, . . . , A_nとして構成されている。

【0030】

例えば、最初にデータベースに対して蓄積されていたドキュメント・ベクトルの集合は、行列Aとして区切られており、Aについて処理を実行する。その後追加されたドキュメント・ベクトルの集合を、例えば1週間といった所定の期間まとめてドキュメント・アトリビュート行列として構成したものが、副行列A₁として示されている。同様にして、順次追加されたドキュメント・ベクトルを所定の期間ごとにまとめたもので、最も新しい副行列が、図1においては副行列A_nとして示されている。

【0031】

上述した所定の期間としては、上述したように日単位、週単位、月単位としてまとめることができるが、特に常に一定の期間ではなく、必要に応じてその時点までに蓄積されたデータをまとめて副行列A_iとすることができる。本発明は、上述した副行列を使用して、QR分解で得られる上三角行列Rまたは共分散行列を算出することで、データが蓄積されていくドキュメント・アトリビュート行列すべてを一割してSVDを実行させる労力を削減することにより、次元削減の計算時間を低減して効率化を達成する。また、本発明においては、重要ではあるがデータ数として数%程度でデータベースに含まれるいわゆるアウトライア・ドキュメントを次元削減プロセスにおいて無視してしまう可能性を可能な限り低減させ、検索精度を向上させることを可能とする。以下、本発明の実施の形態について詳細に説明する。

【0032】

<第1の実施の形態>

本発明の第1の実施の形態においては、上述した副行列と前回までに得られた行列Rからなる行列に対してQR分解を適用して、新たに得られる上三角行列Rに、SVDを適用して特異ベクトルを得、得られた大きなものからk番目までの特異ベクトルを含んで構成された特異行列を使用して、データベースに含まれるドキュメント・アトリビュート行列をk次元へと次元削減させるものである。

【0033】

また、副行列A₁は、図2で示されるように、もとのドキュメント・アトリビュート行列AをQR分解した時の行列Rの底に加えた形式で保持する。

10

20

30

40

50

【 0 0 3 4 】

図 2 の左側の行列 R に副行列 A_i が追加された部分は適当な直交置換行列 P により行列 R の直下に移動させる。

【 0 0 3 5 】

図 3 は、本発明の第 1 の実施の形態として QR 分解法を使用してドキュメント - アトリビュート行列 A の次元削減を実行するプロセスを示したフローチャートである。図 3 に示されるように、プロセスは、ステップ S 1 から開始し、ステップ S 2 において、ドキュメント - アトリビュート行列 A の QR 分解を実行する。ステップ S 3 においては、QR 分解して得られた行列 Q と、行列 R のうち、行列 R のみを使用する。本発明において行列 R のみを使用する理由は、(a) R 行列が上三角行列であり、SVD を行うためにきわめて迅速に計算を実行することができること、(b) 下記式に示されるように行列 R を使用してもドキュメント - アトリビュート行列における特異値または固有値は保存されていること、の理由に基づくものである。

10

【 0 0 3 6 】

【 数 2 】

$$A^T A = (QR)^T QR = R^T Q^T Q R = R^T R$$

【 0 0 3 7 】

次いで、本発明の第 1 の実施の形態におけるプロセスにおいては、ステップ S 4 において追加された副行列 A_i を使用して行列 R を更新する。この際、更新された行列 R の算出は、本発明の好適な実施の形態においては、この更新された行列 R の算出については、より詳細に後述する。ステップ S 5 では、上述したようにして得られた R 行列を使用して SVD を実行する。

20

【 0 0 3 8 】

【 数 3 】

$$R = U \Sigma V^T$$

(上記式中、U、V は、N × N 正規直交行列であり、 Σ は、N × N 対角行列である。)
上述のようにして得られた行列 R の特異値または固有値は、上述したようにドキュメント - アトリビュート行列 A の特異値と同じ特異値を保持している。また、R は、上三角行列であるため、SVD をきわめて容易に行うことができ、従来のプロセスにおいてきわめて計算時間を要した SVD に割り当てられる計算時間を著しく低減することが可能となる。

30

【 0 0 3 9 】

さらに本発明のプロセスにおいては、ステップ S 6 で、得られた特異値から特異ベクトルまたは固有ベクトルを得、得られた特異ベクトルまたは固有ベクトルを、特異値または固有値の大きな方から特異ベクトルまたは固有ベクトルの k 番目までを使用して次元削減した k 次元の特異行列を生成してドキュメント - アトリビュート行列の次元を削減させる。

【 0 0 4 0 】

ステップ S 7 において次元の減少した行列を使用して情報検索を実行し、ステップ S 8 において本発明の第 1 の実施例の情報検索方法を終了する。

40

【 0 0 4 1 】

以下各ステップにおける処理を詳細に説明する。まず、ステップ S 2 で、ドキュメント - アトリビュート行列 A について QR 分解を実行する。A についての QR 分解を下記式に示す。

【 0 0 4 2 】

【 数 4 】

$$A = Q \begin{array}{|c|} \hline R \\ \hline O \\ \hline \end{array}$$

【0043】

上記式中、Oで示された行列は、要素がすべて0の行列を意味する。図4には、一般的な行列Dに対する上述したQR分解を実行させるための擬似コードを示す。

本発明においては、ステップS3において、行列Aについて上述のQR分解により得られた行列Rのみを使用する。図5には、行列AのQR分解により得られる行列の構成要素を概略的に示す。図5において、Mは、ドキュメントの数であり、Nは、アトリビュートの数である。また、図5においては、行列Qの列ベクトルを q_i ($i = 1, \dots, M$)により示し、行列Rの要素がゼロの部分をも0で示している。また、図5に示した実施の形態においては、得られる行列Rは、 $N \times N$ の上三角行列として得られている。

10

【0044】

本発明においてはさらに、ステップS4において、副行列 A_i としてデータが加えられたドキュメント-アトリビュート行列につき、ハウスホルダー変換を実行する。

【0045】

本発明においては、ステップS5において行列Rに対して直接SVDを適用して、R行列の特異値または固有値を求める。この際、SVDに使用することができる方法としては、これまで知られた種々の方法を使用することができ、例えばハウスホルダー変換を使用する方法、またはランチョス法を使用することができる。

20

【0046】

さらに本発明においては、ステップS6においてSVD計算により得られた特異値または固有値から特異ベクトルまたは固有ベクトルを生成し、特異値または固有値の大きな方から所定の数の特異ベクトルまたは固有ベクトルを降順に配置して特異行列を形成させ、これをドキュメント-アトリビュート行列Aに乗じて、下記式(5)にしたがってドキュメント-アトリビュート行列Aの次元を削減させる。

【0047】

【数5】

$$A_k = U_k \Sigma_k V_k^T$$

30

図6に示すように上記式中、 A_k は、k番目までの特異値を使用して得られた次元削減されたドキュメント-アトリビュート行列であり、 Σ_k はk個の特異値から成る $k \times k$ 次元の対角行列であり、 U_k はk個の特異値に対する左特異ベクトルから成る行列であり、 V_k はk個の特異値に対する右特異ベクトルから成る行列である。

【0048】

ステップS7においては、上述したようにして次元削減されたドキュメント-アトリビュート行列を使用して、クエリー・ベクトルとの乗算を実行し、ユーザが所望するクエリー・ベクトルに基づいた情報検索を実行させることになる。

40

【0049】

本発明の第1の実施の形態におけるドキュメント-アトリビュート行列Aの次元削減は、上述したようにQR分解により得られた行列Rを更新しつつ、直接SVDに提供し、精度良くSVDを実行することができると共に、次元削減に対する計算時間を著しく低減させることができ、この結果、メモリ資源を節約しつつ、高精度、高効率の情報検索を実行することが可能となる。

【0050】

<第2の実施の形態>

50

本発明の第2の実施の形態においては、共分散行列を使用した次元削減プロセスにおいて、すでに計算された以前のドキュメント・アトリビュート行列Aに対する共分散行列を使用して新たに加えられた副行列A1を含む新たな共分散行列を生成し、得られた副行列A1を反映した新たな共分散行列をSVDプロセスに提供するものである。本発明の第2の実施の形態は、ドキュメントの追加に対応できるばかりではなく、容易にドキュメントの削除についても適用することができる。

【0051】

以下、本発明の第2の実施の形態について詳細に説明する。本発明の第2の実施の形態を詳細に説明する前に、共分散行列を用いた情報検索について概略的な説明を行う。M×Nの要素からなるドキュメント・アトリビュート行列Aに対して、その共分散行列Cは、下記式により与えられる。

【0052】

【数6】

$$C = \frac{1}{M} \sum_{i=1}^M \mathbf{d}_i \mathbf{d}_i^T - \bar{\mathbf{d}}(M) \bar{\mathbf{d}}(M)^T = \frac{1}{M} \text{SUM}(M)_1 - \text{SUM}(M)_2$$

上式中、SUM(M)₁は、ドキュメント・ベクトルの積和行列であり、SUM(M)₂は、ドキュメント・ベクトルの平均の積行列である。d_i、 $\bar{\mathbf{d}}(M)$ 、 $\bar{\mathbf{d}}(M)_i$ は、それぞれドキュメント・アトリビュート行列Aの要素i、jを使用して下記式で定義される。このうち、SUM(M)₁およびSUM(M)₂は、共にN×Nの対称な正方行列である。

【0053】

【数7】

$$\mathbf{d}_i = [d(i,1), \dots, d(i,N)]^T$$

$$\bar{\mathbf{d}}(M) = [\bar{d}(M)_1, \dots, \bar{d}(M)_N]^T$$

$$\bar{d}(M)_i = \frac{1}{M} \sum_{j=1}^M d(j,i)$$

【0054】

上述したように定義された共分散行列は、N×Nの正方行列として得られる。また、共分散行列の特異値または固有値は、例えば特願2000-175848号に記載されているように、ドキュメント・アトリビュート行列の特異値または固有値を保存し、かつ正方行列であるのでドキュメントが著しく多い場合であっても特異値または固有値、ひいては特異ベクトルまたは固有ベクトルを迅速に計算でき、高い効率の情報検索の実行を可能とする。

【0055】

図7は、本発明の第2の実施の形態の情報検索方法のフローチャートを示した図である。なお、本発明の第2の実施の形態においてもドキュメント・アトリビュート行列構成は、図1に示したようにAおよび逐次的に蓄積されて行く副行列A1、A2、...、Anから構成されているものとして説明する。以下に説明する実施の形態においては、まず最初に行列Aから共分散行列を生成して行くものとして説明する。

【0056】

図7に示されるように、本発明の情報検索方法の第2の実施の形態は、ステップS10から開始し、ステップS11において行列Aから上記式にしたがって共分散行列を生成する。

【0057】

ステップS12においては、H個のドキュメント・ベクトルを含む副行列A1が、行列A

10

20

30

40

50

に追加され、ドキュメント - アトリビュート行列が形成される。ステップ S 1 3 においては、副行列 A 1 のドキュメント - アトリビュート行列から、下記式の計算を実行させる。

【 0 0 5 8 】

【 数 8 】

$$\text{SUM}(H)_1 = \sum_{i=1}^H \mathbf{d}_{M+i} \mathbf{d}_{M+i}^T$$

さらにステップ S 1 4 において、下記式の計算を実行させる。

【 0 0 5 9 】

【 数 9 】

$$\text{SUM-MEAN}(H)_1 = \sum_{j=1}^H d(M+j, i) = H \text{MEAN}(H)_1$$

10

次いで、ステップ S 1 5 において、下記式の計算を実行させる。

【 0 0 6 0 】

【 数 1 0 】

$$\text{SUM}(M+H)_1 = \text{SUM}(M)_1 + \text{SUM}(H)_1$$

【 0 0 6 1 】

ここで、下記式の関係が各構成要素について成り立つ。

20

【 0 0 6 2 】

【 数 1 1 】

$$\bar{d}(M+H)_i = \frac{1}{M+H} \sum_{j=1}^{M+H} d(j, i) = \frac{1}{M+H} (\bar{M}d(M)_i + \bar{H}d(H)_i)$$

【 0 0 6 3 】

次いで、下記式を使用して、 $\text{SUM-MEAN}(M+H)_1$ を更新する。この際の計算は、単なる N 回の加算により実行できるので計算時間を短縮することができる。

【 0 0 6 4 】

【 数 1 2 】

$$\text{SUM-MEAN}(M+H)_1 = \text{SUM-MEAN}(M)_1 + \text{SUM-MEAN}(H)_1$$

30

【 0 0 6 5 】

次いで、ステップ S 1 6 において $\text{SUM}(M+H)_2$ を、下記式にしたがって $\text{SUM-MEAN}(M+H)_1$ を使用して更新する。

【 0 0 6 6 】

【 数 1 3 】

$$(M+H)^2 \text{SUM}(M+H)_2 = \text{SUM-MEAN}(M+H)_1 \text{SUM-MEAN}(M+H)_1^T$$

40

この更新の後、アップデートされた共分散行列 C' を、下記式にしたがって得ることができる。

【 0 0 6 7 】

【 数 1 4 】

$$C' = \frac{1}{(M+H)} \text{SUM}(M+H)_1 - \text{SUM}(M+H)_2$$

この後、ステップ S 1 7 において更新された共分散行列 C' を、SVD に提供して特異値または固有値を得、大きな方から k 個の特異値または固有値に対応する特異ベクトルまたは固有ベクトルを選択して次元の減少された特異行列または固有行列を生成する。ステッ

50

プ S 1 8 において、次元の減少された k 本の特異行列または固有ベクトルまたは特異ベクトルを使用してドキュメント - アトリビュート行列 A の次元削減を行い、ステップ 1 9 で情報検索を実行させ、ステップ S 2 0 で、本発明の第 2 の実施の形態のプロセスを終了する。図 8 には、上述した本発明のプロセスのうち、ステップ S 1 1 において説明した共分散行列の生成のための擬似コードを示す。

【 0 0 6 8 】

また、本発明においては、上述したプロセスにおいて $SUM(M - H)_1$ を生成させて同様に計算を繰り返すことにより、ドキュメント・ベクトルが何らかの理由により削減された場合にも容易に対応することが可能となる。

【 0 0 6 9 】

さらに、本発明の第 2 の実施の形態においては、アトリビュート自体のアップデートおよびダウodateを行うことも可能となる。アトリビュートのアップデートは、新たなアトリビュートが加えられた場合に実行され、アトリビュートの削減は、例えばアトリビュートが検索において非現実的なものとなったとき、または検索する必要が無くなったときに実行される。アトリビュートの追加・削除は、アトリビュート・ハッシュ・テーブルを使用して実行される。

【 0 0 7 0 】

図 9 には、本発明の第 2 の実施の形態に適用される、ドキュメントの追加または削除におけるアトリビュート・ハッシュ・テーブルを変更するプロセスを示した図である。図 9 (a) に示すように、ユーザからのドキュメント (i) の追加 / 削除の要求があると、本発明においては、ドキュメント (i) を追加する場合には、まずバイナリ・モデルといった適切な方法を使用してドキュメント・ベクトルを形成する。また、ユーザがドキュメント (i) を削除しようとする場合には、削除しようとするドキュメント (i) を特定する。

【 0 0 7 1 】

次いで、図 9 (b) に示すようにドキュメント (i) に含まれる非ゼロのアトリビュート A T を特定する。図 9 (b) においては、ドキュメント (i) が含む非ゼロのアトリビュート att_3 , att_{n-1} が、数値 1 に対応することが示されているが、本発明において重み付けをそれぞれのアトリビュート A T に対して適用する場合には、非ゼロの要素は、ウエイト・ファクタに対応した 1 以外の値とされていても良い。

【 0 0 7 2 】

上述したようにドキュメント (i) における非ゼロのアトリビュートが特定されると、本発明においては、図 1 0 に示されたアトリビュート・ハッシュ・テーブルを参照する。

【 0 0 7 3 】

アトリビュート・ハッシュ・テーブルには、それぞれのアトリビュートと、当該アトリビュートを含むドキュメント数とが対応して記憶されており、アトリビュートからドキュメント数を参照することが可能とされている。図 1 0 を使用してドキュメントの追加 / 削除の実施の形態について説明すると、例えば、アトリビュート 3 とアトリビュート n - 1 とを含むドキュメント (i) が追加 / 削除される場合には、もともとアトリビュート 3 を含むドキュメント数 6 が、追加の場合にはドキュメント数 7 とされ、削除の場合にはドキュメント数 5 と変更される。

【 0 0 7 4 】

これに対応して、図 1 0 に示される実施の形態の場合には、アトリビュート n - 1 も非ゼロの要素となっているので、もともとアトリビュート n - 1 に関連するドキュメントのドキュメント数 3 3 が、追加 / 削除に対応して 3 4 または 3 2 へと変更される。上述したアトリビュート・ハッシュ・テーブルは、さらに別の識別子を使用することにより所定のアトリビュートを有する個々のドキュメントの参照を可能とするようにされていてもよい。

【 0 0 7 5 】

図 1 1 は、上述したアトリビュートの削除を実行させる場合の概念図を示す。図 1 1 (a) が、アトリビュート削除前のドキュメント・ベクトルの構成を示したものであり、図 1 1 (b) がアトリビュート削除後のドキュメント・ベクトルの構成を示した実施の形態を

10

20

30

40

50

示した図である。図 1 1 に示した実施の形態においては、アトリビュート 4 が削除される実施の形態を示しているが、本発明において削除されるアトリビュートの示された位置、または一度に削除されるアトリビュートの数は、図 1 1 に示される以外にもいかなるものでも用いることができる。

【 0 0 7 6 】

図 1 2 は、図 1 1 に示されたアトリビュート 4 を削除する場合のドキュメント・ベクトルの積和行列の変更を例示した図である。図 1 2 (a) は、アトリビュート 4 の削除前のドキュメント・ベクトルの積和により得られる積和行列を示し、図 1 2 (b) は、アトリビュート 4 の削除後のドキュメント・ベクトルから得られる積和行列を示す。上述したように、本発明においては、共分散行列を生成する際に積和行列を用いるので、アトリビュートの追加、削除も容易に含めることができる。

10

【 0 0 7 7 】

図 1 3 は、アトリビュートを加える際の処理を説明した概略図である。図 1 3 に示した実施の形態においては、図 1 3 (a) にアトリビュート追加前のドキュメント・ベクトルが示されており、このドキュメント・ベクトルに対して、図 1 3 (b) に示されるように、アトリビュート $n + 1$ がドキュメント・ベクトルに対して加えられる。図 1 4 には、上述したように加えられたアトリビュート $n + 1$ を含んだドキュメント・ベクトルから積和により形成された $(N + 1) \times (N + 1)$ 行列を示す。

【 0 0 7 8 】

上述した従来の共分散行列を作る場合のドキュメント・データは、ドキュメント数とアトリビュート数とがあらかじめわかっていることが前提とされていた。これに対して本発明は、文書数の総数もアトリビュート数も、最初はわかってない状態から計算を逐次実行させることを可能とする。

20

【 0 0 7 9 】

図 1 5 には本発明の情報検索方法を実施するためのコンピュータ・システムの概略図を示す。図 1 5 に示した本発明のコンピュータ・システムは、コンピュータ 1 0 と、このコンピュータ 1 0 とデータの伝送を行うことが可能であるように接続されたデータベース 1 2 とを含んで構成されている。本発明において使用することができるコンピュータ 1 0 としては、本発明の方法を実行することができる中央処理装置 (C P U)、 R A M などのメモリなどのハードウェア資源を含んで構成されている限り、パーソナル・コンピュータ、ワークステーションといったいかなるものでも用いることができる。また、本発明において使用することができるデータベース 1 2 としては、データが追加して書込みができるものであれば、従来知られているいかなるデータベースでも使用することができる。

30

【 0 0 8 0 】

また、図 1 5 に示したコンピュータ・システムは、例えばインターネット、ローカル・エリア・ネットワーク (L A N)、ワイド・エリア・ネットワーク (W A N) といった、これまで知られていかなる T C P / I P といったプロトコルを使用して遠隔的に配置されたネットワーク 1 4 を介して接続されたコンピュータ 1 6 と通信を行うことができるように構成することができる。図 1 5 に示した本発明の実施の形態においては、データベース 1 2 に接続されたコンピュータ 1 0 をサーバとして使用し、このサーバに対して遠隔的に接続されたコンピュータ 1 6 をクライアント・コンピュータとして使用する、いわゆるクライアント・サーバ・システムとして構成することができる。クライアント・コンピュータのユーザは、所望する情報を検索するべく、コンピュータ 1 6 へと例えばキーワードといったアトリビュートの入力を行う。

40

【 0 0 8 1 】

入力されたキーワードは、ネットワーク 1 4 を介してコンピュータ 1 0 へと伝送され、コンピュータ 1 0 において情報検索を行うために使用される。データベース 1 2 には、オリジナルのドキュメント D が保持されている。また、データベース 1 2 には、このドキュメント D から所定のアトリビュートを抽出し、例えばバイナリ・モデルを使用して予め数値化されたドキュメント・アトリビュート行列も保持されている。本発明においては、上述

50

したドキュメント・アトリビュート行列は、一度すでに本発明によるQR分解法または共分散行列法などの方法を使用してすでに次元削減が行われて、実際の情報検索に提供されている。

【0082】

図15に示されたデータベース12には、随時にドキュメントDNが追加されていて、副行列として蓄積されて行く。このドキュメントDNは、本発明においては、例えば日単位、週単位、月単位、あるいはサーバ管理者の規定する間隔で、前回SVDを実行した時点からその時点までに蓄積された分ごとに上述したQR分解法または共分散行列法により、以前に得られているRまたは $SUM(M)_1$ 、 $SUM(M)_2$ を使用することで、効率よく以前の結果を含めて次元削減が実行される。本発明においては、上述したようにしてアップデートまたはダウンデートされたドキュメントに対しての情報検索を行うために使用されるドキュメント・アトリビュート行列の効率的な次元削減を可能とし、高精度、かつ高効率の情報検索システムが提供されている。

10

【0083】

また図15においては、情報検索システムを特にクライアント・サーバ・システムとして説明したが、本発明においては、特にネットワークを介して接続されたクライアント・サーバ・システムばかりではなく、スタンド・アローンのパーソナル・コンピュータまたはワークステーションを使用した情報検索システムとすることができる。

【0084】

以下、本発明をより具体的に実施例をもって説明するが、本発明は後述する実施例によって限定されるものではない。

20

【0085】

【実施例】

(実施例1)

本発明のQR分解法を用いた情報検索について、表1に一部を示したサンプル・データベースを使用して、本発明の効果について検討を加えた。サンプル・データベースは、ドキュメント数300、キーワード数50の300×50のドキュメント・キーワード行列として構成した。

【0086】

表1に示したサンプル・データベースは、現実に使用されるデータベースよりも遙かに小さいので、計算速度においては、従来例で知られた方法に比較して大きな違いは見られなかった。このため、本発明の方法により得られた結果と、後述する比較例により得られた結果との計算精度の比較を行った。

30

【0087】

【表1】

300 50

1 1.65244e-05 7 0.19628 8 3.61624e-08 9 0.00018 10 0.01373 12 2.23141e-05
14 4.64130e-05 15 1.78852e-05 19 2.19083e-10 20 0.00028 23 1.21350e-07 24
9.98935e-10 25 7.85754e-07 28 1.00379e-05 29 2.20989e-05 39 8.26993e-08 40
3.46313e-09 42 3.42076e-10 43 1.40159e-05 45 1.12157e-06 50 5.29161e-07
5 1.95247e-05 6 9.65067e-09 11 4.55447e-07 12 0.00840 14 4.10627e-09 16
2.13934e-10 17 0.00022 22 0.04135 27 3.80890e-09 31 4.73880e-10 32
9.24676e-05 34 6.28455e-06 40 2.17201e-05 41 2.23796e-05 42 1.28397e-06 46
1.50510e-08 47 1.45181e-06 50 9.61504e-09

2 2.20380e-05 3 9.09460e-06 7 0.01776 8 0.02260 10 4.63745e-06 12 0.00147
13 0.00493 14 1.09471e-06 15 0.00366 17 1.82302e-05 18 2.04803e-05 20
2.29817e-05 22 3.91119e-09 25 3.34495e-06 26 0.02291 34 0.00014 37
7.98834e-09 39 1.33136e-07 40 5.08289e-08 41 9.11475e-09 42 3.64340e-05 43
5.91690e-10 44 8.69801e-09 45 0.01932 46 1.31295e-06 47 4.29275e-08 49
1.30851e-10

4 1.62408e-07 7 4.75727e-06 8 2.24351e-05 9 4.66984e-06 10 3.25484e-06 12
2.15089e-05 14 5.29682e-06 17 1.63133e-10 26 0.00766 28 7.18088e-09 31
1.53839e-09 34 8.06182e-10 39 6.13628e-08 41 4.69203e-06 45 2.34841e-10 46
8.81336e-07 47 1.23390e-05 48 3.38329e-09 50 3.24139e-08

2 4.46380e-10 5 5.26576e-09 7 2.09751e-07 12 7.67514e-06 13 2.75356e-05 15
2.15944e-09 18 3.78191e-07 19 8.21337e-09 21 0.18159 27 1.10242e-08 28
2.09905e-05 32 1.38427e-07 34 1.04496e-06 37 2.14666e-08 38 0.00268 40
0.00482 45 6.75967e-09 46 9.43715e-07 47 1.75043e-08

10

20

30

.
.
.

【 0 0 8 8 】

得られた計算精度を実施例 1、比較例 1、比較例 2 の計算結果と共に以下に示す。なお、
特異値計算の精度を比較する簡単な方法として、最大の特異値 λ_1 と最小の特異値 λ_n の
比（以下、この値をコンディション・ナンバーとして参照し c とで示す）を採用した。な
お、比較例 1 は、ドキュメント・アトリビュート行列 A に更新された副行列 A_i を加えた
ものを新たに A として、 A 全体に対して SVD 法を適用して計算を行うもの（以降、ナイ
ーブな SVD 法と呼ぶ）であり、比較例 2 は、 $A^T A$ の逐次アップデート法により SVD
を更新する方法である。また、データベースの更新をシミュレートするため、300 のド
キュメント・キーワード行列を 50×50 チャンクに区切り、チャンク毎に所定の方法を
用いて SVD を実行させた。

40

【 0 0 8 9 】

実施例および比較例とも、Pentium（登録商標）III、クロック周波数 733
MHz（Intel 社、登録商標）を使用した Windows（登録商標）2000（

50

M i c r o s o f t 社、登録商標)をOSとして使用するパーソナル・コンピュータを使用した。また、データは、64ビット浮動小数点精度を使用して計算を実行した。

【0090】

上述した条件の下で得られた結果を、コンディション・ナンバーの比として示す。

【0091】

c₁ = 1 0 7 3 2 . 7 1 4 5 7 0 2 2 3 1 8 3 (実施例 1)
c₂ = 1 0 7 3 2 . 7 1 4 5 7 0 2 2 3 2 1 7 (比較例 1)
c₃ = 1 0 7 3 2 . 7 1 4 5 7 0 4 1 1 2 3 4 (比較例 2)
という結果が得られた。

【0092】

したがって、本発明によるQR分解法を使用してSVDを実行させる場合には、ナイーブなSVD法を使用して得られる結果c₂と本発明による方法c₁との比が、c₂/c₁ = 1 . 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 3 となり、小数点以下9桁まで一致していることが示された。一方で、c₂/c₃の比は、0 . 9 9 9 9 9 9 9 9 9 8 2 4 8 1 9 となって、精度が本発明の方法よりも劣ることが示された。最小の特異値 σ_n についてみれば、実施例1の方法と、比較例1の方法とは、小数点以下15桁まで一致していることが見出された。実施例1の結果から、本発明の方法は、ほぼマシンエプシロンまで精度的に問題ないことが示された。

【0093】

(実施例 2)
実施例2においては、本発明のQR分解法を使用する実施の形態において算出するデータのスケラビリティをあげて、CPUの全占有時間、反復計算時間からみた計算速度およびメモリ資源の必要量の点について検討を加えた。使用したデータは、

- (a) 100x100、
- (b) 1000x100、
- (c) 10000x100、
- (d) 100000x100、
- (e) 1000000x100

の5つの行列を使用した。行列の要素は、ランダムに生成した密行列として作成した。計算にあたって使用したチャンクのサイズは、いずれも100×100とした。したがって、たとえば(e)で示される場合には、1万回反復してアップデートを行いながら計算を実行させることになる。また、本発明の方法を実行するプログラムは、言語Java(登録商標)で実装した。下記表に示すメモリ使用量のうち、概ね4MBがJava(登録商標)(Sun Microsystems社の登録商標)のVirtual Machineが占める。したがって、データが使用する実質のメモリサイズは、下記表4におけるメモリ使用量から4MBを引いた値と見積もることができる。

【0094】

【表2】

| GPU total | (a) | (b) | (c) | (d) | (e) |
|-----------|---------|------------|-------------|---------------|---------------|
| 比較例1 | 454 ミリ秒 | 3 秒 922 ミリ | 92 秒 843 ミリ | 1040 秒 391 ミリ | N/A |
| 比較例2 | 453 ミリ秒 | 1 秒 203 ミリ | 8 秒 781 ミリ | 106 秒 969 ミリ | 764 秒 203 ミリ |
| 実施例2 | 578 ミリ秒 | 1 秒 922 ミリ | 12 秒 297 ミリ | 126 秒 983 ミリ | 1291 秒 532 ミリ |

【0095】

【表3】

| CPU/反復 | (a) | (b) | (c) | (d) | (e) |
|--------|---------|------------|-------------|---------------|----------|
| 比較例1 | 454 ミリ秒 | 3 秒 922 ミリ | 92 秒 843 ミリ | 1040 秒 391 ミリ | N/A |
| 比較例2 | 453 ミリ秒 | 120.3 ミリ | 87.8 ミリ | 106.9 ミリ | 76.4 ミリ |
| 実施例2 | 578 ミリ秒 | 192.2 ミリ | 122.9 ミリ | 126.9 ミリ | 129.1 ミリ |

10

20

30

40

50

【 0 0 9 6 】

【 表 4 】

| メモリ使用量 | (a) | (b) | (c) | (d) | (e) |
|--------|-------|-------|--------|---------|--------|
| 比較例1 | 5.9MB | 8.4MB | 31.0MB | 251.6MB | 割り当て不能 |
| 比較例2 | 6.5MB | 6.5MB | 6.5MB | 6.5MB | 6.5MB |
| 実施例2 | 6.5MB | 6.5MB | 6.5MB | 6.5MB | 6.5MB |

【 0 0 9 7 】

計算時間の点でいえば、比較例2で示される $A^T A$ を反復して特異値を求めるのが最速の結果を与えた。このことは理論的にも明らかで、比較例2では $A^T A$ の計算における行列 - 行列積の計算と行列 - 行列和の計算だけで済むためである。

10

【 0 0 9 8 】

一方、各反復計算の時間についてみれば、実施例2で示される本発明の方法は、QR分解の時間が必要とされる分だけ、反復時間を要することになることが示された。より正確に言えば、実施例2で使用した密行列に対してチャンク・サイズを H としたとき、比較例2では、乗算に $O(HN^2)$ の計算量が必要となる。一方、実施例2で用いるQR分解法は、一回につき、 $O(2N^2((H+N)-N/3))$ の計算量を必要とする(Golub & Van Loan)。

【 0 0 9 9 】

ただし、実施例1においても説明したように、比較例2の $A^T A$ の反復繰り返し計算は、速度を優先し、誤差を犠牲にするものである。なお、データが 100×100 の場合(上述した(a)の行列の場合)は、反復がないので計算手法そのものに要する時間の差を示しているといえる。一回の反復におけるCPU時間は、比較例2で概ね平均100ミリ秒程度、実施例2では、130ミリ秒程度となる。しかしながら、比較例1のナイーブなSVDと本発明の方法を比較すると、遙かに計算時間が早くなっていることが示される。

20

【 0 1 0 0 】

表4に示すメモリ使用量に関してみれば、比較例2も実施例2もほぼ同程度にナイーブなSVDにより得られた比較例1において使用されるメモリよりも著しく小さなメモリですむことが示されている。したがって、本発明の方法は、メモリ資源の消費といった点でも従来の方法と劣るものではないことが示される。

【 0 1 0 1 】

(実施例3)

<アップデート例>

実施例3として、本発明において共分散行列を使用してドキュメントのアップデートを行う場合について、検討を加えた。実施例3においては、図16に示すような7個のデータ・ファイル(その内容は、キーワードとその出現頻度)が、時間順に与えられたものとした。図16に示されるように、データ・ファイルは、タイムス・タンブとデータ・ファイル名とを記したデータ例を示す。この実施例では、“data set”という名前のファイルとして参照する。例えば、図16中、20010701は、2001年7月1日を意味するタイム・スタンプである。

30

【 0 1 0 2 】

図16に現れる個々のファイルの内容を、図17から図23に示す。個々のファイルの内容は、アトリビュートとして使用されるキーワードと、その重み付け(重み付けは、正の実数値で、大きいほどそのキーワードがその文書に含まれる寄与率が高い)とからなるペアとして構成している。また、図17から図23までの文書データ・ファイルは、説明する実施例では20の文書数として構成した。

40

【 0 1 0 3 】

以上のデータが与えられた場合、まず2001年7月1日のタイム・スタンプを含むdata1から順に読み込んでいき、図23のdata7まで本発明の共分散行列法を順次適用する。最後に、data7まで処理したところで共分散行列を使用してSVDを実行して、特異値および特異ベクトルを得た。比較のため、はじめから表5~表10に示す14

50

0 文書のデータ・ファイルを与え、これから共分散行列を生成してSVDを実行して得られた場合と特異値および特異ベクトルの結果を比較した(ベクトルの場合、順序を除いた)。

【0104】

なお、表5の第1行目の140 40は、文書数のトータルが140個で、40個のキーワードからなるものであることを示している。図24には、本発明の実施例3において使用したキーワードを示す。

【0105】

【表5】

| | |
|------------------------------------|----|
| 1行目:140 40 | 10 |
| 2行目:1 1.0 2 1.0 4 0.4 | |
| 3行目:1 1.0 2 1.0 5 0.3 | |
| 4行目:1 1.0 2 1.0 6 0.1 | |
| 5行目:1 1.0 2 1.0 4 0.3 | |
| 6行目:1 1.0 2 1.0 6 0.2 | |
| 7行目:1 1.0 2 1.0 4 0.2 | |
| 8行目:1 1.0 2 1.0 6 0.3 | 20 |
| 9行目:1 1.0 2 1.0 4 0.1 | |
| 10行目:1 1.0 2 1.0 6 0.4 | |
| 11行目:1 1.0 2 1.0 | |
| 12行目:1 1.0 2 1.0 3 0.5 5 0.3 | |
| 13行目:1 1.0 2 1.0 3 0.5 4 0.4 | |
| 14行目:1 1.0 2 1.0 3 0.5 6 0.1 | |
| 15行目:1 1.0 2 1.0 3 0.5 | 30 |
| 16行目:1 1.0 2 1.0 3 0.5 5 0.3 | |
| 17行目:1 1.0 2 0.5 3 1.0 4 0.1 | |
| 18行目:1 1.0 2 0.5 3 1.0 4 0.4 | |
| 19行目:1 1.0 2 0.5 3 1.0 6 0.1 | |
| 20行目:1 1.0 2 0.5 3 1.0 | |
| 21行目:1 1.0 2 0.5 3 1.0 5 0.3 6 0.2 | 40 |
| 22行目:1 1.0 3 1.0 6 0.1 | |
| 23行目:1 1.0 3 1.0 4 0.2 | |
| 24行目:1 1.0 3 1.0 | |
| 25行目:1 1.0 3 1.0 5 0.3 | |
| 26行目:1 1.0 3 1.0 4 0.4 | |
| 27行目:7 1.0 9 1.0 10 0.1 11 0.4 | |
| 【0106】 | 50 |

【表 6】

(続き)

28 行目:7 1.0 9 1.0 11 0.3

29 行目:7 1.0 9 1.0 10 0.2

30 行目:7 1.0 9 1.0

31 行目:7 1.0 9 1.0 10 0.1

32 行目:7 1.0 9 1.0 11 0.2

10

33 行目:7 1.0 9 1.0 10 0.3

34 行目:7 1.0 9 1.0 11 0.4

35 行目:7 1.0 9 1.0 11 0.1

36 行目:7 1.0 9 1.0 10 0.1

37 行目:7 1.0 8 0.5 9 1.0 11 0.2

38 行目:7 1.0 8 0.5 9 1.0 10 0.3

39 行目:7 1.0 8 0.5 9 1.0 10 0.2

20

40 行目:7 1.0 8 0.5 9 1.0 11 0.1

41 行目:7 1.0 8 0.5 9 1.0

42 行目:7 1.0 8 1.0 9 0.5 10 0.1

43 行目:7 1.0 8 1.0 9 0.5 11 0.3

44 行目:7 1.0 8 1.0 9 0.5 11 0.2

45 行目:7 1.0 8 1.0 9 0.5 11 0.1

30

46 行目:7 1.0 8 1.0 9 0.5

47 行目:7 1.0 8 1.0 10 0.1

48 行目:7 1.0 8 1.0

49 行目:7 1.0 8 1.0 11 0.3

50 行目:7 1.0 8 1.0 10 0.2

51 行目:7 1.0 8 1.0 11 0.1

52 行目:12 1.0 18 0.1

40

53 行目:12 1.0 21 0.1

【 0 1 0 7 】

【表 7】

(続き)

| | |
|--|----|
| 54 行目 : 12 1.0 24 0.1 | |
| 55 行目 : 12 1.0 35 0.1 | |
| 56 行目 : 12 1.0 | |
| 57 行目 : 13 1.0 19 0.1 | |
| 58 行目 : 13 1.0 22 0.1 | |
| 59 行目 : 13 1.0 28 0.1 | 10 |
| 60 行目 : 13 1.0 33 0.1 | |
| 61 行目 : 13 1.0 | |
| 62 行目 : 14 1.0 16 0.1 | |
| 63 行目 : 14 1.0 23 0.1 | |
| 64 行目 : 14 1.0 29 0.1 | |
| 65 行目 : 14 1.0 37 0.1 | |
| 66 行目 : 14 1.0 | 20 |
| 67 行目 : 15 1.0 17 0.1 | |
| 68 行目 : 15 1.0 25 0.1 | |
| 69 行目 : 15 1.0 30 0.1 | |
| 70 行目 : 15 1.0 38 0.1 | |
| 71 行目 : 15 1.0 | |
| 72 行目 : 21 0.3 29 0.4 33 0.2 37 1.0 | |
| 73 行目 : 19 0.2 28 0.3 29 0.2 33 0.3 36 0.2 38 1.0 | 30 |
| 74 行目 : 20 0.2 32 0.2 35 0.2 39 1.0 | |
| 75 行目 : 16 0.3 22 0.4 23 0.2 26 0.1 28 0.3 40 1.0 | |
| 76 行目 : 16 1.0 18 0.2 40 0.3 | |
| 77 行目 : 17 1.0 26 0.2 27 0.4 40 0.1 | |
| 78 行目 : 16 0.2 18 1.0 22 0.3 23 0.4 24 0.2 27 0.2 29 0.4 34 0.3 36 0.2 | |
| 79 行目 : 19 1.0 22 0.1 25 0.3 29 0.4 33 0.3 34 0.4 35 0.4 | 40 |

【 0 1 0 8 】

【 表 8 】

(続き)

80 行目:20 1.0 22 0.1 24 0.3 25 0.3

81 行目:21 1.0 29 0.3 33 0.2 35 0.2 36 0.2

82 行目:22 1.0 26 0.4 34 0.4 35 0.2

83 行目:23 1.0 30 0.3 37 0.3

84 行目:20 0.3 24 1.0

85 行目:22 0.3 25 1.0 30 0.2 38 0.4

10

86 行目:22 0.1 26 1.0 29 0.2 31 0.2 34 0.3 36 0.3

87 行目:18 0.2 25 0.3 26 0.1 27 1.0 34 0.3 35 0.3

88 行目:27 0.2 28 1.0

89 行目:17 0.3 23 0.4 29 1.0 37 0.4 40 0.4

90 行目:18 0.4 24 0.3 27 0.3 29 0.3 30 1.0 37 0.2

91 行目:16 0.3 17 0.2 25 0.3 27 0.3 29 0.2 31 1.0 33 0.3 35 0.3

92 行目:20 0.2 25 0.1 30 0.3 32 1.0 36 0.4

20

93 行目:19 0.2 20 0.3 26 0.4 28 0.2 33 1.0 35 0.4

94 行目:17 0.3 34 1.0

95 行目:21 0.4 30 0.3 34 0.3 35 1.0 38 0.4

96 行目:17 0.4 20 0.4 21 0.3 23 0.4 27 0.4 36 1.0

97 行目:18 0.3 22 0.2 29 0.2 37 1.0 39 0.2

98 行目:18 0.2 23 0.3 24 0.3 29 0.2 32 0.3 36 0.2 38 1.0

99 行目:18 0.2 28 0.3 39 1.0

30

100 行目:40 1.0

101 行目:16 1.0 18 0.3 20 0.2 26 0.3 30 0.3 37 0.3 39 0.3

102 行目:17 1.0 20 0.4 24 0.2 26 0.2 27 0.2 30 0.4 36 0.4 38 0.4

103 行目:18 1.0 19 0.3 30 0.2

104 行目:19 1.0 25 0.2 31 0.3 33 0.3 36 0.4

105 行目:18 0.2 20 1.0 22 0.4 23 0.3 33 0.3 37 0.4 38 0.4 40 0.2

40

【 0 1 0 9 】

【 表 9 】

(続き)

106 行目:19 0.2 21 1.0 22 0.2 25 0.1 37 0.2

107 行目:22 1.0 28 0.2 32 0.4 37 0.4

108 行目:21 0.3 23 1.0 34 0.3 40 0.2

109 行目:24 1.0

110 行目:19 0.3 24 0.1 25 1.0 27 0.4 30 0.3

111 行目:20 0.4 26 1.0 28 0.4 30 0.4 35 0.2

10

112 行目:24 0.3 25 0.1 27 1.0 29 0.3

113 行目:26 0.3 28 1.0 30 0.3 35 0.4 40 0.4

114 行目:16 0.4 20 0.4 24 0.3 29 1.0 32 0.2 33 0.1 35 0.4

115 行目:23 0.4 30 1.0 35 0.3

116 行目:31 1.0 34 0.4

117 行目:26 0.3 32 1.0 39 0.4

118 行目:33 1.0 35 0.2 36 0.2 37 0.2 38 0.2

20

119 行目:17 0.2 22 0.2 27 0.2 31 0.2 34 1.0

120 行目:18 0.2 35 1.0

121 行目:28 0.1 33 0.3 35 0.2 36 1.0 38 0.3

122 行目:37 1.0

123 行目:18 0.4 19 0.2 20 0.3 22 0.3 23 0.3 38 1.0 39 0.3

124 行目:20 0.3 29 0.2 34 0.2 39 1.0

125 行目:16 0.3 17 0.3 23 0.3 25 0.4 33 0.4 34 0.3 36 0.2 37 0.4 39 0.2 40 1.0

30

126 行目:16 1.0 17 0.4 19 0.2 24 0.3 27 0.2 31 0.1 32 0.1 36 0.3 39 0.1

127 行目:17 1.0 23 0.2 26 0.4 28 0.2 39 0.3

128 行目:18 1.0 27 0.2 33 0.3 35 0.2

129 行目:19 1.0 24 0.2 26 0.2 30 0.1 33 0.2

130 行目:20 1.0 21 0.2 31 0.2 33 0.2

131 行目:21 1.0 36 0.2 37 0.3

40

【 0 1 1 0 】

【 表 1 0 】

(続き)

132 行目 : 16 0.4 22 1.0 30 0.3 32 0.3 35 0.2 37 0.4

133 行目 : 23 1.0 28 0.3 30 0.3 33 0.3 40 0.2

134 行目 : 24 1.0 33 0.4 35 0.2 38 0.4

135 行目 : 23 0.3 25 1.0 35 0.2 38 0.4 39 0.2 40 0.4

136 行目 : 17 0.2 18 0.3 22 0.3 24 0.2 26 1.0 38 0.4

137 行目 : 16 0.3 18 0.1 21 0.2 25 0.3 27 1.0 36 0.2

10

138 行目 : 23 0.2 26 0.2 28 1.0

139 行目 : 29 1.0 31 0.4 38 0.3 40 0.4

140 行目 : 18 0.1 26 0.2 28 0.3 29 0.3 30 1.0 33 0.2 35 0.1

141 行目 : 24 0.1 26 0.2 31 1.0

【 0 1 1 1 】

実施例 3 では、まず図 1 7 に示した `data 1` が入力されるものとした。図 2 5 には、`data 1` だけを処理した時点での $SUM - MEAN(M)_1$ と、 $SUM(M)_1$ とを示す。なお、最初の 6 行は、新しいキーワードが見つかった文書番号とそのキーワード名とを示す。図 2 5 に示されるように、`data 1` を処理した時点では、キーワード総数は 6 で、したがって $SUM - MEAN(M)_1$ は、6 次元ベクトルであり、 $SUM(M)_1$ は、 6×6 次元の対称行列となる。

20

【 0 1 1 2 】

図 2 5 においては、対称性から要素の半分だけを書き出して示している。また、キーワードは見つかるたびに、キーワードを管理するキーワード・ハッシュ・テーブルに追加される。キーワード・ハッシュ・テーブルにはこの他に、何個の文書がそのキーワードを含んでいるかのカウント数も保持されている。実施例 3 においては、キーワード・ハッシュ・テーブル以外に、総文書数 $M = 20$ と、総キーワード数 $N = 6$ とを保持させた。その後、`data 2` をアップデートすることにより、ドキュメントの追加を実行した。その結果を、図 2 6 に示す。なお、図 2 6 においてもデータは、図 2 5 と同様の順および構成として示されている。`data 3` から `data 6` までを同様に処理した。最後に `data 7` をアップデート処理した。この処理後の $SUM - MEAN(M)_1$ を図 2 7 に示す。また、 $SUM(M)_1$ の内容を表 1 1 から表 1 2 に示す。なお、`data 3` から `data 6` までのアップデート処理で、キーワードが合計 40 個出揃っているため、`data 7` の処理においては、新たなキーワードの追加はされなかった。なお、表 1 1 および表 1 2 中「*」で示された要素は、0.0 を意味する。また、コロンの左側の数字は、行列の行番号を示す。

30

【 0 1 1 3 】

40

【 表 1 1 】

1: 25 17.5 2.5 1.5 1.5 12.5 * * * * *

2: 16.25 1.65 1.05 1.25 5.0 * * * * *

3: 0.83 * * 1.30 * * * * *

4: 0.45 0.06 0.90 * * * * *

5: 0.37 0.45 * * * * *

6: 11.25 * * * * *

7: 25.00 17.50 1.70 2.70 12.50 * * * * *

8: 16.25 1.35 2.00 5.00 * * * * *

9: 0.35 0.04 0.65 * * * * *

10: 0.75 1.15 * * * * *

11: 11.25 * * * * *

12: 5.00 0.10 0.10 0.10 0.10 * * * * *

13: 3.86 0.02 0.44 0.47 * 0.38 0.65 0.09 0.38 * 0.73 0.64 0.65 0.55 * 0.06 0.09 0.79 0.80 0.26
0.38 0.06 0.47 0.43 0.10 0.82 0.36 *

14: 3.52 * 0.60 * 0.20 0.20 * 0.30 * 0.06 0.42 0.42 0.80 * 0.12 0.16 0.12 0.16 0.74 0.32 * * *
0.06 0.32 0.21 0.04

15: 3.73 0.32 * 0.29 0.15 * 0.47 * 0.46 0.17 0.62 0.06 * 0.36 0.22 0.43 0.86 0.27 0.80 0.18
0.03 0.30 * 0.57 0.06 0.13

16: 3.33 * 0.48 0.44 0.61 1.11 * 0.33 0.18 0.71 0.12 * 0.06 0.50 0.96 0.66 0.28 0.40 0.18 0.24
0.61 0.24 0.43 0.63 0.30

17: 5.00 0.10 0.10 0.10 0.10 * * * * *

18: 3.39 0.2 0.10 1.06 * 0.2 0.06 0.44 0.04 * 0.08 0.82 0.25 0.4 0.5 0.12 0.02 0.08 0.28 * 0.16
0.4 0.32

19: 3.84 0.32 0.15 * 0.58 0.41 0.22 1.2 * 0.1 0.38 0.36 0.7 0.09 0.59 0.7 0.13 0.84 0.48 0.1
0.76 0.06

20: 3.75 0.47 * 0.09 0.60 0.15 0.08 * 0.20 * 0.85 0.33 0.16 0.22 0.08 0.36 1.15 0.76 0.20 * *
21: 3.26 * 0.25 0.51 0.54 0.68 * 0.18 0.40 0.31 0.87 0.80 0.84 0.02 0.08 0.48 0.52 0.15 0.24

10

20

30

【 0 1 1 4 】

【 表 1 2 】

40

0.43

22: 5.00 0.10 0.10 0.10 0.10 * * * * * * * * * *

23: 3.73 0.23 0.54 0.58 * 0.55 0.30 0.42 * 0.46 0.36 0.30 0.46 0.33 0.90 0.63 0.15 0.40

24: 4.22 0.62 0.70 * 0.57 0.42 1.00 0.84 0.60 0.55 0.09 0.27 0.14 1.24 0.24 0.51 *

25: 4.09 1.06 * 0.34 0.21 0.60 0.70 0.28 0.46 0.26 0.24 0.26 0.80 0.53 0.38 0.64

26: 4.20 * 0.24 0.18 0.50 0.20 0.18 0.46 0.28 0.37 0.09 0.64 0.06 0.12 *

27: 5.00 0.10 0.10 0.10 0.10 * * * * * * * * * *

10

28: 3.72 0.18 0.40 0.48 0.98 0.56 0.04 0.40 1.00 0.52 0.94 0.59 0.28

29: 3.69 0.53 0.80 0.26 0.32 0.10 0.28 0.03 0.80 1.19 0.33 0.36

30: 4.14 0.36 0.28 0.44 0.39 0.09 0.88 0.18 0.50 0.09 *

31: 4.35 0.90 0.86 0.30 0.38 0.48 0.36 0.08 0.12 0.12

32: 2.98 0.64 0.49 0.07 0.38 0.20 0.78 0.21 0.21

33: 4.12 0.32 0.65 0.66 0.20 0.24 0.06 0.20

34: 2.43 0.61 0.30 * 0.02 * 0.01

20

35: 3.56 0.33 0.28 0.02 0.26 0.01

36: 4.01 0.24 0.22 0.49 0.40

37: 3.86 0.04 0.36 0.16

38: 3.90 0.56 0.36

39: 3.06 0.66

40: 3.38

30

【 0 1 1 5 】

本発明においては、上述したように共分散行列を使用して data 1 から data 7 までのデータのアップデートを終了する。また、本発明において、この時点で共分散行列を作って欲しいとのユーザ・リクエストがある場合には、この時点で SUM (M)₂ を下記式で計算する。これは 4 0 × 4 0 の対称行列となる。

【 0 1 1 6 】

【 数 1 5 】

$$SUM(M)_2 = \bar{d}(M) \bar{d}(M)^T$$

40

その後、SUM (M)₁ と、SUM (M)₂ とから共分散行列を計算させる。得られた共分散行列の結果を表 1 3 から表 1 9 に示す。得られる共分散行列は対称行列なので、表 1 3 ~ 表 1 9 では、要素の半分だけを示している。

【 0 1 1 7 】

【 表 1 3 】

共分散行列 C

〔1行目〕

1.5e-01 1.0e-01 1.5e-02 8.8e-03 8.8e-03 7.3e-02 -3.2e-02 -2.2e-02 -2.2e-03 -
 3.4e-03 -1.6e-02 -6.4e-03 -7.9e-03 -6.1e-03 -7.5e-03 -8.5e-03 -6.4e-03 -6.0e-03
 -7.7e-03 -7.3e-03 -8.2e-03 -6.4e-03 -6.8e-03 -8.7e-03 -8.5e-03 -8.4e-03 -6.4e-03
 -6.9e-03 -7.0e-03 -8.7e-03 -8.5e-03 -6.9e-03 -8.2e-03 -4.5e-03 -6.4e-03 -8.3e-03
 -7.1e-03 -7.7e-03 -6.6e-03 -5.6e-03

10

〔2行目〕

1.0e-01 9.6e-03 6.2e-03 7.6e-03 2.5e-02 -2.2e-02 -1.6e-02 -1.5e-03 -2.4e-03 -
 1.1e-02 -4.5e-03 -5.5e-03 -4.3e-03 -5.3e-03 -6.0e-03 -4.5e-03 -4.2e-03 -5.4e-03
 -5.1e-03 -5.7e-03 -4.5e-03 -4.7e-03 -6.1e-03 -6.0e-03 -5.9e-03 -4.5e-03 -4.8e-03
 -4.9e-03 -6.1e-03 -6.0e-03 -4.8e-03 -5.7e-03 -3.1e-03 -4.5e-03 -5.8e-03 -5.0e-03
 -5.4e-03 -4.6e-03 -3.9e-03

〔3行目〕

5.6e-03 -1.9e-04 -1.9e-04 7.7e-03 -3.2e-03 -2.2e-03 -2.2e-04 -3.4e-04 -1.6e-03
 -6.4e-04 -7.9e-04 -6.1e-04 -7.5e-04 -8.5e-04 -6.4e-04 -6.0e-04 -7.7e-04 -7.3e-04
 -8.2e-04 -6.4e-04 -6.8e-04 -8.7e-04 -8.5e-04 -8.4e-04 -6.4e-04 -6.9e-04 -7.0e-04
 -8.7e-04 -8.5e-04 -6.9e-04 -8.2e-04 -4.5e-04 -6.4e-04 -8.3e-04 -7.1e-04 -7.7e-04
 -6.6e-04 -5.6e-04

20

〔4行目〕

3.1e-03 3.1e-04 5.5e-03 -1.9e-03 -1.3e-03 -1.3e-04 -2.1e-04 -9.6e-04 -3.8e-04
 -4.7e-04 -3.7e-04 -4.5e-04 -5.1e-04 -3.8e-04 -3.6e-04 -4.6e-04 -4.4e-04 -4.9e-04
 -3.8e-04 -4.1e-04 -5.2e-04 -5.1e-04 -5.1e-04 -3.8e-04 -4.1e-04 -4.2e-04 -5.2e-04
 -5.1e-04 -4.1e-04 -4.9e-04 -2.7e-04 -3.8e-04 -5.0e-04 -4.3e-04 -4.6e-04 -4.0e-04
 -3.4e-04

30

【 0 1 1 8 】

【 表 1 4 】

共分散行列 C (続き)

【5行目】

2.5e-03 2.3e-03 -1.9e-03 -1.3e-03 -1.3e-04 -2.1e-04 -9.6e-04 -3.8e-04 -4.7e-04
 -3.7e-04 -4.5e-04 -5.1e-04 -3.8e-04 -3.6e-04 -4.6e-04 -4.4e-04 -4.9e-04 -3.8e-04
 -4.1e-04 -5.2e-04 -5.1e-04 -5.1e-04 -3.8e-04 -4.1e-04 -4.2e-04 -5.2e-04 -5.1e-04
 -4.1e-04 -4.9e-04 -2.7e-04 -3.8e-04 -5.0e-04 -4.3e-04 -4.6e-04 -4.0e-04 -3.4e-04

【6行目】

7.2e-02 -1.6e-02 -1.1e-02 -1.1e-03 -1.7e-03 -8.0e-03 -3.2e-03 -4.0e-03 -3.1e-03
 -3.8e-03 -4.3e-03 -3.2e-03 -3.0e-03 -3.8e-03 -3.6e-03 -4.1e-03 -3.2e-03 -3.4e-03
 -4.3e-03 -4.3e-03 -4.2e-03 -3.2e-03 -3.4e-03 -3.5e-03 -4.3e-03 -4.3e-03 -3.4e-03
 -4.1e-03 -2.2e-03 -3.2e-03 -4.1e-03 -3.6e-03 -3.8e-03 -3.3e-03 -2.8e-03

10

【7行目】

1.5e-01 1.0e-01 1. e-02 1.6e-02 7.3e-02 -6.4e-03 -7.9e-03 -6.1e-03 -7.5e-03 -
 8.5e-03 -6.4e-03 -6.0e-03 -7.7e-03 -7.3e-03 -8.2e-03 -6.4e-03 -6.8e-03 -8.7e-03
 -8.5e-03 -8.4e-03 -6.4e-03 -6.9e-03 -7.0e-03 -8.7e-03 -8.5e-03 -6.9e-03 -8.2e-03
 -4.5e-03 -6.4e-03 -8.3e-03 -7.1e-03 -7.7e-03 -6.6e-03 -5.6e-03

20

【8行目】

1.0e-01 8.1e-03 1.2e-02 2.5e-02 -4.5e-03 -5.5e-03 -4.3e-03 -5.3e-03 -6.0e-03 -
 4.5e-03 -4.2e-03 -5.4e-03 -5.1e-03 -5.7e-03 -4.5e-03 -4.7e-03 -6.1e-03 -6.0e-03
 -5.9e-03 -4.5e-03 -4.8e-03 -4.9e-03 -6.1e-03 -6.0e-03 -4.8e-03 -5.7e-03 -3.1e-03
 -4.5e-03 -5.8e-03 -5.0e-03 -5.4e-03 -4.6e-03 -3.9e-03

30

【9行目】

2.4e-03 5.2e-05 3.6e-03 -4.3e-04 -5.4e-04 -4.2e-04 -5.1e-04 -5.8e-04 -4.3e-04
 -4.1e-04 -5.2e-04 -4.9e-04 -5.6e-04 -4.3e-04 -4.6e-04 -5.9e-04 -5.8e-04 -5.7e-04
 -4.3e-04 -4.7e-04 -4.8e-04 -5.9e-04 -5.8e-04 -4.7e-04 -5.6e-04 -3.0e-04 -4.3e-04
 -5.6e-04 -4.9e-04 -5.2e-04 -4.5e-04 -3.8e-04

【 0 1 1 9 】

40

【表 1 5】

共分散行列 C (続き)

〔10行目〕

5.0e-03 6.5e-03 -6.9e-04 -8.5e-04 -6.6e-04 -8.1e-04 -9.2e-04 -6.9e-04 -6.5e-04
 -8.3e-04 -7.9e-04 -8.8e-04 -6.9e-04 -7.3e-04 -9.4e-04 -9.2e-04 -9.1e-04 -6.9e-04
 -7.4e-04 -7.6e-04 -9.4e-04 -9.2e-04 -7.4e-04 -8.8e-04 -4.8e-04 -6.9e-04 -9.0e-04
 -7.7e-04 -8.3e-04 -7.2e-04 -6.1e-04

〔11行目〕

7.2e-02 -3.2e-03 -4.0e-03 -3.1e-03 -3.8e-03 -4.3e-03 -3.2e-03 -3.0e-03 -3.8e-03
 -3.6e-03 -4.1e-03 -3.2e-03 -3.4e-03 -4.3e-03 -4.3e-03 -4.2e-03 -3.2e-03 -3.4e-03
 -3.5e-03 -4.3e-03 -4.3e-03 -3.4e-03 -4.1e-03 -2.2e-03 -3.2e-03 -4.1e-03 -3.6e-03
 -3.8e-03 -3.3e-03 -2.8e-03

10

〔12行目〕

3.4e-02 -8.7e-04 -5.1e-04 -7.9e-04 -9.9e-04 -1.3e-03 -1.2e-03 -1.5e-03 -1.5e-03
 -1.6e-03 -1.3e-03 -1.4e-03 -1.7e-03 -1.7e-03 -1.7e-03 -1.3e-03 -1.4e-03 -1.4e-03
 -1.7e-03 -1.7e-03 -1.4e-03 -1.6e-03 -8.9e-04 -1.3e-03 -1.7e-03 -1.4e-03 -1.5e-03
 -1.3e-03 -1.1e-03

20

〔13行目〕

2.6e-02 -1.4e-03 1.3e-03 1.2e-03 -1.6e-03 1.2e-03 2.7e-03 -1.2e-03 6.9e-04 -
 1.6e-03 3.5e-03 2.4e-03 2.5e-03 1.8e-03 -1.6e-03 -1.3e-03 -1.1e-03 3.5e-03
 3.6e-03 1.5e-04 6.9e-04 -6.8e-04 1.8e-03 1.0e-03 -1.1e-03 4.0e-03 9.3e-04 -
 1.4e-03

30

〔14行目〕

2.4e-02 -1.4e-03 2.6e-03 -1.2e-03 2.8e-04 -4.1e-05 -1.4e-03 5.8e-04 -1.2e-03 -
 8.7e-04 1.3e-03 1.4e-03 4.1e-03 -1.2e-03 -4.7e-04 -2.0e-04 -8.1e-04 -5.0e-04
 4.0e-03 7.2e-04 -8.6e-04 -1.2e-03 -1.6e-03 -9.4e-04 8.2e-04 2.3e-04 -7.9e-04

【 0 1 2 0 】

【 表 1 6 】

共分散行列 C (続き)

[15行目]

2.5e-02 2.7e-04 -1.5e-03 6.6e-04 -7.3e-04 -1.7e-03 1.4e-03 -1.5e-03 1.7e-03
 -8.3e-04 2.4e-03 -1.6e-03 -1.5e-03 9.5e-04 -8.4e-05 1.0e-03 4.1e-03 3.0e-04
 3.8e-03 2.3e-04 -1.3e-03 1.9e-04 -1.7e-03 2.3e-03 -1.1e-03 -4.0e-04

[16行目]

2.1e-02 -1.7e-03 1.8e-03 1.1e-03 2.4e-03 5.7e-03 -1.7e-03 5.5e-04 -1.0e-03
 2.8e-03 -1.4e-03 -1.7e-03 -1.4e-03 1.7e-03 4.5e-03 2.4e-03 1.5e-04 6.7e-04
 8.9e-05 5.1e-06 2.1e-03 -2.0e-04 1.0e-03 2.7e-03 6.4e-04

10

[17行目]

3.4e-02 -4.8e-04 -8.2e-04 -7.4e-04 -9.2e-04 -1.3e-03 -1.4e-03 -1.7e-03 -1.7e-03
 -1.7e-03 -1.3e-03 -1.4e-03 -1.4e-03 -1.7e-03 -1.7e-03 -1.4e-03 -1.6e-03 -8.9e-04
 -1.3e-03 -1.7e-03 -1.4e-03 -1.5e-03 -1.3e-03 -1.1e-03

[18行目]

20

2.3e-02 -1.0e-05 -6.5e-04 6.0e-03 -1.2e-03 1.6e-04 -1.2e-03 1.5e-03 -1.3e-03
 1.2e-03 -7.2e-04 4.5e-03 1.6e-04 1.3e-03 2.3e-03 -6.8e-04 -7.0e-04 -6.3e-04
 4.4e-04 -1.3e-03 -3.0e-04 1.6e-03 1.2e-03

[20行目]

2.6e-02 5.4e-04 -8.9e-04 -1.5e-03 2.5e-03 8.5e-04 -4.8e-04 6.6e-03 -1.5e-03
 9.4e-04 1.0e-03 4.9e-04 2.9e-03 -1.0e-03 2.3e-03 3.9e-03 -6.0e-04 4.0e-03
 1.7e-03 -1.1e-03 3.8e-03 -9.2e-04

30

[21行目]

2.5e-02 1.5e-03 -1.5e-03 -9.0e-04 2.3e-03 -8.8e-04 -1.3e-03 -1.5e-03 -1.4e-04
 -1.6e-03 4.1e-03 4.1e-04 -4.3e-04 -2.9e-04 -4.5e-04 1.1e-03 6.3e-03 3.8e-03
 3.2e-04 -1.5e-03 -1.3e-03

【 0 1 2 1 】

【 表 1 7 】

共分散行列 C (続き)

【22行目】

2.1e-02 -1.6e-03 5.5e-05 1.4e-03 1.7e-03 2.7e-03 -1.6e-03 -4.8e-04 1.1e-03 -
 6.1e-06 4.0e-03 4.0e-03 3.9e-03 -1.0e-03 -1.1e-03 1.3e-03 1.9e-03 -8.9e-04
 1.6e-05 1.6e-03

【23行目】

3.4e-02 -6.4e-04 -1.0e-03 -9.9e-04 -9.7e-04 -1.3e-03 -1.4e-03 -1.4e-03 -1.7e-03
 -1.7e-03 -1.4e-03 -1.6e-03 -8.9e-04 -1.3e-03 -1.7e-03 -1.4e-03 -1.5e-03 -1.3e-03
 -1.1e-03

10

【24行目】

2.5e-02 -2.0e-04 2.0e-03 2.4e-03 -1.4e-03 2.5e-03 6.6e-04 1.2e-03 -1.8e-03
 1.8e-03 8.4e-04 1.2e-03 1.9e-03 6.0e-04 4.9e-03 2.9e-03 -3.3e-04 1.7e-03

【25行目】

2.8e-02 2.1e-03 2.7e-03 -1.7e-03 2.2e-03 1.1e-03 4.8e-03 3.7e-03 2.4e-03
 1.7e-03 -5.7e-04 1.9e-04 -1.3e-03 6.9e-03 -3.7e-04 1.8e-03 -1.5e-03

20

【26行目】

2.7e-02 5.3e-03 -1.7e-03 5.8e-04 -3.8e-04 2.0e-03 2.7e-03 1.5e-04 1.1e-03
 6.6e-04 5.1e-06 -3.6e-04 3.8e-03 1.7e-03 9.4e-04 3.1e-03

【27行目】

2.8e-02 -1.7e-03 -1.0e-04 -5.7e-04 1.3e-03 -8.3e-04 -5.3e-04 1.1e-03 8.2e-04
 9.6e-04 -1.5e-03 2.7e-03 -1.6e-03 -8.9e-04 -1.5e-03

30

【27行目】

3.4e-02 -6.6e-04 -6.9e-04 -1.0e-03 -9.9e-04 -1.4e-03 -1.6e-03 -8.9e-04 -1.3e-03
 -1.7e-03 -1.4e-03 -1.5e-03 -1.3e-03 -1.1e-03

【28行目】

2.5e-02 -2.3e-04 9.8e-04 1.6e-03 5.5e-03 2.2e-03 -6.8e-04 1.5e-03 5.4e-03
 2.2e-03 5.1e-03 2.8e-03 7.9e-04

40

【 0 1 2 2 】

【表 1 8 】

共分散行列 C (続き)

[29行目]

2.5e-02 1.9e-03 3.8e-03 3.4e-04 4.9e-04 -2.7e-04 6.0e-04 -1.6e-03 4.1e-03
6.8e-03 9.0e-04 1.3e-03

[30行目]

2.7e-02 2.5e-04 1.3e-04 9.2e-04 1.6e-03 -1.1e-03 4.0e-03 -6.6e-04 1.5e-03 ·
1.2e-03 -1.5e-03

10

[31行目]

2.9e-02 4.6e-03 4.0e-03 9.5e-04 1.0e-03 1.2e-03 6.6e-04 -1.5e-03 -9.2e-04 ·
6.5e-04

[32行目]

2.0e-02 2.8e-03 2.5e-03 -8.8e-04 9.2e-04 -1.1e-04 3.9e-03 6.7e-05 2.9e-04

[33行目]

2.7e-02 1.1e-03 3.0e-03 2.6e-03 -4.0e-04 -2.4e-04 -1.3e-03 -8.2e-06

20

[34行目]

1.7e-02 3.5e-03 9.8e-04 -1.0e-03 -9.3e-04 -9.3e-04 -7.1e-04

[35行目]

2.4e-02 7.0e-04 5.7e-04 -1.4e-03 5.3e-04 -1.1e-03

【 0 1 2 3 】

【 表 1 9 】

(続き)

30

[36行目]

2.6e-02 -1.4e-04 -4.2e-04 1.8e-03 1.4e-03

[37行目]

2.6e-02 -1.4e-03 1.1e-03 -1.1e-04

[38行目]

2.6e-02 2.4e-03 1.2e-03

[39行目]

2.0e-02 3.5e-03

40

[40行目]

2.3e-02

【 0 1 2 4 】

上記の逐次アップデートで計算された共分散行列 C ' を用いて特異値分解し、最初の 1 0 個の特異値と最大の特異値に対する特異ベクトルを計算した結果を、図 2 8 に示す。また、特異ベクトル(最大の特異値に対するもの)を表 2 0 に示す。

50

【 0 1 2 5 】

【表 2 0】

〔第 1 特異値に対する特異ベクトル〕 (N = 4 0 次元)

-5.2960e-01 -3.8319e-01 -5.2864e-02 -3.2729e-02 -3.1732e-02
 -2.6187e-01 5.2893e-01 3.8302e-01 3.6178e-02 5.7374e-02
 2.6109e-01 3.1860e-05 4.9484e-05 3.5469e-05 4.4623e-05
 5.2818e-05 3.1851e-05 3.6387e-05 4.8000e-05 4.3140e-05
 5.1465e-05 3.1915e-05 4.2917e-05 5.5013e-05 5.4694e-05
 5.1524e-05 3.1898e-05 4.3605e-05 4.3651e-05 5.4658e-05
 5.5368e-05 4.3439e-05 5.1774e-05 2.6056e-05 3.7219e-05
 5.1642e-05 4.4973e-05 4.7018e-05 3.9580e-05 3.2341e-05

10

【 0 1 2 6 】

比較例 3 として、従来手法（あらかじめ文書総数、キーワード総数がわかっている場合）で表 5 ~ 表 1 0 のようなデータが与えられた場合に、直接共分散行列を計算し、これを特異値分解して特異値（大きい方から 1 0 個）と最大特異値に対する特異ベクトルを求めた結果を図 2 9 に示す。

20

【 0 1 2 7 】

また、第 1 特異値に対する特異ベクトル (N = 4 0 次元) を下記表 2 1 に示す

【 0 1 2 8 】

【表 2 1】

〔第 1 特異値に対する特異ベクトル〕 (N = 4 0 次元)

-5.2960e-01 -3.8319e-01 -2.6187e-01 -5.2864e-02 -3.2729e-02
 -3.1732e-02 5.2893e-01 2.6109e-01 3.8302e-01 3.6178e-02
 5.7374e-02 3.1860e-05 3.1851e-05 3.1915e-05 3.1898e-05
 4.2917e-05 4.3605e-05 4.9484e-05 3.6387e-05 5.1774e-05
 3.5469e-05 4.8000e-05 5.5013e-05 4.4623e-05 4.3651e-05
 5.1642e-05 4.7018e-05 4.3140e-05 5.4694e-05 5.4658e-05
 3.2341e-05 2.6056e-05 5.1465e-05 3.9580e-05 5.2818e-05
 4.3439e-05 5.1524e-05 5.5368e-05 3.7219e-05 4.4973e-05

30

【 0 1 2 9 】

上述した実施例 3 と、比較例 3 とを比較した結果、本発明の方法で得られる特異値と、全体のデータから共分散行列を作って、これを特異値分解して得られる特異値は、一致していることが示された。

40

【 0 1 3 0 】

また、最大特異値についてみても、たとえば、従来手法で得られる最大特異値に対する特異ベクトルの第 3 要素は -2.6187e-01 であり、これは本発明で得られる最大特異値に対する特異ベクトルの第 6 要素に対応していることが示された。すなわち、最大特異値に対する特異ベクトルは、順列(permutation)を除いて、一致していることが示された。

【 0 1 3 1 】

(実施例 4)

50

< ダウンデート例 >

実施例 4 として、実施例 1 で使用したデータから日付が 2001 年 7 月 1 日 (data 1) と 2 日 (data 2) のデータを削除することにより、データベースからのドキュメントのダウンデートについて検討した。ダウンデートのため、まず最初の処理は、SUM-MEAN(M)₁ と、SUM(M)₁ とを更新した。このために、data 1 と data 2 とをまず走査して、削除分の SUM-MEAN(D)₁ と、SUM(D)₁ とをまず作成した。ここで D は、削除する文書数である。SUM-MEAN(D)₁ と、SUM(D)₁ とは、それぞれ、削除する文書の平均ベクトル (N 次元) と削除される文書だけからなる積和行列成分 (N × N の対称行列) となる。

【0132】

具体的には、SUM-MEAN(D)₁ と、SUM(D)₁ とを求めた後、新たな SUM-MEAN(M)₁ と、SUM(M)₁ とを下記式を使用して算出した。

【0133】

【数 16】

$$\text{SUM-MEAN}(M)_1 = \text{SUM-MEAN}(M)_1 - \text{SUM-MEAN}(D)_1$$

$$\text{SUM}(M)_1 = \text{SUM}(M)_1 - \text{SUM}(D)_1$$

【0134】

上記式は、N 次元ベクトルの引き算であり、上記式は、N × N の対称行列の引き算 (要素ごと) である。これで、SUM-MEAN(M)₁ と、SUM(M)₁ との更新の第 1 ステップを終了する。

【0135】

次いで、削除することによってキーワード・ハッシュ・テーブルに保持する必要がなくなるキーワードが出てくる可能性があるため、これを D 個の文書に関してもう一度走査して調査した。実施例 3 では、data 1 と data 2 とを削除するので、これに伴い図 30 に示すキーワードが削除対象として検出された。

【0136】

SUM-MEAN(M)₁ に関しては、上記のキーワード各々に対応するインデックスの位置の要素を削除した (これを (i₁、i₂、i₃、i₄、i₅、i₆) と表記する。) 。この結果、キーワード 6 個分だけ減少した 34 次元のベクトルが得られた。SUM(M)₁ に関しては、この 6 つのインデックス各々に対応する行列の縦方向と横方向の要素を削除した。たとえば、"AlGore" に対するインデックス i₁ に関しては、SUM(M)₁ の中で行番号または列番号が i₁ のものをすべて削除した。この時点で N × N 行列は、(N - 1) × (N - 1) 行列として生成され、最終的には、(N - 6) × (N - 6) の行列として、SUM(M)₁ を生成した。

【0137】

上述した削除を、"Japan" に対応するインデックス i₆ まで繰り返した。このプロセスを終了した時点で、SUM-MEAN(M)₁ と SUM(M)₁ との更新を終了した。この結果として、元の 40 次元から 6 次元少ない 34 次元ベクトルの SUM-MEAN(M)₁ と SUM(M)₁ とを得た。

【0138】

最後に SUM(M)₂ をアップデートの場合と同じ式で計算し、共分散行列を求め、これを特異値分解して特異値 (大きい方から 10 個) と最大特異値に対する特異ベクトルを求めた。その結果を、図 31 に示す。また、表 22 には、第 1 特異値に対する特異ベクトル (N = 34 次元) を示す

【0139】

【表 22】

10

20

30

40

〔第1特異値に対する特異ベクトル〕 (N = 34次元)

6.6945e-01 1.7290e-01 2.6265e-02 6.7805e-02 6.6945e-01 -2.1384e-02
 -5.6159e-02 -3.5503e-02 -4.6750e-02 -5.8561e-02 -2.1388e-02 -3.9300e-02
 -5.4406e-02 -4.5086e-02 -5.8817e-02 -2.1630e-02 -4.9441e-02 -6.3746e-02
 -6.3778e-02 -5.6538e-02 -2.1591e-02 -5.0169e-02 -4.8973e-02 -6.2525e-02
 -6.6006e-02 -4.9784e-02 -5.9844e-02 -2.6616e-02 -3.7788e-02 -5.7980e-02
 -5.1621e-02 -5.1754e-02 -4.1644e-02 -3.2110e-02

10

【0140】

また、比較例4としてdata3からdata7までをまとめた全体の行列を表5～表10と同様なフォーマットとして一度に共分散行列を計算してSVDを実行させることにより得た特異値(大きい方から10個)を図32に示す。また、最大特異値に対する特異ベクトルを求めた結果を表23に示す。

【0141】

【表23】

〔第1特異値に対する特異ベクトル〕 (N = 34次元)

6.6945e-01 6.6945e-01 1.7290e-01 2.6265e-02 6.7805e-02 -2.1384e-02
 -2.1388e-02 -2.1630e-02 -2.1591e-02 -4.9441e-02 -5.0169e-02 -5.6159e-02
 -3.9300e-02 -5.9844e-02 -3.5503e-02 -5.4406e-02 -6.3746e-02 -4.6750e-02
 -4.8973e-02 -5.7980e-02 -5.1754e-02 -4.5086e-02 -6.3778e-02 -6.2525e-02
 -3.2110e-02 -2.6616e-02 -5.8817e-02 -4.1644e-02 -5.8561e-02 -4.9784e-02
 -5.6538e-02 -6.6006e-02 -3.7788e-02 -5.1621e-02

20

【0142】

実施例4および比較例4の結果を比較すると、最大から10個の特異値に関しては、従来手法と一致していることが示される。また、本手法で得られる最大特異値に対する特異ベクトルと従来手法で得られる特異ベクトルとを比較すると、順序を除き、それぞれの特異ベクトルは、ほぼ一致していることが示された。たとえば、実施例4により得られた特異ベクトルの2番目のベクトル要素1.7290e-01は、従来手法(一括手法)の3番目のベクトルの要素に対応していることが示される。

30

【0143】

本発明の上述した情報検索を実行させるためのプログラムは、これまで知られたいかなる言語を使用しても記述することができ、例えばC言語、C++言語、Java(登録商標)といった言語を使用して本発明の方法を実行させるためのプログラムとすることができる。また、本発明の方法を実行させるためのプログラムは、フロッピー(登録商標)ディスク、ハードディスク、磁気テープといったコンピュータ可読な磁気記録媒体、光磁気ディスク、CD-ROM、DVDといったコンピュータ可読な記録媒体に記憶させることができる。

40

【0144】

上述したように、本発明によれば、ドキュメントが逐次アップデートされていく大規模なデータベースにおける以前に計算された次元削減に関する結果を使用して、アップデートされたドキュメントを含むドキュメント-アトリビュート行列の特異値分解を効率的に実行することが可能となる。また、本発明により得られた特異ベクトルを使用してアップデートされたドキュメント-アトリビュート行列の次元削減を実行させた後、情報検索を行

50

うことにより、効率的で、かつ高精度の情報検索を実行することが可能となる。

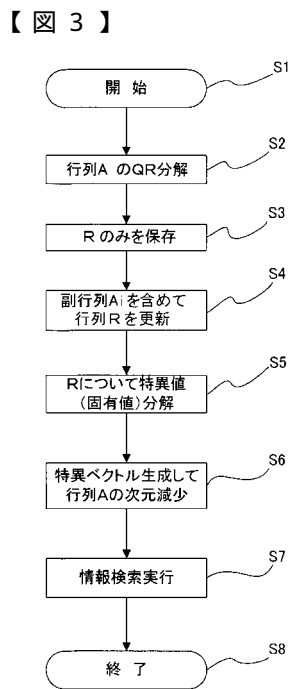
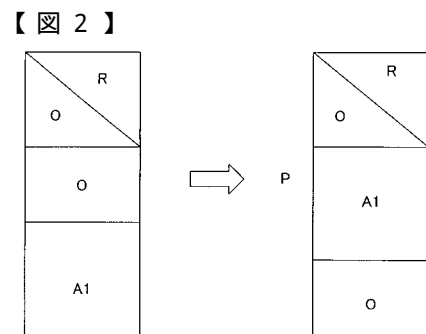
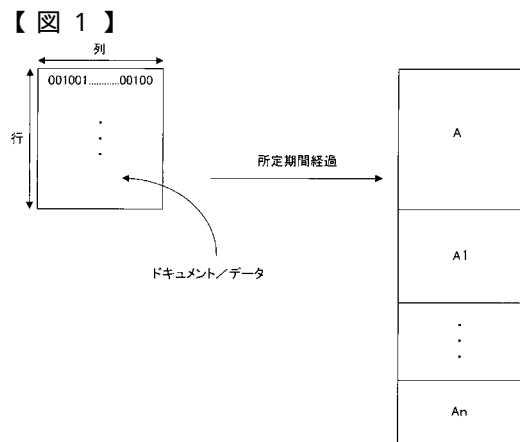
【0145】

これまで本発明を図面に示した実施の形態に基づいて詳細に説明してきたが、本発明は、図面に示した実施の形態に限定されるものではなく、種々の変更、別の実施の形態を採用することが可能である。例えば、本発明においては、ドキュメントを文書ドキュメントとして説明してきたが、本発明において使用することができるドキュメントは、文書ドキュメントに限定されるものではなく、オーディオ、グラフィックス、動画といったドキュメントを含むことができる。

【図面の簡単な説明】

- 【図1】 本発明において使用するデータベースの概略構成を示した図。 10
- 【図2】 副行列A1が加えられた時の更新を示した図。
- 【図3】 本発明のQR分解方を用いた情報検索方法のフローチャート。
- 【図4】 一般的な行列Dに対するQR分解の実行のための擬似コードを示した図。
- 【図5】 行列AのQR分解により得られる行列の構成要素を示した概略図。
- 【図6】 特異値分解してk次元への次元削減を示した図。
- 【図7】 本発明の第2の実施形態の情報検索方法のフローチャートを示した図。
- 【図8】 図7に示したステップS11において使用する共分散行列生成のための擬似コードを示した図。
- 【図9】 アトリビュート・ハッシュ・テーブルを変更するプロセスを示した概略図。
- 【図10】 本発明において使用するアトリビュート・ハッシュ・テーブルを示した図。 20
- 【図11】 アトリビュートの削除を実行する場合の概略図。
- 【図12】 アトリビュートの削除と、生成される行列との関係を示した図。
- 【図13】 アトリビュートの追加と、生成される行列との関係を示した図。
- 【図14】 アトリビュートの追加と、生成される行列との関係を示した図。
- 【図15】 本発明の情報検索システムを示した概略図。
- 【図16】 タイム・スタンプとデータ・ファイルメイトからなるデータ例を示した図。
- 【図17】 図16で説明した個々のファイルを示した図。
- 【図18】 図16で説明した個々のファイルを示した図。
- 【図19】 図16で説明した個々のファイルを示した図。
- 【図20】 図16で説明した個々のファイルを示した図。 30
- 【図21】 図16で説明した個々のファイルを示した図。
- 【図22】 図16で説明した個々のファイルを示した図。
- 【図23】 図16で説明した個々のファイルを示した図。
- 【図24】 本発明において使用したキーワードを例示した図。
- 【図25】 data1のみを処理した段階でのSUM-MEAN(M)₁と、SUM(M)₁とを示した図。
- 【図26】 data2をアップデートしてデータを追加して得られた結果を示した図。
- 【図27】 data7までアップデートして得られたSUM-MEAN(M)₁を示した図。
- 【図28】 アップデートされた共分散行列C'を使用して得られた特異ベクトルを示した図。 40
- 【図29】 はじめから表5～表10のデータとして与えられた行列を直接特異値分解して得た特異値を示した図。
- 【図30】 削除の対象となったキーワードを示した図。
- 【図31】 本発明によりダウンデートされた共分散行列により得られた特異値を示した図。
- 【図32】 表5～表10のデータから予めデータをダウンデートして直接共分散行列を生成して得られた特異値を示した図。
- 【符号の説明】
- 10...コンピュータ

- 1 2 ... データベース
- 1 4 ... ネットワーク
- 1 6 ... コンピュータ



【 図 4 】

Pseudo Code for the QR-decomposition of $D \in R^{m \times n}$

$$D = \begin{bmatrix} d_{11}^T \\ d_{21}^T \\ \vdots \\ d_{m1}^T \end{bmatrix} = \begin{bmatrix} d_{21} & d_{22} & \dots & d_{2n} \\ d_{31} & d_{32} & \dots & d_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{m1} & d_{m2} & \dots & d_{mn} \end{bmatrix}$$

(when information about Q does not need to be stored).

for $k = 1, 2, \dots, n$

determine a reflector $Q_k = I - \gamma_k u^{(k)} u^{(k)T}$ such that

$$Q_k [d_{k1} \dots d_{mk}]^T = [-\sigma_k \ 0 \ \dots \ 0]^T$$

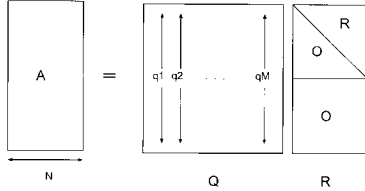
set $d_{kk} = -\sigma_k$

set $[d_{k+1k} \dots d_{mk}]^T$ to zero

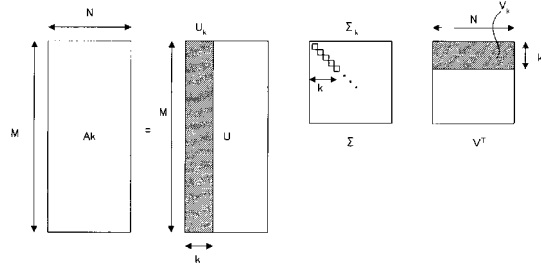
for $j = k + 1, k + 2, \dots, n$

$$[d_{kj} \dots d_{mj}]^T = Q_k [d_{kj} \dots d_{mj}]^T$$

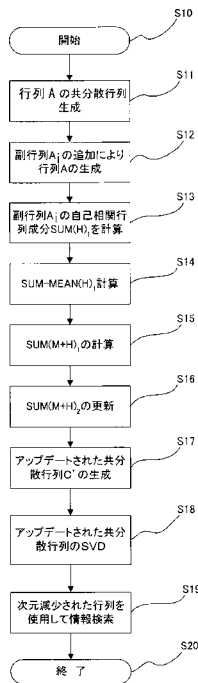
【 図 5 】



【 図 6 】



【 図 7 】



【 図 8 】

Pseudo Code for Constructing Covariance Matrix (basic case, no up/downdating)

Code to determine the covariance matrix for a set of n -dimensional vectors

$$\{d_i \mid d_i = (d(i, 1), d(i, 2), \dots, d(i, N))^T \in R^N\}_{i=1}^M$$

where:

M is the cardinality of the set $\{d_i\}$ (unknown - to be determined),

$D(\cdot)$ is the component-wise sum of the document vectors,

$C(\cdot, \cdot)$ is an N -by- N covariance matrix,

temp is a parameter for storing intermediate results during computations.

routine for initialization

$M = 0$

for $i = 1$ to N

$D(i) = 0$

for $j = 1$ to N

$C(i, j) = 0$

routine to compute component-wise sum of the document vectors D and cardinality M

for $i = 1$ to end-of-stack

for $j = 1$ to N

$D(j) = D(j) + d(i, j)$

$M = M + 1$

routine to compute $M \times C(\cdot, \cdot)$

for $i = 1$ to N

for $j = 1$ to N

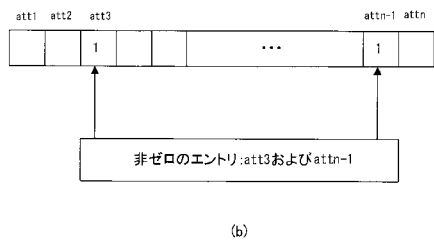
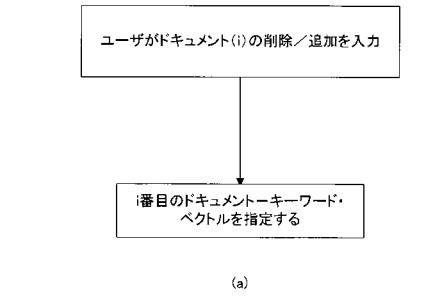
temp = 0

for $k = 1$ to M

temp = temp + $d(k, i) \cdot d(k, j)$

$C(i, j) = [\text{temp} - D(i) \cdot D(j) / M] / M$

【 図 9 】



【 図 10 】

| 属性 | 属性を含むドキュメントの数 |
|-------|---------------|
| 属性1 | 1 |
| 属性2 | 25 |
| 属性3 | 6 |
| ... | ... |
| 属性n-1 | 33 |
| 属性n | 12 |

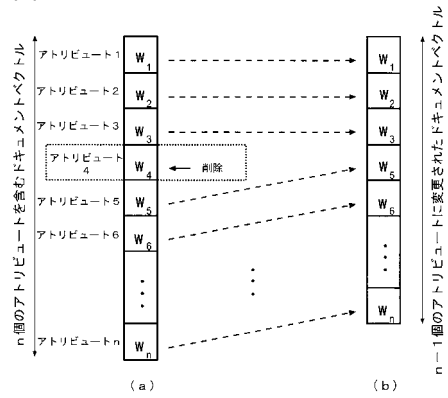
5 ← 削除 (属性1, 2, 3)

7 → 追加 (属性1)

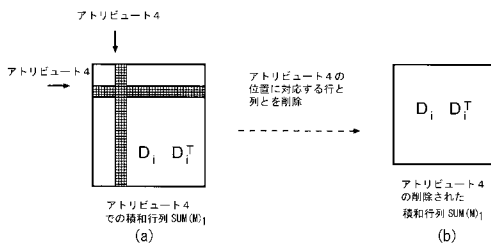
32 ← 削除 (属性n-1)

34 → 追加 (属性n)

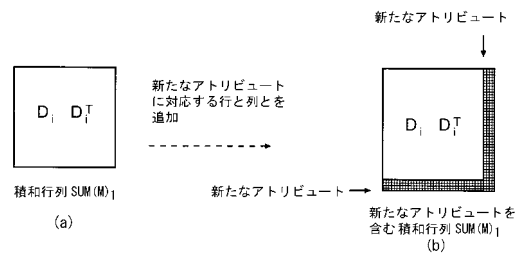
【 図 11 】



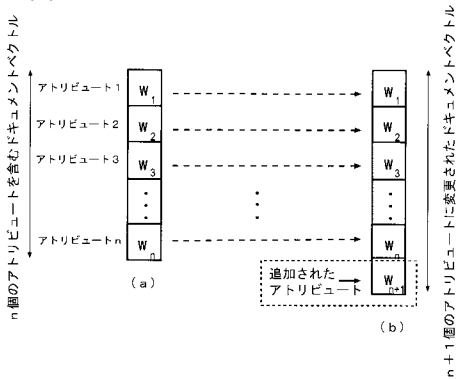
【 図 12 】



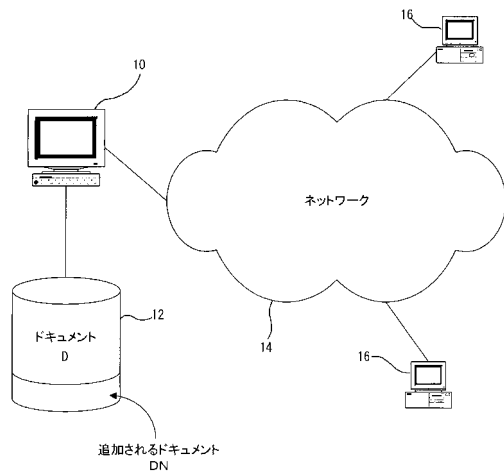
【 図 14 】



【 図 13 】



【 図 15 】



【 16 】

- 20010701 data1
 - 20010702 data2
 - 20010703 data3
 - 20010704 data4
 - 20010705 data5
 - 20010706 data6
 - 20010707 data7
- dataset

【 17 】

- 1| Clinton 1.0 AlGore 1.0 Chelsea 0.4 ↓
- 2| Clinton 1.0 AlGore 1.0 Bush 0.3 ↓
- 3| Clinton 1.0 AlGore 1.0 Japan 0.1 ↓
- 4| Clinton 1.0 AlGore 1.0 Chelsea 0.3 ↓
- 5| Clinton 1.0 AlGore 1.0 Japan 0.2 ↓
- 6| Clinton 1.0 AlGore 1.0 Chelsea 0.2 ↓
- 7| Clinton 1.0 AlGore 1.0 Japan 0.3 ↓
- 8| Clinton 1.0 AlGore 1.0 Chelsea 0.1 ↓
- 9| Clinton 1.0 AlGore 1.0 Japan 0.4 ↓
- 10| Clinton 1.0 AlGore 1.0 ↓
- 11| Clinton 1.0 AlGore 1.0 Hillary 0.5 Bush 0.3 ↓
- 12| Clinton 1.0 AlGore 1.0 Hillary 0.5 Chelsea 0.4 ↓
- 13| Clinton 1.0 AlGore 1.0 Hillary 0.5 Japan 0.1 ↓
- 14| Clinton 1.0 AlGore 1.0 Hillary 0.5 ↓
- 15| Clinton 1.0 AlGore 1.0 Hillary 0.5 Bush 0.3 ↓
- 16| Clinton 1.0 AlGore 0.5 Hillary 1.0 Chelsea 0.1 ↓
- 17| Clinton 1.0 AlGore 0.5 Hillary 1.0 Chelsea 0.4 ↓
- 18| Clinton 1.0 AlGore 0.5 Hillary 1.0 Japan 0.1 ↓
- 19| Clinton 1.0 AlGore 0.5 Hillary 1.0 ↓
- 20| Clinton 1.0 AlGore 0.5 Hillary 1.0 Bush 0.3 Japan 0.2 ↓

【 18 】

- 1| Clinton 1.0 Hillary 1.0 Japan 0.1 ↓
- 2| Clinton 1.0 Hillary 1.0 Chelsea 0.2 ↓
- 3| Clinton 1.0 Hillary 1.0 ↓
- 4| Clinton 1.0 Hillary 1.0 Bush 0.3 ↓
- 5| Clinton 1.0 Hillary 1.0 Chelsea 0.4 ↓
- 6| Java 1.0 JSP 1.0 apache 0.1 servlet 0.4 ↓
- 7| Java 1.0 JSP 1.0 servlet 0.3 ↓
- 8| Java 1.0 JSP 1.0 apache 0.2 ↓
- 9| Java 1.0 JSP 1.0 ↓
- 10| Java 1.0 JSP 1.0 apache 0.1 ↓
- 11| Java 1.0 JSP 1.0 servlet 0.2 ↓
- 12| Java 1.0 JSP 1.0 apache 0.3 ↓
- 13| Java 1.0 JSP 1.0 servlet 0.4 ↓
- 14| Java 1.0 JSP 1.0 servlet 0.1 ↓
- 15| Java 1.0 JSP 1.0 apache 0.1 ↓
- 16| Java 1.0 applet 0.5 JSP 1.0 servlet 0.2 ↓
- 17| Java 1.0 applet 0.5 JSP 1.0 apache 0.3 ↓
- 18| Java 1.0 applet 0.5 JSP 1.0 apache 0.2 ↓
- 19| Java 1.0 applet 0.5 JSP 1.0 servlet 0.1 ↓
- 20| Java 1.0 applet 0.5 JSP 1.0 ↓

【 19 】

- 1| Java 1.0 applet 1.0 JSP 0.5 apache 0.1 ↓
- 2| Java 1.0 applet 1.0 JSP 0.5 servlet 0.3 ↓
- 3| Java 1.0 applet 1.0 JSP 0.5 servlet 0.2 ↓
- 4| Java 1.0 applet 1.0 JSP 0.5 servlet 0.1 ↓
- 5| Java 1.0 applet 1.0 JSP 0.5 ↓
- 6| Java 1.0 applet 1.0 apache 0.1 ↓
- 7| Java 1.0 applet 1.0 ↓
- 8| Java 1.0 applet 1.0 servlet 0.3 ↓
- 9| Java 1.0 applet 1.0 apache 0.2 ↓
- 10| Java 1.0 applet 1.0 servlet 0.1 ↓
- 11| Bluetooth 1.0 dream 0.1 ↓
- 12| Bluetooth 1.0 cat 0.1 ↓
- 13| Bluetooth 1.0 god 0.1 ↓
- 14| Bluetooth 1.0 museum 0.1 ↓
- 15| Bluetooth 1.0 ↓
- 16| soccer 1.0 center 0.1 ↓
- 17| soccer 1.0 spring 0.1 ↓
- 18| soccer 1.0 silk 0.1 ↓
- 19| soccer 1.0 lecture 0.1 ↓
- 20| soccer 1.0 ↓

【 20 】

- 1| matrix 1.0 heat 0.1 ↓
- 2| matrix 1.0 power 0.1 ↓
- 3| matrix 1.0 magnet 0.1 ↓
- 4| matrix 1.0 gourmet 0.1 ↓
- 5| matrix 1.0 ↓
- 6| DNA 1.0 music 0.1 ↓
- 7| DNA 1.0 shark 0.1 ↓
- 8| DNA 1.0 snail 0.1 ↓
- 9| DNA 1.0 montage 0.1 ↓
- 10| DNA 1.0 ↓
- 11| cat 0.3 magnet 0.4 lecture 0.2 gourmet 1.0 ↓
- 12| center 0.2 silk 0.3 magnet 0.2 lecture 0.3 pattern 0.2 montage 1.0 ↓
- 13| whisper 0.2 saga 0.2 museum 0.2 torch 1.0 ↓
- 14| heat 0.3 spring 0.4 power 0.2 stock 0.1 silk 0.3 flag 1.0 ↓
- 15| heat 1.0 dream 0.2 flag 0.3 ↓
- 16| music 1.0 stock 0.2 insect 0.4 flag 0.1 ↓
- 17| heat 0.2 dream 1.0 spring 0.3 power 0.4 god 0.2 insect 0.2 magnet 0.4 pavillion 0.3 pattern 0.2 ↓
- 18| center 1.0 spring 0.1 shark 0.3 magnet 0.4 lecture 0.3 pavillion 0.4 museum 0.4 ↓
- 19| whisper 1.0 spring 0.1 god 0.3 shark 0.3 ↓
- 20| cat 1.0 magnet 0.3 lecture 0.2 museum 0.2 pattern 0.2 ↓

【 22 】

- 1| music 1.0 whisper 0.4 god 0.2 stock 0.2 insect 0.2 snail 0.4 pattern 0.4 montage 0.4 ↓
- 2| dream 1.0 center 0.3 snail 0.2 ↓
- 3| center 1.0 shark 0.2 latitude 0.3 lecture 0.3 pattern 0.4 ↓
- 4| dream 0.2 whisper 1.0 spring 0.4 power 0.3 lecture 0.3 gourmet 0.4 montage 0.4 flag 0.2 ↓
- 5| center 0.2 cat 1.0 spring 0.2 shark 0.1 gourmet 0.2 ↓
- 6| spring 1.0 silk 0.2 saga 0.4 gourmet 0.4 ↓
- 7| cat 0.3 power 1.0 pavillion 0.3 flag 0.2 ↓
- 8| god 1.0 ↓
- 9| center 0.3 god 0.1 shark 1.0 insect 0.4 snail 0.3 ↓
- 10| whisper 0.4 stock 1.0 silk 0.4 snail 0.4 museum 0.2 ↓
- 11| god 0.3 shark 0.1 insect 1.0 magnet 0.3 ↓
- 12| stock 0.3 silk 1.0 snail 0.3 museum 0.4 flag 0.4 ↓
- 13| heat 0.4 whisper 0.4 god 0.3 magnet 1.0 saga 0.2 lecture 0.1 museum 0.4 ↓
- 14| power 0.4 snail 1.0 museum 0.3 ↓
- 15| latitude 1.0 pavillion 0.4 ↓
- 16| stock 0.3 saga 1.0 torch 0.4 ↓
- 17| lecture 1.0 museum 0.2 pattern 0.2 gourmet 0.2 montage 0.2 ↓
- 18| music 0.2 spring 0.2 insect 0.2 latitude 0.2 pavillion 1.0 ↓
- 19| dream 0.2 museum 1.0 ↓
- 20| silk_0.1 lecture 0.3 museum 0.2 pattern 1.0 montage 0.3 ↓

【 21 】

- 1| spring 1.0 stock 0.4 pavillion 0.4 museum 0.2 ↓
- 2| power 1.0 snail 0.3 gourmet 0.3 ↓
- 3| whisper 0.3 god 1.0 ↓
- 4| spring 0.3 shark 1.0 snail 0.2 montage 0.4 ↓
- 5| spring 0.1 stock 1.0 magnet 0.2 latitude 0.2 pavillion 0.3 pattern 0.3 ↓
- 6| dream 0.2 shark 0.3 stock 0.1 insect 1.0 pavillion 0.3 museum 0.3 ↓
- 7| insect 0.2 silk 1.0 ↓
- 8| music 0.3 power 0.4 magnet 1.0 gourmet 0.4 flag 0.4 ↓
- 9| dress 0.4 god 0.3 insect 0.3 magnet 0.3 snail 1.0 gourmet 0.2 ↓
- 10| heat 0.3 music 0.2 shark 0.3 insect 0.3 magnet 0.2 latitude 1.0 lecture 0.3 museum 0.3 ↓
- 11| whisper 0.2 shark 0.1 snail 0.3 saga 1.0 pattern 0.4 ↓
- 12| center 0.2 whisper 0.3 stock 0.4 silk 0.2 lecture 1.0 museum 0.4 ↓
- 13| music 0.3 pavillion 1.0 ↓
- 14| cat 0.4 snail 0.3 pavillion 0.3 museum 1.0 montage 0.4 ↓
- 15| music 0.4 whisper 0.4 cat 0.3 power 0.4 insect 0.4 pattern 1.0 ↓
- 16| dream 0.3 spring 0.2 magnet 0.2 gourmet 1.0 torch 0.2 ↓
- 17| dream 0.2 power 0.3 god 0.3 magnet 0.2 saga 0.3 pattern 0.2 montage 1.0 ↓
- 18| dream 0.2 silk 0.3 torch 1.0 ↓
- 19| flag 1.0 ↓
- 20| heat 1.0 dream 0.3 whisper 0.2 stock 0.3 snail 0.3 gourmet 0.3 torch 0.3 ↓

【 23 】

- 1| gourmet 1.0 ↓
- 2| dream 0.4 center 0.2 whisper 0.3 spring 0.3 power 0.3 montage 1.0 torch 0.3 ↓
- 3| whisper 0.3 magnet 0.2 pavillion 0.2 torch 1.0 ↓
- 4| heat 0.3 music 0.3 power 0.3 shark 0.4 lecture 0.4 pavillion 0.3 pattern 0.2 gourmet 0.4 torch 0.2 flag 1.0 ↓
- 5| heat 1.0 music 0.4 center 0.2 god 0.3 insect 0.2 latitude 0.1 saga 0.1 pattern 0.3 torch 0.1 ↓
- 6| music 1.0 power 0.2 stock 0.4 silk 0.2 torch 0.3 ↓
- 7| dream 1.0 insect 0.2 lecture 0.3 museum 0.2 ↓
- 8| center 1.0 god 0.2 stock 0.2 snail 0.1 lecture 0.2 ↓
- 9| whisper 1.0 cat 0.2 latitude 0.2 lecture 0.2 ↓
- 10| cat 1.0 pattern 0.2 gourmet 0.3 ↓
- 11| heat 0.4 spring 1.0 snail 0.3 saga 0.3 museum 0.2 gourmet 0.4 ↓
- 12| power 1.0 silk 0.3 snail 0.3 lecture 0.3 flag 0.2 ↓
- 13| god 1.0 lecture 0.4 museum 0.2 montage 0.4 ↓
- 14| power 0.3 shark 1.0 museum 0.2 montage 0.4 torch 0.2 flag 0.4 ↓
- 15| music 0.2 dream 0.3 spring 0.3 god 0.2 stock 1.0 montage 0.4 ↓
- 16| heat 0.3 dream 0.1 cat 0.2 shark 0.3 insect 1.0 pattern 0.2 ↓
- 17| power 0.2 stock 0.2 silk 1.0 ↓
- 18| magnet 1.0 latitude 0.4 montage 0.3 flag 0.4 ↓
- 19| dream 0.1 stock 0.2 silk 0.3 magnet 0.3 snail 1.0 lecture 0.2 museum 0.1 ↓
- 20| god 0.1 stock 0.2 latitude 1.0 ↓

【 図 2 4 】

1 Clinton 2 AlGore 3 Hillary 4 Chelsea 5 Bush 6 Japan 7 Java 8 applet 9 JSP 10 apache
 11 servlet 12 Bluebooth 13 soccer 14 matrix 15 DNA 16 heat 17 music 18 dream 19 critic 20 whisper
 21 cat 22 spring 23 power 24 god 25 shark 26 stock 27 insect 28 silk 29 magnet 30 snail
 31 latitude 32 saga 33 lecture pavilion 35 museum 36 pattern 37 gourmet 38 montage 39 torch 40 flag

【 図 2 5 】

document 1: new keyword = Clinton
 document 1: new keyword = AlGore
 document 1: new keyword = Chelsea
 document 2: new keyword = Bush
 document 3: new keyword = Japan
 document 11: new keyword = Hillary

data1 で追加される新しいキーワード。
 これらがハッシュテーブルに追加される。

SUM-MEAN(M)_i (N = 6)
 [20.0 17.5 1.9 1.2 1.4 7.5]

data1 処理後のSUM-MEAN(M)_i

SUM(M)_i (N = 6)
 20.00 17.50 1.90 1.20 1.40 7.50
 16.25 1.65 1.05 1.25 5.00
 0.63 0.00 0.00 0.70
 0.36 0.06 0.60
 0.36 0.35
 6.25

data1 が処理された後のSUM(M)_i。
 対称行列なので半分だけ書き出した。

【 図 2 6 】

document 26: new keyword = Java
 document 26: new keyword = JSP
 document 26: new keyword = apache
 document 26: new keyword = servlet
 document 36: new keyword = applet

data2 で追加された新しいキーワード。
 これら5つがハッシュテーブルに追加され、
 合計で 11 個のキーワードとなる。

SUM-MEAN(M)_i (N = 11)
 [25.0 17.5 2.5 1.5 1.5 12.5 15.0 15.0 1.3 1.7 2.5]

data2 処理後のSUM-MEAN(M)_i

SUM(M)_i (N = 11)
 25.00 17.50 2.50 1.50 1.50 12.50 0.00 0.00 0.00 0.00 0.00
 16.25 1.65 1.05 1.25 5.00 0.00 0.00 0.00 0.00 0.00
 0.83 0.00 0.00 1.30 0.00 0.00 0.00 0.00 0.00
 0.45 0.06 0.90 0.00 0.00 0.00 0.00 0.00
 0.37 0.45 0.00 0.00 0.00 0.00 0.00
 11.25 0.00 0.00 0.00 0.00 0.00
 15.00 15.00 1.30 1.70 2.50
 15.00 1.30 1.70 2.50
 0.29 0.04 0.25
 0.51 0.15
 1.25

data2 が処理された後のSUM(M)_i。
 11x11 の対称行列となる。

【 図 2 8 】

特異値 (Largest 10)

0.3161665995009
 0.2324900894427
 0.0590495015116
 0.0589660332610
 0.0487117233259
 0.0404724327529
 0.0383981484518
 0.0367731637155
 0.0359210981307
 0.0358693365225

data2 で追加された後のSUM(M)_i。
 11x11 の対称行列となる。

【 図 2 7 】

SUM-MEAN(M)_i (N=40)
 [25.0 17.5 2.5 1.5 1.5 12.5 25.0 17.5 1.727 12.5 5.0 6.2 4.8 5.9 6.7 5.0 4.7 6.0
 5.7 6.4 5.0 5.3 6.8 6.7 6.6 5.0 5.4 5.5 6.8 6.7 5.4 6.4 3.5 5.0 6.5 5.6 6.0 5.2 4.4]

data7 処理後のSUM-MEAN(M)_i

【 図 2 9 】

特異値 (Largest 10)

0.3161665995009
0.2324900894427
0.0590495015116
0.0589660332610
0.0487117233259
0.0404724327529
0.0383981484518
0.0367731637155
0.0359210981307
0.0358693365225

【 図 3 1 】

特異値 (Largest 10)

0.1987144620508
0.0690170884491
0.0566614854733
0.0537575326572
0.0514825460617
0.0502895403078
0.0502170789069
0.0497783735464
0.0480041373927
0.0467372731528

【 図 3 0 】

AlGore
Clinton
Hillary
Chelsea
Bush
Japan

data1 と data2 の削除で同時に削除されるキーワード。

【 図 3 2 】

特異値 (Largest 10)

0.1987144620508
0.0690170884491
0.0566614854733
0.0537575326572
0.0514825460617
0.0502895403078
0.0502170789069
0.0497783735463
0.0480041373927
0.0467372731528

フロントページの続き

- (72)発明者 寒川 光
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 小林 メイ
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 青野 雅樹
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 竹内 広宜
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内

審査官 深津 始

- (56)参考文献 北研二、ほか、意味および感性に基づくマルチメディア・クロス・コンテンツ検索に関する研究
，財団法人電気通信普及財団 調査研究報告書，第19号，411 - 418頁

(58)調査した分野(Int.Cl.，DB名)

G06F 17/30