

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4967133号  
(P4967133)

(45) 発行日 平成24年7月4日(2012.7.4)

(24) 登録日 平成24年4月13日(2012.4.13)

(51) Int.Cl. F I  
**G06F 17/30 (2006.01)**  
 G06F 17/30 350C  
 G06F 17/30 340Z  
 G06F 17/30 170A

請求項の数 4 (全 31 頁)

(21) 出願番号 特願2007-85469(P2007-85469)  
 (22) 出願日 平成19年3月28日(2007.3.28)  
 (65) 公開番号 特開2008-243024(P2008-243024A)  
 (43) 公開日 平成20年10月9日(2008.10.9)  
 審査請求日 平成21年9月25日(2009.9.25)

(73) 特許権者 504174135  
 国立大学法人九州工業大学  
 福岡県北九州市戸畑区仙水町1番1号  
 (74) 代理人 100099634  
 弁理士 平井 安雄  
 (72) 発明者 野村 浩郷  
 福岡県飯塚市大字川津680-4 九州工業大学内  
 審査官 鈴木 和樹

最終頁に続く

(54) 【発明の名称】 情報取得装置、そのプログラム及び方法

(57) 【特許請求の範囲】

【請求項1】

少なくとも1つの検索情報を取得している情報取得装置において、  
 重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成手段と、  
 全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算手段と、  
 前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算手段と、

前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成手段と、

前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成手段と、  
 前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定手段と、

検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析手段と、  
各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出手段と、

前記抽出された名詞の単体の名詞間の類似度  $S_1$  及び名詞集合の類似度  $S_2$  を計算する名詞集合間類似度比較計算手段と、

抽出された名詞の表示の一致する割合の類似度  $S_3$  を計算する名詞表示一致割合計算手段と、

前記類似度  $S_2$  に類似度  $S_3$  を加えて文類似度  $S$  を計算する文類似度計算手段と、  
検索情報の文タイプによる選定を行う文タイプ選定手段と、

10

20

前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及びノ又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力手段と、

前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正手段とを備え、

前記検索順位決定手段が、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力手段が、前記検索順位決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力することを特徴とする情報取得装置。

10

【請求項2】

前記請求項1に記載された情報取得装置において、  
前記特徴ベクトル作成手段は、  
検索情報の文の形態素解析を行う形態素解析手段と、  
情報毎に単語とその単語の出現回数TFを計算するTF計算手段と、  
全単語について文書頻度DF及びそのIDFを計算するIDF計算手段と、  
各情報の各単語についてTF-IDF法を用いて単語重み付けを計算する単語重み計算手段と、

前記単語重み付けから各文書の特徴ベクトルを作成する特徴ベクトル作成手段とを備えることを特徴とする情報取得装置。

20

【請求項3】

少なくとも1つの検索情報を取得している情報取得装置としてコンピュータを機能させる情報取得プログラムにおいて、

重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成手段、  
全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算手段、  
前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算手段、

前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成手段、

前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成手段、  
前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定手段、

30

検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析手段、  
各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出手段、  
前記抽出された名詞の単体の名詞間の類似度 $S_1$ 及び名詞集合の類似度 $S_2$ を計算する名詞集合間類似度比較計算手段、

抽出された名詞の表示の一致する割合の類似度 $S_3$ を計算する名詞表示一致割合計算手段、

前記類似度 $S_2$ に類似度 $S_3$ を加えて文類似度 $S$ を計算する文類似度計算手段、  
検索情報の文タイプによる選定を行う文タイプ選定手段、

40

前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及びノ又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力手段、

前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正手段としてコンピュータを機能させ、

前記検索順位決定手段が、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力手段が、前記検索順位決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力することを特徴とする情報取得プログラム。

50

## 【請求項4】

少なくとも1つの検索情報を取得している情報取得装置のコンピュータが、  
 重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成ステップと、  
 全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算ステップと、

前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算ステップと、

前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成ステップと、

前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成ステップと

10

、  
 前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定ステップと、

検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析ステップと、

各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出ステップと、

前記抽出された名詞の単体の名詞間の類似度 $S_1$ 及び名詞集合の類似度 $S_2$ を計算する名詞集合間類似度比較計算ステップと、

抽出された名詞の表示の一致する割合の類似度 $S_3$ を計算する名詞表示一致割合計算ステップと、

20

前記類似度 $S_2$ に類似度 $S_3$ を加えて文類似度 $S$ を計算する文類似度計算ステップと、

検索情報の文タイプによる選定を行う文タイプ選定ステップと、

前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及び/又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力ステップと、

前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正ステップとを実行し、

前記検索順位決定ステップが、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力ステップが、前記検索順位決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力することを特徴とする情報取得方法。

30

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

関連内容の情報の検索とそれらを集約する情報取得装置に関する。

## 【背景技術】

## 【0002】

Web検索を含めた多くの情報検索システムでは、キーワード検索を元にしており、ユーザは検索キーワードをシステムに与えることで検索結果を得る。単純な単語のマッチングのみを条件として検索を行うために、検索結果が膨大になることが多く、またノイズも多い。現状ではユーザは膨大な量で、しかも玉石混淆の検索結果から要求に合致したテキストを探さなければならない。そのため、ユーザが検索結果から合致した情報を得るためには多大な労力を必要とする。また、キーワード検索の途中で関連する情報を発見したいという状況も頻繁に発生することがある。さらに、近年情報機器の普及により様々な情報が電子化されており、大量の情報がいつでも閲覧できるようになった現在、その中から必要な情報を効率よく選ぶ作業は、情報の電子化が急速に進んでいる中、困難になっている。

40

そこで、複数のデータベースを検索して所望の情報を取得し、その情報をユーザの望む形式に編集・加工する情報編集・加工方法が、特開平9-185632号公報に開示され

50

ている。

【0003】

背景技術の情報編集・加工方法は、遠隔にある少なくとも1つのデータベースが保有していると推定される目的情報の検索指示、検索した情報の編集加工指示、編集加工した情報の出力形態決定指示とを受け付ける第1の過程と、前記検索指示に基づいて前記少なくとも1つのデータベースを検索するコマンドを生成し、前記データベースを検索する第2の過程と、前記第2の過程により取得した少なくとも1つの目的情報に対して前記第1の過程の指示に従って、編集・加工を施す第3の過程と、前記第3の過程によって、編集・加工された目的情報に対して前記情報の出力形態決定指示に従って、所定の出力形態に変換し、視覚、聴覚または他の感覚に捉え得る方法によって出力する第4の過程よりなることを特徴とする。前記第3の過程は、前記第4の過程において表示する検索結果をユーザが指定した出力順にソートする。前記出力順として、情報の関連度順、情報発生時間順、または検索順のいずれかを用いるものである。

10

【特許文献1】特開平9-185632号公報

【発明の開示】

【発明が解決しようとする課題】

【0004】

上述のように背景技術の情報検索・編集方法及び装置によれば、得られた情報間の関連を意識するので、関連のある情報同士を近接してユーザに提供することができる。また、複数のデータベースから得られた情報をユニフォームに扱うため、異なるデータベースから得られた情報の提供を時間順であっても関連度順であっても適切に行うことができる。

20

【0005】

しかしながら、提供された情報がトピックスについては条件を満たしていても、情報内容についての質あるいは量が、ユーザの要求を必ずしも十分に満たしているとは限らない場合がある。そのため、ユーザの希望する情報が不足している場合に改めて情報検索を行わなければならないという課題を有する。

【0006】

また、あるトピックスに対して複数の観点から作成された文書等の情報に関して、それらの情報を比較し、理解を深めるということも可能になってきてはいるが、その作業もまた情報量の増大につれて困難になってきている。しかも、関連情報に関しては情報の重複個所を何度も繰り返し取得するためにユーザの負担が大きくなるという課題もある。

30

【0007】

本発明は、前記課題を解決するためになされたものであり、続報情報を発見する検索装置において関連した情報を比較、整理して効率よく必要な情報を取得し、利用者の要望に沿った形式で出力を行うことができる情報取得装置の提供を目的とする。

【課題を解決するための手段】

【0008】

本発明に係る情報取得装置は、少なくとも1つの検索情報を取得している情報取得装置において、重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成手段と、全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算手段と、前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算手段と、前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成手段と、前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成手段と、前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定手段と、検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析手段と、各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出手段と、前記抽出された名詞の単体の名詞間の類似度 $S_1$ 及び名詞集合の類似度 $S_2$ を計算する名詞集合間類似度比較計算手段と、抽出された名詞の表示の一致する割合の類似度 $S_3$ を計算する名詞表示一致割合計算手段と、前記類似度 $S_2$ に類似度 $S_3$ を加えて文類似度 $S$ を計算する文類似度計算手段と、検

40

50

索情報の文タイプによる選定を行う文タイプ選定手段と、前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及び/又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力手段と、前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正手段とを備え、前記検索順位決定手段が、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力手段が、前記検索順位決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力するものである。ここで、「情報」には、例えば、文、文書、記事、画像、音声等を含む。また、「文タイプ」とは、例えば、重複個所、固有個所、補足説明等である。さらに、要旨、予定、理由、分析、補足説明、様態・伝聞、比況・推量等を含むものとする。

10

## 【0009】

これにより、重み付けされた検索情報の特徴ベクトルを作成し、全検索情報の特徴ベクトルの組み合わせの類似度を計算し、前記情報間類似度計算によって得られた数値の類似度行列を計算し、前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求め、前記検索情報の問い合わせ内容の質問ベクトルを作成し、前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定し、前記検索順位決定された関連情報の内容を統合して出力するので、すでに検索された情報から、その情報と類似度が高い情報を優先的に選択し、類似度に応じた確率で読み進めていくという仮想的なユーザを考えると、無限時間後に定常状態になった時点で、どの情報に行きつくかという確率に相当することを判断しながら関連情報を検索し、ユーザにとって必要な情報を取得することができる。

20

また、検索された情報の文中に含まれる品詞の係り受け関係を解析し、各文中の動詞を含む文節に係る文節中の名詞を抽出し、前記抽出された名詞の単体の名詞間の類似度  $S_1$  及び名詞集合の類似度  $S_2$  を計算し、抽出された名詞の表示の一致する割合の類似度  $S_3$  を計算し、前記類似度  $S_2$  に類似度  $S_3$  を加えて文類似度  $S$  を計算し、検索情報の文タイプによる選定を行い、前記文類似度計算及び文タイプ選定された関連情報の内容を統合して出力するので、検索された関連情報について、重複個所、固有個所、補足説明等の情報内容を整理した状態で、関連情報を取得することができる。また、膨大な量の関連情報に含まれる重複情報が何度も繰り返し表示されることによるユーザの負担を軽減でき、整理された固有個所や補足説明の情報を効率よく利用者の要望に沿った形式により取得することができる。

30

さらに、検索順位決定によりスコアリングされた検索結果を開示し、前記検索結果の適否を入力し、特徴及び質問ベクトルを修正して、検索情報の内容を出力するので、特徴ベクトル及び質問ベクトルをユーザが適合していると判断した情報に近づけ、不適合であると判断した情報から遠ざけるように特徴ベクトル及び質問ベクトルを生成していくことができる。これを繰り返し適用することにより、確実にユーザの望む検索結果を得ることができる。

## 【0014】

40

本発明に係る情報取得装置は必要に応じて、前記特徴ベクトル作成手段は、検索情報の文の形態素解析を行う形態素解析手段と、情報毎に単語とその単語の出現回数  $TF$  を計算する  $TF$  計算手段と、全単語について文書頻度  $DF$  及びその  $IDF$  を計算する  $IDF$  計算手段と、各情報の各単語について  $TF - IDF$  法を用いて単語重み付けを計算する単語重み計算手段と、前記単語重み付けから各文書の特徴ベクトルを作成する特徴ベクトル作成手段とを備えるものである。

## 【0015】

これにより、重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成手段は、検索情報の文の形態素解析を行い、情報毎に単語とその単語の出現回数  $TF$  を計算し、全単語について文書頻度  $DF$  及びその  $IDF$  を計算し、各情報の各単語について  $TF -$

50

IDF法を用いて単語重み付けを計算し、前記単語重み付けから各文書の特徴ベクトルを作成するので、特徴ベクトル作成手段情報検索において目的の情報を探すために、関連情報についての重要度を判断するためのひとつの指標とすることができ、文書と単語の関連性の数値演算を行い、その値の高いものを特徴ベクトルに反映することにより、よりユーザにとって重要な関連情報を取得することができる。

【0016】

本発明に係る情報取得プログラムは、少なくとも1つの検索情報を取得している情報取得装置としてコンピュータを機能させる情報取得プログラムにおいて、重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成手段、全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算手段、前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算手段、前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成手段、前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成手段、前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定手段、検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析手段、各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出手段、前記抽出された名詞の単体の名詞間の類似度 $S_1$ 及び名詞集合の類似度 $S_2$ を計算する名詞集合間類似度比較計算手段、抽出された名詞の表示の一致する割合の類似度 $S_3$ を計算する名詞表示一致割合計算手段、前記類似度 $S_2$ に類似度 $S_3$ を加えて文類似度 $S$ を計算する文類似度計算手段、検索情報の文タイプによる選定を行う文タイプ選定手段、前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及び/又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力手段、前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正手段としてコンピュータを機能させ、前記検索順位決定手段が、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力手段が、前記検索順位決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力するものである。

【0018】

本発明に係る情報取得方法は、少なくとも1つの検索情報を取得している情報取得装置のコンピュータが、重み付けされた検索情報の特徴ベクトルを作成する特徴ベクトル作成ステップと、全検索情報の特徴ベクトルの組み合わせの類似度を計算する情報間類似度計算ステップと、前記情報間類似度計算によって得られた数値の類似度行列を計算する類似度行列計算ステップと、前記類似度計算結果を数値解析し、特徴ベクトルの最大固有値の固有ベクトルを求める固有ベクトル作成ステップと、前記検索情報の問い合わせ内容の質問ベクトルを作成する質問ベクトル作成ステップと、前記特徴ベクトルと質問ベクトルの余弦の計算値に固有ベクトルの数値を乗じて求められる関連情報の検索順位を決定する検索順位決定ステップと、検索された情報の文中に含まれる品詞の係り受け関係を解析する係り受け解析ステップと、各文中の動詞を含む文節に係る文節中の名詞を抽出する名詞抽出ステップと、前記抽出された名詞の単体の名詞間の類似度 $S_1$ 及び名詞集合の類似度 $S_2$ を計算する名詞集合間類似度比較計算ステップと、抽出された名詞の表示の一致する割合の類似度 $S_3$ を計算する名詞表示一致割合計算ステップと、前記類似度 $S_2$ に類似度 $S_3$ を加えて文類似度 $S$ を計算する文類似度計算ステップと、検索情報の文タイプによる選定を行う文タイプ選定ステップと、前記文類似度計算及び文タイプ選定された関連情報の内容を統合したもの、並びに、前記検索順位決定手段によりスコアリングされた検索結果を出力すると共に、前記検索結果の適否及び/又はパラメータの重み付けの度合いを入力するための入力フォームを出力する出力ステップと、前記入力フォームに入力された情報に基づいて、前記特徴ベクトル及び質問ベクトルを修正する修正ステップとを実行し、前記検索順位決定ステップが、前記修正手段にて修正された前記特徴ベクトル及び質問ベクトルに基づいて、再度前記関連情報の検索順位を決定し、前記出力ステップが、前記検索順位

10

20

30

40

50

決定手段によりスコアリングされた検索結果を、前記検索情報間の経時的な関連性を含めて出力するものである。

【発明を実施するための最良の形態】

【0020】

ここで、本発明は多くの異なる形態で実施可能である。したがって、下記の実施形態の記載内容のみで解釈すべきではない。実施形態では、主に装置について説明するが、所謂当業者であれば明らかな通り、本発明は、コンピュータで使用可能なプログラムとしても実施できる。また、本発明では、ハードウェア、ソフトウェア、または、ソフトウェア及びハードウェアの実施形態で実施可能である。プログラムは、ハードディスク、CD ROM、DVD-ROM、光記憶装置または磁気記憶装置等の任意のコンピュータ可読媒体に記録できる。さらに、プログラムはネットワークを介した他のコンピュータに記録することが出来る。

10

【0021】

[1. ハードウェア構成]

図1に本発明の実施形態における情報取得装置のハードウェア構成図を示す。コンピュータ1は、例えば、CPU(Central Processing Unit)2、メインメモリ3、HDD(Hard Disk Drive)4、ビデオカード5、マウス6、キーボード7、光学ディスク8等を含む。なお、必要に応じて、データベース等を接続することもできる。

【0022】

20

[2. ブロック構成]

図2に本発明の実施形態に係る情報取得装置のブロック構成図を示す。本発明は、主として、入力部10、続報情報検索部20、情報内容統合部30、出力部40を含む。ここで、続報情報検索部20は、ベクトル作成部21、検索順位決定部22、適合・非適合判定部23を含む。さらに、ベクトル作成部21は、形態素解析部211、TF計算部212、IDF計算部213、単語重み計算部214、特徴ベクトル作成部215、情報間類似度計算部216、類似度行列計算部217、固有ベクトル作成部218、質問ベクトル作成部219を含む。また、情報内容統合部30は、係り受け解析部31、名詞抽出部32、名詞集合間類似度比較計算部33、名詞表示一致割合計算部34、文類似度計算部35、文タイプ選定部36、要約文作成部37、記事集約部38を含む。

30

まず、入力部10により入力された記事データは、続報情報検索部20におけるベクトル作成部21送られて処理されることになる。ここで、本発明の実施形態の例として、検索対象に新聞記事を一例に挙げて、各構成の内容について以下に詳説する。

【0023】

[2.1 続報情報検索]

[2.1.1 ベクトル空間モデル]

ベクトル作成部21において、まず、記事データは形態素解析部211、TF計算部212、IDF計算部213、単語重み計算部214、特徴ベクトル作成部215で処理される。ここで、ベクトル空間モデル(vector-space model)は検索対象となる個々のデータの性質を表現するための特徴量として、多次元ベクトルを個々のデータに対応づける。この間に類似度(Similarity)を定義することにより、問い合わせ(質問)と類似したものを探し出す方法である。いま、検索対象の特徴としてn個の属性が備わっており、i番目の属性を $w_i$ とする。そしてj番目のデータに(数式1)のベクトルを対応させることを考える。これらのベクトルが線形独立であれば、n次元のベクトル空間が定義される。このように定義されたベクトル空間において、j番目データの特徴ベクトルは

40

【0024】

【数 1】

$$\vec{D}_j = (w_{dj1}, w_{dj2}, w_{dj3}, \Lambda, w_{djn})$$

のように表すことができる ( $d_{ji}$  はの  $w_i$  に対する値)。

ベクトル空間モデルにおける検索システムへの問い合わせ (質問) もベクトルで表される。n 次元のベクトル空間に対するその質問ベクトルは

【0 0 2 5】

【数 2】

$$\vec{Q} = (w_{q1}, w_{q2}, w_{q3}, \Lambda, w_{qn})$$

のように表すことができる ( $q_i$  は質問ベクトルの  $w_i$  に対する値)。

【0 0 2 6】

検索は検索対象の (数式 1) の特徴ベクトルと (数式 2) の質問ベクトルの類似度を計算することにより行われる。この特徴ベクトル (feature vector) を得る方法は、検索の目的や、その対象であるデータの種類などによって異なる。例えば検索対象が文献データならば単語の出現頻度を基にベクトルの各要素の重み付けを行い、画像であれば画素ごとの濃淡や色のデータなどを用いることができる。なお、新聞記事では、1 記事を 1 つのベクトルに割り当て、記事中の単語の TF - IDF をベクトルの重みづけに利用している。

【0 0 2 7】

[ 2 . 1 . 2 単語の重みづけ ]

情報検索において目的の文書を探すために、文書と単語の関連性の数値演算を行い、その値の高いものを候補とする。そこで用いられる評価値は文書中には重要な単語がどれくらい多く含まれているかを表している。文書中の単語がどの程度重要であるか重み付けに用いられているのが以下に述べる TF - IDF 法である。この手法は次の 2 つのキーワードの性質に注目している。

( 1 ) 文書に数多く、高い頻度で現れる単語は重要である

( 2 ) 少ない数の文書にしか現れない単語は重要である

単語出現頻度 (Term Frequency : TF) 単語  $t$  が文書  $d$  に高い頻度で現れるなら、 $t$  は  $d$  を良く特徴付ける。この考えによる尺度が単語出現頻度、 $tf$  (Term Frequency) である。ある文書  $d$  における単語  $t$  の出現頻度  $tf(d, t)$  は次式で定義され、TF 計算部 2 1 2 において計算が行われる。

【0 0 2 8】

【数 3】

$$tf(d, t) = freq(d, t)$$

$freq(d, t)$  : 文書  $d$  における単語  $t$  の出現頻度。

【0 0 2 9】

文書出現頻度 (Document Frequency :  $df$ )  $tf$  が大きいというのは重要な性質だが、それだけでは十分に文書の特徴付けることはできない。例えば、日本語文書で「は」という助詞はどんな文書でも高い頻度で現れるが、特定の文書の特徴付けにないことは明白である。そこで、単語  $t$  が検索対象となる文書集合のうちの少数の文書にしか現れないという性質が重要である。単語  $t$  の出現する文書数を文書出現頻度  $df$  (document frequency) は、次式で定義される。

【0 0 3 0】

10

20

30

40

50



【数4】

$$df(t) = dfreq(t)$$

$dfreq(t)$  : 単語  $t$  が出現する文書数

$df$  が小さいことが単語  $t$  の文書を特徴付ける能力が高いことを表すので、実際にはこの逆数を  $\log$  と文書集合中の文書総数  $N$  により正規化した  $idf$  (inverse document frequency) を用いる。

【0031】

【数5】

10

$$idf(t) = \log \frac{N}{dfreq(t)} + 1$$

$N$  : 文書の数

$freq(d, t)$  : 文書  $d$  における単語  $t$  の出現頻度

【0032】

なお、IDF 計算部 213 では、まず、 $df$  を求めた後に、 $idf$  を計算することになる。TF-IDF による重み付け単語  $t$  について、その単語が文書内に出てくる回数とそれが全文書内に占める割合の積を計算することで、その単語の重要性和、その出現頻度によって文書の重要性を表すことが目的である。単語  $t$  が  $tf$  と  $idf$  の両者の性質を併せ持つ、すなわち  $tf$  が大きく、 $df$  が小さいならば、単語  $t$  は文書  $d$  を真に特徴付けるといえる。この考え方を数値の尺度として表現したのが TF-IDF による重み付けである。文書  $d$  におけるキーワード  $t$  の重み  $w(t, d)$  は次のように定義され、単語重み計算部 214 で計算される。

20

【0033】

【数6】

$$w(t, d) = tf(d, t) \cdot idf(t)$$

30

$dfreq(t)$  : 単語  $t$  が出現する文書数

そして、特徴ベクトル作成部 215 において、これらの求められた数値を利用して特徴ベクトルを作成する。

【0034】

[2.1.3 類似度]

特徴ベクトル作成部 215、質問ベクトル作成部 219 で処理されたデータは検索順位決定部 22 に送られる。ここで、ベクトル空間モデルにおいて、検索を行うためにはベクトル間の類似度を定義しなければならない。類似度の尺度としては様々なものがあるが、ここではベクトル間の余弦を用いる。

類似度として2つのベクトル間の余弦の値を利用する方法である。特徴ベクトル  $D$  と質問ベクトル  $Q$  の類似度  $sim(D, Q)$  は以下ようになる。

40

【0035】

【数7】

$$sim(D, Q) = \frac{\overrightarrow{D} \cdot \overrightarrow{Q}}{\|\overrightarrow{D}\| \|\overrightarrow{Q}\|} = \frac{\sum_{i=1}^n (w_i q_i)}{\sqrt{\sum_{i=1}^n (w_i)^2} \sqrt{\sum_{i=1}^n (q_i)^2}}$$

50

$\text{sim}(D, Q)$ の値は0以上1以下であり、1に近づくほど類似度が高くなる。検索順位決定部22では、余弦、いわゆるコサイン相関値を用いた類似度評価を行う。

【0036】

[2.1.4 ベクトル空間モデルにおける関連性フィードバック]

検索順位決定部22において検索結果が得られた場合、出力部40で処理される。一度の検索で最終的な結果を得るのではなく、結果に対するユーザのフィードバックを元に新たな質問を生成し、繰り返し検索を行い、徐々に検索結果をユーザの求める結果に近づけていくフィードバック検索を行う。つまり、改めてユーザの検索結果に対する適否データを入力部10において入力する。

【0037】

ベクトル空間モデルにおける、関連性フィードバック(Relevance Feedback)では、装置への質問式も質問ベクトル作成部219により作成された多次元ベクトルで表現される。質問の結果については、質問ベクトル作成部219で作成された質問ベクトルと特徴ベクトル作成部215で作成されたデータの特徴ベクトルの類似度を計算した結果のデータの集合として求める。この類似度が高いデータほど、質問の答えとしてふさわしいものであると考え、検索結果に含まれるデータに、それがどれだけ質問に適合していたかという順位をつけてユーザに提示する。ユーザは提示された検索結果からフィードバックを返す。

【0038】

[2.1.5 ユーザからのフィードバック]

入力された検索結果の適否データは、適合・非適合判定部23に送られ、質問ベクトル作成部219及び単語重み計算部214に送られる。

具体的なユーザからのフィードバックとして、最も多いのは結果の正例(positive example)、負例(negative example)の提示である。また、正例のみをフィードバックするもの、それぞれの妥当性の度合いをランクづけてフィードバックするものなど、様々なものがある。また、ユーザからのフィードバックを検索に反映させる方法としては、大きく以下の二つに分けることができる。

(1) 質問ベクトル修正(Query Vector Movement)は、検索質問のベクトルを修正・変換して、正例の特徴に近づけ、負例から遠ざける。

(2) 再重みづけ(Feature Re-weighting)は、特徴ベクトルに対応するための重みをユーザのニーズにあわせて調節する。すなわち、正例を検索するのに好都合な次元を強調し、負な例のものとの影響を減らすように重みづけを動的に変更する。本発明では、この両方のフィードバックを利用する。

【0039】

[2.1.6 質問ベクトルの修正]

ユーザのフィードバックした結果から、検索結果をユーザの求めるものに近づける手法として、質問ベクトルをユーザが適合していると判断した記事に近づけ、不適合であると判断した記事から遠ざけるように質問ベクトルを生成していく。これを繰り返し適用することにより、徐々にユーザの望む検索結果を得ることができる。このために良く利用されるのはRocchioフィードバック手法であり、Rocchioの式は以下のように与えられる。

【0040】

【数8】

$$Q_{i+1} = Q_i + \alpha \frac{1}{Rn} \sum_{j \in R} \vec{D}_j - \beta \frac{1}{Nn} \sum_{j \in N} \vec{D}_j$$

10

20

30

40

## 【 0 0 4 1 】

$Q_i$  は前回の検索時に用いられた質問ベクトルであり、 $Q_{i+1}$  が新しく生成された質問ベクトルである。 $R$  は適合だと判断された文書  $D_j$  に対する特徴ベクトルであり、 $N$  は不適合であると判断された文書に対する特徴ベクトルである。 $R_n$ 、 $N_n$  はそれぞれ適合文献数、不適合文献数である。 $\alpha$ 、 $\beta$  はそれぞれ適合文献、不適合文献に対する変数であり、 $\alpha$  の値が高いと適合文献による変更が重要視され、 $\beta$  の値が高いと不適合文献による変更が重要視される。適合フィードバックの結果として、問合せ位置は  $Q_i$  から  $Q_{i+1}$  に移動するととらえることができる。ここで、特徴ベクトル作成部 2 1 5 におけるデータは、質問ベクトル作成部においても処理される。

## 【 0 0 4 2 】

## [ 2 . 1 . 7 状態遷移確率を考慮に入れた重要度評価 ]

情報間類似度計算部 2 1 6、類似度行列計算部 2 1 7、固有ベクトル作成部 2 1 8 は、特徴ベクトル作成部 2 1 5 からのデータを以下の内容で処理する。

図 3 は本発明の実施形態に係る情報取得装置の記事間の類似度による記事の重要度評価の説明図である。Page Rank は、www 上のハイパーリンクによって結ばれた Web ページ群において、「多くの良質なページからリンクされているページは、やはり良質なページである」、という再帰的な関係をもとに、Web ページの重要度を評価する理論、およびそれによって求められるページの重要度である。Page Rank を用いることで、ハイパーリンク構造のような相互参照関係があるときに、どのページがもっとも重要であるかを定量的に求めることができる。

## 【 0 0 4 3 】

図 3 ( a ) は Page Rank の概念図を示す。この図を例に基本的な Page Rank の計算方法を説明すると、まず全ての Web ページはそれぞれ Page Rank の値を持っている。そしてこの値はそのページがリンクしている先のページへ均等に分配されることになる。図 3 ( a ) を例にとると、図中にある 1 0 0 の値を持ったページは 2 つのページへのリンクを持っているので、このページの持つ 1 0 0 の値は 2 つに分割されてリンク先へ与えられる。つまり、リンク先のページはそれぞれ 5 0 ずつの値を得ることになる。

## 【 0 0 4 4 】

## [ 2 . 1 . 8 記事間の類似度による記事の重要度評価 ]

Page Rank がページ間のリンクの重みを平等に扱っているのに対し、本発明では各記事との類似度で重み付けを行う。これによって新聞記事群を関連度の強さに応じたリンクによって結ばれたグラフ構造と考える。そのなかから、より関連性が高いとしてリンクされている記事を、Page Rank 同様、遷移確率の最大固有値における固有ベクトルを算出することで求める。図 3 ( b ) は新聞記事間の類似度を示すものであり、その算出方法を以下に説示する。

まず、記事数を  $N$  とするとき、情報間類似度計算部 2 1 6 が、 $N \times N$  の  $N$  次正方行列、要素に各記事間の類似度をそれぞれ計算し、類似度行列計算部 2 1 7 が類似度行列を作成する。図 3 ( b ) について、類似度行列を求めた結果である行列  $A$  を以下に示す。

## 【 0 0 4 5 】

## 【数 9】

$$A = \begin{pmatrix} 0 & 4 & 1 & 3 \\ 4 & 0 & 2 & 5 \\ 1 & 2 & 0 & 6 \\ 3 & 5 & 6 & 0 \end{pmatrix}$$

次に、各記事、すなわち各列について合計が 1 になるように正規化し行列  $A$  を状態遷移

10

20

30

40

50

確率行列 M とする。図 3 ( C ) は、新聞記事間の類似度による重み付けを行った遷移確率を示す。このときの記事間の関係は図 3 ( C ) のように示される。

【 0 0 4 6 】

【 数 1 0 】

$$M = \begin{pmatrix} 0 & 4/1 & 1 & 1/9 & 3/1 & 4 \\ 4/8 & 0 & 2/9 & 5/1 & 4 \\ 1/8 & 2/1 & 1 & 0 & 6/1 & 4 \\ 3/8 & 5/1 & 1 & 6/9 & 0 \end{pmatrix}$$

10

行列 M の状態遷移確率行列から、固有ベクトル作成部 2 1 8 が、最大固有値の固有ベクトルを計算した結果を図 4 に示す。

【 0 0 4 7 】

図 4 は本発明の実施形態に係る情報取得装置の記事の重要度計算例である。図 4 は、より多くの記事から高い重みで参照されている記事ほどスコアが高くなっていることを示している。この図 4 のスコアは、現在見ている記事から、その記事と類似度が高い記事を優先的に選択し、類似度に応じた確率で読み進めていくという仮想的なユーザを考えると、無限時間後に定常状態になった時点で、どの記事に行きつくかという確率に相当する。すなわち、その記事が類似性があるとしてユーザが興味を持ち、辿り着きやすいかというスコアであり、また、記事群の中でどの記事が多くの記事から類似性を持っているとして高い重みでリンクされているか、というのを示すスコアであるともいえる。

20

【 0 0 4 8 】

以上によって求められたスコアを、ベクトル空間モデル上の類似度を計算したスコアに併用することにより、質問ベクトルとの類似度でユーザの興味を考慮に入れつつ、そのなかで代表らしい記事を結果として示すことができる。これによりユーザが検索結果の判断に用いるのに適している記事を得て、効率良くフィードバック検索を行おうとするものである。

30

【 0 0 4 9 】

[ 2 . 2 情報内容統合システム ]

続報情報検索部 2 0 で得られた続報記事の情報データは、情報内容統合部 3 0 において情報データの整理・分類処理される。その際に行われる重複箇所、固有箇所、補足説明の各カテゴリの設定、及びカテゴリ分けを行う類似度、文タイプによる判定について以下に説示する。

【 0 0 5 0 】

[ 2 . 2 . 1 カテゴリ設定及び分類 ]

内容統合において、複数新聞記事を文カテゴリに分類し、それらの組み合わせにより、利用者の要望に沿った形式の出力を目指す。よって、その際の各カテゴリは、ユーザの情報取得の選択肢を広げ、複数新聞記事を比較する際の利点に沿ったものでなければならない。そこで、各記事に共通の箇所である重複箇所、各記事に固有の箇所である固有箇所と、記事中における補足的な内容である補足説明という合計 3 つのカテゴリを設定する。なお、重複箇所中の文の対応の定義として、一方の文に比較対象の文の話題が、完全にまたは部分的に含まれていることとする。

40

対象記事を「重複箇所」、「固有箇所」、「補足説明」の 3 つのカテゴリに分類するために、その判定基準として、「文単位の類似度」、「文タイプ」の 2 つを用いる。

【 0 0 5 1 】

「文単位の類似度」では、記事データが、係り受け解析部 3 1 で解析処理される。そして、名詞抽出部 3 2、名詞集合間類似度比較計算部 3 3 により各文中の動詞をキーとした

50

名詞集合中の名詞単体の概念間の距離と表記を利用して求めた値と、名詞表示一致割合計算部34により求めたそれらの結果を利用した名詞単語中の表記が同じ名詞の割合の合計を、文類似度計算部35により算出された結果の値とする。

また「文タイプ」では、文タイプ選定部36が、各文に対して文のタイプ付けを行う。以下に、それぞれの判定に関する詳細な説明を述べる。

#### 【0052】

##### [2.2.2 複数新聞記事間における文単位の類似度]

重複箇所、固有箇所の選定の一基準として、文単位の類似度を採用している。以下に類似度の算出方法について述べる。一般に、文の類似度の指標には、構文構造の類似度と意味的な類似度が考えられる。類似文検索では、構文構造の類似度を求めるために「動詞への係り受け」を使用する。また、意味的な類似度を求めるために「動詞に直接係る文節中の名詞の意味属性」、「名詞表記の一致の割合」を利用する。類似文の検索は、次の4つのステップで行われる。

(1) 動詞を含む文節に係る文節中の名詞の検出

(2) (1)で抽出した動詞をキーとする名詞集合毎の類似度の比較

(3) (2)の結果を利用した名詞表記の一致の割合

(4) (2)と(3)の結果を利用した類似度の算出

#### 【0053】

##### [2.2.2 動詞を含む文節に係る文節中の名詞の抽出]

図5は、本発明の実施形態に係る情報取得装置の動詞を含む文節に係る文節中の名詞の抽出の例である。図5の場合には網掛け部分の名詞A, B, C, D, G, H, I, 7, 8, 9を抽出する。かかる処理は、文中に含まれる動詞に関する係り受けを利用することから、例えば日本語係り受け解析器cabochaを用いて行うことができる。

#### 【0054】

##### [2.2.3 動詞を含む文節に係る文節中の名詞の概念関係を利用した比較]

前述の日本語形態素解析器cabochaにより部分的な重複を文間の類似度の情報に入れるため、各文中の動詞を含む文節に係る文節中の名詞を抽出し、各動詞に対する名詞集合を作成する。その際の集合中の各名詞間の類似度は、表記が異なるものはEDR電子化辞書により、概念間の距離からその値を求め、表記が同一のものはその値を最大値にする。そして、各名詞単体同士の類似度から名詞集合同士の類似度を算出し、その中で最も類似度が高い値をとる。名詞の類似度を測る方法としては、意味属性体系上での共通親属性の位置や、両意味属性間のパスの長さから類似度を求める方法が考えられる。しかし、一般に名詞には複数の意味属性を割り当てることができる。そのため、名詞の類似度を求めるために、その名詞がどの意味属性の名詞として使われているのかを、文脈情報などから一意に決定しなければならない。本発明においては、この多義性の問題には立ち入らずに、「EDR電子化辞書を用いた単語類似度計算法」[参考文献：崔ら、情報処理学会報告NL-93-1, pp1-6]で提案されている手法である名詞に割り当てられた複数の意味属性から総合的に名詞の類似度を求める。また、動詞に係る文節中の名詞を類似度の指標として特に取り上げているのは、一文に含まれる話題の数の違いを考慮したことによる。

#### 【0055】

##### [2.2.4 EDR電子化辞書について]

図6は、本発明の実施形態に係る情報取得装置のEDR辞書の構造図である。EDR電子化辞書は、コンピュータによる先進的な言語処理のために開発され、単語辞書などのいくつかの大規模な個別辞書から構成されている。辞書は、単語辞書中で定義した概念の類義を記述する概念体系(シソーラス)、辞書記述の典拠としてのコーパスDB(例文集)を統合した日本語と英語の語彙知識総目録と呼ぶにふさわしい機械処理用の電子化辞書である。言語学的偏向を極力排除し、各種応用へのチューンアップの容易さを保持することを開発方針として採用してあるものである。EDR電子化辞書は単語辞書、対訳辞書、概念辞書、共起辞書、専門用語辞書とEDRコーパスから構成されている。

## 【 0 0 5 6 】

本発明では、名詞の概念間の距離を調べるために、概念辞書、及び日本語単語辞書を利用する。日本語単語辞書は約26万語の語彙を持つ単語辞書である。基本的役割は、単語と概念（意味）との対応関係を記述し、この対応関係が成り立つときの文法的特性を与えることである。概念辞書は、単語辞書に語義として導入された約41万の概念についての知識が記述され、情報の種類によって、概念体系辞書と概念記述辞書に分けられる。概念体系辞書は約41万の概念に対して、それらの間の上位下位関係を記述したものである。上位下位関係とは概念間の包含関係であり、一種のシソーラスと見なすことができる。概念記述辞書は文中に共起する概念間（2項）の意味的關係（動作主、道具、場所、等）を整理したものを記述したものである。

10

## 【 0 0 5 7 】

## [ 2 . 2 . 5 名詞集合間の類似度を算出する処理 ]

名詞集合間の類似度を算出するための処理について以下に述べる。名詞同士の比較を行い、表記が同じものは類似度を最大値の1として算出する。それ以外の表記が異なるものがある場合には、概念辞書を利用した比較を行う。概念辞書による名詞の比較の手順を以下に示す。まず、名詞の概念を表す概念識別子を日本語単語辞書からとりだし、それを利用して概念辞書から意味属性のリストを得る。次に両名詞の持つ意味属性から名詞間の関係を「類似文の比較による省略可能な格要素の認定」[参考文献：篠原ら，情報処理学会研究報告，NL - 139 - 14，pp101 - 108]の提案による同義関係と類似関係とに分類する。

20

## 【 0 0 5 8 】

図7は、本発明の実施形態に係る情報取得装置の名詞間の同義・類似関係図である。図7(a)は、同義関係を示す。また、図7(b)は類似関係を示す。この2つの関係に基づき、概念辞書を利用した表記が異なる名詞間の類似度を求める。同義関係の類似度aと類似関係の類似度bはそれぞれ次式により求める。各式については、篠原らの名詞間の概念関係の式を採用する。また篠原らは、他に同一関係という概念識別子が同一であるという関係を定義しているが、EDR電子化辞書においてはかなり詳細に概念が定義されているので、同一関係というものは採用していない。

同義関係の類似度 a

## 【 0 0 5 9 】

## 【数11】

$$a = \frac{2D_a}{A_1 + A_2} \quad (0 \leq a \leq 1)$$

A<sub>n</sub> : 名詞nの意味属性数 (n = 1, 2)D<sub>a</sub> : 重複する意味属性数

類似関係の類似度b

## 【 0 0 6 0 】

## 【数12】

$$b = \frac{1}{N_1 N_2} \sum_{i,j=0}^{N_1, N_2} \frac{2D_{ij}}{N_{1i} + N_{2j}} \quad (0 \leq a \leq 1)$$

N<sub>n</sub> : 名詞nの意味属性数N<sub>ni</sub> : 名詞nの意味属性iの上位概念数D<sub>ij</sub> : 意味属性i, jの上位概念の重複数

求めた類似度a、bを使用し、次式により概念間の距離による名詞同士の類似度S<sub>1</sub>を求める。

50

概念間の距離による名詞同士の類似度  $S_1$

【 0 0 6 1 】

【 数 1 3 】

$$S_1 = 1 - e^{-(a+b)} \quad (0 \leq S_1 \leq 1)$$

以上より、単体の名詞間の類似度を求める。そして、以下にそれらを利用した動詞をキーとした名詞集合間の類似度の算出方法を述べる。

10

【 0 0 6 2 】

図 8 は、本発明の実施形態に係る情報取得装置の名詞集合間の類似度算出の例である。c a b o c h a により得られた係り受け情報から、動詞が含まれる文節に係る文節の中の名詞句を動詞をキーとした組として取り出す。図 8 中の名詞集合 1 と名詞集合 3 との類似度を算出する際には、記事 1 を主体と考えた場合に、名詞 A と名詞 F、G 間で類似度が高い方を名詞 A に対する類似した名詞とし、ここでは名詞 F とする。同様に名詞 B も名詞 F、G 間で類似度が高い方を名詞 B に対する類似した名詞とし、ここでは名詞 G とする。そして、主体側の名詞の数を  $n$ 、名詞 A と名詞 F の類似度を  $S_{AF}$ 、名詞と名詞 G の類似度を  $S_{BG}$  した場合には、名詞集合間の類似度を  $S_2$  とすると、 $S_2$  は以下のようになる。

20

【 0 0 6 3 】

【 数 1 4 】

$$S_2 = \frac{S_{AF} + S_{BG}}{n} \quad (0 \leq S_2 \leq 1)$$

同様に、名詞集合 2 と名詞集合 3 を比較し、集合間の類似度を求める。そこで名詞集合 1 と名詞集合 3、名詞集合 2 と名詞集合 3 の類似度をそれぞれ比較し、値が高い方を動詞をキーとする名詞集合間の類似度とする。ここでは名詞集合 1 と名詞集合 3 の類似度  $S_2$  とする。

30

【 0 0 6 4 】

[ 2 . 2 . 6 名詞表記の一致 ]

前工程では名詞の概念間の距離を利用して最も類似度が高い動詞をキーとした名詞集合を各文で選んだ。ここでは、そこで選んだ名詞集合以外の文中の名詞単語中の表記が同じ名詞の割合を算出する。以下にその類似度  $S_3$  を示す。

表記の一致の割合による  $S_3$

【 0 0 6 5 】

【 数 1 5 】

$$S_3 = \frac{2D_{ij}}{A_i + A_j} \quad (0 \leq S_3 \leq 1)$$

40

$D_{ij}$  : 文  $i$  と文  $j$  の動詞に係る文節以外の部分の名詞の内的一致した数

$A_i$  : 文  $i$  中の動詞に係る文節以外の部分の名詞の数

$A_j$  : 文  $j$  中の動詞に係る文節以外の部分の名詞の数

【 0 0 6 6 】

[ 2 . 2 . 7 類似度の算出 ]

文の類似度  $S$  は前述の  $S_2$  と  $S_3$  により以下のようになる。

文の類似度  $S$

50

【 0 0 6 7 】

【数 1 6】

$$S = S_2 + S_3 \quad (0 \leq S \leq 1)$$

【 0 0 6 8 】

[ 2 . 2 . 8 文タイプによる選定 ]

より新聞記事の特色を利用した重複箇所の選定方法として、各文に新聞記事の特徴を考慮した文タイプを設定し、それに基づいた重複文・固有文・補足説明の選定を行う新しい手法を提案する。この手法により、新たに新聞記事特有の言い回し、表現というものを選定の指標として採り入れることが可能となる。

10

【 0 0 6 9 】

[ 2 . 2 . 9 文タイプの種類 ]

また、各文タイプは従来の要約処理において定義されていた多くの文タイプの中から、新聞記事の特徴から要旨、予定、理由、分析、補足説明の5つの文タイプを、また「日本語のシンクタンクスと意味2」[参考文献：寺村秀夫：くろしお出版]の記載による概言のムードと上記データ解析から様態・伝聞、比況・推量の2つの文タイプを本発明の実施の一例とする。以下に要旨、予定、理由、分析、補足説明の5つの文タイプと様態・伝聞、比況・推量の2つの文タイプの特徴、判断基準等について述べる。

20

【 0 0 7 0 】

[ 2 . 2 . 1 0 文タイプ：要旨、予定、理由、分析、補足説明について ]

要旨、予定、理由、分析、補足説明の5種類の文タイプについて、判断基準と特徴について述べる。また判断基準に際し、断定的表現、日時を表す表現に関しては判断基準中においてはそれぞれ#[dantei], #[nitizi?](?=1or2or3)としている。断定的表現に関しては実験データ記事を解析した結果以下のように設定する。なお、以下の判断基準の表記形式はrubyの正規表現の表現形式に準ずる。

【 0 0 7 1 】

【表 1】

である|だとしている|となっている|だ|です

30

また、日時を表す表現については、時間を表す部分(nitizi1)と季節や日付を表す語句(nitizi2)、そしてそれらに付随する語(nitizi3)に大きく分けて設定する。

【 0 0 7 2 】

【表 2】

●nitizi1

((1|2|3|4|5|6|7|8|9|0|1|2|3|4|5|6|7|8|9|0|一|二|三|四|五|六|七|八|九|〇)+(時|日|分|秒))

40

●nitizi2

(昨|今|来)(春|夏|秋|冬|日|月|年)

●nitizi3

(において|にて|に当たって|にかけて|(に際し|に際して)|(のうち|のうちに)|  
|(の際|の際に)||(の点|の点で)||(の時|の時に)|までに|頃|中|から|, |まで|  
後|を指し|に)

【 0 0 7 3 】

これらのうち時間を表す部分(nitizi1)と季節や日付を表す語句(nitizi

50



i 2) は一般的に考えられるものと実験データ 200 記事から設定したものである。また、それらに付随する語 (nitizi3) は実験データの解析と、「自然言語処理の基礎技術」[参考文献: 野村浩郷, 社団法人電子情報通信学会, pp 246] に記載されている図 7.7 の格助詞総当語における時空関係群を参考に設定する。

要旨は新聞記事中で第一文としてある全体の要約が述べられていると考えられる文である。判断基準は記事中の第一文を要旨の文タイプとしている。

【0074】

予定はその文がこれから行われる出来事の日時等を述べられているなど、その文が出来事の予定を表す際につけられる文タイプである。判断基準は前述のトレーニングデータを人手で分類した結果から以下のようにになっている。

【0075】

【表 3】

. \*今後. \*|. \*(方針|見込み|計画|見通し|予定)#{dantei}.  
|. \*({nitizi1}|#{nitizi2})#{nitizi3}. \*

理由はその文が理由を述べている場合に付けられる文タイプである。判断基準は前述のトレーニングデータを人手で分類した結果から以下のように作成する。

【0076】

【表 4】

. \*(という|が|を)(計画|狙い|目的|理由|動機|考え)  
(#{dantei}|.). \*|. \*(と|という)もの.\*

分析はその文が記者の観点から見た出来事に対する分析、意見の場合につけられる文タイプである。判断基準は前述のトレーニングデータを人手で分類した結果から以下のように作成する。

【0077】

【表 5】

\*見通し#{dantei}. |. \*見込(み|んでいる|まれている). |  
. \*(期待|予想|と)(している|されている|みている|みられている). |  
. \*(未知数|程度|确实|初めて|初|とみられるから|形|実情)  
#{dantei}. |. \*と(みられる|いえそうだ).

【0078】

補足説明はその文が記事の補足的な説明の場合に付けられる文タイプである。補足説明には前述のトレーニングデータ記事を分析した結果、以下のような種類がある。

- (1) 記事における登場人物の素性の紹介
- (2) 記事内の出来事についての識者、関係者の話
- (3) 記事内容に対する補足的説明(専門用語等)

補足説明の際には記事中においてその箇所に特殊な記号が記事中で使用されていること

10

20

30

40

50

に着目した。そこで、その記号が出現した後の部分を捕足説明として文タイプを設定する。補足説明の判断に用いた特殊記号は、「 $\llcorner$ 」、「 $\lrcorner$ 」、「 $\triangleright$ 」、「 $\ast$ 」、「 $\llcorner$ 」、「 $\lrcorner$ 」である。

【0079】

[ 2.2.11 文タイプ: 様態・伝聞、比況・推量について ]

後述の様態・伝聞、比況・推量という2つの文タイプについてその採用理由と判断基準について述べる。様態・伝聞は様態と伝聞が合わさった文タイプである。様態とは物の存在や行動のありさまを伝える文タイプである。また、伝聞とは直接にはなく人から伝え聞いているような文タイプである。判断基準を以下に示す。

【0080】

【表6】

. \*という. |. \*によれば. \*|. \*によると. \*|. \*と(話している|  
述べている|している|説明している|いうこと#{dantei}). \*|. \*そうだ. \*

比況・推量は比況と推量が合わさった文タイプである。比況とは動作・状態などをほかのものにたとえて表すような文タイプである。また推量とはある根拠・理由や、確かな論理的要請などに基づいて、込み入った事情や人の心の中などをおしはかっているような文タイプである。判断基準を以下に示す。

【0081】

【表7】

. \*らしい. |. \*(よう|みたい)#{dantei}. |. \*かもしれない. |. \*だろう。  
これら2つの文タイプは、前述の「日本語のシンクタンクと意味2」に記載された二次的

【0082】

これら2つの文タイプは、前述の「日本語のシンクタンクと意味2」に記載された二次的ムードの助動詞中の概言のムードより抜粋する。前述の「日本語のシンクタンクと意味2」によると、現実のいろいろな場で、話し手が、コトを相手の前にもち出すもち出し方、態度を表す部分を「ムード」という構文要素としている。前述の「日本語のシンクタンクと意味2」の記載ではムードとして確言のムード、概言のムード、説明のムードをあげている。ここでは、その中で新聞記事の文タイプとして話し手のいろいろな主観を表すという点で有効と考えられる概言のムード中の様態・伝聞、比況・推量を文タイプとする。様態、伝聞、比況、推量は従来の文法書において「助動詞」という項目の中で、他の形式、たとえば、ナイ、(ラ)レル、(サ)セル、タイ、タなどと並んで個別的にその用法が記述されてきたものである。すなわち、動詞連用形、形容詞語幹につく「ラシイ」は「(根拠のある)推量」、「ヨウダ」は「伝聞」を表す、とされる。文法書によっては、「ヨウダ」あるいは伝聞の「ソウダ」も、形式体言に形式体言「ダ」がついたものとし、助動詞とは認めないものもある。また様態の「ソウダ」を、接尾語に「ダ」がついたものとする見方もある。寺村は、これらを、一定の統語的特徴と、一定の(最大公約数的な)意味を共有するものとして統語的に扱っている。少数の形式については、統語的特徴から外れるが、意味的な特徴から見て、この中に入れる。また、ふつうの文法書では助動詞としては扱われない「カモシレナイ」「カモワカラナイ」「ニチガイナイ」「トイウ」なども、もともと助詞や動詞や助動詞であったものが結びつき、その結びつきが強くなって、一語化したものと見て、前述の統語的、意味的特徴から、やはり概言の助動詞の中を含める。このような理由から概言のムードである様態・伝聞、比況・推量を文タイプの例とする。

【0083】

[ 2.2.12 文タイプの適用優先順位 ]

10

20

30

40

50

以上のように7種類の文タイプを設定したが、文によっては複数の文タイプを兼ねるものも多数存在する。その際に文タイプを設定する優先順位というものを考慮に入れる必要が出てくる。そこで、前述の実験データ200記事に対して文タイプの優先順位が未実装である装置を試作し、各文に対し文タイプを設定した結果を示す。図9は、本発明の実施形態に係る情報取得装置の文タイプごとの割合である。このうち、要旨は記事全体の要約であるという性質上、各文に一程度設定されていると考えられるので、優先順位は最上位とする。また補足説明も補足的な説明を表すという性質上優先順位を最上位とする。それ以外の文タイプ(予定、理由、分析、様態・伝聞、比況・推量)を出現数が少ないものから優先する。結果、優先順位については以下ようになる。

「要旨 = 補足説明 > 比況・推量 > 理由 > 分析 > 予定 > 様態・伝聞」

10

【0084】

[2.2.13 文タイプによる重複文、固有文、補足説明の出力]

図10は、本発明の実施形態に係る情報取得装置の重複箇所、固有箇所、補足説明のカテゴリ分けである。文タイプによる重複箇所、固有箇所、補足説明の出力までの流れである。図10のように同じ文タイプがない場合には前の工程で算出した類似度を利用した重複文選定を行っている。また要旨、補足説明については先に述べた性質から前もって対象文タイプを抜き取る。

【0085】

[2.2.14 複数記事間の集約]

文タイプ選定部36から送られたデータは要約文作成部37又は記事集約部38に送られる。

20

【0086】

[2.2.14.1 2記事間の集約]

2記事間の集約は、1つめの記事の重複箇所、固有箇所の文章に、2つめの記事の固有箇所を合わせた文章を2記事の集約としている。これにより、2記事間で重複しているものや、補足的なものをカットした集約ができる。

【0087】

[2.2.14.2 3記事間の集約]

3つの記事がある場合、2記事同士の集約を3つ作成する。3つの記事を記事A、記事B、記事Cとする。続報記事を分類するので、時系列順でみた初めの2記事A、Bを集約したものを基本とする。この集約と、記事B、Cを集約したものを1文ずつ比較し、含まれていない文を集約として追加する。これにより、記事A、B、Cについての内容についての集約ができる。記事A、Cの集約において、C、の重複箇所と判断されたものが、もし集約に含まれていた場合、それを削除することによって、新しい集約とする。もちろん、この手法を応用することで、3つ以上の記事を一度に集約することもできる。

30

【0088】

[2.2.14.3 記事間の関連度を考慮に入れた集約]

上記の方法に、各記事間の関連度を考慮に入れて集約を行う。記事間の関連度は、各記事を1つのベクトルで表し、そのベクトルを比較して求める。検索を行う際に記事間の関連度を得るが、続報記事であるため、1つめの記事と2つめの記事、2つめの記事と3つ目の記事は繋がりがあっても、1つめの記事と3つめの記事の繋がりが弱い場合がある。記事同士の関連度が低い場合は、記事A、Cの比較を行わないようにすることで、無駄な処理をしなくて済む。

40

【0089】

[3.動作]

ユーザは結果を見ながら記事の適合、非適合の選択、あるいはパラメータを修正することで記事の重要度の変更を行う。システムはそれをフィードバックとして得て、質問ベクトルの修正を行い、順位付けをして再びユーザに検索結果を返す。

図11は、本発明の実施形態に係る情報取得装置の続報記事検索のフローシートである。データが入力される(S100)。形態素解析部211が形態素解析により品詞分解し

50

、名詞（普通、固有、サ変）のみを取り出す（S110）。ここで、形態素解析器茶筌を利用することもできる。TF計算部212が、各記事毎に単語とその単語の出現回数（tf）を記事データベースに登録し、TFを計算する（S120）。IDF計算部213が記事データベースに登録された全単語について文書頻度（df）及びそのIDFを計算する（S130）。単語重み計算部214が各記事の各単語についてTF-IDF法を用いて評価値を求め、単語重みを計算する（S140）。特徴ベクトル作成部215が評価値から各文書の特徴ベクトルを作成する（S150）。つまり、各記事はその記事に出現する全名詞のTF-IDF値を要素にもつベクトルである。情報間類似度計算部216が、全記事ベクトルの組み合わせについて類似度を計算する（S160）。類似度行列計算部217が類似度行列を計算する（S170）。固有ベクトル作成部218がその結果を数値解析プログラムOctaveに渡し、そのベクトルの最大固有値の固有ベクトルを求めることで各記事の参照重要度を計算する（S180）。質問ベクトル作成部219が質問ベクトルを作成する（S190）。検索順位決定部22が検索結果のための類似度計算により検索順位を決定する（S200）。ここで、ある記事 $D_i$ のスコアは、ユーザからの質問ベクトル $Q$ との類似度によるスコア $sim(D_i, Q)$ と固有ベクトルにより求められた重要度を掛け合わせることで求められる。

【0090】

【数17】

$$score(i) = \mu_i sim(D_i, Q)$$

$\mu_i$ ：固有ベクトルにより求められた重要度

【0091】

出力部40が順位付け決定部でスコアリングされた結果をユーザに開示する（S210）。ユーザは検索結果をフィードバックするために検索結果の適否を入力する（S220）。適合・非適合判定部23が検索結果は適合か否かを判定する（S230）。検索結果が適合でないと判定された場合に特徴及び質問ベクトルを修正する（S240）。検索が適合であると判定された場合に情報内容を統合する（S250）。なお、情報内容統合については以下に説示する。

【0092】

図12は、本発明の実施形態に係る情報取得装置の記事内容統合のフローシート（2）である。検索された記事が入力される（S251）。係り受け解析部31が文中に含まれる動詞に関する係り受けを利用するために係り受け解析を行う（S252）。名詞抽出部32が各文中の動詞を含む文節に係る文節中の名詞を抽出する（S253）。名詞集合間類似度比較計算部33が単体の名詞間の類似度 $S_1$ を計算する（S254）。また、名詞集合間類似度比較計算部33が名詞集合間の類似度 $S_2$ を計算する（S255）。名詞表示一致割合計算部34が名詞表示一致割合の類似度 $S_3$ を計算する（S256）。文類似度計算部35が類似度 $S_2$ 及び類似度 $S_3$ から文類似度 $S$ を計算する（S257）。文タイプ選定部36が文タイプによる選定を行う（S258）。要約文作成部37が記事内容を要約する（S259）。記事集約部38が記事内容を集約する（S260）。ここで、記事内容は、要約されたかどうかに関わらず集約することができる。出力部40が内容統合記事の出力を行う（S261）。なお、検索結果及び内容統合記事及び出力の具体的な内容を以下に説示する。

【0093】

[4. 可視化手法]

本発明の実施の形態に係る情報取得装置は、情報を分類する機能を用いて続報記事の情報をユーザに見やすいように整理し、その機能に応じたインターフェースを表示する。そこで、以下にその内容を詳説する。

図13は、本発明の実施形態に係る情報取得装置の実行図である。情報検索においてイ

10

20

30

40

50

インタラク션을促進する関連技術として、情報の可視化は欠くことのできない存在である。可視化によって、システムが提示するデータを効率的にユーザに伝えることができるだけでなく、ユーザのより柔軟なデータへのアクセスが可能になる。Aは関連記事のタイトル表示、Bは記事内容の表示、Cはレーダーチャート、Dは各種コマンド、Eは記事間の関連表示である。

【0094】

図14は、本発明の実施形態に係る情報取得装置の検索式拡張・質問ベクトル選定のためのインターフェースである。図14(a)は、図13のAの関連記事のタイトル表示であり、検索式拡張のためのインターフェースを示す。検索式拡張は、一旦検索した結果に対し利用者が適合文書であったか非適合文書であったかをフィードバックされた結果に基づいて検索式を拡張して、再度検索する。そのために利用者が検索された結果に対し、フィードバックを行うことができる。検索結果一覧表示をした画面に対して適合・非適合の入力を受け付けるように各々の記事の適合・非適合を選択可能とする。図14(a)に示すように、文書検索結果としてユーザ画面には検索結果の記事のタイトルがリスト表示される。ユーザがボタンを押すと、図13のBの本文表示用領域に本文が表示され、確認しながら適合・非適合を選択することができる。

10

【0095】

図14(b)に図13のCのレーダーチャートである質問ベクトル選定のためのインターフェースを示す。質問ベクトルの選定は質問ベクトルとその元となった記事ベクトルとの類似度をレーダーチャートで示すことで行う。レーダーチャートの中心に向かう程、質問ベクトルとの類似性が低く、逆に外側に向かう程、類似性が高くなるように配置している。このように表示することで、質問ベクトルを発散させている記事は凹型になってあらわれる。そういった記事を質問ベクトルから外す、もしくはレーダーチャートの頂点をマウスでドラッグしてその記事に対する重みを補正することで、クエリベクトルの洗練を行う。なお、関連関係と関連度の関連計算は、TF/IDF、ベクトル空間での類似度の判定、統計的手法による類似度の判定、PageRankによる記事の重要度の判定、などを総合して行う。それぞれの関連計算をコントロールするためのものがレーダーチャートである。スクリーンの広さの制約から、レーダーチャートは二つのみ表示されている。これらのレーダーチャートは、上記の判定計算のいずれにも入れ替えることができる。また、サイズを小さくして、前記の四つを表示することもできる。関連計算のコントロールは、総合判定のとき、TF/IDF、ベクトル空間での類似度の判定、統計的手法による類似度の判定、PageRankによる記事の重要度の判定でそれぞれの重み付けを変えたり、それぞれの関連計算においてそれらの計算要素の重み付けを変えたりするものである。これらのコントロールは、スクリーン上でマウスなどの入力手段を使って行う。レーダーチャートの軸は、TF/IDF、ベクトル空間での類似度の判定、統計的手法による類似度の判定、PageRankによる記事の重要度の判定計算の計算要素である。軸上の値が円周に近いほど「重みの値が大きく」その計算要素が重要視される。軸上の値が円の中心に近いほどその計算要素の重要度が小さくされる。総合判定に関するレーダーチャートもある。このときの軸は、TF/IDF、ベクトル空間での類似度の判定、統計的手法による類似度の判定、PageRankによる記事の重要度の判定計算である。軸上の値は、それぞれの重要度を制御するための重みである。以上のインターフェースを提供することで絞りこみ等を行うことができる。

20

30

40

【0096】

図15は、本発明の実施形態に係る情報取得装置の検索結果表示、続報記事発見のためのインターフェースである。

図13のEには記事間の関連表示である検索結果をグラフィカルに表示する領域をそなえている。X軸に時間、Y軸に検索ベクトルに対するスコアをとり、その空間上に記事を表す点を配置している。また、記事を表す点をクリックすることで、その記事本文を図13のBの本文表示用領域に表示する。各記事からはその記事に対する関連性の強さに応じたリンクが結ばれている。各リンクの関連性の強さを線の太さや種類等で識別することが

50

できる。図15では関連性の強いものを実線で表し、関連性の弱いものを点線で表している。ここで、各リンクは強さに応じて色分けをすることもできる。また、ユーザの全体的な興味に対する指標をY軸の座標で、ユーザの局所的な興味に対する指標を記事間を結ぶリンクで表現しているため、例えば、ユーザはこの中から、できるだけY座標が大きく、また、現在見ている記事とのリンクが関連の強い色のリンクで結ばれた記事を読み進めることで、関連記事の中から続報性を持つ記事を発見することができる。

#### 【0097】

なお、記事内容の時間的経緯を考慮することもできる。図13のスクリーンショットのグラフでは、新聞記事を例としているため、横軸が日付となっている。縦軸は、記事の重要度を示し、上部に表示されているほど重要度が高く、下部に表示されているほど重要度が低い。グラフ上の点は記事を表し、点と点を結ぶ線は「関連」があることを示している。点をクリックすると、記事内容が図13のBのウィンドに表示される。点と点を結ぶ線は実線は最も関連度が高いことを示し、点線は最も関連度が低いことを示す。なお、関連性が高い線を赤色に、関連性が低い線を黄色に着色することもできる。さらに、この線の彩色は、より顕示性を上げるために、赤色から青色へのスペクトル変化に対応させることもできる。

#### 【0098】

図16(a)は本発明の実施形態に係る情報取得装置の複数の記事本文表示のためのインターフェース及び(b)新たに記事の本文を表示する場合の表示方法の図である。図13のBの記事本文表示領域には指定した複数の記事本文を同時に表示するためのインターフェースがある。ユーザが本文を表示したい記事を左クリックすると、複数の表示領域の中で、記事本文が表示されていない領域に指定した記事の本文が表示される。もしも3つの表示領域が埋まっていた場合、左側の領域を初期化しその後右側の領域の記事本文を表示する。そして指定した記事本文を右側の領域に表示する。

#### 【0099】

図17は、本発明の実施形態に係る情報取得装置の記事選択補助のための見出し表示領域である。Aによってリストアップされた複数の新聞記事について、図13のDには記事の見出しを表示する領域と複数の記事を比較したり、関連度を計算させたり、要約させたり、集約させたり、再検索させたり、要約または集約した状態から元の記事に戻したりするための各種コマンドのボタンを用意する。なお、図17では、一例として3つのウィンドの場合を示すが、ウィンドの数はいくつであってもよい。また、ユーザが記事を右クリックすると記事の見出しが表示される。これによって表示している記事の本文を変えずに見たい記事を探すことができる。

#### 【0100】

図18は、本発明の実施形態に係る情報取得装置の記事を集約した結果を表示するインターフェースである。記事を集約する方法は、複数の記事同士を1つの文書とみなして、文章構文解析と要約文を作成する。文章構文解析では、意味段落を作成し、意味段落の接続関係を作成しながら、接続関係を崩さずに文章構成を再編し、意味段落の飛地構造解析も行う。文章構造解析された意味段落から、陳述形式による重要句を評価し、重み付けした句を抽出し、さらに、語の類似度を考慮した句の抽出を行い、語の補完を行うことで、集約した文章を作成する。なお、集約を行わず、各文書の要約のみを行い、必要ならば分類タグを付け、文書の分類整理に使うこともできる。分類タグは、例えば、TF/IDFの計算で得られた重要語から作成したり、ベクトル空間の計算のときのベクトルから作成したり、統計的手法による類似度計算に使った重みの値が高い重要要素から作成したりできる。

#### 【0101】

図19は、本発明の実施形態に係る情報取得装置の関連記事検索結果の一例である。Aの記事リストの一番上のものを指定して関連記事を検索したものである。二番目以下の記事が検索された関連記事のリストである。関連記事の検索開始時に指定するものは、記事でもよいし(Aの記事リストの一番上のもの)、自由に記述した文章でもよいし、キーワ

10

20

30

40

50

ードの組み合わせでもよい。関連関係と関連度の計算は、下に述べるユーザ・コントロールが行われていないときは、システムのデフォルト値を使って行われる。Eには、記事間の関連関係と関連度のグラフが表示されている。Bには、Aで表示指定するか、またはEのグラフ内の点をクリックしたときに、それらの記事内容が表示される。Bのサブウィンドの数はいくつでもよい。スクリーンショットでは、3つのサブウィンドが表示されている。Bの上部に選択ボタンが示されているように、現在は、サブウィンドの数は、1つ、2つ、3つの3種類を選択できるようになっている。サブウィンドの数をさらに増やすと、サブウィンドが細い縦長になり、読みづらくなる。レーダーチャートは、デフォルト値が表示されている。

【0102】

図20は、本発明の実施形態に係る情報取得装置の各記事を要約した一例である。Bの各記事を要約したものを表示している。関連関係と関連度を計算するとき、記事の長さが短いほうが計算速度が速くなるため、要約する。また、要約により記事の重要な情報に絞られているため、集約処理の品質がよくなる。

【0103】

図21は、本発明の実施形態に係る情報取得装置の3つの要約文書を1つに集約した一例である。Bの3つの要約文書を1つの文書に集約したものを独立なウィンドに表示している。集約は、要約をしない記事についておこなうこともできる。

【0104】

図22は、本発明の実施形態に係る情報取得装置の集約結果の検討の一例である。集約文書をBの右端のウィンドに移し、集約結果を検討し、必要に応じて、真中や左端の記事を新しく選択し表示する。これらの記事について、関連関係と関連度を再計算し、その結果に基づいて再検索し、新しい関連関係と関連度をグラフ表示する。その結果として、Eのグラフが更新される。必要に応じて、このような操作を繰り返し、最終的な集約を得る。この集約が情報検索の結果である。すなわち、いわゆる情報検索の結果は、一つの文書として出力される。これは、現在の多くの情報検索システムがURLのリストを情報検索の出力としているのとはまったく異なるものである。

【0105】

以上の前記実施形態により本発明を説明したが、本発明の技術的範囲は実施形態に記載の範囲には限定されず、これら各実施形態に多様な変更又は改良を加えることが可能である。そして、かような変更又は改良を加えた実施の形態も本発明の技術的範囲に含まれる。このことは、特許請求の範囲及び課題を解決する手段からも明らかなことである。

【図面の簡単な説明】

【0106】

【図1】本発明の実施形態に係る情報取得装置のハードウェア構成図である。

【図2】本発明の実施形態に係る情報取得装置のブロック構成図である。

【図3】本発明の実施形態に係る情報取得装置の記事間の類似度による記事の重要度評価の説明図である。

【図4】本発明の実施形態に係る情報取得装置の記事の重要度計算例である。

【図5】本発明の実施形態に係る情報取得装置の動詞を含む文節に係る文節中の名詞の抽出の例である。

【図6】本発明の実施形態に係る情報取得装置のEDR辞書の構造図である。

【図7】本発明の実施形態に係る情報取得装置の名詞間の同義・類似関係図である。

【図8】本発明の実施形態に係る情報取得装置の名詞集合間の類似度算出の例である。

【図9】本発明の実施形態に係る情報取得装置の文タイプごとの割合である。

【図10】本発明の実施形態に係る情報取得装置の重複箇所、固有箇所、補足説明のカテゴリ分けである。

【図11】本発明の実施形態に係る情報取得装置の続報記事検索のフローシートである。

【図12】本発明の実施形態に係る情報取得装置の記事内容統合のフローシートである。

【図13】本発明の実施形態に係る情報取得装置の実行図である。

10

20

30

40

50

【図14】本発明の実施形態に係る情報取得装置の検索式拡張・質問ベクトル選定のためのインターフェースである。

【図15】本発明の実施形態に係る情報取得装置の検索結果表示、続報記事発見のためのインターフェースである。

【図16】本発明の実施形態に係る情報取得装置の複数の記事本文表示のためのインターフェース及び新たに記事の本文を表示する場合の表示方法の図である。

【図17】本発明の実施形態に係る情報取得装置の記事選択補助のための見出し表示領域である。

【図18】本発明の実施形態に係る情報取得装置の記事を集約した結果を表示するインターフェースである。

10

【図19】本発明の実施形態に係る情報取得装置の関連記事検索結果の一例である。

【図20】本発明の実施形態に係る情報取得装置の各記事を要約した一例である。

【図21】本発明の実施形態に係る情報取得装置の3つの要約文書を1つに集約した一例である。

【図22】本発明の実施形態に係る情報取得装置の集約結果の検討の一例である。

【符号の説明】

【0107】

1 コンピュータ

2 CPU

3 メインメモリ

4 HDD

5 ビデオカード

6 マウス

7 キーボード

8 光学ディスク

10 入力部

20 続報情報検索部

21 ベクトル作成部

22 検索順位決定部

23 適合・非適合判定部

30 情報内容統合部

31 係り受け解析部

32 名詞抽出部

33 名詞集合間類似度比較計算部

34 名詞表示一致割合計算部

35 文類似度計算部

36 文タイプ選定部

37 要約文作成部

38 記事集約部

40 出力部

20

30

40

211 形態素解析部

212 TF計算部

213 IDF計算部

214 単語重み計算部

215 特徴ベクトル作成部

216 情報間類似度計算部

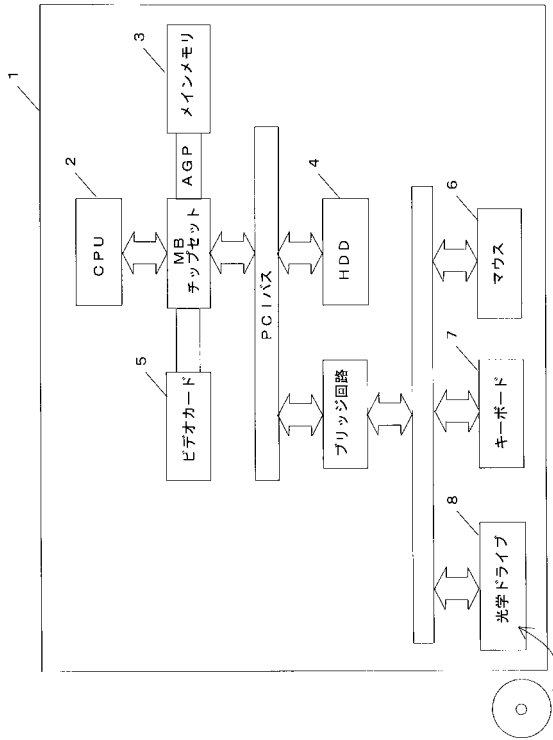
217 類似度行列計算部

218 固有ベクトル作成部

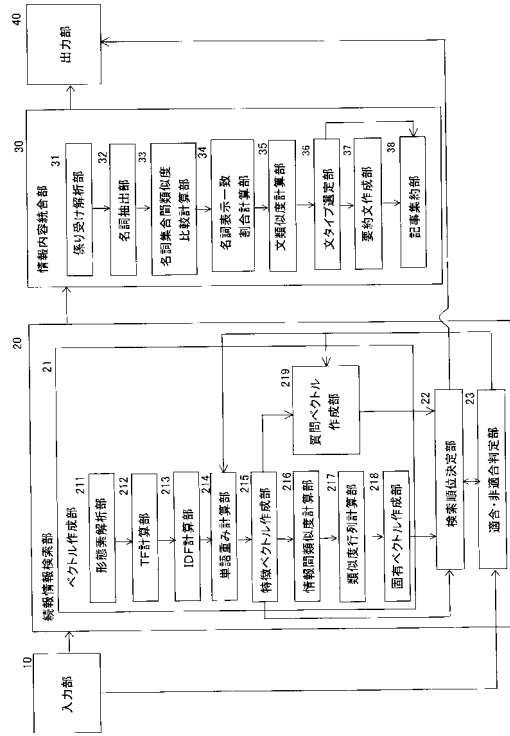
219 質問ベクトル作成部



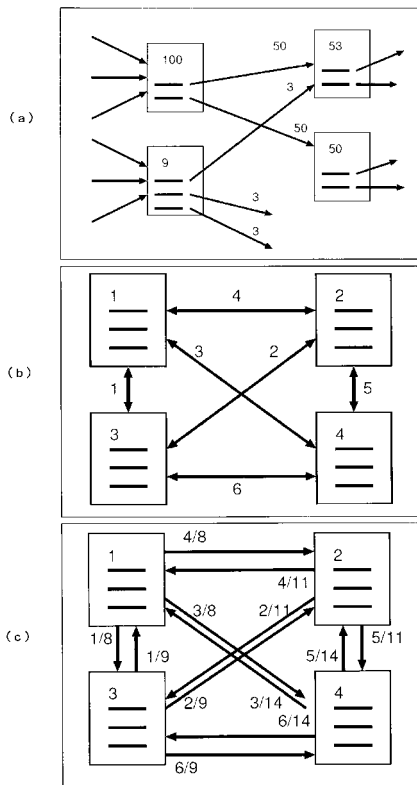
【図1】



【図2】



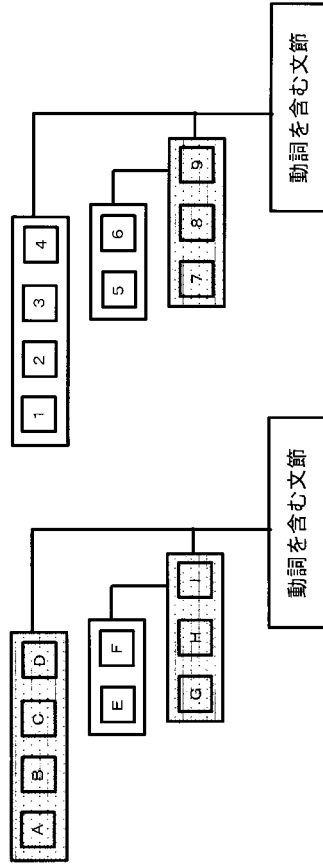
【図3】



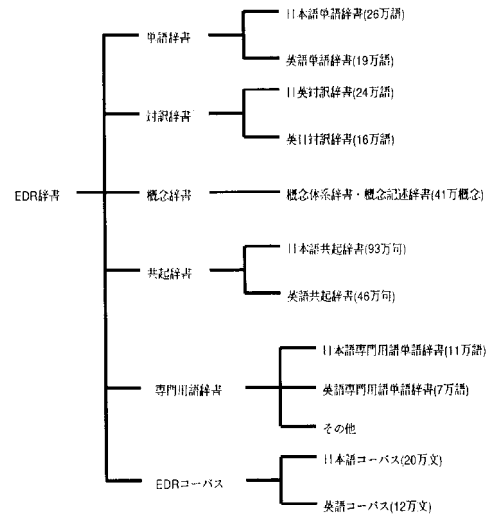
【図4】

文書番号	発リンク	被リンク (合計)	PageRank	順位
1	1/8, 1/8, 3/8	1/11, 1/9, 3/14 (0.6688)	0.19773	1
2	4/11, 2/11, 5/11	1/8, 2/9, 5/14 (1.0793)	0.26097	2
3	1/9, 2/9, 6/9	1/8, 2/11, 6/14 (0.7353)	0.21689	3
4	3/11, 5/14, 6/14	3/8, 5/11, 6/9 (1.4961)	0.32431	1

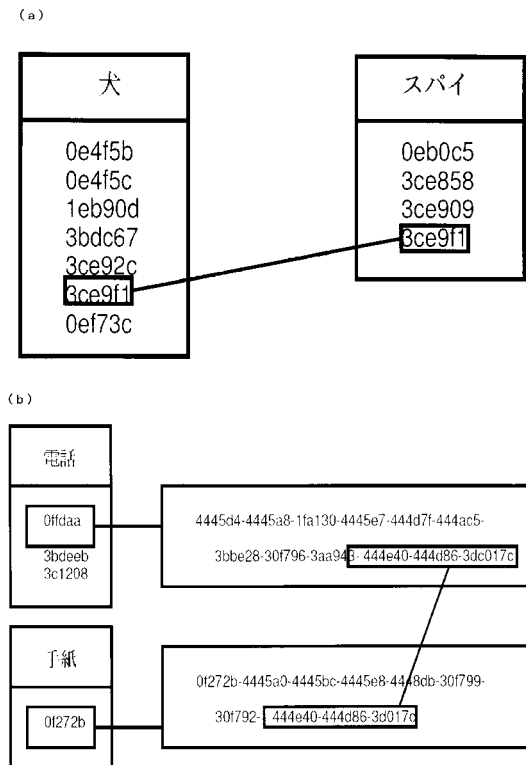
【図5】



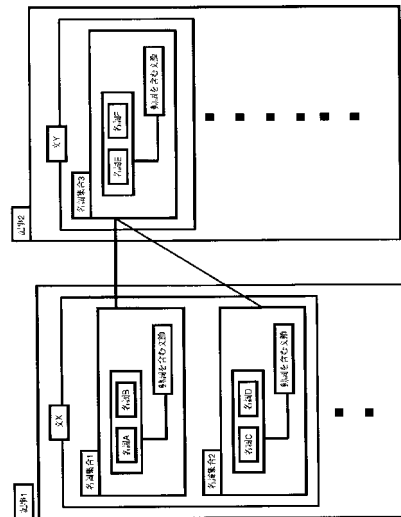
【図6】



【図7】



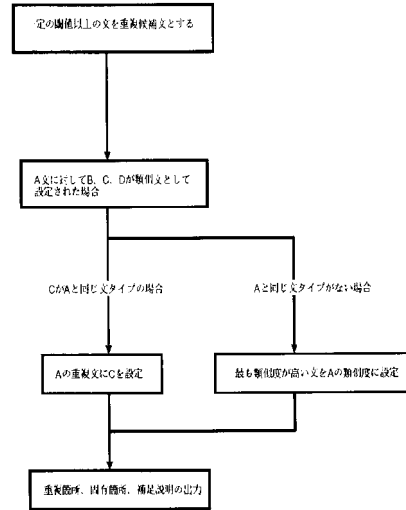
【図8】



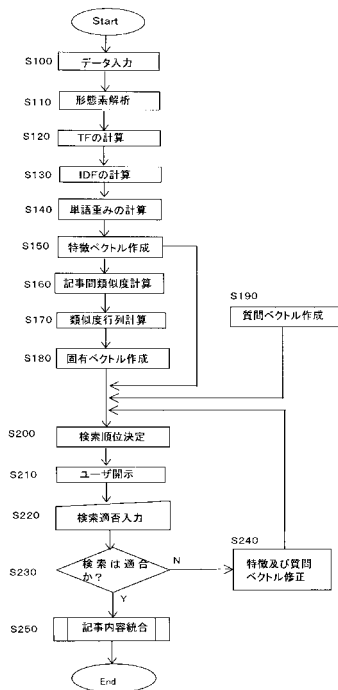
【図9】

文タイプ名	数	割合(%)
要旨	199/1396	14.26
予定	109/1396	7.81
理由	11/1396	0.79
分析	69/1396	4.94
補足説明	105/1396	7.52
様態・伝聞	116/1396	8.31
比況・推量	3/1396	0.21
その他	784/1396	56.16

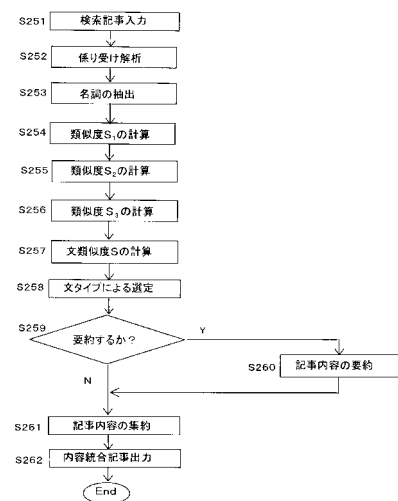
【図10】



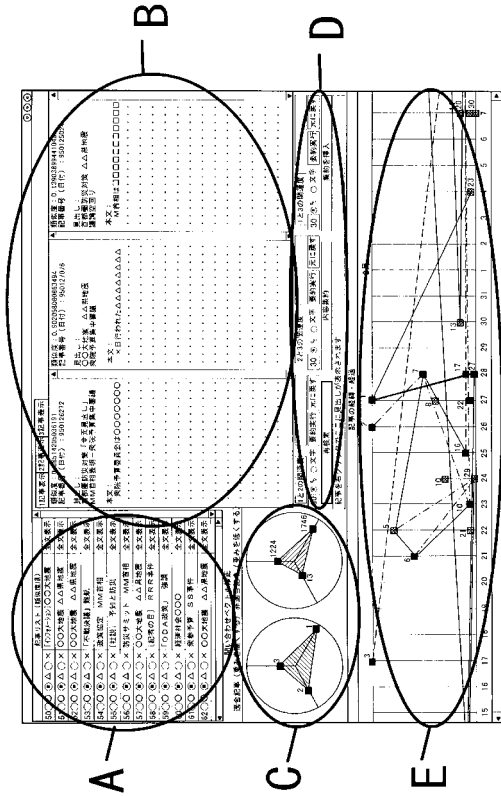
【図11】



【図12】



【図13】

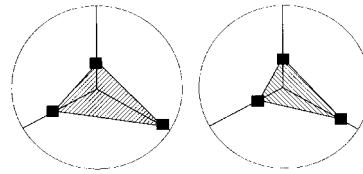


【図14】

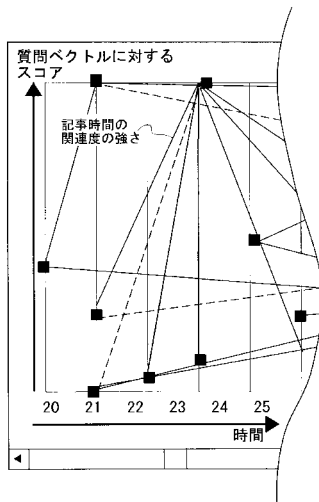
(a)

記事リスト (類似度順)		
50	○●△×	【カメフラシ】○○大地震 全文表示
51	○●△×	○○大地震 △△県地震 全文表示
52	○●△×	○○大地震 △△県地震 全文表示
53	○●△×	「不戦決議」難航 全文表示
54	○●△×	政策協定 MM首相 全文表示
55	○●△×	【社説】 予知と防災 全文表示
56	○●△×	防災サミット MM首相 全文表示
57	○●△×	○○大地震 △△県地震 全文表示
58	○●△×	【記者の目】 RRR事件 全文表示
59	○●△×	「ODA政策」 強調 全文表示
60	○●△×	経済社会○○○ 全文表示
61	○●△×	衆参予算 SS事件 全文表示
62	○●△×	○○大地震 △△県地震 全文表示

(b)



【図15】

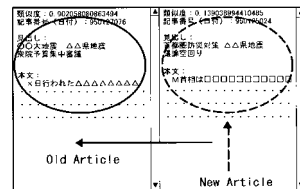


【図16】

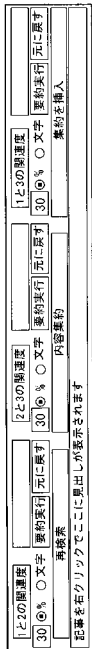
(a)

類似度: 0.5814220503181 記事番号 (日付): 95017877	類似度: 0.90205808063494 記事番号 (日付): 950127016	類似度: 0.1392899410485 記事番号 (日付): 950125024
題名: 政府機関が「全面禁止」 MM首相参拜-英領下野党中絶	題名: ○○大地震 △△県地震 衆参予算案中絶	題名: 防衛増強計画 △△県地震 議員辞任
本文: 高松参事官が○○○○○○	本文: ○行われた△△△△△△	本文: ○行われた○○○○○○○○

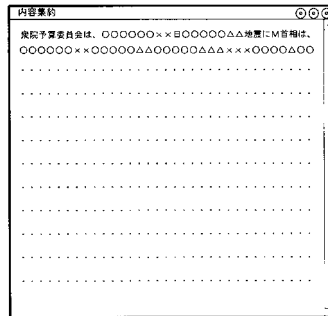
(b)



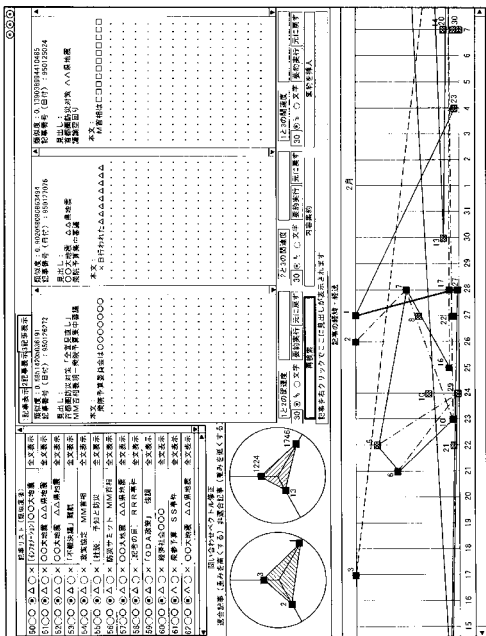
【図 17】



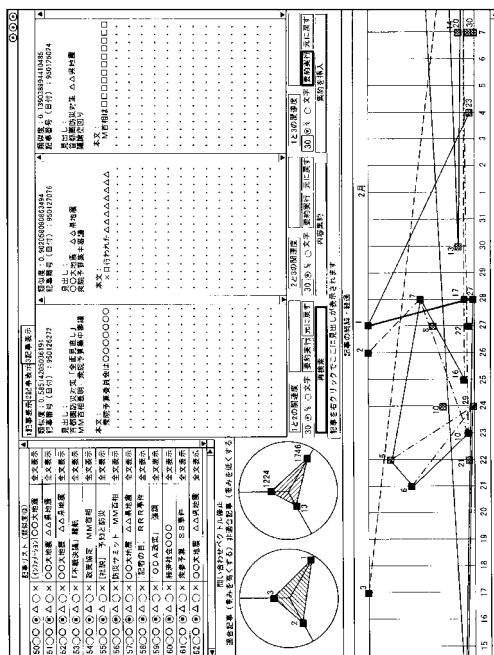
【図 18】



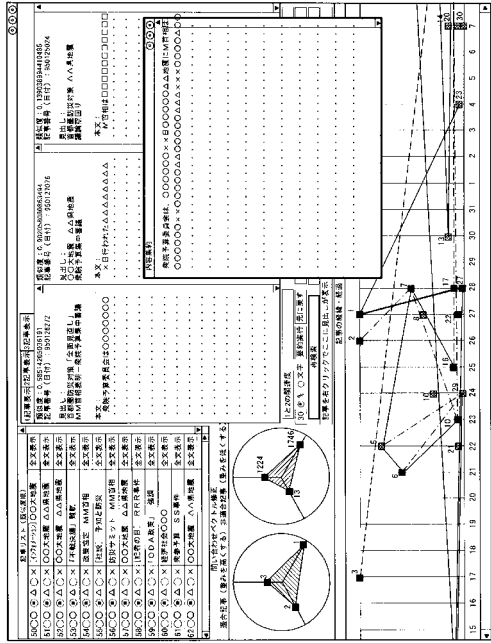
【図 19】



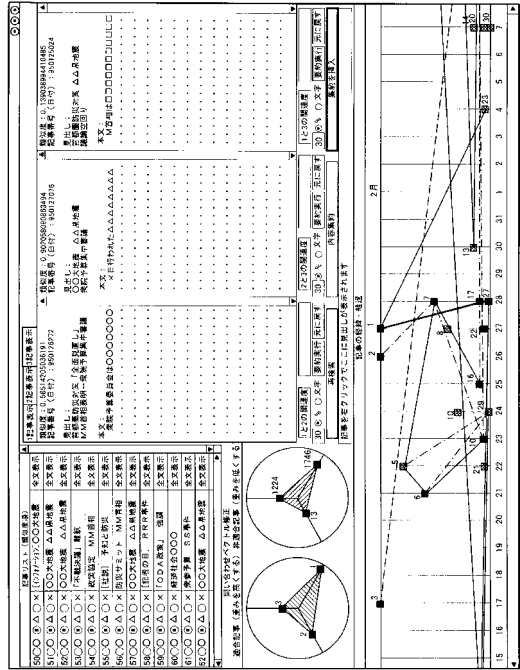
【図 20】



【 2 1 】



【 2 2 】



## フロントページの続き

- (56)参考文献 特開2005-327225(JP,A)  
特開2004-185515(JP,A)  
特開平10-134066(JP,A)  
特開2003-228581(JP,A)  
荒谷寛和、外2名、ウェブページ間類似度に基づく推薦リンクを用いたウェブ検索システムの設計、電子情報通信学会技術研究報告(AI2004-12~18)、日本、社団法人電子情報通信学会、2004年7月22日、第104巻、第233号、p.7-12  
中村貞吾、外2名、文タイプと文間関係に基づく要約処理、言語処理学会第4回年次大会ワークショップ論文集、日本、言語処理学会、1998年3月27日、p.50-55  
中山聡、外4名、EDRコーパスを利用した動詞の語義分類、電子情報通信学会技術研究報告(NLC95-40~45)、日本、社団法人電子情報通信学会、1995年10月20日、第95巻、第321号、p.23-30  
平田陽一、外3名、Web検索における意味的適合フィードバック機構、情報処理学会研究報告(2000-DBS-122)、日本、社団法人情報処理学会、2000年7月28日、第2000巻、第69号、p.137-144  
篠原直嗣、外2名、類似文の比較による省略可能な格要素の認定、情報処理学会研究報告(2000-NL-139)、日本、社団法人情報処理学会、2000年9月22日、第2000巻、第86号、p.101-108  
金田晃征、外1名、情報検索と情報集約による情報取得システム、情報処理学会研究報告(2007-NL-179)、日本、社団法人情報処理学会、2007年5月24日、第2007巻、第47号、p.31-36

## (58)調査した分野(Int.Cl., DB名)

G06F 17/30