

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5869948号
(P5869948)

(45) 発行日 平成28年2月24日 (2016. 2. 24)

(24) 登録日 平成28年1月15日 (2016. 1. 15)

(51) Int. Cl. F 1
G 0 6 F 17/27 (2006.01) G 0 6 F 17/27 6 6 5

請求項の数 14 (全 14 頁)

(21) 出願番号	特願2012-95344 (P2012-95344)	(73) 特許権者	000005108
(22) 出願日	平成24年4月19日 (2012. 4. 19)		株式会社日立製作所
(65) 公開番号	特開2013-222418 (P2013-222418A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成25年10月28日 (2013. 10. 28)	(74) 代理人	110001689
審査請求日	平成27年1月8日 (2015. 1. 8)		青稜特許業務法人
		(74) 代理人	110000350
			ポレール特許業務法人
		(72) 発明者	柿下 容弓
			東京都国分寺市東恋ヶ窪一丁目280番地
			株式会社日立製作所 中央研究所内
		(72) 発明者	服部 英春
			東京都国分寺市東恋ヶ窪一丁目280番地
			株式会社日立製作所 中央研究所内

最終頁に続く

(54) 【発明の名称】 パッセージ分割方法、装置、及びプログラム

(57) 【特許請求の範囲】

【請求項 1】

処理部により、ドキュメントをパッセージに分割するパッセージ分割方法であって、前記処理部は、
 前記ドキュメントを文単位に分割し、
 分割した前記文をクエリとして、予め記憶されている複数のドキュメントから、関連するドキュメントを抽出して、特徴量を作成し、
 作成した前記特徴量の内の二つの特徴量の類似度が所定の閾値以上である、当該二つの特徴量の共通要素を用いて特徴量を更新し、
 前記特徴量の類似度が所定の閾値以上である、当該二つの特徴量に対応する前記文、あるいはパッセージ候補を連結して新たなパッセージ候補とし、
 前記特徴量の更新は前記新たなパッセージ候補の特徴量を得るものであることを特徴とするパッセージ分割方法。

【請求項 2】

請求項 1 に記載のパッセージ分割方法であって、
 前記処理部は、
 前記特徴量として、ドキュメントベクトルを用いる、
 ことを特徴とするパッセージ分割方法。

【請求項 3】

請求項 2 に記載のパッセージ分割方法であって、

前記処理部は、

前記二つの特徴量である、二つのドキュメントベクトル V_i 、 V_j の類似度が所定の閾値以上である場合、二つの前記ドキュメントベクトル V_i 、 V_j の共通要素 V_{ij} を抜き出し、検索クエリを生成し、

生成した前記検索クエリを用いて、新たなドキュメントベクトル V'_{ij} を得る、
ことを特徴とするパッセージ分割方法。

【請求項 4】

請求項 3 に記載のパッセージ分割方法であって、

前記処理部は、

前記新たなドキュメントベクトル V'_{ij} が、前記共通要素 V_{ij} の要素を含む度合い
10 に対応して、前記新たなドキュメントベクトル V'_{ij} のベクトルサイズを修正する、
ことを特徴とするパッセージ分割方法。

【請求項 5】

請求項 3 に記載のパッセージ分割方法であって、

前記処理部は、

前記新たなドキュメントベクトル V'_{ij} に対応する前記文、あるいはパッセージ候補
として、二つの前記ドキュメントベクトル V_i 、 V_j に対応する前記文、あるいはパッセージ候補を連結して、新たなパッセージ候補とする、
20 ことを特徴とするパッセージ分割方法。

【請求項 6】

請求項 1 に記載のパッセージ分割方法であって、

前記処理部は、

前記特徴量として、単語ベクトルを用いる、

ことを特徴とするパッセージ分割方法。

【請求項 7】

請求項 6 に記載のパッセージ分割方法であって、

前記二つの特徴量である、二つの単語ベクトル V_i 、 V_j の類似度が所定の閾値以上
ある場合、二つの前記単語ベクトル V_i 、 V_j の共通要素 V_{ij} を抜き出し、検索クエリを生成し、

生成した前記検索クエリを用いて、新たな単語ベクトル V'_{ij} を得る、
30 ことを特徴とするパッセージ分割方法。

【請求項 8】

請求項 7 に記載のパッセージ分割方法であって、

前記処理部は、

前記新たな単語ベクトル V'_{ij} が、前記共通要素 V_{ij} の要素を含む度合いに対応し
て、前記新たな単語ベクトル V'_{ij} のベクトルサイズを修正する、
40 ことを特徴とするパッセージ分割方法。

【請求項 9】

請求項 8 に記載のパッセージ分割方法であって、

前記処理部は、

前記新たな単語ベクトル V'_{ij} に対応する前記文、あるいはパッセージ候補として、
二つの前記単語ベクトル V_i 、 V_j に対応する前記文、あるいはパッセージ候補を連結して、新たなパッセージ候補とする、
40 ことを特徴とするパッセージ分割方法。

【請求項 10】

入力されるドキュメントをパッセージに分割するパッセージ分割装置であって、

処理部と記憶部とを備え、

前記処理部は、

前記ドキュメントを文単位に分割し、

分割した前記文をクエリとして、予め前記記憶部に記憶されている複数のドキュメント
50

から、関連するドキュメントを抽出して、特徴量を作成し、

作成した前記特徴量の内の二つの類似度が所定の閾値以上である、当該特徴量の共通要素を用いて特徴量を更新し、

前記特徴量の類似度が所定の閾値以上である、当該二つの特徴量に対応する文またはパッセージ候補を連結して新たなパッセージ候補とし、

前記特徴量の更新は前記新たなパッセージ候補の特徴量を得るものであることを特徴とするパッセージ分割装置。

【請求項 1 1】

請求項 1 0 に記載のパッセージ分割装置であって、

前記処理部は、

前記特徴量として、関連する前記ドキュメントに基づく、ドキュメントベクトルあるいは単語ベクトルを用いる、

ことを特徴とするパッセージ分割装置。

10

【請求項 1 2】

請求項 1 1 に記載のパッセージ分割装置であって、

前記処理部は、

前記二つの特徴量である、二つのドキュメントベクトル、或いは単語ベクトル V_i 、 V_j の類似度が所定の閾値以上である場合、二つの前記ドキュメントベクトル、或いは単語ベクトル V_i 、 V_j の共通要素 V_{ij} を抜き出し、検索クエリを生成し、

生成した前記検索クエリを用いて、新たなドキュメントベクトル、或いは単語ベクトル V'_{ij} を得、

前記新たなドキュメントベクトル、或いは単語ベクトル V'_{ij} が、前記共通要素 V_{ij} の要素を含む度合いに対応して、前記新たなドキュメントベクトル、或いは単語ベクトル V'_{ij} のベクトルサイズを修正する、

ことを特徴とするパッセージ分割装置。

20

【請求項 1 3】

請求項 1 2 に記載のパッセージ分割装置であって、

前記処理部は、

前記新たなドキュメントベクトル V'_{ij} に対応する前記文、あるいはパッセージ候補として、二つの前記ドキュメントベクトル V_i 、 V_j に対応する前記文、あるいはパッセージ候補を連結し、新たに連結されたパッセージ候補を前記記憶部に記憶する、

ことを特徴とするパッセージ分割装置。

30

【請求項 1 4】

処理部と記憶部とを備え、入力されるドキュメントをパッセージに分割するパッセージ分割装置の処理部で実行されるパッセージ分割プログラムであって、

前記処理部を、

前記ドキュメントを文単位に分割し、

分割した前記文をクエリとして、予め前記記憶部に記憶されている複数のドキュメントから、関連するドキュメントを抽出し、

抽出した前記関連するドキュメントを用いて特徴量を作成し、

作成した前記特徴量の内の二つの類似度が所定の閾値以上である、当該特徴量の共通要素を用いて特徴量を更新し、

前記特徴量の類似度が所定の閾値以上である、当該二つの特徴量に対応する文またはパッセージ候補を連結して新たなパッセージ候補とし、

前記特徴量の更新は前記新たなパッセージ候補の特徴量を得るものであるよう動作させる、

ことを特徴とするパッセージ分割プログラム。

40

【発明の詳細な説明】

【技術分野】

【0001】

50

本発明は、電子化された文書の処理に係り、特に電子化書類のパッセージ分割技術に関する。

【背景技術】

【0002】

近年、文書の電子化やデータベース化が進んだことで、自然言語処理技術も大きく発展し、例えば文書の自動要約や文書検索のための自動キーワード抽出などの研究が多くなされてきた。しかしこれらの技術の対象となる文書はパッセージ毎、すなわち、話題、あるいは内容的、意味的なまとまり単位毎に分割されている、または単一のパッセージしか含まない文書を想定していることが多い。そのため、複数のパッセージを含む文書に対しては、予めパッセージを分割することが有効である。従来、このようなパッセージ分割手法

10

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開2009-15795号公報

【特許文献2】特開2004-145790号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

20

しかし、従来のパッセージ分割、テキストセグメンテーションに関する手法は意味の近い文、すなわちその特徴量が似た文を含む複数のパッセージが、一つの文書に含まれる場合、パッセージを正しく分割することが難しい。その結果、文書の自動要約や文書検索のための自動キーワード抽出などを効率的に進めることができない。

【0005】

本発明の目的は、上記課題に鑑みてなされたものであり、複数のパッセージを含む文書を有効に分割するパッセージ分割方法、装置、及びプログラムを提供することにある。

【課題を解決するための手段】

【0006】

上記の目的を達成するため、本発明においては、処理部により、ドキュメントをパッセージに分割するパッセージ分割方法であって、処理部は、ドキュメントを文単位に分割し、分割した文をクエリとして、予め記憶されている複数のドキュメントから、関連するドキュメントを抽出して、特徴量を作成し、作成した特徴量の内の二つの特徴量の類似度が所定の閾値以上である、当該二つの特徴量の共通要素を用いて特徴量を更新するパッセージ分割方法を提供する。

30

【0007】

又、上記の目的を達成するため、本発明においては、入力されるドキュメントをパッセージに分割するパッセージ分割装置であって、処理部と記憶部とを備え、処理部は、ドキュメントを文単位に分割し、分割、記憶した文をクエリとして、予め記憶部に記憶されている複数のドキュメントから、関連するドキュメントを抽出して、特徴量を作成し、作成した特徴量の内の二つの類似度が所定の閾値以上である、当該特徴量の共通要素を用いて特徴量を更新する構成のパッセージ分割装置を提供する。

40

【0008】

更に、上記の目的を達成するため、本発明においては、処理部と記憶部とを備え、入力されるドキュメントをパッセージに分割するパッセージ分割装置の処理部で実行されるパッセージ分割プログラムであって、処理部を、ドキュメントを文単位に分割し、分割した文をクエリとして、予め記憶部に記憶されている複数のドキュメントから、関連するドキュメントを抽出し、抽出した関連するドキュメントを用いて特徴量を作成し、作成した特徴量の内の二つの類似度が所定の閾値以上である、当該特徴量の共通要素を用いて特徴量を更新するよう動作させるパッセージ分割プログラムを提供する。

50

【発明の効果】

【0009】

本発明によれば、意味の近い文、すなわち特徴量が似た文を含む、複数のパッセージが一つの文書に含まれる場合でも、パッセージを正しく分割することが可能となる。

【図面の簡単な説明】

【0010】

【図1A】第1の実施例のパッセージ分割装置の一機能構成を示す図である。

【図1B】第1の実施例のパッセージ分割装置の一ハードウェア構成を示す図である。

【図2】第1の実施例に係る、パッセージ分割プログラムの動作の一例を示す図である。

【図3】第1の実施例に係る、ドキュメントベクトルの類似度に応じて文が連結される様子を示す図である。

【図4】第2の実施例のパッセージ分割装置の一機能構成を示す図である。

【図5】第2の実施例に係る、パッセージ分割プログラムの動作の一例を示す図である。

【図6】各実施例に係る、ドキュメントベクトルの一例を説明するための図である。

【図7】各実施例に係る、単語ベクトルの一例を説明するための図である。

【発明を実施するための形態】

【0011】

以下、本発明の実施例を図面に従い説明するが、本発明は以下に説明する実施例に限定されるものではない。本明細書において、「文書」と「ドキュメント」とは、同義であることとする。また、「パッセージ」とは、話題、あるいは内容的、意味的なまとまりのある単位を意味する。更に、ドキュメントベクトルとは、蓄積されたドキュメントを次元とするベクトルを意味し、単語ベクトルとは、全ドキュメント中に出現する全ての単語を次元とするベクトルを意味するものとする。そして、本明細書において、文の「特徴量」とは、文の意味を定量的に示すものであり、例えば、ドキュメントベクトル、あるいは単語ベクトルはその一例として説明する。

【実施例1】

【0012】

第1の実施例は、類似度計算にドキュメントベクトルを、類似文書検索に単語ベクトルを用いるパッセージ分割方法、装置、及びプログラムの実施例である。本実施例において、ドキュメントベクトルとは、分割装置のコーパス部に含まれる全てのドキュメントを次元とするベクトルである。

【0013】

本実施例の詳細を説明するに先立ち、ドキュメントベクトルと単語ベクトルの一例を説明する。

図6にドキュメントベクトルの一例を示す。図6において、コーパス部に含まれるドキュメントの総数を10として例示した。そして、検索の結果得られるドキュメントが、1、3、4、8である場合、ドキュメントベクトルは、同図の(a)に示すドキュメントベクトル601のように表わすことができる。同様に、検索の結果、検索スコアが得られる場合、得られた検索スコアを用いて、同図の(b)に示すようなドキュメントベクトル602として表わすことができる。

【0014】

図7に単語ベクトルの一例を示した。単語ベクトルとは、全文書中に出現する全ての単語を次元とするベクトルであり、図7の単語ベクトルでは、全てのドキュメントに出現する単語の種類を10として例示した。そして、あるドキュメントに含まれる単語が、3、6、7、8であり、出願頻度がそれぞれ、1、5、3、9である場合、該当する要素に出現頻度を代入することで、同図に示す単語ベクトル701を得る。

【0015】

図1Aは、実施例1に係るパッセージ分割装置の機能ブロックの一例を示す図である。図1Bは、実施例1のパッセージ分割装置を実現するハードウェア構成の一例を示す図である。図1Bのハードウェア構成は、通常の処理部である中央処理部(Central

10

20

30

40

50

Processing Unit: CPU) 11、メモリ、RAM、ROM、ハードディスクドライブ(HDD)、記憶装置等の記憶部12、入出力部13、ネットワークインタフェースである通信部14からなり、これらの各ブロックは、内部バス15によって相互に接続されているコンピュータを示している。

【0016】

図1Aにおいて、パッセージ分割装置100は、制御部101と、入力部102と、文分割部103と、特徴量算出部104と、類似度計算部105と、検索クエリ生成部106と、特徴量更新部107と、パッセージ更新部108と、出力部109と、文記憶部110と、コーパス部111と、特徴量記憶部112と、パッセージ記憶部113と、形態素解析部114とを有する。前提として、コーパス部111には、例えば新聞記事のような文書、ドキュメントが S_D 個記憶されているものとする。

10

【0017】

この内、入力部102、出力部109が入出力部13や通信部14に対応し、文記憶部110と、コーパス部111と、特徴量記憶部112と、パッセージ記憶部113が記憶部12のメモリや記憶装置に対応している。その余の制御部101、文分割部103と、特徴量算出部104と、類似度計算部105と、検索クエリ生成部106と、特徴量更新部107と、パッセージ更新部108と、形態素解析部114は、CPU11における、オペレーティングシステム(OS)や、ROM等の記憶部に記憶された各種のプログラムの処理で実現できる。

【0018】

図1Aに示した実施例1のパッセージ分割装置の各機能ブロックの動きを順次説明する。

20

まず、パッセージ分割の対象となるドキュメントが入力部102から装置に入力される。文分割部103は、処理部であるCPU11の所定プログラムの実行により、入力されたドキュメントを文単位に分割し、文記憶部110に分割結果である複数の文を記憶する。

【0019】

同様に、特徴量算出部104は、文記憶部110から読み込んだ文各々を用いて、コーパス部111から関連するドキュメントを取得し、得られた複数の関連ドキュメントを、ドキュメントベクトル化して特徴量記憶部112に記憶する。すなわち、特徴量算出部104は、取得した関連ドキュメントに対応する次元に値を代入することで、図6で例示したようなドキュメントベクトルを生成する。

30

【0020】

検索クエリ生成部106は、検索クエリを生成し、制御部101に送る機能を持つ。

【0021】

特徴量算出部104は、制御部101を介して、検索クエリが与えられた場合、当該検索クエリに関連するドキュメントを文記憶部110から取得し、得られた複数の関連ドキュメントをドキュメントベクトル化し、特徴量として、特徴量記憶部に112に記憶すると共に、制御部101を介して、特徴量更新部107に出力する。

【0022】

類似度計算部105は、制御部101の指定に基づいて、二つのドキュメントベクトルを特徴量記憶部112から読み出し、二つのドキュメントベクトルの類似度を計算する機能を有する。本実施例における類似度の計算方法については後述する。更に、類似度計算部105は、計算して得られた類似度が所定の閾値以上か否かを判断する。

40

【0023】

検索クエリ生成部106は、制御部101の指定に基づいて、二つのドキュメントベクトルを特徴量記憶部112から読み出し、二つのドキュメントベクトルに共通するドキュメント群をコーパス部111から抽出する。抽出された共通するドキュメント群から検索クエリを生成し、制御部101へ出力する。この検索クエリの生成方法については後述する。

50

【 0 0 2 4 】

特徴量更新部 1 0 7 は、制御部 1 0 1 の指定に基づいて二つのドキュメントベクトル V_i , V_j を特徴量記憶部 1 1 2 から読み出す。また制御部 1 0 1 から一つのドキュメントベクトル V_k が特徴量更新部 1 0 7 に入力される。入力された三つのドキュメントベクトル V_k , V_i , V_j から信頼度を計算し、信頼度に基づいて V_k を修正する。この信頼度については後述する。その後、 V_i , V_j を特徴量記憶部 1 1 2 から削除し、 V_k を特徴量記憶部 1 1 2 に記憶する。

【 0 0 2 5 】

パッセージ更新部 1 0 8 は、制御部 1 0 1 の指定に基づいて、文記憶部 1 1 0 またはパッセージ記憶部 1 1 3 の中から二つの文またはパッセージ候補を読み出す。読み出された文またはパッセージ候補を文記憶部 1 1 0 またはパッセージ記憶部 1 1 3 の中から削除し、読み出された文またはパッセージ候補を連結して、その連結結果を、パッセージ候補としてパッセージ記憶部 1 1 3 に記憶する。

10

【 0 0 2 6 】

出力部 1 0 9 は文記憶部 1 1 0 とパッセージ記憶部 1 1 3 からそれぞれ文、パッセージ候補を読み出し、不明パッセージか否かを判定した上で、その判定結果に基づき、パッセージにラベルを付与して出力する。ここで不明パッセージとは、どのパッセージと連結するか判定できなかった文またはパッセージ候補を指す。不明パッセージの判定方法については後述する。

【 0 0 2 7 】

図 2 は本実施例に係るパッセージ分割装置で実行されるパッセージ分割プログラムの動作を示すフロー図である。以下、図 2 を用いてパッセージ分割プログラムの動作の一例について説明する。

20

ここでは例として、二つのパッセージを含むドキュメントが入力された場合について述べるが、入力されるドキュメント中のパッセージ数は二つ以上であっても良く、以後の処理は同じであるので、二つのパッセージを含むドキュメントを例にして説明する。

【 0 0 2 8 】

第一のパッセージに含まれる文を a_1 , a_2 , ... , a_N 、第二のパッセージに含まれる文を b_1 , b_2 , ... , b_M と定義する。ここで N は第一のパッセージに含まれる文の数（自然数）、 M は第二のパッセージに含まれる文の数（自然数）である。

30

【 0 0 2 9 】

まず、ステップ 2 0 1 で入力部 1 0 2 からドキュメントが入力される。

ステップ 2 0 2 では入力されたドキュメントが、文分割部 1 0 3 により文単位に分割され、文記憶部 1 1 0 に記憶される。

【 0 0 3 0 】

ステップ 2 0 3 では文記憶部 1 1 0 に記憶された全ての文 a_1 , a_2 , ... , a_N 、 b_1 , b_2 , ... , b_M を特徴量算出部 1 0 4 に入力し、先に説明した通り、ドキュメントベクトルを得る。ドキュメントベクトルの算出方法としては、例えば、コサイン尺度を用いる方法が挙げられる。コサイン尺度とは二つのベクトルの類似度を計る手法の一つとして用いられるものである。二つのベクトル Q 、 P のコサイン尺度は以下の式 1 で計算される。

40

【 0 0 3 1 】

【 数 1 】

$$\frac{\sum_i q_i p_i}{\sqrt{\sum_i q_i^2} \sqrt{\sum_i p_i^2}} \quad (q_i \in Q, p_i \in P) \quad \text{----- (式1)}$$

本実施例においては、上述の通り、類似するドキュメントの検索に単語ベクトルを用いる。そこで、例えば、コーパス部 1 1 1 に記憶された各ドキュメントに対して、含まれる単語の出現頻度を要素とする単語ベクトル W_i ($0 \leq i < S_D$) を作成しておく。入力さ

50

れた文についても同様に単語ベクトル化し、 $W_{current}$ とする。単語ベクトル $W_{current}$ と、単語ベクトル W_i ($0 \leq i < S_D$)のコサイン尺度を計算し、得られた類似度が高いドキュメントから L 番目 (L は所定の自然数)までのドキュメントを得て、ドキュメントベクトル化し、特徴量記憶部112に蓄積する。

【0032】

尚、ここでは類似度計算の例として、コサイン尺度を用いたが、その他の尺度を用いて、類似度を計算しても良い。ドキュメントベクトルの各要素の値としては、図6の(a)、(b)で説明したように、選定されたドキュメントは1、その他のドキュメントは0としても良いし、算出された類似度を用いるなど、なんらかの重み付けを行っても良い。

【0033】

次にステップ204では、特徴量記憶部112に蓄積されているドキュメントベクトルを二つ読み出し、類似度計算部105を用いて、最も類似度の高いドキュメントベクトルの組 V_i, V_j を見つける。この場合における類似度の計算方法としては、上述したコサイン尺度等を用いても良いし、二つのドキュメントベクトルの両方に存在する要素、すなわち共通要素の数などを用いても良い。

【0034】

ステップ205では、類似度計算部105が、ステップ204で算出した最大類似度が、予め設定した閾値以上か否かを判定する。閾値は予め設定した固定値でも良いし、ステップ204で類似度を計算した際に、計算した類似度の平均や分散を計算しておき、これを用いても良い。

【0035】

ステップ206およびステップ207は検索クエリ生成部106にて行われる。ステップ206では、ステップ204で算出された最大類似度が閾値以上である場合、ドキュメントベクトルの組 V_i, V_j の共通要素を抽出し、これをドキュメントベクトルの共通要素 V_{ij} とする。

【0036】

ステップ207では、ステップ206で得られた共通要素 V_{ij} から検索クエリを生成する。検索クエリの生成方法としては、例えばTFIDFを用いた方法が挙げられる。TFIDFとは単語に関する重みの一種である。TF (Term Frequency) と IDF (Inverse Document Frequency) はそれぞれ次の式で表され、TFIDFはTFとIDFの積で求められる。

【0037】

【数2】

$$tf_i = \frac{n_i}{\sum_k n_k}, \quad idf_i = \log \frac{|D|}{|\{d: t_i \in d\}|} \quad \text{----- (式2)}$$

ここで n_i はドキュメント d における単語 i の出現回数、 $|D|$ は総ドキュメント数、 $|\{d: t_i \in d\}|$ は単語 t_i を含むドキュメント数である。本実施例においては、総ドキュメント数 D はコーパス部111に記憶されている全ドキュメント数に相当する。

【0038】

ドキュメント d に対して、形態素解析部114を用いて形態素解析を行い、TFIDFが大きい順に S_w 個の単語を抽出し、これを検索クエリとする。TFIDF以外にも、例えば出現頻度の多さで重要度を決めても良いし、ドキュメントのタイトルをクエリとしても良いし、その他の方法で検索クエリを生成しても良い。

【0039】

ステップ208では、ステップ207で生成された検索クエリを、制御部101を介して特徴量算出部104に入力し、特徴量算出部104において、新たなドキュメントベクトル V'_{ij} を得る。

【0040】

10

20

30

40

50

続いて、新たに得られたドキュメントベクトル V'_{ij} の信頼度の計算等を実行するステップ 209 およびステップ 210 を実行する。これらのステップ 209 およびステップ 210 は、図 1 に示した特徴量更新部 107 にて実行される。まず、ステップ 209 では、ステップ 208 で得られたドキュメントベクトル V'_{ij} の信頼度を計算し、その結果に応じて、ドキュメントベクトルのベクトルサイズを修正する。

【0041】

本実施例において信頼度とは、ドキュメントベクトル V'_{ij} に共通要素 V_{ij} の要素がどれだけ含まれているかを数値化した指標である。信頼度の算出としては、例えばドキュメントベクトル V'_{ij} がドキュメントベクトルの組 V_i, V_j の共通要素 V_{ij} の要素をいくつ含んでいるかを数え上げ、共通要素 V_{ij} の要素数で割る方法が挙げられる。その他にも、共通要素 V_{ij} の要素が重要度によって重み付けされている場合、重み付けされた重要度の高さに応じて信頼度を算出しても良い。何れにしろ、この信頼度が、所定の値より低い場合、得られたドキュメントベクトル V'_{ij} のベクトルサイズを増減する等の信頼度のフィードバックを行う。

10

【0042】

ステップ 210 では、共通要素 V_{ij} を生成した際のドキュメントベクトル V_i, V_j を、特徴量記憶部 112 から削除し、新たに得られたドキュメントベクトル V'_{ij} を特徴量記憶部 112 に記憶させる。

【0043】

ステップ 211 では、本実施例のパスセージ分割方法のために、パスセージ更新部 108 にて、 V_i, V_j に対応する二つの文またはパスセージ候補を連結する。一度も連結されていない文は文記憶部 110 に記憶されている。文が連結された場合、連結前の文を文記憶部 110 から削除する。パスセージ候補と文が連結された場合、あるいはパスセージ候補同士が連結された場合には、連結前の文の削除のみならず、連結前のパスセージ候補をパスセージ記憶部 113 から削除する。連結された文またはパスセージ候補は新たなパスセージ候補としてパスセージ記憶部 113 に記憶する。

20

【0044】

本実施例のパスセージ分割方法、装置においては、図 2 のフローにおいて、ステップ 204 からステップ 211 を繰り返すことで、目的とするパスセージを作成する。そして、ステップ 205 において、二つのドキュメントベクトルの最大類似度が所定の閾値未満の場合、パスセージの作成を終了するため、ステップ 212 を実行する。

30

【0045】

ステップ 212 は、出力部 109 にて実行され、不明パスセージの判定とパスセージの出力を行うステップである。不明パスセージの判定方法の一例として、文またはパスセージ候補の中に含まれる形態素数を調べる方法がある。文またはパスセージ候補の中に含まれる形態素数が少ない場合、ドキュメントベクトルが適切に作成されず、連結が難しい場合がある。よって、ステップ 21 において、残された文またはパスセージ候補に含まれる形態素数がある閾値以下の場合、出力部 409 は、不明パスセージのラベルをつけて出力し、処理フローを終了する。

【0046】

図 3 は本実施例において、ドキュメントベクトルの類似度に応じて、文が連結されていく様子を模式的に示した一例である。図 2 のステップ 205 における閾値は“10”とする。

40

一度目の類似度算出結果が 301 である。結果 301 の中で最も類似度が高いのは、 a_2 と a_3 の組の類似度 40 である。

【0047】

よってこの組に対して図 2 のステップ 205 からステップ 211 の処理を行い、再度図 2 のステップ 204 に戻る。連結された結果を a_{23} と表す。同様に結果 302 では b_1 と b_2 、結果 303 では a_1 と a_{23} が類似度の最も高い組として選定され、図 2 のステップ 205 から図 2 のステップ 211 の処理が行われる。閾値を 10 と設定したので、結

50

果 3 0 4 で選ばれる組はなく、パッセージの作成が完了する。

【 0 0 4 8 】

以上詳述した実施例 1 によれば、意味の近い文、すなわち、特徴量が似た文を含む複数のパッセージが、一つの文書に含まれる場合でも、複数のパッセージを正しく分割することが可能となり、更には、文書の自動要約や文書検索のための自動キーワード抽出など。

【 実施例 2 】

【 0 0 4 9 】

実施例 2 は類似度計算に単語ベクトルを、類似文書検索にも単語ベクトルを用いたパッセージ分割方法、装置、及びプログラムの実施例である。

図 4 は実施例 2 に係るパッセージ分割装置の機能ブロック図である。同図のパッセージ分割装置のハードウェア構成も、実施例 1 の図 1 A の装置同様、図 1 B に示したコンピュータ等で実現できることは言うまでもなく、ここではハードウェア構成の図示説明を省略する。

10

【 0 0 5 0 】

入力部 4 0 2 と、文分割部 4 0 3 と、パッセージ更新部 4 0 8 と、出力部 4 0 9 と、文記憶部 4 1 0 と、特徴量記憶部 4 1 2 と、パッセージ記憶部 4 1 3 と、形態素解析部 4 1 4 とは実施例 1 の対応するブロックと共通であるので、実施例 1 と異なる、コーパス部 4 1 1 と、特徴量算出部 4 0 4 と、類似度計算部 4 0 5 と、検索クエリ生成部 4 0 6 と、特徴量更新部 4 0 7 についてのみ説明する。なお、形態素解析部 4 1 4 は特徴量算出部 4 0 4 に接続される。

20

【 0 0 5 1 】

コーパス部 4 1 1 には、例えば新聞記事などのドキュメントの集合やシソーラス、あるいはその両方を用いる。

【 0 0 5 2 】

特徴量算出部 4 0 4 は、文記憶部 4 1 0 から読み込んだ文に対し、形態素解析部 4 1 4 を用いて形態素解析を行い、文を単語ベクトルへ変換する。単語ベクトルの要素数が十分でない場合にはコーパス部 4 1 1 を使用して要素数を増やす方法が有効である。例えばコーパスとしてシソーラスを用いた場合、入力文から得られた各単語をクエリとして類義語を検索し、結果として得られた類義語を単語ベクトルに追加する。またコーパスとしてドキュメントの集合を用いた場合、入力文から得られた単語ベクトルに、コーパス内の各ドキュメントから抽出した単語ベクトルを追加することができる。

30

【 0 0 5 3 】

単語ベクトルの要素を追加する方法の他の例として、上位数件のドキュメントから T F I D F 等を用いて重要語を抜き出し、単語ベクトルに追加する方法が挙げられる。これに限らず、他の方法で文に関連する単語を得て追加して、単語ベクトルの要素数を十分にしてもよい。そして、得られた単語ベクトルを特徴量記憶部 4 1 2 に記憶する。また検索クエリ生成部 4 0 6 から、制御部 4 0 1 を介して単語ベクトルが特徴量算出部 4 0 4 に与えられた場合も、同様の方法で単語ベクトルの要素数を拡充し、特徴量記憶部 1 1 2 に記憶すると共に、制御部 4 0 1 を介して特徴量更新部 4 0 7 へ単語ベクトルを出力する。

40

【 0 0 5 4 】

本実施例の類似度計算部 4 0 5 は、制御部 4 0 1 の指定に基づいて、二つの単語ベクトルを特徴量記憶部 4 1 2 から読み出し、二つの単語ベクトルの類似度を計算する。類似度の計算方法としては、例えば、上述したコサイン尺度等が挙げられる。

【 0 0 5 5 】

本実施例の検索クエリ生成部 4 0 6 は、制御部 4 0 1 の指定に基づいて、二つの単語ベクトルを特徴量記憶部 4 1 2 から読み出し、二つの単語ベクトルに共通する単語群をコーパス 4 1 1 から抽出する。抽出された共通する単語群から単語ベクトルを作成し、制御部 4 0 1 を介して特徴量算出部 4 0 4 に出力する。

【 0 0 5 6 】

特徴量更新部 4 0 7 は、制御部 4 0 1 の指定に基づいて二つの単語ベクトル V_i , V_j

50

を特徴量記憶部 4 1 2 から読み出す。また制御部 4 0 1 から一つの単語ベクトル V_k が入力される。入力された三つの単語ベクトル V_k, V_i, V_j から信頼度を計算し、信頼度に基づいて V_k のベクトルサイズを修正する。その後 V_i, V_j を特徴量記憶部 4 1 2 から削除し、 V_k を特徴量記憶部 4 1 2 に記憶する。

【 0 0 5 7 】

図 5 は実施例 2 に係るプログラムの動作を示した処理フロー図である。実施例 1 では、類似度計算としてドキュメントベクトルを用いているが、実施例 2 では上述の通り、単語ベクトルを用いており、その点が実施例 1 と異なるが、それ以外の動作は実施例 1 と同様である。

【 0 0 5 8 】

実施例 2 によれば、意味の近い文、即ち、特徴量が似た文を含む複数のパッセージが、一つの文書に含まれる場合でも、パッセージを正しく分割することが可能となる。

【 0 0 5 9 】

なお、本発明は上記した実施例に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施例は本発明を分かりやすく説明するために詳細に説明したのであり、必ずしも説明の全ての構成を備えるものに限定されものではない。また、ある実施例の構成に他の実施例の構成を加えることが可能である。また、各実施例の構成の一部について、他の構成の追加・削除・置換をすることが可能である。

【 0 0 6 0 】

上記の各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、上記の各構成、機能等は、それぞれの機能を実現するプログラムを実行することによりソフトウェアで実現する場合を例示して説明したが、各機能を実現するプログラム、テーブル、ファイル等の情報はメモリのみならず、ハードディスク、SSD (Solid State Drive) 等の記憶装置、または、ICカード、SDカード、DVD等の記録媒体におくことができるし、必要に応じてネットワーク等を介してダウンロード、インストールすることも可能である。

【 符号の説明 】

【 0 0 6 1 】

- 1 1 CPU
- 1 2 記憶部
- 1 3 入出力部
- 1 4 通信部
- 1 0 0、4 0 0 パッセージ分割装置
- 1 0 1、4 0 1 制御部
- 1 0 2、4 0 2 入力部
- 1 0 3、4 0 3 文分割部
- 1 0 4、4 0 4 特徴量算出部
- 1 0 5、4 0 5 類似度計算部
- 1 0 6、4 0 6 検索クエリ生成部
- 1 0 7、4 0 7 特徴量更新部
- 1 0 8、4 0 8 パッセージ更新部
- 1 0 9、4 0 9 出力部
- 1 1 0、4 1 0 文記憶部
- 1 1 1、4 1 1 コーパス部
- 1 1 2、4 1 2 特徴量記憶部
- 1 1 3、4 1 3 パッセージ記憶部
- 1 1 4、4 1 4 形態素解析部

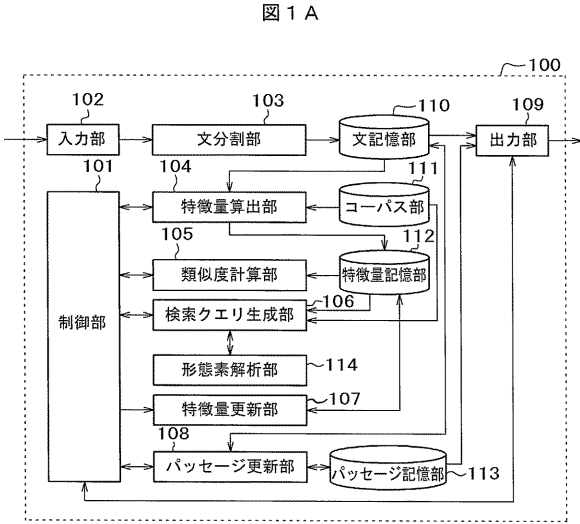
10

20

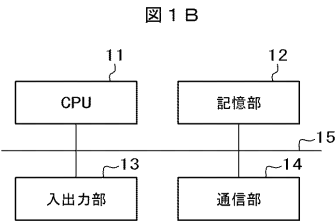
30

40

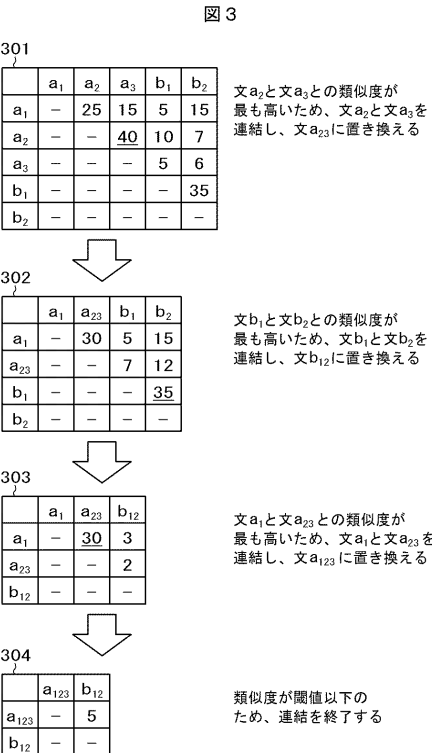
【図1A】



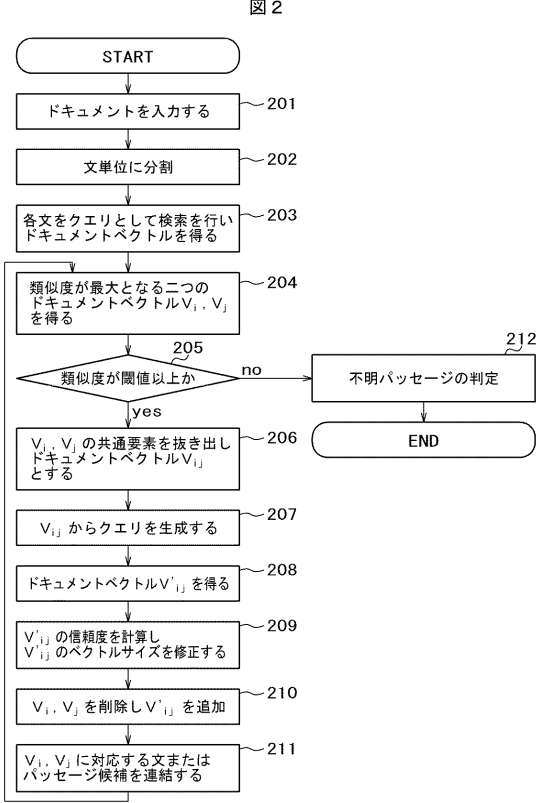
【図1B】



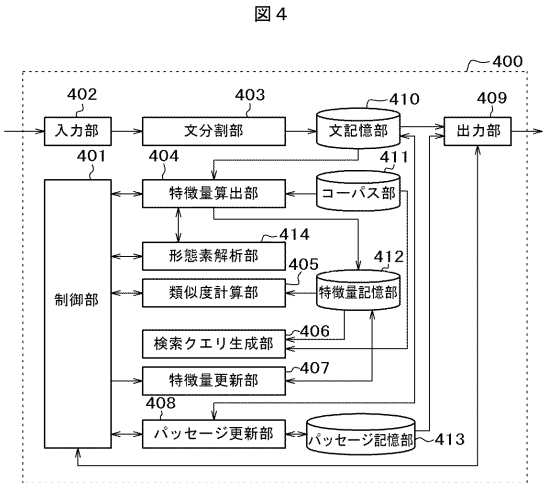
【図3】



【図2】

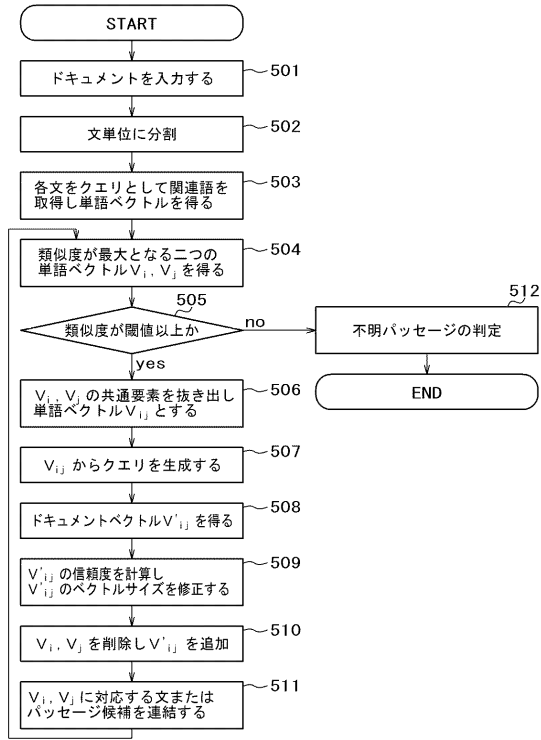


【図4】



【図5】

図5



【図6】

図6

コーパスに含まれるドキュメントの総数=10とする。
 検索の結果得られたドキュメントがドキュメント1, 3, 4, 8である場合 (a)のように該当する要素に1を代入することでドキュメントベクトルを得る。
 検索スコアが得られる場合 (b)のようにスコアを代入してもよい。ここでは以下のようにスコアが算出されたとする。
 ドキュメント1のスコア=0.1
 ドキュメント3のスコア=0.5
 ドキュメント4のスコア=0.4
 ドキュメント8のスコア=0.9

601		602	
0	ドキュメント0	0	ドキュメント0
1	ドキュメント1	0.1	ドキュメント1
0	ドキュメント2	0	ドキュメント2
1	ドキュメント3	0.5	ドキュメント3
1	ドキュメント4	0.4	ドキュメント4
0	ドキュメント5	0	ドキュメント5
0	ドキュメント6	0	ドキュメント6
0	ドキュメント7	0	ドキュメント7
1	ドキュメント8	0.9	ドキュメント8
0	ドキュメント9	0	ドキュメント9

【図7】

図7

全てのドキュメントに出現する単語の種類を10とする。
 あるドキュメントに含まれる単語が単語3, 6, 7, 8であり、以下のように出現頻度が求められたとする。
 単語3の出現頻度=1
 単語6の出現頻度=5
 単語7の出現頻度=3
 単語8の出現頻度=9
 該当する要素に出現頻度を代入することで単語ベクトルを得る。

701	
0	単語0
0	単語1
0	単語2
1	単語3
0	単語4
0	単語5
5	単語6
3	単語7
9	単語8
0	単語9

フロントページの続き

(72)発明者 村上 智一

東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内

(72)発明者 今一 修

東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内

審査官 早川 学

(56)参考文献 特開2004-145790(JP,A)

望月源、外2名、語彙的連鎖に基づくパッセージ検索、自然言語処理、言語処理学会、1999年4月10日、第6巻、第3号、pp.101~126

(58)調査した分野(Int.Cl., DB名)

G06F 17/27

G06F 17/21

G06F 17/30