

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5143057号
(P5143057)

(45) 発行日 平成25年2月13日(2013.2.13)

(24) 登録日 平成24年11月30日(2012.11.30)

(51) Int.Cl.

F I

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 2 1 0 A

G 0 6 F 17/30 1 7 0 A

請求項の数 9 (全 28 頁)

(21) 出願番号 特願2009-48550 (P2009-48550)
 (22) 出願日 平成21年3月2日(2009.3.2)
 (65) 公開番号 特開2010-204866 (P2010-204866A)
 (43) 公開日 平成22年9月16日(2010.9.16)
 審査請求日 平成22年5月28日(2010.5.28)

(73) 特許権者 000004226
 日本電信電話株式会社
 東京都千代田区大手町二丁目3番1号
 (74) 代理人 100070150
 弁理士 伊東 忠彦
 (74) 代理人 100124844
 弁理士 石原 隆治
 (72) 発明者 近藤 光正
 東京都千代田区大手町二丁目3番1号 日
 本電信電話株式会社内
 (72) 発明者 中辻 真
 東京都千代田区大手町二丁目3番1号 日
 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 重要キーワード抽出装置及び方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

Web 文書に含まれるテキストから重要なキーワードを抽出する重要キーワード抽出装置であって、

前記 Web 文書を取得して該 Web 文書の特徴量を抽出し、主要コンテンツを抽出する主要コンテンツ抽出手段と、

オンライン百科辞典に代表される、文書集合内において一意の見出し語を持つ文書で、かつ、文書集合内において引用関係もしくは参照関係を持つ文書集合（以下、辞典文書集合と記す）の見出し語を形態素解析用の辞書として登録しているキーワード辞書を参照して、前記主要コンテンツから重要キーワード候補を抽出する重要キーワード候補抽出手段と、

前記 Web 文書内の前記重要キーワード候補に重みを付ける出現頻度算出手段と、

知名度や話題性が高く、キーワードについて語る際に様々な話題が挙がる情報量が豊富なキーワードの重要度が高くなるように前記重要キーワード候補の固有重要度を算出し、第1の記憶手段に格納するキーワード重要度算出手段と、

前記重要キーワード候補の前記 Web 文書中の位置に基づいて出現位置キーワード重要度を求め、第2の記憶手段に格納する位置情報算出手段と、

前記出現頻度算出手段による重要キーワード候補の重み、前記第1の記憶手段に格納されている重要キーワード候補の固有重要度及び前記第2の記憶手段に格納されている前記出現位置キーワード重要度を乗算した値に基づいて、最終重要度付きのキーワード集合を

出力するキーワード出力手段と、を有し、

前記主要コンテンツ抽出手段は、

前記Web文書を所定の分割規則に基づいてセグメントに分割するWeb文書分割手段と、

前記セグメント毎に主要コンテンツ判定のための特徴量を抽出する特徴量抽出手段と、

前記セグメント毎の特徴量に基づいて機械学習アルゴリズムを用いて主要コンテンツか否かの判定を行う主要コンテンツ判定手段と、

主要コンテンツと判断された部位を結合して前記主要コンテンツとして出力する主要コンテンツ出力手段を含み、

前記特徴量抽出手段は、

前記Web文書で表示される文字列の特徴量、タグ情報の特徴量、アンカーリンク情報の特徴量を抽出し、各特徴量間の比率を用いて最終的な特徴量とする手段を含み、

前記出現頻度算出手段は、

検索エンジンが収集したWeb文書集合における出現確率を考慮した重みであるWeb IDFを用いて前記重要キーワード候補に重みを付ける手段、

または、

前記重要キーワード候補の前記Web文書内の出現頻度を求め、該出現頻度と文書集合内での出現分布に基づいて各重要キーワード候補に重み付けを行うBM25を用いて該重要キーワード候補に重みを付ける手段、

のいずれかを有することを特徴とする重要キーワード抽出装置。

【請求項2】

前記出現頻度算出手段は、

形態素解析手段の形態素解析結果から得られた形態素の固有名詞を含む前記重要キーワード候補の出現頻度を用いて、該重要キーワード候補に重みを付ける手段を含む

請求項1記載の重要キーワード抽出装置。

【請求項3】

前記キーワード重要度算出手段は、

前記辞典文書集合内のリンク構造を用いて、前記重要キーワード候補の固有重要度を求める手段を含む請求項1または2記載の重要キーワード抽出装置。

【請求項4】

前記キーワード重要度算出手段は、

検索エンジンに投入された回数が多いキーワードほど、重要なキーワードとなるような値を用いて、前記重要キーワード候補の固有重要度を求める手段を含む請求項1乃至3のいずれか1項記載の重要キーワード抽出装置。

【請求項5】

Web文書中に含まれるテキストから重要なキーワードを抽出する重要キーワード抽出方法であって、

主要コンテンツ抽出手段が、前記Web文書を取得して該Web文書の特徴量を抽出し、主要コンテンツを抽出する主要コンテンツ抽出ステップと、

重要キーワード候補抽出手段が、オンライン百科辞典に代表される、文書集合内において一意の見出し語を持つ文書で、かつ、文書集合内において引用関係もしくは参照関係を持つ文書集合（以下、辞典文書集合と記す）の見出し語を形態素解析用の辞書として登録しているキーワード辞書を参照して、前記主要コンテンツから重要キーワード候補を抽出する重要キーワード候補抽出ステップと、

出現頻度算出手段が、前記Web文書内の前記重要キーワード候補に重みを付ける出現頻度算出ステップと、

キーワード重要度算出手段が、知名度や話題性が高く、キーワードについて語る際に様々な話題が挙がる情報量が豊富なキーワードの重要度が高くなるように前記重要キーワード候補の固有重要度を算出し、第1の記憶手段に格納するキーワード重要度算出ステップと、

10

20

30

40

50

位置情報算出手段が、前記重要キーワード候補の前記W e b 文書中の位置に基づいて出現位置キーワード重要度を求め、第2の記憶手段に格納する位置情報算出ステップと、

キーワード出力手段が、前記出現頻度算出ステップによる重要キーワード候補の重み、前記第1の記憶手段に格納されている重要キーワード候補の固有重要度及び前記第2の記憶手段に格納されている前記出現位置キーワード重要度を乗算した値に基づいて、最終重要度付きのキーワード集合を出力するキーワード出力ステップと、を行い、

前記主要コンテンツ抽出ステップでは、

前記W e b 文書を所定の分割規則に基づいてセグメントに分割するW e b 文書分割ステップと、

前記セグメント毎に主要コンテンツ判定のための特徴量を抽出する特徴量抽出ステップと、

前記セグメント毎の特徴量に基づいて機械学習アルゴリズムを用いて主要コンテンツか否かを判定する主要コンテンツ判定ステップと、

主要コンテンツと判断された部位を結合して前記主要コンテンツとして出力する主要コンテンツ出力ステップと、を含み、

前記特徴量抽出ステップでは、

前記W e b 文書で表示される文字列の特徴量、タグ情報の特徴量、アンカーリンク情報の特徴量を抽出し、各特徴量間の比率を用いて最終的な特徴量とするステップを含み、

前記出現頻度算出ステップでは、

検索エンジンが収集したW e b 文書集合における出現確率を考慮した重みであるW e b IDFを用いて前記重要キーワード候補に重みを付けるステップ、

または、

前記重要キーワード候補の前記W e b 文書内の出現頻度を求め、該出現頻度と文書集合内での出現分布に基づいて各重要キーワード候補に重み付けを行うBM25を用いて該重要キーワード候補に重みを付けるステップ、

のいずれかを行うことを特徴とする重要キーワード抽出方法。

【請求項6】

前記出現頻度算出ステップは、

形態素解析手段の形態素解析結果から得られた形態素の固有名詞を含む前記重要キーワード候補の出現頻度を用いて、該重要キーワード候補に重みを付ける

請求項5記載の重要キーワード抽出方法。

【請求項7】

前記キーワード重要度算出ステップは、

前記辞典文書集合内のリンク構造を用いて、前記重要キーワード候補の固有重要度を求める

請求項5または6記載の重要キーワード抽出方法。

【請求項8】

前記キーワード重要度算出ステップは、

検索エンジンに投入された回数が多いキーワードほど、重要なキーワードとなるような値を用いて、前記重要キーワード候補の固有重要度を求める

請求項5乃至7のいずれか1項記載の重要キーワード抽出方法。

【請求項9】

請求項1乃至4のいずれか1項記載の重要キーワード抽出装置を構成する各手段としてコンピュータを機能させるための重要キーワード抽出プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、重要キーワード抽出装置及び方法及びプログラムに係り、特に、文書に含まれるテキストから重要なキーワードを抽出する重要キーワード抽出装置及び方法及びプログラムに関する。

10

20

30

40

50

【背景技術】

【0002】

従来の重要キーワード抽出法は、キーワードの出現頻度や複合語を構成する形態素の連接頻度等の"頻度"を用いて抽出する手法が主流であった(例えば、非特許文献1参照)。

【先行技術文献】

【非特許文献】

【0003】

【非特許文献1】中川裕志、森辰則、湯本紘彰、「出現頻度と連接頻度に基づく専門用語抽出。自然言語処理」Vol.10 No.1, pp.27-45, 2003年1月

【発明の概要】

10

【発明が解決しようとする課題】

【0004】

従来の重要キーワード抽出法は、キーワードの出現頻度や連接頻度等の"頻度"を用いる手法が主流であったため、事前に比較的量の多い文書集合を要する必要があった。また、頻度が高いからといって、重要であるという判定を行う手法は場合によっては上手く作用しない場合が多く、精度的にもあまりよい結果を得られていない。

【0005】

本発明は、上記の点に鑑みなされたもので、キーワードそのものの固有の重要性を算出でき、文書内の出現頻度等に左右されずに重要キーワードを頑健に抽出することが可能な重要キーワード抽出装置及び方法及びプログラムを提供することを目的とする。

20

【課題を解決するための手段】

【0006】

図1は、本発明の原理構成図である。

【0007】

本発明(請求項1)は、Web文書に含まれるテキストから重要なキーワードを抽出する重要キーワード抽出装置であって、

Web文書を取得して該Web文書の特徴量を抽出し、主要コンテンツを抽出する主要コンテンツ抽出手段10と、

オンライン百科辞典(例えば、Wikipedia(登録商標))に代表される文書集合内において一意の見出し語を持つ文書で、かつ、文書集合内において引用関係もしくは参照関係を持つ文書集合(以下、辞典文書集合と記す)の見出し語を形態素解析用の辞書として登録しているキーワード辞書101を参照して、主要コンテンツから重要キーワード候補を抽出する重要キーワード候補抽出手段2と、

30

Web文書内の重要キーワード候補に重みを付ける出現頻度算出手段3と、

知名度や話題性が高く、内容的に興味深く、キーワードについて語る際に様々な話題が挙がる情報量が豊富なキーワードの重要度が高くなるように重要キーワード候補の固有重要度を算出し、第1の記憶手段105に格納するキーワード重要度算出手段4と、

重要キーワード候補のWeb文書中の位置に基づいて出現位置キーワード重要度を求め、第2の記憶手段106に格納する位置情報算出手段5と、

出現頻度算出手段3による重要キーワード候補の重み、第1の記憶手段105に格納されている重要キーワード候補の固有重要度及び第2の記憶手段106に格納されている出現位置キーワード重要度を乗算した値に基づいて、最終重要度付きのキーワード集合を出力するキーワード出力手段6と、を有し、

40

主要コンテンツ抽出手段10は、

Web文書を所定の分割規則に基づいてセグメントに分割するWeb文書分割手段12と、

セグメント毎に主要コンテンツ判定のための特徴量を抽出する特徴量抽出手段13と、

セグメント毎の特徴量に基づいて機械学習アルゴリズムを用いて主要コンテンツか否かの判定を行う主要コンテンツ判定手段14と、

主要コンテンツと判断された部位を結合して主要コンテンツとして出力する主要コンテ

50

ンツ出力手段 15 を含み、

主要コンテンツ抽出手段 10 の特徴量抽出手段 13 は、

Web 文書で表示される文字列の特徴量、タグ情報の特徴量、アンカーリンク情報の特徴量を抽出し、各特徴量間の比率を用いて最終的な特徴量とする手段を含み、

出現頻度算出手段 3 は、

検索エンジンが収集した Web 文書集合における出現確率を考慮した重みである Web IDF を用いて重要キーワード候補に重みを付ける手段、

または、

重要キーワード候補の Web 文書内の出現頻度を求め、該出現頻度と文書集合内での出現分布に基づいて各重要キーワード候補に重み付けを行う BM25 を用いて該重要キーワード候補に重みを付ける手段、のいずれかを有する。

10

【0008】

また、本発明（請求項 2）は、上記請求項 1 の出現頻度算出手段 3 が、

形態素解析手段の形態素解析結果から得られた形態素の固有名詞を含む重要キーワード候補の出現頻度を用いて、該重要キーワード候補に重みを付ける手段を含む。

【0009】

また、本発明（請求項 3）は、上記請求項 1 または 2 のキーワード重要度算出手段 4 が、

辞典文書集合内のリンク構造を用いて、重要キーワード候補の固有重要度を求める手段を含む。

20

【0010】

また、本発明（請求項 4）は、上記請求項 1 乃至 3 のキーワード重要度算出手段 4 が、

検索エンジンに投入された回数が多いキーワードほど、重要なキーワードとなるような値を用いて、重要キーワード候補の固有重要度を求める手段を含む。

【0017】

図 2 は、本発明の原理を説明するための図である。

【0018】

本発明（請求項 5）は、Web 文書中に含まれるテキストから重要なキーワードを抽出する重要キーワード抽出方法であって、

主要コンテンツ抽出手段が、Web 文書を取得して該 Web 文書の特徴量を抽出し、主要コンテンツを抽出する主要コンテンツ抽出ステップ（ステップ 1）と、

30

重要キーワード候補抽出手段が、オンライン百科辞典（例えば、Wikipedia（登録商標））に代表される文書集合内において一意の見出し語を持つ文書で、かつ、文書集合内において引用関係もしくは参照関係を持つ文書集合（以下、辞典文書集合と記す）の見出し語を形態素解析用の辞書として登録しているキーワード辞書を参照して、主要コンテンツから重要キーワード候補を抽出する重要キーワード候補抽出ステップ（ステップ 2）と、

出現頻度算出手段が、Web 文書内の重要キーワード候補に重みを付ける出現頻度算出ステップ（ステップ 3）と、

キーワード重要度算出手段が、知名度や話題性が高く、キーワードについて語る際に様々な話題が挙がる情報量が豊富なキーワードの重要度が高くなるように重要キーワード候補の固有重要度を算出し、第 1 の記憶手段に格納するキーワード重要度算出ステップ（ステップ 4）と、

40

位置情報算出手段が、重要キーワード候補の Web 文書中の位置に基づいて出現位置キーワード重要度を求め、第 2 の記憶手段に格納する位置情報算出ステップ（ステップ 5）と、

キーワード出力手段が、出現頻度算出ステップによる重要キーワードの重み、第 1 の記憶手段に格納されている重要キーワード候補の固有重要度及び第 2 の記憶手段に格納されている出現位置キーワード重要度を乗算した値に基づいて、最終重要度付きのキーワード集合を出力するキーワード出力ステップ（ステップ 6）と、を行い、

主要コンテンツ抽出ステップ（ステップ 1）では、

50

Web文書を所定の分割規則に基づいてセグメントに分割するWeb文書分割ステップと、

セグメント毎に主要コンテンツ判定のための特徴量を抽出する特徴量抽出ステップと、セグメント毎の特徴量に基づいて機械学習アルゴリズムを用いて主要コンテンツが否かを判定する主要コンテンツ判定ステップと、

主要コンテンツと判断された部位を結合して主要コンテンツとして出力する主要コンテンツ出力ステップと、を含み、

特徴量抽出ステップでは、

Web文書で表示される文字列の特徴量、タグ情報の特徴量、アンカーリンク情報の特徴量を抽出し、各特徴量間の比率を用いて最終的な特徴量とするステップを含み、

出現頻度算出ステップ(ステップ3)では、

検索エンジンが収集したWeb文書集合における出現確率を考慮した重みであるWeb IDFを用いて重要キーワード候補に重みを付けるステップ、

または、

重要キーワード候補のWeb文書内の出現頻度を求め、該出現頻度と文書集合内での出現分布に基づいて各重要キーワード候補に重み付けを行うBM25を用いて該重要キーワード候補に重みを付けるステップ、

のいずれかを行う。

【0019】

また、本発明(請求項6)は、請求項11の出現頻度算出ステップ(ステップ3)において、

形態素解析手段の形態素解析結果から得られた形態素の固有名詞を含む重要キーワード候補の出現頻度を用いて、該重要キーワード候補に重みを付ける。

【0020】

また、本発明(請求項7)は、請求項11または12のキーワード重要度算出ステップ(ステップ4)において、

辞典文書集合内のリンク構造を用いて、重要キーワード候補の固有重要度を求める。

【0021】

また、本発明(請求項8)は、請求項11乃至13のキーワード重要度算出ステップ(ステップ4)において、

検索エンジンに投入された回数が多いキーワードほど、重要なキーワードとなるような値を用いて、重要キーワード候補の固有重要度を求める。

【0023】

本発明(請求項9)は、請求項1乃至4のいずれか1項記載の重要キーワード抽出装置を構成する各手段としてコンピュータを機能させるための重要キーワード抽出プログラムである。

【発明の効果】

【0024】

本発明によれば、文書中の重要キーワードの抽出が可能となり、それらを用いた文書管理システムや、重要キーワードからの広告配信等での応用が可能である。さらに、従来の技術では、文書中におけるキーワードの出現頻度のみから重要キーワードの特定を行っていたが、オンライン百科事典(例えば、Wikipedia(登録商標))のリンク構造や実際の検索クエリの投入回数を用いることで、従来の技術では実現できなかったキーワードそのものの固有の重要性を算出でき、文書内の出現頻度等に左右されない重要キーワードの頑健な抽出が可能になる。

【図面の簡単な説明】

【0025】

【図1】本発明の原理構成図である。

【図2】本発明の原理を説明するための図である。

【図3】本発明の一実施の形態における重要キーワード抽出装置の構成図である。

10

20

30

40

50

【図4】本発明の処理の概要を示す図である。

【図5】本発明の一実施の形態における主要コンテンツ抽出部の構成図である。

【図6】本発明の一実施の形態におけるWeb文書取得・入力部の構成図である。

【図7】本発明の一実施の形態におけるWeb文書分割部の構成図である。

【図8】本発明の一実施の形態における特徴量抽出部の構成図である。

【図9】本発明の一実施の形態における主要コンテンツ判定部の構成図である。

【図10】本発明の一実施の形態における主要コンテンツ出力部の構成図である。

【図11】本発明の一実施の形態における特徴量のパラメータ推定方法のフローチャートである。

【図12】本発明の一実施の形態における主要コンテンツ例(その1)である。

10

【図13】本発明の一実施の形態における主要コンテンツ例(その2)である。

【発明を実施するための形態】

【0026】

以下、図面と共に本発明の実施の形態を説明する。

【0027】

本発明は、文書中における重要キーワードを抽出するシステムである。従来の技術においては比較的量の多い文書集合を事前に用意する必要があったが、本発明では文書が少ない場合においても精度が落ちにくい手法を提案する。

【0028】

本発明では、「文書集合内において一意の見出し語を持つ文書で、かつ文書集合内において引用関係もしくは参照関係を持つ文書集合」を用いて、重要キーワードを抽出する。例をあげると、オンライン百科事典である「Wikipedia(登録商標)」や、「はてなブックマーク(登録商標)」、「マイペディア(登録商標)」等の文書集合がそれに該当する。キーワード候補は、これらの文書集合における文書の見出し語を用いる。百科事典のような一意の見出し語を有する文書集合の見出し語を用いることで事象を一意に示す複数形態素から構成される複合語の切り出しが可能になる。さらに、このような見出し語は一般的に重要だと思われるキーワードを網羅しているため、掲載されている見出し語をキーワード候補として用いることで、重要なキーワードの絞り込みが可能になる。

20

【0029】

本実施の形態において、対象とする文書集合は「文書集合内において一意の見出し語を持つ文書で、かつ文書集合内において引用関係もしくは参照関係を持つ文書集合」とするが、以下では、Wikipedia(登録商標)を例に挙げて説明する。

30

【0030】

図3は、本発明の一実施の形態における重要キーワード抽出装置の構成図である。

【0031】

同図に示す重要キーワード抽出装置は、主要コンテンツ抽出部10を有する文書入力部1、重要キーワード候補抽出部2、出現頻度算出部3、キーワード重要度算出部4、位置情報算出部5、重要キーワード出力部6、キーワード辞書101、WebIDF記憶部102、事前に求められているWikipedia(登録商標)のリンク構造からのキーワード重要度を格納するキーワード重要度1記憶部103、事前に求められている検索エンジンの検索クエリ投入回数を用いたキーワード重要度を格納するキーワード重要度2記憶部104、固有重要度記憶部105、位置重要度記憶部106から構成される。

40

【0032】

図4は、本発明の処理の概要を示す図である。

【0033】

ステップ100) 文書入力部1の主要コンテンツ抽出部10は、入力された対象文書がWeb文書であるかを判定し、Web文書である場合はステップ200に移行し、Web文書でない場合はステップ300に移行する。

【0034】

ステップ200) Web文書である場合は、主要コンテンツの抽出を行う(主要コン

50

テンツの抽出処理)。

【0035】

ステップ300) 重要キーワード候補抽出部2は、文書内のテキストからキーワード候補を抽出する(キーワード候補の抽出処理)。

【0036】

ステップ400) 出現頻度算出部3において文書中に出現するキーワード候補の出現頻度を求め、キーワード重要度算出部4ではキーワード候補の重要度(固有重要度)を求め、位置情報算出部5は、キーワード候補が文書中に出現する位置により重要度(位置重要度)を求め、重要キーワード出力部6において、固有重要度と位置重要度に基づいて、キーワードの順序付けを行う(キーワードの重要度算出処理)。

10

【0037】

以上が、文書中における重要キーワードを抽出する処理の主な流れである。

【0038】

[1]主要コンテンツの抽出処理(ステップ200)：

当該処理は、文書入力部1の主要コンテンツ抽出部10が行う処理である。

【0039】

入力文書がWeb文書である場合、ナビゲーションリンクや広告テキスト等のWeb文書の内容とは関係のないテキストが存在する。そのため、それら不要テキストを除去する主要コンテンツ部分の抽出を行う必要がある。Web文書か否かの判定は、ファイルの拡張子を用いて行うものとする。

20

【0040】

[1-1]処理の流れ：

文書入力部1の主要コンテンツ抽出部10は、図5に示すように、Web文書取得・入力部11、Web文書分割部12、特徴量抽出部13、主要コンテンツ判定部14、主要コンテンツ出力部15から構成される。ここで、特徴量抽出部13と主要コンテンツ判定部14は、Web文書分割部12で分割されたWeb文書毎に処理を行う。

【0041】

Web文書取得・入力部11は、処理するWeb文書(データ)の入力を行うもので、図6に示すように、データ入力部111、Web文書ファイル入力部112、URL入力部113、Web文書取得部114、文字コード変換部115から構成される。

30

【0042】

Web文書分割部12は、取得したWeb文書を分割するものであり、図7に示すように、広告対象領域部121、ノイズとなるタグや領域除去部122、Web文書の分割部123を有する。広告対象領域部121では、Web文書に広告対象領域が存在する場合は、当該広告対象領域を抽出する。ノイズとなるタグや領域除去部122は、広告対象領域または入力されたWeb文書のノイズとなるタグや領域を除去し、Web文書分割部123において文書を分割する。

【0043】

特徴量抽出部13は、Web文書分割部12で分割されたWeb文書毎に、主要コンテンツ判定のための特徴量を抽出するものであり、図8に示すように、アンカーリンク情報特徴量抽出部131、タグ情報特徴量抽出部132、Web文書で表示される文字列特徴量抽出部133、特徴量の正規化部134、特徴量の比率特徴量抽出部135から構成される。

40

【0044】

主要コンテンツ判定部14は、特徴量抽出部13で抽出された特徴量から主要コンテンツか否かの判定を行うもので、図9に示すように、特徴量入力部141、テキスト判定部142、主要コンテンツ判定部143から構成される。

【0045】

主要コンテンツ出力部15は、図10に示すように、タグ付きテキストを要求された場合に用いられるタグ付テキスト出力部151、タグなしテキストを要求された場合に用い

50

られるタグなしテキスト出力部 1 5 2、タグ付テキスト出力部 1 5 1 またはタグなしテキスト出力部 1 5 2 から出力された主要コンテンツと判定されたセグメントを結合して出力するデータ出力部 1 5 3 から構成される。

【 0 0 4 6 】

主要コンテンツ抽出部 1 0 の処理は、概略以下のようになる。

【 0 0 4 7 】

(1) W e b 文書取得・入力部 1 1 が、処理する W e b 文書 (データ) の入力を行う。

【 0 0 4 8 】

(2) W e b 文書分割部 1 2 が、(1) で取得した W e b 文書を分割する。

【 0 0 4 9 】

以下の処理は、分割した W e b 文書毎に行う。

【 0 0 5 0 】

(3) 特徴量抽出部 1 3 が主要コンテンツ判定のための特徴量を抽出する。

【 0 0 5 1 】

(4) 主要コンテンツ判定部 1 4 が、(3) で抽出した特徴量から主要コンテンツか否かの判定を行う。

【 0 0 5 2 】

(5) 主要コンテンツ出力部 1 5 が、(4) で主要コンテンツと判定された部位を結合して最終出力とする。

【 0 0 5 3 】

以下に主要コンテンツ抽出部 1 0 の動作を詳細に説明する。

【 0 0 5 4 】

[1 - 2] W e b 文書の入力 :

W e b 文書の入力は、W e b 文書取得・入力部 1 1 で行われる処理である。

【 0 0 5 5 】

W e b 文書取得・入力部 1 1 のデータ入力部 1 1 1 は、ユーザから入力された主要コンテンツを抽出した W e b 文書の U R L、もしくは、ファイルそのものを取得する。入力が U R L の場合は U R L 入力部 1 1 3 及び W e b 文書取得部 1 1 4 において、その U R L 先の W e b 文書を取得し、ファイルが直接入力された場合は W e b 文書ファイル入力部 1 1 2 がそのファイルを取得する。文字コード変換部 1 1 5 は、W e b 文書ファイル入力部 1 1 2 及び W e b 文書取得部 1 1 4 から取得した W e b 文書の文字コードを U T F - 8 に変換し統一する。

【 0 0 5 6 】

[1 - 3] W e b 文書の分割 :

W e b 文書の分割は、W e b 文書分割部 1 2 で行われる処理である。

【 0 0 5 7 】

最初に、W e b 文書分割部 1 2 の広告対象領域抽出部 1 2 1 は、インターネット広告等のコンテンツタグを含む領域がある場合、その領域を抽出する。ここで、インターネット広告とは、google (登録商標) や、Overture (登録商標) 等の広告会社が広告配信のための主要コンテンツ絞込みに用いるタグである。google (登録商標) の広告の場合、<!--google_ad_section_start--> から、<!--google_ad_section_end--> までがその領域に該当する。これらのタグは W e b 文書によって文字列が少々異なったり大文字で表記されるので、大文字と小文字を区別しない正規表現を用いたり、ワイルドカードの正規表現を用いる等により、多少の文字列表記の違いを吸収する処理を行う。以下、正規表現を用いる処理の説明の際には、多少の違いを吸収する処理を行っているものとする。

【 0 0 5 8 】

以降の処理は、インターネット広告の領域が存在する場合、上記で述べた領域を抽出する処理を行い、インターネット広告の領域がない場合は、最初に入力された W e b 文書に対して処理を行う。

【 0 0 5 9 】

10

20

30

40

50

次に、ノイズとなるタグや領域除去部 1 2 2 は、余計なタグや領域、特定の文字列を除去する処理を行う。除去されるタグや領域は、Web 文書の HTML を説明するコメントタグであったり、JavaScript であったり、form タグであったりする。除去するタグと領域を以下に記載する。

【 0 0 6 0 】

- ・ "<!-->" で始まり、"-->" で終わるコメントタグ；
- ・ "<script>" タグから、"</script>" タグで囲まれる領域；
- ・ "<style>" タグから、"</style>" タグで囲まれる領域；
- ・ "<select>" タグから、"</select>" タグで囲まれる領域；
- ・ "<noscript>" タグから、"</noscript >" タグで囲まれる領域；
- ・ "<form>" タグから、"</form>" タグで囲まれる領域；
- ・ 連続した空白文字列（単一の空白は除く）；
- ・ 連続したタブ文字列（単一のタブは除く）；

ノイズとなるタグや領域除去部 1 2 2 は、以上のタグ、領域、文字列を正規表現を用いて除去する。タグ内に alt 属性や class 属性が存在する場合も考えられるため、その場合はそれらを含めたタグを考慮した正規表現を用いて分割を行う（例：<style class="hoge">）。

【 0 0 6 1 】

次に、Web 文書の分割部 1 2 3 は、Web 文書の分割を行う。分割の規則は、以下のタグを用いて分割を行う。

【 0 0 6 2 】

- ・ <div>
- ・ </div>
- ・ <td>
- ・ </td>

タグ内に alt 属性や class 属性が存在する場合も考えられるため、その場合は、それらを含めたタグを考慮した正規表現を用いて分割を行う（例：<div class="hoge".）。

【 0 0 6 3 】

Web 文書分割部 1 2 で分割された Web 文書の一つ一つを『セグメント』と呼び、これらはメモリ（図示せず）に格納される。以降、特徴量抽出部 1 3 と主要コンテンツ判定部 1 4 の処理は、当該セグメント毎に行う。

【 0 0 6 4 】

[1 - 4] 特徴量の抽出：

ここでは、特徴量抽出部 1 3 が抽出する特徴量について述べる。特徴量抽出部 1 3 では、メモリ（図示せず）に格納されたセグメント毎に、以降で述べる特徴量を抽出し、Web 文書の主要コンテンツ部分の判定を行う。

【 0 0 6 5 】

[1 - 4 - 1] Web 文書で表示される文字列

当該特徴量は、Web 文書で表示される文字列特徴量抽出部 1 3 3 で抽出される特徴量である。ここで述べる文字列とは、HTML タグ等の Web ブラウザで表示されない文字列を含まないものとする。

【 0 0 6 6 】

Web 文書で表示される文字列として、文字列の量、句読点の数がある。

【 0 0 6 7 】

[1 - 4 - 1 - 1] 文字列の量：

一般的に Web 文書の主要コンテンツ部分は、主要コンテンツでない部分と比較すると多くの文字列が含まれている。また、全体的に文字列の少ない Web 文書においても同様のことがいえる。そのため Web 文書で表示される文字列特徴量抽出部 1 3 3 では、分割された Web 文書に含まれる文字列の数を特徴量とする。そして、特徴量の正規化部 1 3 4 では、文字列の量を正規化して特徴とする手法と、文字列の絶対値を用いて特徴量とす

10

20

30

40

50

る手法の2つを実行し、最終的な文字列の量の特徴量とする。

【0068】

・文字列の量の正規化を行い特徴量とする手法：

特徴量の正規化部134では、全てのセグメントにおいて最大の文字列の量をもつセグメントの特徴量を1とする正規化を行う。例えば、全てのセグメントにおいて最大の文字列の量が200で、あるセグメント内の文字列の量が100だった場合には、そのセグメントの文字列の量の特徴量は0.5となる。このような正規化を行うことで、全体的に文字列の少ないWeb文書においても主要コンテンツの抽出が可能になる。

【0069】

・文字列の量の絶対値を用いて特徴量とする手法：

上記で述べた正規化を行い特徴量とする手法は、全体的に文字列の少ないWeb文書において有効であったが、正規化を行うことで、全体的に文字列の量が多く、主要コンテンツ部分のセグメント間の文字列の量の差が大きい場合に不都合が生じる。例えば、全てのセグメントにおいて最大の文字列の量が1000で、あるセグメント内の文字列の量が100だった場合、そのセグメントの文字列の量の特徴量は0.1になる。文字列の量としては多いはずだが、正規化を行うことで、このような弊害が生じる。そこで、文字列の絶対値を用いて特徴量とする手法を行う必要がある。具体的には、特徴量の正規化部134では、ある特定の値を超えた場合にその文字列の特徴量を1とする手法を用いる。例えば、あるセグメント内の文字列の量が100の場合、文字列の量が5以上の場合の特徴量が1となり、文字列の量が10以上の場合の特徴量が1となり、...、文字列の量が105以上の特徴量は0となり、...、文字列の量が200以上の特徴量は0となるように特徴量を作成する。このように、ある特定の文字列量を超えた場合に特徴量を1とする手法を用いることで、特徴量の最大値は1のままで文字列量の絶対値を特徴量とすることができる。また、例における文字列の量の絶対値の特徴量の間隔は5としたが、場合において適切な間隔を用いるのが好ましい。8, 16, 32, 64といった2の乗数を用いて特徴量の間隔とする手法も考えられる。文字列の量がx以上の...の最大のxも同様に、場合において適切な値に変更する。主要コンテンツ判定における計算量を減らしたい場合にはxの値を小さくすればよい。

【0070】

また、セグメント内に文字列が全くない場合も考えられる。その場合は、以降で説明する特徴量を抽出するまでもなく、主要コンテンツ判定部14内のテキスト判定部142で主要コンテンツでないと判断できる。そのため、実際の実行時には以下の特徴量抽出は行わず、該当セグメントを非主要コンテンツとして判別する。この処理は、特徴量の学習を行う際には行わない。

【0071】

[1-4-1-2] 句読点の数：

Web広告等のノイズとなりやすいセグメントは、文字列の量が多いが句読点の数が少ない傾向にある。そのため、句読点の数を特徴とする。具体的には、セグメント内の文字列に含まれる、『、』、『, 』、『。』、『. 』、『! 』、『. 』、『? 』、『... 』の数を特徴量としてカウントする。この特徴量も、文字列の量で述べた正規化による特徴量と、絶対値による特徴量の二通りを算出する。算出方法においては、[1-4-1-1]の文字列の量で述べた手法と同じものを用いる。

【0072】

[1-4-2] タグ情報

ここでは、タグ情報特徴量抽出部132で扱うHTMLタグ等のタグ情報に関する特徴量について述べる。タグ情報には、

- ・テキスト系のHTMLタグの数；
- ・テキスト系のHTMLタグの連続出現数；
- ・リンクリストタグの数、Web文書で表示される文字列を含まない文字列の量；

がある。

10

20

30

40

50

【 0 0 7 3 】

[1 - 4 - 2 - 1] テキスト系の HTML タグの数 :

あるセグメント内において、Web 文書で表示される文字列が多い場合、テキストに関する HTML タグが多く含まれる。また、ブログ等の CGM (Consumer Generated Media) においては、Web 文書で表示される文字列は少ないが、ユーザが改行タグを多用する事例が多くみられる。そこで、タグ情報特徴量抽出部 132 は、テキストに関する HTML タグの数を特徴量として用いる。この特徴量も、Web 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。例えば、全てのセグメントにおいて最大の HTML タグの量が 10 で、あるセグメント内の HTML タグの量が 5 だった場合には、そのセグメントの HTML タグの量の特徴量は 0 . 5 となる。そして、実際に使用するテキスト系の HTML タグは、以下のタグを対象とする。

10

【 0 0 7 4 】

- ・ <p>
- ・ </p>
- ・

- ・ </br>
- ・
- ・

タグ内に size 属性や class 属性が存在する場合も考えられるため、その場合は、それらを含めたタグを考慮した正規表現を用いてカウントを行う (例 :)。

20

【 0 0 7 5 】

[1 - 4 - 2 - 2] テキスト系の HTML タグの連続出現数 :

Web 文書で表示される文字列が集中して記述してあるセグメントは、テキスト系のタグが多く存在すると同時に、テキスト系の HTML タグが連続して出現する。ここでいう、連続して出現するというのは、他のアンカーリンク等の HTML タグが間に出てこないということである。そこで、[1 - 4 - 2 - 1] で述べたテキスト系の HTML タグの連続出現数を特徴量とする。この特徴量も、Web 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。例えば、全てのセグメントにおいて最大の HTML タグの連続量が 10 で、あるセグメント内の HTML タグの連続量が 5 だった場合には、そのセグメントの文字列の量の特徴量は 0 . 5 となる。

30

【 0 0 7 6 】

[1 - 4 - 2 - 3] リンクリストタグの数 :

あるセグメント内においてリンクリストタグが多い場合、そのセグメントにはナビゲーションリンク等の多くのアンカーリンクが存在し、そのセグメントは主要コンテンツとならない可能性が高い。そこで、リンクリストタグの数を特徴量とする。この特徴量も、Web 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。例えば、全てのセグメントにおいて最大のリンクリストタグの量が 10 で、あるセグメント内のリンクリストタグの量が 5 だった場合には、そのセグメントのリンクリストタグの量の特徴量は、0 . 5 となる。そして、具体的に使用するリンクリストタグは、以下のタグを対象とする。

40

【 0 0 7 7 】

- ・
- ・
- ・ <dl>
- ・ <dd>
- ・

タグ内に alt 属性や class 属性が存在する場合も考えられるため、その場合は、それらを含めたタグを考慮した正規表現を用いてカウントを行う (例 : <font class="hoge")。

50

【 0 0 7 8 】

[1 - 4 - 2 - 4] W e b 文書で表示される文字列を含まない文字列 (H T M L タグを含む) の量 :

あるセグメント内において、W e b で表示されない文字列 (H T M L タグを含む) が多い場合、そのセグメントは広告等の主要コンテンツでない可能性が高い。そこで、W e b 文書で表示される文字列以外の文字列 (H T M L タグを含む) 量を特徴量とする。この特徴量も、W e b 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。例えば、全てのセグメントにおいて最大の W e b で表示されない文字列の量が 1 0 0 で、あるセグメント内の W e b で表示されない文字列の量が 5 0 だった場合には、そのセグメントの W e b で表示されない文字列の量の

10

特徴量は 0 . 5 となる。

【 0 0 7 9 】

[1 - 4 - 3] アンカーリンク情報 :

以下では、特徴量抽出部 1 3 のアンカーリンク情報特徴量抽出部 1 3 1 で扱うアンカーリンクに関する特徴量の抽出方法について述べる。アンカーリンク情報には、

- ・アンカーリンクの数 ;
- ・各アンカーリンクの文字列の平均量 ;
- ・すべてのアンカーリンク文字列の合計値 ;
- ・最大文字列のアンカーリンク U R L の量 ;
- ・広告に関するアンカーリンクを含むか ;

20

がある。

【 0 0 8 0 】

[1 - 4 - 3 - 1] アンカーリンクの数 :

あるセグメントにおいて、アンカーリンクの数が多数含まれているセグメントは主要コンテンツでない可能性が高い。そこで、アンカーリンクの数を特徴量として用いる。具体的には、` ... ` タグで表されるアンカーリンクの数を特徴量とする。この特徴量も、特徴量の正規化部 1 3 4 において、W e b 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。例えば、全てのセグメントにおいて最大のアンカーリンクの数が 1 0 で、あるセグメント内のアンカーリンクの数が 5 だった場合には、そのセグメントのアンカーリンクの数の特徴

30

量は 0 . 5 となる。アンカーリンクタグには、class 属性や alt 属性が含まれる場合もあるので、アンカーリンクタグの数は正規表現を用いてカウントする。

【 0 0 8 1 】

[1 - 4 - 3 - 2] 各アンカーリンクの文字列の平均量 :

各アンカーリンクの文字列が平均して多い場合、そのセグメントは、関連記事等のナビゲーションリンクである可能性が高い。また、アンカーリンクの文字列が平均して少ない場合、主要コンテンツ内に含まれるキーワード検索リンクである可能性が高い。そこで、セグメントに含まれるアンカーリンクの文字列の平均量を特徴量として用いる。アンカーリンクの文字列とは、` ` の 部分に該当する。この特徴量も特徴量の正規化部 1 3 4 において W e b 文書で表示される文字列で用いた特徴量の

40

正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。

【 0 0 8 2 】

[1 - 4 - 3 - 3] すべてのアンカーリンク文字列の合計値 :

セグメント内に含まれるアンカーリンクの文字列の合計量が多い場合、そのセグメントはナビゲーションリンクである可能性が高い。そこで、セグメント内に含まれるアンカーリンクの文字列の合計量を特徴量として用いる。アンカーリンクの文字列とは、` ` の 部分に該当する。この特徴量も特徴量の正規化部 1 3 4 で W e b 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ (図示せず) に格納する。

【 0 0 8 3 】

50

[1 - 4 - 3 - 4] 最大文字列のアンカーリンクURLの量 :

アンカーリンク先のURL文字列が非常に長い場合、そのセグメントは広告である可能性が高い。そこで、セグメント内で最大長のアンカーリンク先のURLの文字列を特徴量として用いる。ここで述べるアンカーリンク先のURL文字列とは、` ... ` の 部分に該当する。この特徴量も特徴量の正規化部 1 3 4 で Web 文書で表示される文字列で用いた特徴量の正規化を同様に言い、最終的な特徴量とし、メモリ (図示せず) に格納する。

【 0 0 8 4 】

[1 - 4 - 3 - 5] 広告に関するアンカーリンクを含むか :

広告に関するURLを含むアンカーリンクは特徴的な文字列を含む可能性が高い。例えば、「adclick」、「adnet」、「banner」等がそれにあたる。そこで、アンカーリンク情報特徴量抽出部 1 3 1 では、このような広告となりやすい文字列を含んだURLを含むアンカーリンクが存在する場合、特徴量を 1 とし、存在しない場合を 0 とする特徴量を抽出し、メモリ (図示せず) に格納する。広告になりやすい文字列は、広告除去用のFirefox用アドインであるadblock plugin等のサイトに記載されているため、それを用いる。

【 0 0 8 5 】

[1 - 4 - 4] 特徴量の比率 :

ここでは、特徴量の比率特徴量抽出部 1 3 5 で扱う、[1 - 4 - 1] ~ [1 - 4 - 3] で述べた特徴量の比率を用いた特徴量について述べる。

【 0 0 8 6 】

[1 - 4 - 4 - 1] テキスト系のタグ数とテキスト系のタグの連続出現数の比率 :

テキスト系のタグが多数あり、またテキスト系のタグの連続出現数が多いセグメントは、テキストが密に書かれているため、主要コンテンツである可能性が高い。しかしながら、テキスト系のタグは多数あるが、テキスト系のタグの連続出現数が少ないセグメントは、セグメントのサイズが大きいただけで、主要コンテンツでない可能性が高いといえる。そこで、特徴量の比率特徴量抽出部 1 3 5 では、テキスト系のタグ数とテキスト系のタグの連続出現数の比率を特徴量として用いる。具体的には、テキスト系のタグ数を分母とし、テキスト系のタグの連続出現数を分子とした値を特徴量として用いる。ここで、テキスト系のタグ数が 0 の場合は、分母が 0 となってしまうため、この場合のテキスト系のタグ数とテキスト系のタグの連続出現数の比率の特徴量は 0 とする。本特徴量も特徴量の正規化部 1 3 4 において Web 文書で表示される文字列で用いた特徴量の正規化を同様に言い、最終的な特徴量とし、メモリ (図示せず) に格納する。この特徴量が大きければ大きいほど主要コンテンツである可能性が高い。

【 0 0 8 7 】

[1 - 4 - 4 - 2] Web で表示される文字列とタグの比率 :

あるセグメント内において、Web で表示される文字列が多い場合は、主要コンテンツとなる可能性が高いが、同じセグメント内において、HTML タグ等のタグが多い場合もある。この場合、前述の [1 - 4 - 4 - 1] で述べたようにセグメントのサイズが大きいただけで主要コンテンツでない可能性がある。そこで、特徴量の比率特徴量抽出部 1 3 5 では、Web で表示される文字列とタグの比率を特徴量として用いることでこの場合に対処する。具体的には、Web で表示される文字列の数を分子とし、タグの数を分母とした値を特徴量とする。この特徴量が大きければ大きいほど、主要コンテンツである可能性が高い。本特徴量も、特徴量の正規化部 1 3 4 において、Web 文書で表示される文字列で用いた特徴量の正規化を同様に言い、最終的な特徴量とし、メモリ (図示せず) に格納する。タグの数が 0 の場合は分母が 0 となってしまうため、特徴量は 1 とする。

【 0 0 8 8 】

[1 - 4 - 4 - 3] アンカーリンクの数とリンクリストタグの数の比率 :

あるセグメント内において、アンカーリンクの数や、リンクリストタグの数が多ければ多いほど、そのセグメントは主要コンテンツでない可能性が高いが、セグメントが広いためにこれらの特徴量が偶然高くなってしまう場合も考えられる。そこで、特徴量の比率特

10

20

30

40

50

微量抽出部 135 では、アンカーリンクの数とリンクリストタグの数の比率を特徴量として用いる。具体的にはアンカーリンクの数を分母とし、リンクリストタグの数を分子とし、特徴量とする。この特徴量が大きければ大きいほど、セグメントの面積に対し密度の高いリンク数が存在することになり、主要コンテンツでない可能性が高い。本特徴量も、特徴量の正規化部 134 で、Web 文書で表示される文字列で用いた特徴量の正規化を同様に行い、最終的な特徴量とし、メモリ（図示せず）に格納する。アンカーリンクの数が 0 の場合は分母が 0 となってしまうため、特徴量は 0 とする。

【0089】

[1-5] 主要コンテンツの判定：

ここで、主要コンテンツの判定部 14 が、特徴量抽出部 13 で求められ、メモリ（図示せず）に格納されている [1-4] で述べた特徴量を用いて、主要コンテンツか否かの判定を用いて行う手法について述べる。判定には、Support Vector Machine(SVM)や最大エントロピー法、ナイーブベイズ法等の機械学習アルゴリズムを用いて判定を行う。

10

【0090】

図 11 は、本発明の一実施の形態における特徴量のパラメータ推定方法のフローチャートである。最初に、人手で主要コンテンツか否かを特徴量を抽出したセグメント毎に判定し、訓練データを作成する（ステップ 301～303）。ここで、[1-4-1-1] で、Web で表示される文字列が存在しない場合は主要コンテンツと見做さないと記述したが、機械学習を用いた手法において負例として学習に有効であるため、訓練データにはそのようなデータも採用する。そして、そのセグメントの特徴量を用いて学習を行い（ステップ 304）、特徴量毎の重みを算出する（ステップ 305）。速度を重視する場合は、最大エントロピー法で学習し、精度を重視する場合には、二次の多項式カーネルを用いた Support Vector Machine を用いて学習を行う。そしてこれらの学習したパラメータを用いて、特徴量入力部 141 がセグメントの特徴量を主要コンテンツ判定部 143 に入力し、主要コンテンツ判定部 143 がセグメント毎に主要コンテンツか否かを判定する。

20

【0091】

[1-6] 主要コンテンツ部分の出力：

上記の [1-5] で説明した主要コンテンツ判定部 14 が主要コンテンツか否かの判定を行った後、主要コンテンツ出力部 15 のデータ出力部 153 は、学習器によって主要コンテンツと判断されたセグメントのみを、結合して最終出力とする。出力の例を図 12、図 13 に示す。

30

【0092】

ここで、情報検索の事前処理として本装置を用いたい場合には、タグ付きテキスト出力部 151 を用いて HTML タグ等のタグを残して出力する。また、情報推薦等で Web 文書の内容を解析したい場合には、タグなしテキスト出力部 152 を用いて HTML タグ等のタグを削除して出力する。

【0093】

[1-7] 精度向上のための処理：

[1-4] で述べた特徴量を抽出する事前処理として、不要文字列等を除去する手法が有効である。以下に記述する不要文字列を事前に除去しておくことで、主要コンテンツの判定精度を高める。

40

【0094】

- ・
- ・<
- ・>
- ・&
- ・«
- ・»

これらの文字列は、HTML タグ等で用いる記号を Web ブラウザ上で表示する際に用いる特殊文字である。また、上記で挙げた特殊文字以外の HTML 特殊文字も削除の対象

50

とする。特殊文字は実際表示される文字列に対して、文字列の量が少ないため、学習の際のノイズとなりやすい。

【 0 0 9 5 】

[1 - 8] 実装理由による特徴量の選択：

主要コンテンツ抽出手法をユーザPC等に組み込む場合、[1 - 4]で述べた全ての特徴量を用いて処理することは処理量的に難しい。そのため、抽出する特徴量を絞り込むことで処理量を削減する。ここで、機械学習による学習モデルは、絞り込んだ特徴量モデル毎に学習モデルを作成する。

【 0 0 9 6 】

[2] キーワード候補の抽出処理 (ステップ 3 0 0) :

以下では、重要キーワード候補抽出部 2 について説明する。

【 0 0 9 7 】

従来の一般的なキーワード抽出手法は形態素解析や固有表現抽出を用いた手法であったが、形態素解析手法では、複合名詞の抽出に関する問題や重要キーワードの絞込みに関する問題があり、また、固有表現抽出においては、人名、組織名、地名といった狭い範囲のキーワード抽出しかできないため、ユーザの興味を網羅するキーワード候補の抽出ができなかった。そこで、本発明では、ユーザ参加型オンライン百科事典であるWikipedia (登録商標) の見出し語をキーワード候補として用いることで、これらの解決を試みた。Wikipedia (登録商標) の見出し語のため、体系的にまとめられており、かつ実世界の事象を一意に表す特徴を持つため、ユーザの興味対象を幅広く網羅した言語資源であるといえる。また、Wikipedia (登録商標) に記載されているキーワードは、重要で、キーワードについて語る際に様々な話題が挙がる情報量が豊富なキーワードのみが登録されているため、重要なキーワードの絞込みが可能となる。具体的なキーワード抽出手法は、Wikipedia (登録商標) の見出し語を形態素解析用の辞書としてキーワード辞書 1 0 1 に登録し、そのキーワード辞書 1 0 1 を用いた形態素解析結果からキーワードを抽出する。なお、キーワード辞書 1 0 1 に登録する際、最も長いキーワードを抽出するよう辞書の重み付けを行う。

【 0 0 9 8 】

[3] キーワードの重要度算出処理 (ステップ 4 0 0) :

前述の [2] では、Wikipedia (登録商標) の見出し語を用いて重要だと思われるキーワード候補の絞込みを実現したが、[3] では、出現頻度算出部 3、キーワード重要度算出部 4 により、その中からさらに重要なキーワードを上位に位置づけるために、キーワードの重要度を算出する手法について述べる。

【 0 0 9 9 】

[3 - 1] WebIDF :

検索エンジンが収集した大規模なWeb文書集合における出現確率を考慮した重みがWebIDFであり、WebIDF記憶部 1 0 2 に格納されているものとする。具体的には、出現頻度算出部 3 は、検索エンジンにキーワードを入力し、その結果得られたWeb文書のヒット数からIDF値を算出し、これをWebIDFとする。Web文書集合中で多数出現するキーワードは、一般的なキーワードで特徴的なキーワードではないと判断する重みである。以下に、キーワード k のWebIDF算出式を示す。

【 0 1 0 0 】

【数 1】

$$WebIDF(k) = \log_2 \left(\frac{N}{n_k + 1} + 1 \right) \quad (1)$$

10

20

30

40

ここで、 N はキーワードの閾値で、 n_k はキーワード k のヒット数である。ここでキーワード k の検索エンジンのヒット数の閾値 N は、 n_k を降順にソートし、不必要なキーワードが少なくなってきた辺りの n_k の値を用いる。その理由として、 n_k が高ければ高いほど一般的なキーワードである可能性が高く、また、閾値 N を設定すると閾値 N 以上の n_k を持つキーワード k のWebIDF値は常に負となり、常にスコアとして低い値となるためである。なお、 N の推奨値は20000000である。

【0101】

[3-2] BM25 :

ここでは、キーワード重要度算出部4による文書内の出現頻度と文書集合内での出現分布を用いたキーワード重み付けの手法について述べる。これらの重み付けには既存の技術であるBM25を用いる。

10

【0102】

【数2】

$$BM25(d, k) = \frac{(k_1 + 1) \cdot tf(d, k)}{k_1 \left((1 - b) + b \cdot \frac{dl}{avdl} \right) + tf(d, k)} \cdot \log \left(\frac{N - df(k) + 0.5}{df(k) + 0.5} \right) \quad (2)$$

20

ここで、 d は文書であり、 k は文書に含まれるキーワードである。そして、 $tf(d, k)$ は d 内のキーワード k の出現頻度、 dl は文書長、 $avdl$ は文書集合内における平均文書長、 N は総文書数である。 k_1 と b は自由パラメータで、 k_1 は $tf(d, k)$ に関するパラメータ、 b は dl に関するパラメータである。

【0103】

また、文書集合を用意できなく、文書1枚からキーワードを抽出したい場合には、以下の算出式を用いる。

【0104】

【数3】

$$BM25(d, k) = \frac{(k_1 + 1) \cdot tf(d, k)}{k_1 + tf(d, k)} \quad (3)$$

30

なお、式に関するノテーションは(2)式と同様である。

【0105】

[3-3] タイトルと本文のキーワード重要度の算出 :

40

上記の[3-2]では、出現頻度算出部3によるBM25を用いた出現頻度による重み付け手法を述べた。本節では、キーワード重要度算出部4での文書のタイトル中での出現頻度と、本文中での出現頻度を用いたキーワード重要度の算出手法について述べる。

【0106】

具体的には、キーワード重要度算出部4は、文書のタイトル(もしくはファイル名)から算出したBM25の値と、本文から算出したBM25の値の線形和を用いてスコアを算出する。以下に算出式を示す。

【0107】

【数4】

$$BM25_{all}(d, k) = (1 - \alpha) \cdot BM25_{title}(d, k) + \alpha \cdot BM25_{body}(d, k) \quad (4)$$

ここで、 $BM25_{all}$ はタイトルと本文の出現頻度を用いたBM25スコアで、 $BM25_{title}$ はタイトルでの出現頻度を用いたBM25スコア、 $BM25_{body}$ は本文でのBM25スコアとなる。 α はパラメータである。 α の推奨値は0.7とする。

10

【0108】

[3-4] 形態素の固有名詞を用いたキーワード重要度算出：

一般的なニュース記事内において、人名や組織名等の固有名詞は最初に正式名称で記述された後に、省略形で記述される場合が多い。例えば、「麻生次郎」と最初に書かれた後に「麻生」と書かれる場合がそれに該当する。省略形で記述された場合、記事の主題に関するキーワードでも他の重要でないキーワードより出現回数が少ない場合が生じる問題があった。そこで、形態素解析結果を用いて簡易的に省略形に対応する方法を提案する。一般的な形態素解析器は、品詞として、人名(姓、名)、地名、組織名の出力が可能である。そのため、形態素解析結果から得られた形態素の固有名詞を含むキーワード候補の出現頻度を、その形態素の出現頻度と置き換えることで、この問題に対処する。例えば、「麻生次郎」というキーワードが1回出現し、「麻生」という形態素の固有名詞が3回出てきた場合、「麻生次郎」の出現頻度を3回とする。

20

【0109】

逆に、地名等は省略されずに記述されるため、キーワードの出現回数が大きくなりすぎてしまう問題がある。そのため、形態素解析結果から得られた地名の名詞の出現回数を提案する。以下に、キーワード重要度算出部4で用いる、上記で述べた手法の算出式を示す。

【0110】

$$tf(d, k) = \text{person} \cdot \text{match}(d, k, m_k)$$

$$tf(d, k) = \text{location} \cdot \text{match}(d, k, m_k)$$

$$tf(d, k) = \text{organization} \cdot \text{match}(d, k, m_k)$$

(5)

30

ここで、 $tf(d, k)$ は文書 d における出現頻度によるキーワード k のスコア、 person 、 location 、 organization はそれぞれ、人名、地名、組織名に関する係数である。そして、 $\text{match}(d, k, m_k)$ は文書 d におけるキーワード k に含まれる形態素 m_k の出現頻度である。上記の式(5)で算出した $tf(d, k)$ は、[3-2]、[3-3]で述べた出現頻度を用いたキーワード重要度の算出式で用いる。

【0111】

[3-5] キーワードの重要性の算出(キーワード固有重要度)：

本節では、キーワード重要度算出部4でのキーワードの重要性の算出について述べる。キーワードの重要性とは、知名度や話題性の高い、キーワードについて語る際に様々な話題が挙がる情報量が豊富な(内容の深い)キーワード程重要であると定義し、このキーワードの重要度をキーワード固有重要度と呼ぶ。当該キーワード重要度算出部4で算出されたキーワード固有重要度は固有重要度記憶部105に格納する。

40

【0112】

[3-5-1] Wikipedia(登録商標)内のリンク構造を用いた手法：

以下では、キーワード重要度算出部4において、Wikipedia(登録商標)内のリンク構造を用いてキーワード固有重要度を算出する手法について説明する。

【0113】

HITSやPageRankといった一般的なWeb文書のランキング手法は、Webページのリンク構造を用いて、Web文書のランキングを行っている。しかしながら、Wikipedia(登

50

録商標)の文書には、1つの文書に付き1つの見出し語(キーワード)がついているため、Wikipedia(登録商標)のリンク構造から得られた文書のランキングをキーワードのランキングと見做すことができる。そのため、本装置では、キーワード重要算出部4において、Wikipedia(登録商標)に特化したランキング手法を適用することで、キーワード重要度1記憶部103に格納されているキーワード重要度を用いてキーワード固有重要度を算出し、固有重要度記憶部105に格納する。ベースとなるアルゴリズムには、HITSアルゴリズムを用いる。

【0114】

HITSアルゴリズムには、全てのWeb文書をauthority(コンテンツ)とhub(リンク集)の2つから構成されると定義する。そして、良いhubから多数リンクされるauthority程良いauthorityであるという仮説と、良いauthorityに多数リンクしているhub程良いhubであるという二つの仮説を繰り返し実行することでWeb文書のランキングを行う。しかしながら、HITSアルゴリズムはWeb世界におけるWeb文書のリンク構造をモデルにしたアルゴリズムのため、リンク構造が非常に密なWikipedia(登録商標)にそのまま適用した場合、やや難がある。そこで、本装置では、Wikipedia(登録商標)の特徴的な構造と密なリンク構造に対応させた手法を提案する。そして、本アルゴリズムから算出したauthorityの値による順位を、本手法が提案する減衰関数に近似させ、最終的なキーワード固有重要度とする。

10

【0115】

テキスト量の考慮：

Wikipedia(登録商標)の見出し語は、知名度が高く話題性の高い見出し語ほど、テキストの記述量が多い傾向がある。そこで、authority値の算出の際に、自文書のテキスト量が多ければ多いほどその文書は重要であるといった重み、 $text(k)$ を考慮する。

20

【0116】

自リンクと被リンクの比率：

一般的にWikipedia(登録商標)の見出し語は、有名なキーワードほど、自リンクと被リンクの数が多くなっている。しかしながら、地名やジャンル名のような広い概念を持つキーワードは、引用しやすいキーワードのため、自リンク数に比べて圧倒的に被リンクの数が多い傾向がある。通常のHITSアルゴリズムは良いhubから多数リンクされているauthorityは良いauthorityであるといった仮説を用いるが、圧倒的に被リンクが多い場合においては、これらの仮説は成り立たないと予想される。また、その一方で、最近知名度が高くなってきている新人俳優や話題語等の見出し語は、誕生してから日が浅いため引用数は少ないが自リンクは多い傾向にある。そのため少ない被リンク数においても、authorityを高める必要がある。これらの問題を解決するために、authority値の算出の際に、 $flink(k)/blink(k)$ を考慮する。ここで、 $flink(k)$ はキーワードkの文書内に含まれる自リンクの数を表す、 $blink(k)$ はキーワードkの文書にリンクしている被リンクの数を表す。

30

【0117】

明らかにauthority算出とならない見出し語の扱い：

Wikipedia(登録商標)の見出し語には「～年」や「～一覧」といった明らかにauthorityとならない見出し語が存在する。これらの見出し語は自リンクが非常に多く、被リンクも非常に多い場合があるためノイズとなりやすい。そこで、明らかにauthorityとならない見出し語のauthority値は常に変更しないことで、この問題に対処する。

40

【0118】

hubの平均的なリンクの質：

Wikipedia(登録商標)の文書には、自リンクが多数あるが、hubとして質の悪い文書がある。そこで、リンク先キーワードのauthorityが平均的に高いhubは重要であるといった仮説に変更することで、自リンクは多いがhubとして質の低い文書のhub値を下げる重み

【0119】

【数5】

$$\sum_{k'} \log(a(k')+1) / K'$$

を考慮する。

【0120】

10

リダイレクトの扱い：

Wikipedia（登録商標）の文書には、見出し語の異表記を解消するために、redirect（リダイレクト）が存在する。例えば、「イチロー」には「鈴木イチロー」、「ICHIRO」のredirectがある。Redirectは異表記のキーワードを一意に纏める効果だけでなく、キーワードの被リンク構造に大きな影響を持つため、redirectキーワードを親ノードに纏めることで、異表記のキーワード固有重要度を算出し、被リンクの問題も解決する。

【0121】

そして、最終的なWikipedia（登録商標）ランキングアルゴリズムは以下の式で定義される。

【0122】

20

【数6】

$$a(k) = \frac{\log(\text{flink}(k)+1)}{\log(\text{blink}(k)+1)} \cdot \text{text}(k) \cdot \sum_{k'} h(k')$$

$$h(k) = \frac{\sum_{k'} \log(a(k')+1)}{K'} \cdot \sum_{k'} a(k') \quad (6)$$

30

ここで、 $a(k)$ は、キーワード k のauthority値で、 $h(k)$ はhub値である。そして、 $\text{flink}(k)$ はキーワード k から自リンク数、 $\text{blink}(k)$ はキーワード k からの被リンク数である。 $\text{text}(k)$ は、キーワード k が見出し語になっているWikipedia（登録商標）文書の文字数（アンカーリンク対象の文字列は除く）であり、 K' はキーワード k が見出し語になっているWikipedia（登録商標）文書内に含まれるリンク数の総数となる。なお、式（6）は交互に繰り返し計算を行うものとし、10回ほど繰り返し行うとよい。

【0123】

そして、上記の式（6）で算出したauthority値を用いて、降順にキーワードを順位付けする。そして、以下のキーワード固有重要度算出式を用いてスコアの近似を行い、Wikipedia（登録商標）内のリンク構造を用いたキーワード固有重要度WKIS(k)とし、固有重要度記憶部105に格納する。

40

【0124】

【数7】

$$WKIS(k) = \exp\left(\frac{\log(y_1 - y_0 + 1)(K - k_r + 1)^a}{K^a}\right) \div y_0 - 1 \quad (7)$$

ここで、 y_1 はキーワード固有重要度の上界であり、 y_0 はキーワード固有重要度の下界である。そして、 k_r はキーワード k のauthorityの値による順位、 K はキーワードの総数、 a はスコアの勾配係数で、 a の値が大きくなればなるほどスコアの勾配が急になる。本関数の特徴は、キーワードの候補数（ x 軸の要素数）に左右されることなく、また、上界と下界を設定でき、またキーワードの候補数の上位10%以内において、最大値と最小値の差が30%～70%以内に収まる減衰曲線を描き、そしてパラメータ係数にて減衰度合いを整数値を用いることで容易に設定できる特徴をもつスコア関数である（勿論、減衰度合いの設定は実数値でも設定可能である）。 y_1 、 y_0 の推奨値は、それぞれ1, 0.1で、勾配係数の推奨値は3から7である。

【0125】

Wikipedia（登録商標）内のリンク構造から算出したキーワード固有重要度は、ユーザ参加型のオンライン百科事典のWikipedia（登録商標）内においての重要なキーワードが上位に位置付けられるため、一般に知名度は低い但实际上には内容が深く重要なキーワードが上位に位置づけられる。

【0126】

[3-5-2] 検索エンジンの検索クエリ投入回数を用いた手法：

この手法は、実際の検索エンジンに投入された回数が多い検索クエリほど、重要なキーワードであるとみなす手法である。検索クエリの投入回数は、ポータルサイト上で投入された検索クエリのような大規模データであることが好ましい。ここで、対象となる検索クエリは、Wikipedia（登録商標）に存在するキーワードであるとする。この手法もWikipedia（登録商標）内のリンク構造を用いた手法と同じく式（7）のキーワード固有重要度算出式を用いてスコアの近似を行い、検索クエリを投入回数を用いたキーワード固有重要度 $QKIS(k)$ とし、固有重要度記憶部105に格納する。

【0127】

【数8】

$$QKIS(k) = \exp\left(\frac{\log(y_1 - y_0 + 1)(K - k_r + 1)^a}{K^a}\right) + y_0 - 1 \quad (8)$$

ここで、 y_1 はキーワード固有重要度の上界であり、 y_0 はキーワード固有重要度の下界である。そして、 k_r はキーワード k の検索回数による順位、 K はキーワードの総数、 a はスコアの勾配係数で、 a の値が大きくなればなるほどスコアの勾配が急になる。本関数の特徴は、キーワードの候補数（ x 軸の要素数）に左右されることなく、また、上界と下界を設定でき、またキーワードの候補数の上位10%以内において、最大値と最小値の差が30%～70%以内に収まる減衰曲線を描き、そしてパラメータ係数にて減衰度合いを整数値を用いることで容易に設定できる特徴をもつスコア関数である（勿論、減衰度合いの設定は実数値でも設定可能である）。

【 0 1 2 8 】

また、Wikipedia（登録商標）には見出し語として存在するが、検索エンジンの検索クエリには存在しない場合がある。その場合は y_0 を該当キーワードのスコアとする。 y_1, y_0 の推奨値は、それぞれ、1, 0.1で、勾配係数の推奨値は3から7である。

【 0 1 2 9 】

キーワード重要度2記憶部104に格納されている検索エンジンの検索クエリ投入回数から算出したキーワード固有重要度は、実際の検索クエリ投入回数によって重要度が決まるため、検索クエリとして投入されやすく、重要なキーワードが上位に来る傾向にある。

【 0 1 3 0 】

[3 - 6] 最終的なキーワード固有重要度算出：

前述の[3 - 5]では、Wikipedia（登録商標）内のリンク構造から算出したキーワード固有重要度WKIS(k)と、検索エンジンの検索クエリ投入回数から算出したキーワード固有重要度QKIS(k)について述べた。しかしながら、QKIS(k)は検索クエリデータの収集期間が短い場合、その間にインターネット上で起きた話題に強く影響されてしまう問題があり、さらにインターネットサイト名等の生活的クエリが多く含まれる傾向にある。そのため、キーワード重要度算出部4では、固有重要度記憶部105に格納されているWKIS(k)とQKIS(k)の線形和を、最終的なキーワード固有重要度とすることで、話題性が高く内容が深く、かつ、検索クエリとして投入されやすいキーワードを上位に位置づける重要度を算出する。

【 0 1 3 1 】

$$\text{Keyword_score}(k) = \alpha \cdot \text{WKIS}(k) + (1 - \alpha) \cdot \text{QKIS}(k) \quad (9)$$

QKIS(k)よりも、WKIS(k)のキーワード重要の方が一般的に精度が良いことが実験により確認できたため、 α の値は0.5~0.8辺りの範囲で調整する。

【 0 1 3 2 】

[3 - 7] 文書中での出現位置におけるキーワード重要度算出：

以下では、位置情報算出部5における処理を説明する。

【 0 1 3 3 】

ニュース記事やコラム等の一般的な記述がなされる記事においては、文の先頭に来れば来るほど、重要なキーワードが含まれている可能性が高い。そのため、文の先頭であればあるほど、そのキーワードは重要であるという重要度を算出する。この手法も、Wikipedia（登録商標）内のリンク構造を用いた手法と同じく式(7)のキーワード固有重要度算出式を用いて、出現位置を用いたキーワード重要度Pos(k)の算出を行い、位置重要度記憶部106に格納する。

【 0 1 3 4 】

【 数 9 】

$$\text{Pos}(k) = \exp\left(\frac{\log(y_1 - y_0 + 1)(P - k_p + 1)^a}{P^a}\right) + y_0 - 1 \quad (10)$$

ここで、 y_1 は出現位置を用いたキーワード重要度の上界であり、 y_0 は出現位置を用いたキーワード重要度の下界である。そして k_p はキーワードkの文位置（文の位置は、文の先頭から1文ずつ（句点までを1文とする）数え上げたものを用い、さらにそのキーワードにおいて、最も先頭の位置を用いる）、Pは最後尾の文の位置、aはスコアの勾配係数で、aの値が大きくなればなるほどスコアの勾配が急になる。本関数の特徴は、キーワードの候補数（x軸の要素数）に左右されることなく、また、上界と下界を設定でき、またキーワードの候補数の上位10%以内において、最大値と最小値の差が30%~70%以内に

10

20

30

40

50

収まる減衰曲線を描き、そしてパラメータ係数にて減衰度合いを整数値を用いることで容易に設定できる特徴をもつスコア関数である（勿論、減衰度合いの設定は実数値でも設定可能である）。 y_1 、 y_0 の推奨値は、それぞれ1、0.5で勾配係数の推奨値は1から5である。ブログ記事のような一般的な記述がなされていない文書を多く処理する場合には、 y_0 の値を大きくし、勾配係数の値も小さくするのがよい。

【0135】

[3-8] 文書における最終的なキーワード重要度：

以下に、重要キーワード出力部6の動作を説明する。

【0136】

文書d内における最終的なキーワード重要度FS(d,k)は、[3-7]節までに述べた手法により求められ、固有重要度記憶部105と位置重要度記憶部106に格納されている各値を乗算することによって算出する。FS(d,k)が高いほど重要なキーワードと見做す。

【0137】

$$FS(d,k) = BM25_{a_{11}}(d,k) \cdot WebIDF(k) \cdot Keyword_score(k) \cdot Pos(k) \quad (11)$$

FD(d,k)は全てのキーワード候補に対して、キーワードの重要度を算出するため、応用アプリケーション等で少数のキーワードしか表示できない場合、上位3～5位のキーワードをその文書内における重要キーワードとして表示する。

【0138】

なお、上記の重要キーワード抽出装置の各構成要素の動作をプログラムとして構築し、重要キーワード抽出装置として利用されるコンピュータにインストールして実行させる、または、ネットワークを介して流通させることが可能である。

【0139】

また、構築されたプログラムをハードディスクや、フレキシブルディスク・CD-ROM等の可搬記憶媒体に格納し、コンピュータにインストールする、または、配布することが可能である。

【0140】

なお、本発明は、上記の実施の形態に限定されることなく、特許請求の範囲内において種々変更・応用が可能である。

【産業上の利用可能性】

【0141】

本発明は、Web文書のキーワードを抽出する技術、例えば、情報検索の事前処理や情報推薦等に適用可能である。

【符号の説明】

【0142】

- 1 文書入力部
- 2 重要キーワード候補抽出手段、重要キーワード候補抽出部
- 3 出現頻度算出手段、出現頻度算出部
- 4 キーワード重要度算出手段、キーワード重要度算出部
- 5 位置情報算出手段、位置情報算出部
- 6 キーワード出力手段、重要キーワード出力部
- 10 主要コンテンツ抽出手段、主要コンテンツ抽出部
- 11 Web文書取得・入力部
- 12 Web文書分割手段、Web文書分割部
- 13 特徴量抽出手段、特徴量抽出部
- 14 主要コンテンツ判定手段、主要コンテンツ判定部
- 15 主要コンテンツ出力手段、主要コンテンツ出力部
- 101 キーワード辞書
- 102 WebIDF記憶部
- 103 キーワード重要度1記憶部
- 104 キーワード重要度2記憶部

10

20

30

40

50

- 1 0 5 第1の記憶手段、固有重要度記憶部
- 1 0 6 第2の記憶手段、位置重要度記憶部
- 1 1 1 データ入力部
- 1 1 2 Web文書ファイル入力部
- 1 1 3 URL入力部
- 1 1 4 Web文書取得部
- 1 1 5 文字コード変換部
- 1 2 1 広告対象領域抽出部
- 1 2 2 ノイズとなるタグや領域除去部
- 1 2 3 Web文書の分割部
- 1 3 1 アンカーリンク情報特徴量抽出部
- 1 3 2 タグ情報特徴量抽出部
- 1 3 3 Web文書で表示される文字列特徴量抽出部
- 1 3 4 特徴量の正規化部
- 1 3 5 特徴量の比率特徴量抽出部
- 1 4 1 特徴量入力部
- 1 4 2 テキスト判定部
- 1 4 3 主要コンテンツ判定部
- 1 5 1 タグ付テキスト出力部
- 1 5 2 タグなしテキスト出力部
- 1 5 3 データ出力部

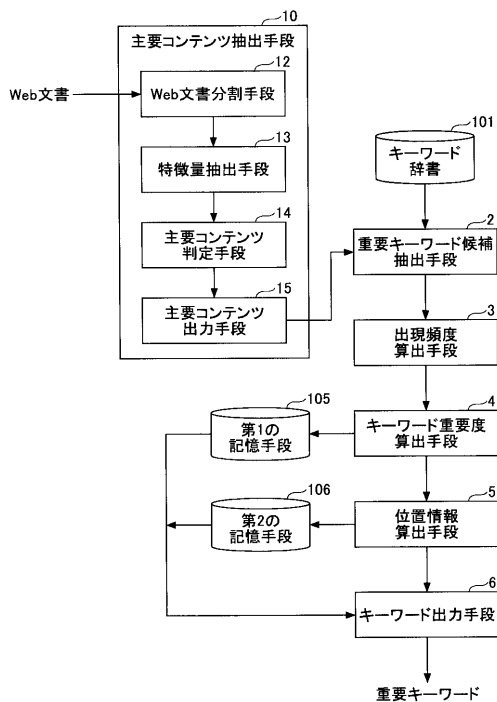
10

20

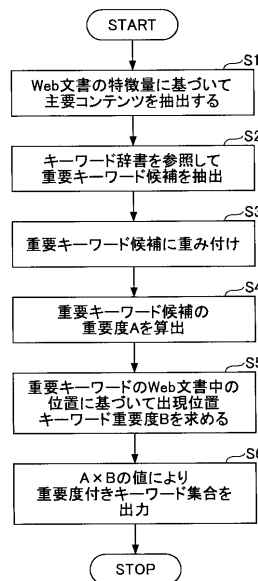
【図1】

【図2】

本発明の原理構成図

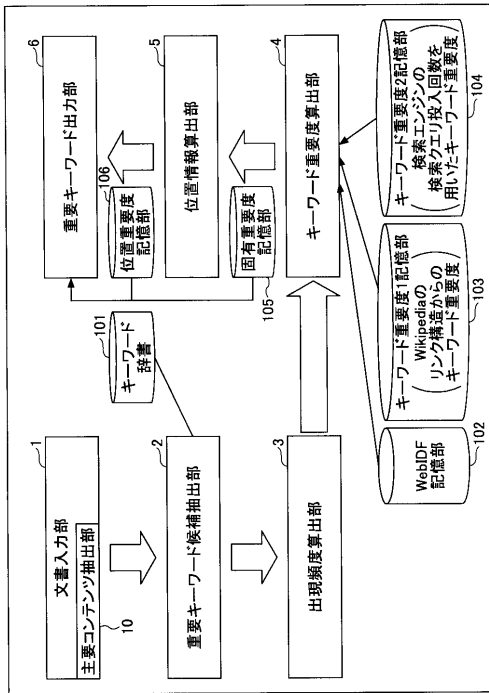


本発明の原理を説明するための図



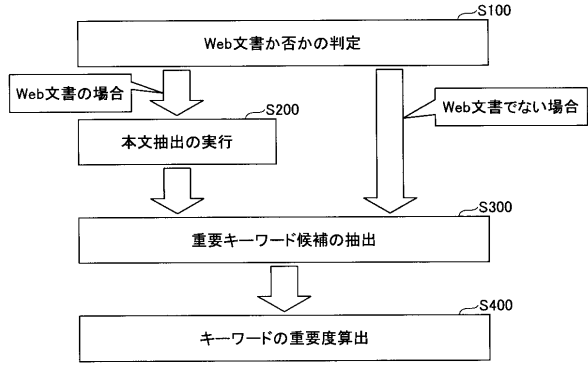
【図3】

本発明の一実施の形態における重要キーワード抽出装置の構成図



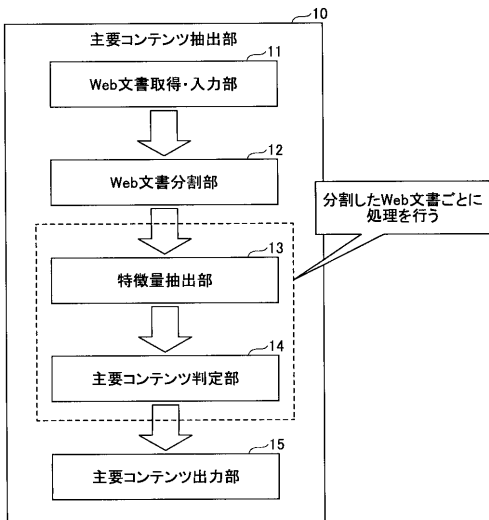
【図4】

本発明の処理の概要を示す図



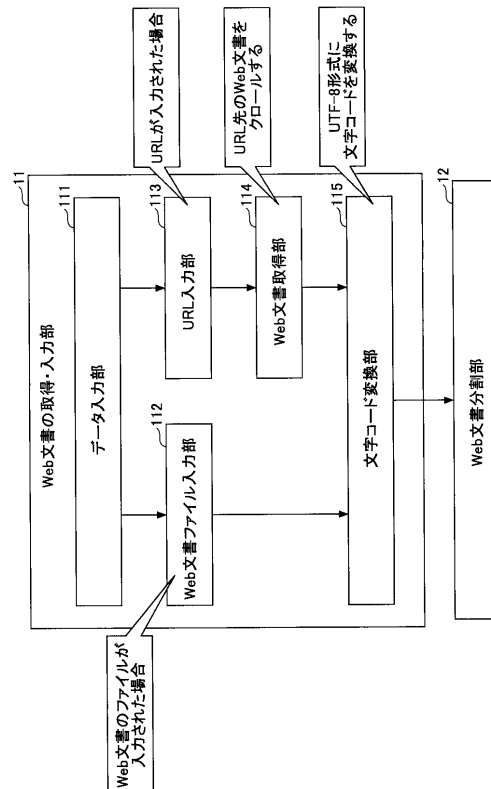
【図5】

本発明の一実施の形態における主要コンテンツ抽出部の構成図



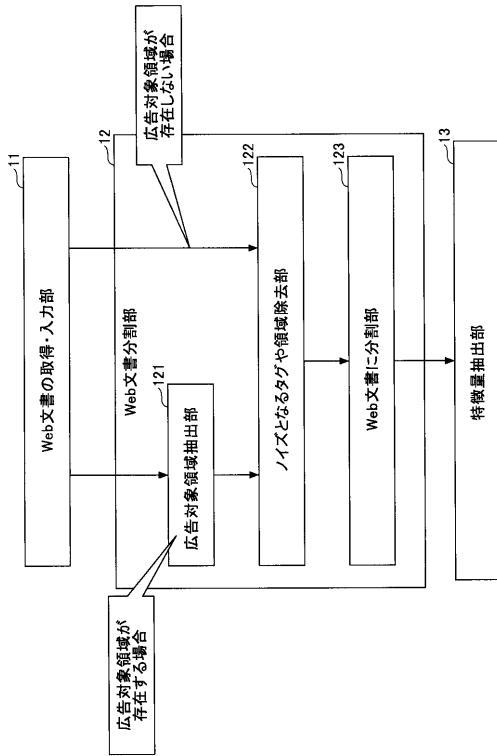
【図6】

本発明の一実施の形態におけるWeb文書取得・入力部の構成図



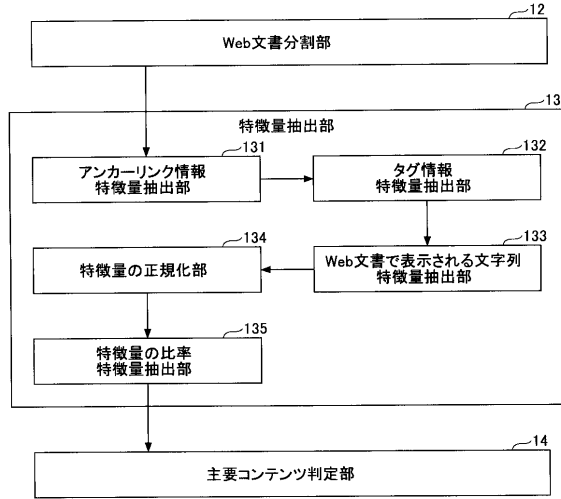
【図7】

本発明の一実施の形態におけるWeb文書分割部の構成図



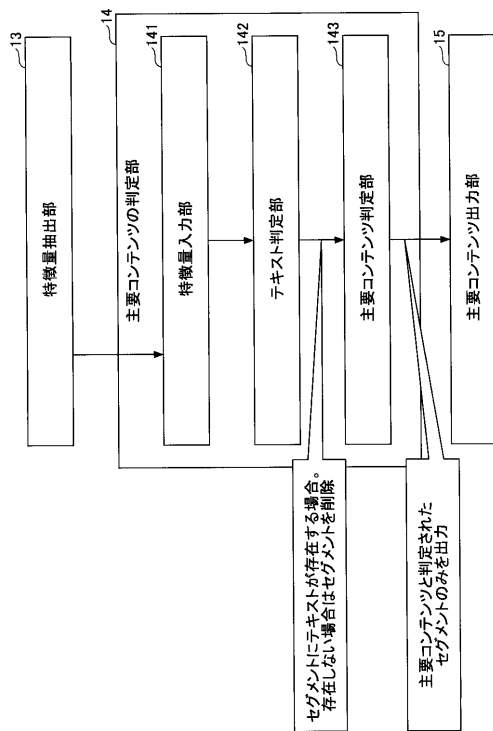
【図8】

本発明の一実施の形態における特徴量抽出部の構成図



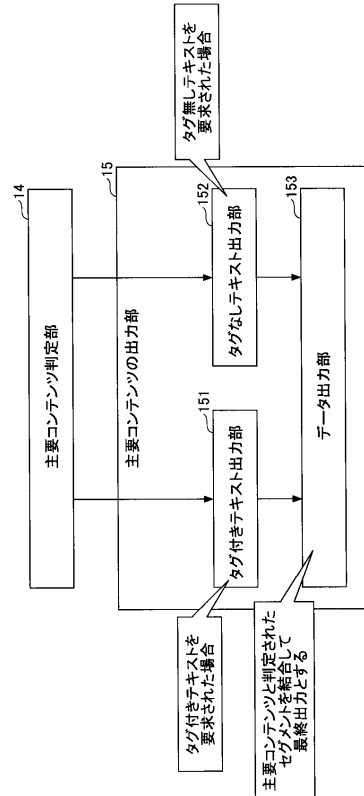
【図9】

本発明の一実施の形態における主要コンテンツ判定部の構成図



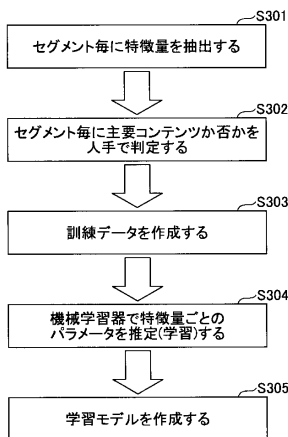
【図10】

本発明の一実施の形態における主要コンテンツ出力部の構成図



【図 1 1】

本発明の一実施の形態における特徴量のパラメータ推定方法のフローチャート



【図 1 2】

本発明の一実施の形態における主要コンテンツ例(その1)

この破線内が主要コンテンツ

▲全体
▲お祝い
▲お祝い
▲お祝い
▲お祝い

▲贈る言葉ブログ

【図 1 3】

本発明の一実施の形態における主要コンテンツ例(その2)

この破線内が主要コンテンツ

▲全体
▲お祝い
▲お祝い
▲お祝い
▲お祝い

▲贈る言葉ブログ

フロントページの続き

(72)発明者 内山 匡

東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内

審査官 鈴木 和樹

- (56)参考文献 近藤光正、外3名、PC上のWeb閲覧履歴からのクエリ抽出技術を用いたモバイル情報検索システム、2008年度人工知能学会全国大会(第22回)論文集 [CD-ROM]、2008年6月11日、p.1-4(2P2-10)
近藤光正、外3名、HITSに基づくWikipediaランキングアルゴリズムとユーザ履歴を用いた個人適応型クエリ推薦、電子情報通信学会 第19回データ工学ワークショップ論文集 [online]、日本、電子情報通信学会データ工学研究専門委員会、2008年4月7日、p.1-8(B2-4)、Internet<URL:http://www.ieice.org/iss/de/DEWS/DEWS2008/proceedings/files/b2/b2-4.pdf>

(58)調査した分野(Int.Cl., DB名)

G06F 17/30