

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第5990178号
(P5990178)

(45) 発行日 平成28年9月7日(2016.9.7)

(24) 登録日 平成28年8月19日(2016.8.19)

(51) Int.Cl.

F I

GO6F 17/30 (2006.01)

GO6F 17/27 (2006.01)

GO6F 17/30 210A

GO6F 17/30 220Z

GO6F 17/30 170A

GO6F 17/27 665

請求項の数 20 (全 18 頁)

(21) 出願番号	特願2013-537776 (P2013-537776)	(73) 特許権者	399037405
(86) (22) 出願日	平成23年11月2日 (2011.11.2)		楽天株式会社
(65) 公表番号	特表2013-544397 (P2013-544397A)		東京都世田谷区玉川一丁目14番1号
(43) 公表日	平成25年12月12日 (2013.12.12)	(74) 代理人	100088155
(86) 国際出願番号	PCT/US2011/058899		弁理士 長谷川 芳樹
(87) 国際公開番号	W02012/061462	(74) 代理人	100113435
(87) 国際公開日	平成24年5月10日 (2012.5.10)		弁理士 黒木 義樹
審査請求日	平成26年8月7日 (2014.8.7)	(74) 代理人	100144440
(31) 優先権主張番号	61/410,392		弁理士 保坂 一之
(32) 優先日	平成22年11月5日 (2010.11.5)	(72) 発明者	スタンキーウィッチ, ゴフィア
(33) 優先権主張国	米国 (US)		アメリカ合衆国, ニューヨーク州, フ
			ォレスト ヒルズ, イエローストーン
			ブルーヴァード 68-37, アパート
			メント シー65

最終頁に続く

(54) 【発明の名称】 キーワード抽出に関するシステム及び方法

(57) 【特許請求の範囲】

【請求項1】

(a) ウェブページからテキストを抽出して少なくとも候補キーワードの第1のセットを生成し、言語処理を適用して少なくとも候補キーワードの第2のセットを生成し、前記候補キーワードの第1及び第2のセットを第1の候補プールに結合する前処理部と、

(b) 少なくとも前記第1の候補プールを記述するデータを前記前処理部から受信して第2の候補プールを生成する候補抽出部と、

(c) 前記前処理部および前記候補抽出部に接続され、少なくとも前記第2の候補プールを記述するデータを受信し、一般的特徴及び言語的特徴について前記第2の候補プールを分析する特徴抽出部と、

(d) 少なくとも前記第2の候補プールを記述する前記データと関連データとを前記特徴抽出部から受信し、前記第2の候補プール内の各候補が1次又は2次キーワードである尤度を判定する分類部と

を備え、

各候補が1次又は2次キーワードである尤度の前記判定が、複数の注釈者からの注釈入力を結合することによって作成されたトレーニングデータに基づき、各注釈が1次キーワードと2次キーワードとの間の区別を含む、

コンピュータシステム。

【請求項2】

前記言語処理の少なくとも一部がトークナイザ及びパーサによって実行される、

請求項 1 に記載のコンピュータシステム。

【請求項 3】

前記言語処理の少なくとも一部がトークナイザ、パーサ、品詞タグ付けプログラム、及び固有表現タグ付けプログラムによって実行される、
請求項 1 に記載のコンピュータシステム。

【請求項 4】

前記言語処理の少なくとも一部がトークナイザによって実行される、
請求項 1 に記載のコンピュータシステム。

【請求項 5】

前記言語処理の少なくとも一部がパーサによって実行される、
請求項 1 に記載のコンピュータシステム。

10

【請求項 6】

前記言語処理の少なくとも一部が品詞タグ付けプログラムによって実行される、
請求項 1 に記載のコンピュータシステム。

【請求項 7】

前記言語処理の少なくとも一部が固有表現タグ付けプログラムによって実行される、
請求項 1 に記載のコンピュータシステム。

【請求項 8】

前記候補キーワードの第 1 のセットがメタデータテキストを含む、
請求項 1 ～ 7 のいずれか一項に記載のコンピュータシステム。

20

【請求項 9】

前記第 2 の候補プールが名詞句及び名詞列を含む、
請求項 1 ～ 8 のいずれか一項に記載のコンピュータシステム。

【請求項 10】

前記第 2 の候補プールが名詞句、名詞列、及び n グラムを含む、
請求項 1 ～ 8 のいずれか一項に記載のコンピュータシステム。

【請求項 11】

前記一般的特徴が頻度、文書中の位置、及び大文字使用のうちの一又は複数を含む、
請求項 1 ～ 10 のいずれか一項に記載のコンピュータシステム。

【請求項 12】

前記言語的特徴が品詞、語句構造、及び固有表現情報のうちの一又は複数に関連する、
請求項 1 ～ 11 のいずれか一項に記載のコンピュータシステム。

30

【請求項 13】

前記一般的特徴が頻度特徴を含み、前記頻度特徴が前記ウェブページ内の関連語出現頻度と語出現頻度のログとのうちの一又は複数を含む、
請求項 1 ～ 10 のいずれか一項に記載のコンピュータシステム。

【請求項 14】

前記一般的特徴が頻度、文書中の位置、及び大文字使用を含み、前記言語的特徴が品詞、語句構造、及び固有表現情報に関連する、
請求項 1 ～ 10 のいずれか一項に記載のコンピュータシステム。

40

【請求項 15】

前記一般的特徴が頻度特徴を含み、前記頻度特徴が前記ウェブページ内の関連語出現頻度と語出現頻度のログとのうちの一又は複数を含み、前記言語的特徴が品詞、語句構造、及び固有表現情報に関連する、
請求項 1 ～ 10 のいずれか一項に記載のコンピュータシステム。

【請求項 16】

前記注釈が、少なくとも一人の注釈者によって前記 1 次キーワードとしてマークを付けられたことを示す第 1 の注釈と、前記複数の注釈者によって前記 2 次キーワードとして選択されたことを示す第 2 の注釈と、一人の注釈者により前記 2 次キーワードとして選択され、かつ他の注釈者により選択された前記 1 次キーワードに部分的に一致することを示す

50

第 3 の注釈とから選択される、

請求項 1 ～ 15 のいずれか一項に記載のコンピュータシステム。

【請求項 17】

コンピュータ処理システムによって実装されるステップを含む方法であって、

(a) ウェブページからテキストを抽出して少なくとも候補キーワードの第 1 のセットを生成し、言語処理を適用して少なくとも候補キーワードの第 2 のセットを生成し、前記候補キーワードの第 1 及び第 2 のセットを第 1 の候補プールに結合するステップと、

(b) 少なくとも前記第 1 の候補プールを記述するデータを受信して第 2 の候補プールを生成するステップと、

(c) 少なくとも前記第 2 の候補プールを記述するデータを受信し、一般的特徴及び言語的特徴について前記第 2 の候補プールを分析するステップと、

(d) 少なくとも前記第 2 の候補プールを記述する前記データと関連データとを特徴抽出部から受信し、前記第 2 の候補プール内の各候補が 1 次又は 2 次キーワードである尤度を判定するステップと

を含み、

各候補が 1 次又は 2 次キーワードである尤度の前記判定が、複数の注釈者からの注釈入力を結合することによって作成されたトレーニングデータに基づき、各注釈が 1 次キーワードと 2 次キーワードとの間の区別を含む、

方法。

【請求項 18】

前記注釈が、少なくとも一人の注釈者によって前記 1 次キーワードとしてマークを付けられたことを示す第 1 の注釈と、前記複数の注釈者によって前記 2 次キーワードとして選択されたことを示す第 2 の注釈と、一人の注釈者により前記 2 次キーワードとして選択され、かつ他の注釈者により選択された前記 1 次キーワードに部分的に一致することを示す第 3 の注釈とから選択される、

請求項 17 に記載の方法。

【請求項 19】

(a) ウェブページからテキストを抽出して少なくとも候補キーワードの第 1 のセットを生成し、言語処理を適用して少なくとも候補キーワードの第 2 のセットを生成し、前記候補キーワードの第 1 及び第 2 のセットを第 1 の候補プールに結合するステップと、

(b) 少なくとも前記第 1 の候補プールを記述するデータを受信して第 2 の候補プールを生成するステップと、

(c) 少なくとも前記第 2 の候補プールを記述するデータを受信し、一般的特徴及び言語的特徴について前記第 2 の候補プールを分析するステップと、

(d) 少なくとも前記第 2 の候補プールを記述する前記データと関連データとを特徴抽出部から受信し、前記第 2 の候補プール内の各候補が 1 次又は 2 次キーワードである尤度を判定するステップと

をコンピュータシステムに実行させ、

各候補が 1 次又は 2 次キーワードである尤度の前記判定が、複数の注釈者からの注釈入力を結合することによって作成されたトレーニングデータに基づき、各注釈が 1 次キーワードと 2 次キーワードとの間の区別を含む、

プログラム。

【請求項 20】

前記注釈が、少なくとも一人の注釈者によって前記 1 次キーワードとしてマークを付けられたことを示す第 1 の注釈と、前記複数の注釈者によって前記 2 次キーワードとして選択されたことを示す第 2 の注釈と、一人の注釈者により前記 2 次キーワードとして選択され、かつ他の注釈者により選択された前記 1 次キーワードに部分的に一致することを示す第 3 の注釈とから選択される、

請求項 19 に記載のプログラム。

【発明の詳細な説明】

10

20

30

40

50

【序文】

【0001】

キーワード抽出は、通常は、ウェブページの内容に合致する広告がページテキストから自動的に選択されたキーワードに基づいて選ばれる、文脈的広告システムの中核的な構成要素としての機能を果たす。そのウェブページに関連し、ユーザにとって興味深いと思われる広告を表示するためには、そのテキスト中に存在する多くの特徴が評価されて、どのキーワードがそのページの内容を的確に反映するかに関する決定を行う必要がある。

【0002】

本明細書に記載の例示的な一実施形態では、キーワード抽出システムはページのurlを入力として取り、上位キーワード候補としてシステムによってランク付けされた10個のキーワード語句を返す。本システムはまずウェブページテキストを処理し、その構造を使用して、キーワード候補プールの役割を果たす語句を抽出する。次いで、各語句はウェブページ上の頻度、テキスト中の位置、大文字使用、及びその言語構造（たとえば、それが名詞句を構成するかどうか）などの特徴のセットによって説明され得る。人間が注釈を付けたキーワードを有するサンプルウェブページのコレクションに基づき、本システムは、候補語句が「良い」キーワードである可能性が高いかどうかの決定にこれらの特徴がいかに寄与するかを学習する。この方式でトレーニングされた後には、本システムは、前に見ていない（すなわち、トレーニングセットになかった）ウェブページ上のキーワードを識別するために使用することができる。

【0003】

既存のキーワード抽出システムの大多数は、tf-idfなどの統計的頻度測定を使用する情報検索モデルに依存する。例示的なシステム実施形態は、性能を改善するために、自然言語処理技法を使用することによってこの手法を改善する。一又は複数の例示的な実施形態では、語句構造に敏感な新しいキーワード候補抽出方法を使用し、より優れた機械学習結果をもたらす追加の言語的特徴を含み得る。

【0004】

ここで、tf-idf重み（語出現頻度 逆文書頻度）は、情報検索及びテキストマイニングでしばしば使用される重みである。この重みは、単語がコレクション又はコーパス中の文書にとってどの程度重要かを評価するために使用される統計的測定である。その重要性は、その単語が文書中に現れる回数に比例して増すが、コーパス中のその単語の頻度によってオフセットされる。

【0005】

例示的な一態様は、（a）ウェブページからテキストを抽出して少なくとも候補キーワードの第1のセットを生成し、言語処理を適用して少なくとも候補キーワードの第2のセットを生成し、候補キーワードの第1及び第2のセットを第1の候補プールに結合する前処理部と、（b）少なくとも第1の候補プールを記述するデータを前処理部から受信して第2の候補プールを生成する候補抽出部と、（c）少なくとも第2の候補プールを記述するデータを受信し、一般的特徴及び言語的特徴について第2の候補プールを分析する特徴抽出部と、（d）少なくとも第2の候補プールを記述するデータと関連データとを特徴抽出部から受信し、第2の候補プール中の各候補が1次又は2次キーワードである尤度を判定する分類部とを備えるコンピュータシステムを含む。

【0006】

一又は複数の例示的な実施形態で、及びそれらの組合せで、（1）言語処理の少なくとも一部はトークナイザ及びパーサによって実行され、（2）言語処理の少なくとも一部はトークナイザ、パーサ、品詞タグ付けプログラム、及び固有表現タグ付けプログラムによって実行され、（3）言語処理の少なくとも一部はトークナイザによって実行され、（4）言語処理の少なくとも一部はパーサによって実行され、（5）言語処理の少なくとも一部は品詞タグ付けプログラムによって実行され、（6）言語処理の少なくとも一部は固有表現タグ付けプログラムによって実行され、（7）候補キーワードの第1のセットはメタデータテキストを含み、（8）第2の候補プールは名詞句及び名詞列を含み、（9）第2

の候補プールは名詞句、名詞列、及びnグラムを含み、(10)一般的特徴は頻度、文書中の位置、及び大文字使用のうちの一又は複数を含み、(11)言語的特徴は品詞、語句構造、及び固有表現情報のうちの一又は複数に関連し、(12)一般的特徴は頻度特徴を含み、頻度特徴はウェブページ内の関連語出現頻度と語出現頻度のログとのうちの一又は複数を含み、(13)各候補が1次又は2次キーワードである尤度の判定は注釈付きトレーニングデータに基づき、(14)各候補が1次又は2次キーワードである尤度の判定は、複数の注釈者からの注釈入力を結合することによって作成されるトレーニングデータに基づき、各注釈は1次キーワードと2次キーワードの区別を含み、(15)一般的特徴は頻度、文書中の位置、及び大文字使用を含み、言語的特徴は品詞、語句構造、及び固有表現情報に関連し、そして/あるいは、(16)一般的特徴は頻度特徴を含み、頻度特徴はウェブページ内の関連語出現頻度と語出現頻度のログとのうちの一又は複数を含み、言語的特徴は品詞、語句構造、及び固有表現情報に関連する。

10

【0007】

もう一つの態様は、コンピュータ処理システムによって実装されるステップを含む方法を含み、当該ステップは、(a)ウェブページからテキストを抽出して少なくとも候補キーワードの第1のセットを生成し、言語処理を適用して少なくとも候補キーワードの第2のセットを生成し、候補キーワードの第1及び第2のセットを第1の候補プールに結合するステップと、(b)少なくとも第1の候補プールについて説明するデータを受信して第2の候補プールを生成するステップと、(c)少なくとも第2の候補プールについて説明するデータを受信し、一般的特徴及び言語的特徴について第2の候補プールを分析するステップと、(d)少なくとも第2の候補プールについて説明するデータと関連データとを特徴抽出部から受信し、第2の候補プール内の各候補が1次又は2次キーワードである尤度を判定するステップとを含む。

20

【0008】

もう一つの態様は、(a)ウェブページからテキストを抽出して少なくとも候補キーワードの第1のセットを生成し、言語処理を適用して少なくとも候補キーワードの第2のセットを生成し、候補キーワードの第1及び第2のセットを第1の候補プールに結合するステップと、(b)少なくとも第1の候補プールについて説明するデータを受信して第2の候補プールを生成するステップと、(c)少なくとも第2の候補プールについて説明するデータを受信し、一般的特徴及び言語的特徴について第2の候補プールを分析するステップと、(d)少なくとも第2の候補プールについて説明するデータと関連データとを特徴抽出部から受信し、第2の候補プール内の各候補が1次又は2次キーワードである尤度を判定するステップとを含むステップ群を実行するように動作可能なソフトウェアを格納した、コンピュータ読み取り可能な記録媒体を備える。

30

【0009】

本明細書で提供される説明及び図面から、他の態様及び実施形態が当業者には明らかとなる。

【図面の簡単な説明】

【0010】

【図1】例示的な一実施形態の処理の概要を示す図である。

40

【図2】例示的な一実施形態が実装され得るコンピュータシステムを示す図である。

【選ばれた例示的な実施形態の詳細な説明】

【0011】

コンピュータで実装される例示的な実施形態の概要を図1に示す。各構成要素は本明細書の残りの節でさらに詳細に説明される。

【0012】

例示的な前処理部

例示的な一実施形態では、潜在的なキーワード語句がそのページから選択され得る前に、そのページのプレーンテキストがHTML形式から抽出されてもよい。順に、このテキストをさらに処理することで、キーワード抽出システムに有用であり得るその構造に関す

50

る情報を取得することができる。本システムの前処理部は、好ましくは、ウェブページテキストの抽出並びにタグ付け及び書式付けを実行して、後に続く候補語句選択及び特徴抽出の段階のための適切な入力を提供する。

【0013】

その前処理段階で、まず、定型的な内容を除去してそのページの主要なテキスト本体のみを保存するBoiler Pipe（たとえば[9]を参照）を使用して、ウェブページから内容テキストを抽出してもよい。本体テキストに加えて、タイトル、メタ記述、メタキーワードなどのヘッダ情報を抽出し、Boiler Pipe出力と結合することで、さらなる処理のためのプレーンテキスト入力を作成してもよい。

【0014】

次いで、ページテキストはトークン化され、そのトークナイザ出力は品詞タグ付けプログラム（たとえば[18]を参照）及びパーサ（たとえば[13]を参照）に渡される。キーワードは名詞句を構成する傾向があるので、そのテキスト中で名詞句を見つけるためにパーサ出力を使用してもよい。チャンカではなくパーサの使用は、キーワード候補抽出を改良するために、基本的名詞句のかたまりに対するものとして、階層的語句構造のより肌理の細かい情報を取得したいという要望によって動機付けられ得る。

【0015】

個人名又は組織名などの固有表現（「NE」）は有用なキーワードとなり得るので、固有表現もウェブページテキストから抽出してもよい。二つの異なるNEシステム（たとえば[18]、[4]を参照）を、エンティティタイプの大きなセットを網羅するために用いるのが好ましい。

【0016】

例示的な候補抽出部

候補抽出は、潜在的なキーワードであり、且つ所与の語句がキーワードである尤度を推定する分類器の入力として使用することができる語句を選択するために用いてもよい。加えて、例示的な分類器が語句をトレーニングしている間に、候補抽出の精度の向上は、キーワードではありそうにない単語の組合せをフィルタリングするのに役立ち、良くないトレーニングサンプルの量を減らし、ひいては、良いトレーニングデータと悪いトレーニングデータとの比率を改善する（キーワード抽出タスクは、良いラベルデータが殆どなく、良いサンプルと悪いサンプルの間の不均衡を有する）。

【0017】

例示的な一実施形態で、キーワード抽出方法は以下のように行われる。まず、基底候補セットが、構文解析されたテキストからすべての名詞句を再帰的に抽出することによって形成される。次いで、名詞のみから成るすべての候補部分列（左から右に抽出された）が候補セットに追加される（たとえば、「最高のニクソン・カメラ・アクセサリ」が候補である場合、「ニクソン・カメラ・アクセサリ」、「カメラ・アクセサリ」及び「アクセサリ」がその候補セットに追加されることになる）。そして、候補セットは、候補語句から抽出されたすべてのユニグラム、バイグラム、及びトライグラムにより増補される。

【0018】

候補セットについて、最頻出の英単語のストップリストに対するフィルタをかけてもよい。ストップワードを含むユニグラム又はバイグラムを候補セットから取り除くのが好ましい。しかし、その真ん中にストップリスト内の単語を含むより長い語句は保持してもよい。

【0019】

例示的なキーワード分類部

どの候補語句がキーワードであるかを識別するために、例示的な一実施形態では、入力（候補語句の特徴）を使用してその語句がキーワードである確率を推定する分類器を使用し、出力ラベル（キーワード又は非キーワード）をその語句に割り当てる。特徴入力をキーワードラベルにマップする分類機能は、監視された機械学習を使用して取得してもよい。すなわち、本マッピングは、「正しい」出力ラベルが注釈者によって提供されたデータ

10

20

30

40

50

セットに基づく分類器システムによって学習され得る。

【 0 0 2 0 】

例示的なシステムの分類器をトレーニングするために、最大エントロピー（ M E ）モデルを用いてもよい（これはときにロジスティック回帰モデルと呼ばれる。その紹介については [1 1] を参照）。 M E モデルは、トレーニングデータから制約条件を導出し、そのトレーニングセットによって網羅されない場合に最大エントロピーの分配を仮定する。その M E 分類器の入力は、各特徴に関連付けられた重みを学習するためにそのモデルによって使用される、各キーワード候補の値のベクトルで構成される。新しい入力データが与えられ、トレーニングされた分類器は、語句が、候補語句の入力値を与えられたキーワードである確率を計算することができる。

10

【 0 0 2 1 】

入力値

【数 1】

\vec{x}

を与えられたラベル c の確率は以下の公式により計算することができる。

【数 2】

$$P(c|\vec{x}) = \frac{\exp(\sum_i \alpha_i f_i(\vec{x}, c))}{1 + \exp(\sum_i \alpha_i f_i(\vec{x}, c))}$$

20

ここで、 f は共同特徴（入力ベクトル及びラベルの関数）であり、 α_i はその特徴に割り当てられた重みである。

【 0 0 2 2 】

最大エントロピー分類器をトレーニングするために、自然言語ツールキット（ N a t u r a l L a n g u a g e T o o l k i t （ N L T K ））。 [1] を参照）で提供される P y t h o n ライブラリを使用することができる。 C G はトレーニング方法として用いることができる。しかし、そのアルゴリズムは一意解に収束するので、トレーニング方法の選択は分類器の性能に影響せず、他のトレーニング方法が本明細書に記載の本発明の範囲を逸脱することなく使用可能であることが当業者に理解されよう。たとえば、サポート・ベクトル・マシン（ r b f カーネル）（ [8] を参照）が使用可能であるが、 M E モデルを使用して得られる結果に比した改善は本発明者によって発見されなかった。

30

【 0 0 2 3 】

ここで、 C G は共役勾配法のことであり、これは、分類器ライブラリ内のトレーニング方法の一つとして提供されるスパース 1 次方程式系を解くための標準反復法である。 C G は P y t h o n 及び N L T K でインストールされるサイパイ（ s c i p y ）パッケージ（ <http://www.scipy.org/> ）を必要とする。

【 0 0 2 4 】

良いトレーニングデータと悪いトレーニングデータとの不均衡（すなわち、トレーニングデータ内の候補語句の大部分が通常はキーワードではない）のために、分類器によって割り当てられたラベル（キーワード又は非キーワード）を使用しないことを選択してもよいが、その代わりに、たとえば所与のウェブページ内で最も高い確率を有する 1 0 個の候補を選択して、確率のスコアに直接基づいて候補をランク付けすることができる。

40

【 0 0 2 5 】

例示的な特徴

特徴値のセットは各キーワード候補について計算されて分類器への入力として使用することができる。特徴の選択は分類器の性能で重要な役割を果たす。その特徴は、（ a ）一般的、非言語的特徴と、（ b ）言語的特徴の二つのタイプに分けることができる。一般的特徴は、 [1 7] に記載のシステムによって使用される特徴と同様でもよく、頻度、文書

50

中の位置、及び大文字使用などの情報を含む。言語的特徴は品詞、語句構造、及び固有表現情報を使用する。二つのタイプの特徴について以下にさらに詳しく説明する。

【 0 0 2 6 】

【表 1】

例示的な分類特徴

特徴名	定義	タイプ	[17]での使用
TF	キーワード候補が文書中で出現する回数／文書中の候補の総数	実数値	はい
TFLog	Log(TF+1)	実数値	はい
DF	Google Ngramコーパスのユニグラム及びバイグラムデータに基づくキーワード候補の相対的頻度。候補語句が2単語より長い場合、その候補語句内のすべてのバイグラムの平均頻度が使用される。	実数値	はい(異なる文書コレクションに基づく)
DFLog	Log(DF+1)	実数値	はい
Title	キーワード候補が文書のタイトル中にある場合は1、そうでない場合は0。	2値	はい
IsCap	キーワード候補内のすべての単語が大文字で書かれている場合は1、そうでない場合は0。	2値	はい
HasCap	キーワード候補内の少なくとも1つの単語が大文字で書かれている場合は1、そうでない場合は0。	2値	はい
Location	文書中の単語の総数に対する、当該文書内のキーワード候補の最初の出現の相対的な位置。	実数値	はい
LocationLog	Log(Location)	実数値	はい
Length	キーワード候補内の単語の数	実数値	はい
Url	キーワード候補が文書のurlに含まれる場合は1、そうでない場合は0。	2値	はい
IsNoun	キーワード候補中のすべての単語が名詞である場合は1、そうでない場合は0。	2値	はい(しかし、固有名詞と総称的名詞の区別で異なるように定義される)
hasNoun	キーワード候補中の少なくとも1つの単語が名詞である場合は1、そうでない場合は0。	2値	はい
isNP	キーワード候補が名詞句である場合は1、そうでない場合は0。	2値	はい
hasNP	キーワード候補が名詞句を含む場合は1、そうでない場合は0。	2値	いいえ
POS sequence	キーワード候補の品詞タグの列	実数値	いいえ
IsNE_oak	キーワード候補がOAKシステムによって見つけられた固有表現である場合は1、そうでない場合は0。	2値	いいえ
HasNE_oak	キーワード候補がOAKシステムによって見つけられた固有表現を含む場合は1、そうでない場合は0。	2値	いいえ
TagNE_oak	OAKシステムによってキーワード候補に割り当てられたNEタグ	実数値	いいえ
IsNE_Stanford	キーワード候補がスタンフォードNERシステムによって見つけられた固有表現である場合は1、そうでない場合は0。	2値	いいえ
HasNE_Stanford	キーワード候補がスタンフォードNERシステムによって見つけられた固有表現を含む場合は1、そうでない場合は0。	2値	いいえ
TagNE_Stanford	スタンフォードNERシステムによってキーワード候補に割り当てられたタグ	実数値	いいえ
Pmi	候補語句の自己相互情報量のスコア	実数値	いいえ
Iswiki	キーワード候補がウィキペディアのタイトルである場合は1、そうでない場合は0。	2値	いいえ(しかし、[14]及び[17]で使用する探索照会特徴と同様)
WikiFreq	キーワードがウィキペディアのタイトルである場合はウィキペディアトラフィック頻度、そうでない場合は0。	実数値	いいえ(しかし、[14]及び[17]で使用する探索照会特徴と同様)

【 0 0 2 7 】

例示的な一般的特徴

10

20

30

40

50

例示的な頻度特徴

頻度特徴は、 $TF \times IDF$ と同様の情報を提供する。頻度特徴は、文書内の関連語出現頻度と語出現頻度のログとの他に DF （文書コレクション中の頻度）と $\log DF$ 値とを含んでもよい。 DF 値は、Google Ngramコーパスからの頻度を使用して概算してもよい。好ましくは、ユニグラム及びバイグラムの頻度情報のみが、 DF を計算するために使用される。2単語よりも長い候補語句については、その語句内のすべてのバイグラムの DF の平均値を DF 値として用いてもよい。平均値は、異なる長さの語句についての値の似た範囲を得るために使用することができる。また、ブログのコレクション全体について計算された DF 値をGoogle Ngramコーパスからの頻度の代わりに用いてもよい。

10

【0028】

ここで、 $TF \times IDF$ は語出現頻度 - 逆文書頻度のことをいい、用語の相対的重要性を評価するために情報検索で使われる標準スコアである。これは、文書のコレクション内のその全体的頻度によって相殺された所与の文書内の用語の頻度に基づく。 tf 及び idf の標準的公式は以下の通りである。

$$tf_{i,j} = n_{i,j} / \sum_k n_{k,j}$$

ここで、 $n_{i,j}$ は、文書 j 内の考察される用語 i の出現回数である。

$$idf_i = \log(|D| / |d : t_i|)$$

これは、コレクション内のすべての文書の数のログを、用語 i を含む文書の数で割ったものである。

20

【0029】

2. タイトル

候補語句がその文書のタイトル中にあるかどうか。

【0030】

3. 大文字使用

大文字で書かれた単語は、所与の文書中の重要な用語としてマークされた固有名詞又は単語を含む。例示的な大文字使用の特徴は、キーワード候補内のすべての単語が大文字で書かれているかどうか、及び、候補語句内の少なくとも一つの単語が大文字で書かれているかどうか、である。

【0031】

4. 位置

文書内でキーワード候補が最初に出現する相対的な位置であり、単語の数で数える。たとえば、文書が20個の単語を有し、候補語句の最初の単語がその文書中の5番目の単語である場合には、位置 = $5 / 20 = 0.25$ である。

30

【0032】

5. 長さ

候補語句中の単語の数。

【0033】

6. URL

候補語句がページurl内にあるかどうか。

40

【0034】

7. Wikiトラフィック

ウィキペディア・トラフィックの統計値に基づく特徴を、頻出する探索/照会項目としてキーワード候補の人気を反映するために用いてもよい。この特徴のセットは、候補語句がウィキペディアのタイトル（リダイレクトを含む）であるかどうかと、その候補語句のトラフィック数字（その候補がウィキペディアのタイトルでない場合には0）とを含み得る。トラフィック統計値は、一定期間（たとえば、2010年6月中の20日間）に亘って集められた1時間当たりのウィキペディアのログに基づくものであってもよい。

【0035】

例示的な言語的特徴

50

1. 名詞句

その候補が名詞句である又は名詞句を含むかどうか。

2. 名詞

その候補語句が少なくとも一つの名詞を含むかどうか、及び、その候補語句が名詞のみから成るかどうか。

3. POS タグ

その候補語句に割り当てられた品詞タグの列。

4. 固有表現

キーワード候補が固有表現であるかどうか、キーワード候補が、固有表現及びその候補に割り当てられた固有表現タグを含むかどうか（その候補語句が NE でない場合には「0」）。

10

5. PMI

自己相互情報量 (Pointwise mutual information (PMI)) は、語句がコロケーションである可能性が高いかどうかを反映する。候補語句の PMI スコアは以下のように計算することができる。

バイグラムについては、

$$PMI(w_1, w_2) = \log \{ P(w_1, w_2) / P(w_1) * P(w_2) \}$$

ただし、 $P(w)$ は単語又は語句の相対的頻度である。

単一の単語については、

$$PMI = \log \{ 1 / P(w) \}$$

20

2 単語よりも長い候補語句については、PMI は、その語句内のすべてのバイグラムの PMI スコアの平均値に設定してもよい。

【0036】

例示的な評価及びデータ

例示的なトレーニングデータ

トレーニングデータは、たとえば 500 個のウェブページ（ブログページのコーパスから無作為に選択されたもの。[3] を参照）を含んでもよい。注釈者は、ブログページから抽出されたプレーンテキストを提示され、そのページの内容を最もよく表すキーワードを選択するように指示され得る。ヘッダからのメタ情報は注釈付きテキストに含まれないのが好ましい。単一のページについて選択することができるキーワードの数には制限を設けないのが好ましい。追加のページについても、注釈を付け、トレーニングに使用されないテストセットとして除外することができる。

30

【0037】

注釈者間の合意及び最も理想的な基準 (golden standard)

各ページについて、キーワードは 2 人の注釈者によって選ばれるのが好ましい。しかし、このタスクへの注釈者間の合意は高くなくてもよい（たとえば、一実装形態では、注釈者のカップスコアは 0.49 であった）。低いカップスコアの要因としては次のことが考えられる。第 1 に、注釈者が、部分的にのみ一致する同様の語句にタグを付けることがある。第 2 に、選択可能なキーワードの最大数が指定されていないときに、一人の注釈者が、所与のテキストについてもう一人よりも多くのキーワードを選択することを選ぶことがある。

40

【0038】

ここで、コーエンのカッパ係数は、分類タスクでの注釈者間の合意を測定するために一般に使用される統計的測定である。カッパは、 $\{ P(A) - P(E) \} / \{ 1 - P(E) \}$ で計算され、ここで、 $P(A)$ はコーダ間で観測された合意であり、 $P(E)$ はコーダが偶然合意した確率である。0.6 ~ 0.7 を上回るカッパスコアは「実質的合意」として考えられる。

【0039】

単独の注釈者に依存しない最も理想的な基準 (Golden Standard (GS)) を作成するために、両方の注釈者からの出力を結合してもよい。キーワードに注釈を

50

付けるとき、注釈者はそのキーワードが「１次キーワード」であるか「２次キーワード」であるかも選択するように指示され得る。１次キーワードは、文書の主題又は中心思想をとらえたキーワードとして規定することができる。２次キーワードは、その文書に関する追加のキー情報（たとえば、その事象が起きた位置、追加のものではあるが重要な、記述された数字など）を提供する重要な語句として規定することができる。両方の注釈者の選択を正確に反映するキーワードセットを作成するために、ＧＳで以下のキーワードを保持することができる。

１．（一人の注釈者又は双方によって）１次としてマークを付けられたすべてのキーワード。

２．両方の注釈者によって選択された２次キーワード。

３．一人のみにによって選択されたが、他方の注釈者によって選択された１次キーワードに部分的に一致する、２次キーワード。

【００４０】

ＧＳを使用する例示的な一実施形態では、各注釈者と標準との間のカップスコアは、注釈者１については０．７５で注釈者２については０．７４であった。１次及び２次キーワードの詳細な合意統計値を下記の表２に示す。

【００４１】

【表２】

注釈者1/注釈者2	1次	2次	キーワードではない
1次	1652	319	1796
2次	437	264	1777
キーワードではない	1069	264	////////

【００４２】

例示的な候補抽出部

上記のように、例示的な一実施形態では、基底候補セットとして名詞句を使用するが、その名詞句から抽出された名詞列とユニグラム、バイグラム、及びトライグラムで候補プールを増補する。

【００４３】

起こり得るすべての候補語句をテキストから取得する一つの従来の手法では、その候補セット内の長さ n （通常は３～５）までのすべての n グラムを含む。この n グラム方法の重大な欠点は、それが、意味のある語句でない及び／又は潜在的キーワードでなさそうな単語列の形で、かなりのノイズをもたらすことである。したがって n グラム方法は低い精度に悩まされる。

【００４４】

従来の一代替方法では、候補を抽出するために言語構造のキューを使用する。キーワードは名詞句である傾向があるので、テキストからのすべての名詞句が候補プールを形成するために使用され得る。しかし、この方法では再現率が n グラム抽出方法よりも著しく低く、これは多くの潜在的なキーワードが候補セットに含まれないことを意味する。

【００４５】

n グラム及び名詞句戦略の精度、再現率、及び F 測定が本発明者により例示的な一実施形態の抽出方法と比較された。言い換えれば、本発明者は、さらなる分類段階なしに、キーワードを選択するための唯一の方法としてそれぞれの手法が使用された場合にそれがどの程度効果的であるかを評価した。結果は以下の表３に要約される。

【００４６】

10

20

30

40

【表 3】

候補抽出方法の比較

方法	総数	候補にな いキー	候補にあ るキー	精度	再現度%	Fスコア
Nグラム	365,779	786	4839	1.3	85.9	2.6
名詞句	14,441	4160	1465	10.4	26	14.6
例示的实施 形態	85,059	1008	4617	5.4	81.95	10.2

【0047】

10

表3に示すように、nグラム手法の再現率は80%を超えるが、それはまた、その3つの方法のうちで精度が最も低い（すなわち、候補セットがかなりの量のノイズを含む）。候補として名詞句を抽出することは精度が上がる点で有利だが、この方法では再現率が非常に低く（わずか26%）、潜在的なキーワードを見落とす蓋然性が高い。

【0048】

対照的に、本発明の方法の例示的な実施形態は、名詞句を抽出することに比べて再現率を改善する。この手法の再現率はnグラム方法と同程度であるが、精度はより高くなる。異なる手法が分類器の性能とどのように結合するかの評価結果を以下に述べる。

【0049】

分類器の性能

20

システム性能全体を評価するために、本発明のシステムによって達成される結果が[17]に基づいて基準値と比較された。基準値システムで、候補抽出方法はnグラム方法であり、特徴は(NP/名詞特徴の単純なセットを加えた)一般的非言語的特徴で構成される。(a)結合された候補抽出方法の使用と(b)分類段階での言語的特徴の追加とでシステム性能がどのように変わったかを分析した。

【0050】

基準値に対して本発明のシステムを比較する際には以下の二つの評価測定を用いた。

1. R精度（上位n個の結果のうちいくつの候補がキーワードであるか。ここで、nはページ上に有り得るキーワードの総数）。

2. 上位10個のスコア（R精度に似ているが、上位10個の結果で切り捨てる。すなわち、すべての $n > 10$ が10にセットされる）。

30

【0051】

上位10個のスコアを有する候補がキーワード出力として選択されるときにどのように分類器が抽出システムとして機能するかの推定を上位10個の測定が提供するので、その測定を評価に用いた。システム性能は、分類器トレーニングでは決して使用されなかった100個のウェブページの提出されたテストセットでテストされ（表4を参照）、そして、相互検証テストは、500ページのトレーニングセットで行われた（それぞれが約50個の文書の10倍。表5を参照）。

【0052】

【表 4】

40

提出されたセットの上位10個のスコア結果

方法	一般的特徴	一般的+言語的特徴
Nグラム	43.71	47.24
本発明	48.28	49.84

【0053】

【表 5】

相互検証テストの上位10個のスコア結果

方法	一般的特徴	一般的+言語的特徴
Nグラム	45.97	49.01
本発明	48.21	51.74

【0054】

基準値と本発明のシステムとの結果の差は統計的に重要である（相互検証結果への対応のある両側 t 検定によれば、 $p = 0.0001$ ）。基準値に対する相対的な改善は 12.55%である。

10

【0055】

関連実施形態

所与のウェブページについて文脈的に適切な広告を選択する二つの好ましい手法がある。一つの手法は、広告プールへのウェブページテキストの直接の突き合わせを含む。もう一方の手法では、そのページと広告の間の突き合わせが中間キーワード抽出ステップを含む。各手法の例を以下に示す。

【0056】

キーワード抽出

1. KEA [5]
2. GenEx [15]、[14]
3. Yih 他 [17]
4. Hult h [7]、[6]
5. その他：[10]、[16]

20

【0057】

文脈的広告

1. Broder 他 [2]
2. Ribeiro - Neto 他 [12]。

【0058】

本明細書に記載のいくつかの例示的なシステム及び方法の実施形態では、キーワード抽出は (a) ウェブページからのテキスト抽出、並びに、品詞タグ付け及び構文解析などの言語的処理を含む前処理と、(b) キーワード候補語句の抽出と、(c) 監視された機械学習を使用する候補分類とを含むのが好ましい。

30

【0059】

本発明のシステム及び方法は、候補選択及び特徴抽出段階の両方で、言語的情報の使用による性能の改善を達成することができる。例示的な一実施形態は、階層的語句構造を使用する候補選択を含み、よりノイズの少ない候補プールをもたらす。分類に使用できる特徴は品詞や固有表現情報などの言語的特徴も含み、分類器の性能の改善をもたらす。

【0060】

実施形態は、当業者には明らかであろうコンピュータ構成要素及びコンピュータ実装ステップを含む。たとえば、計算及び通信は電子的に実行でき、結果をグラフィカルユーザインターフェースを用いて表示することができる。

40

【0061】

そのような例示的なシステムを図 2 に示す。コンピュータ 100 はネットワーク 110 を介してサーバ 130 と通信する。複数のデータソース 120、121 もネットワーク 110 を介してサーバ 130、プロセッサ 150、及び / 又は、情報を計算及び / 又は送信するように動作可能な他の構成要素と通信する。一又は複数のサーバ 130 は、一又は複数の記憶装置 140、一又は複数のプロセッサ 150、及びソフトウェア 160 に結合され得る。

【0062】

本明細書に記載の計算及び同等のものは、一実施形態では、完全に電子的に実行される

50

。当業者には明らかなように、他の構成要素及び構成要素の組合せもまた、本明細書に記載の処理データ又は他の計算をサポートするために使用できる。サーバ130は、一又は複数のプロセッサ150とのデータ通信と、記憶装置140からのデータ通信と、コンピュータ100への通信とを円滑に進めることができる。プロセッサ150は、一時的な情報又は他の情報を記憶するために使用することができるローカル又はネットワークストレージ（図示せず）を任意で含んでもよいし、それと通信してもよい。ソフトウェア160は、コンピュータ100またはプロセッサ150でローカルにインストールされてもよいし、及び/又は、計算と適用とを容易にするために中央でサポートされてもよい。

【0063】

説明を容易にするために、本発明のあらゆるステップ又は要素がコンピュータシステムの部分として本明細書に記載されているわけではないが、各ステップ又は要素は対応するコンピュータシステム又はソフトウェア構成要素を有してもよいことは当業者に理解されよう。したがって、そのようなコンピュータシステム及び/又はソフトウェア構成要素は、それらの対応するステップ又は要素（すなわち、それらの機能性）について説明することによって可能にされ、本発明の範囲内にある。

【0064】

さらに、コンピュータシステムが特定の機能を実行するためのプロセッサを有するものとして説明又は特許請求される場合、そのような使用は、単一のプロセッサが、たとえば、さまざまなプロセッサに委託されたタスクのいくつか又はすべてを実行するシステムを除外するものとして解釈されるべきではないことが当業者に理解されよう。すなわち、本明細書及び/又は特許請求の範囲で指定されるプロセッサの任意の組合せ又はすべては同一のプロセッサでもよい。そのようなすべての組合せは本発明の範囲内である。

【0065】

別の方法として、又は組み合わせて、処理及び意思決定はデジタル信号プロセッサ回路又は特定用途向け集積回路などの機能的に同等の回路によって実行可能である。

【0066】

ループ及び変数の初期化と一時的数値変数の使用などの多数のルーチンプログラム要素は本明細書には記載されていない。さらに、特段の指示のない限り、記載されるステップの特定の順番は説明のみを目的とし、一般に、本発明の範囲から逸脱することなしに変更することができるが、当業者に理解されよう。特段の記述のない限り、本明細書に記載のプロセスは順序付けされていない。すなわち、そのプロセスは任意の妥当な順番で実行することができる。

【0067】

本明細書に記載のすべてのステップは、実行可能な場合、ソフトウェアによる実装が可能であることが当業者に理解されよう。さらに、そのようなソフトウェアは非一時的なコンピュータ可読媒体に格納可能であり、一又は複数のコンピュータプロセッサによって実行可能であることが、当業者に理解されよう。

【0068】

本発明は、本明細書で概説される例示的な態様の実施形態とともに説明されているが、多数の代替、修正、及び変更が当業者には明らかとなる。したがって、本明細書に記載するような本発明の例示的な態様及び実施形態は例示的なであって限定ではない。さまざまな変更本発明の趣旨及び範囲を逸脱することなしに行われ得る。

【0069】

参考文献

[1] Bird, Steven, Edward Loper及びEwan Klein. Natural Language Processing with Python. O'Reilly Media Inc., 2009.

[2] Broder, Andrei及びFontoura, Marcus及びJosifovski, Vanja及びRiedel, Lance. A semantic approach to contextual advertising. SIGIR '07

10

20

30

40

50

: Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval、ページ559~566、New York, NY, USA、2007。ACM。

[3] Kevin Burton及びAkshay Java及びIan Soboroff。ICWSM 2009 Spinn3r Dataset。San Jose, CA、2009。

[4] Finkel, Jenny Rose及びGrenager, Trond及びManning, Christopher。Incorporating non-local information into information extraction systems by Gibbs sampling。ACL '05: Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics、ページ363~370、Morristown, NJ, USA、2005。Association for Computational Linguistics。

[5] Frank, Eibe及びPaynter, Gordon W.及びWitten, Ian H.及びGutwin, Carl及びNevill-Manning, Craig G.。Domain-specific keyphrase extraction。IJCAI '99: Proceedings of the 16th international joint conference on Artificial intelligence、ページ668~673、San Francisco, CA, USA、1999。Morgan Kaufmann Publishers Inc.。

[6] Hulth, Anette。Automatic Keyword Extraction。Combining Machine Learning and Natural Language Processing。Verlag Dr. Muller、2008。

[7] Hulth, Anette。Improved automatic keyword extraction given more linguistic knowledge。Proceedings of the 2003 conference on Empirical methods in natural language processing、ページ216~223、Morristown, NJ, USA、2003。Association for Computational Linguistics。

[8] Joachims, Thorsten。Making large-scale support vector machine learning practical。: 169~184、1999。

[9] Kohlschütter, Christian及びFankhauser, Peter及びNejdl, Wolfgang。Boilerplate detection using shallow text features。WSDM '10: Proceedings of the third ACM international conference on Web search and data mining、ページ441~450、New York, NY, USA、2010。ACM。

[10] Matsuo, Y.及びIshizuka, M.。Keyword Extraction from a Document using Word Co-occurrence Statistical Information。Transactions of the Japanese Society for Artificial Intelligence、17: 217~223、2002。

[11] Adwait Ratnaparkhi。A Simple Introduction to Maximum Entropy Models for Natur

10

20

30

40

50

al Language Processing. Technical report, IRCS, 1997.

[12] Ribeiro-Neto, Berthier及びCristo, Marco及びGolgher, Paulo B.及びSilva de Moura, Edleno. Impedance coupling in content-targeted advertising. SIGIR '05: Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval、ページ496~503、New York, NY, USA、2005。ACM。

10

[13] Sekine, Satoshi及びRalph Grishman. A corpus based probabilistic grammar with only two non-terminals. Fourth International Workshop on Parsing Technology, 1995.

[14] Turney, Peter D.. Coherent keyphrase extraction via web mining. IJCAI '03: Proceedings of the 18th international joint conference on Artificial intelligence、ページ434~439、San Francisco, CA, USA、2003。Morgan Kaufmann Publishers Inc..

20

[15] Turney, Peter D.. Learning Algorithms for Keyphrase Extraction. Inf. Retr.、2(4):303~336、2000。

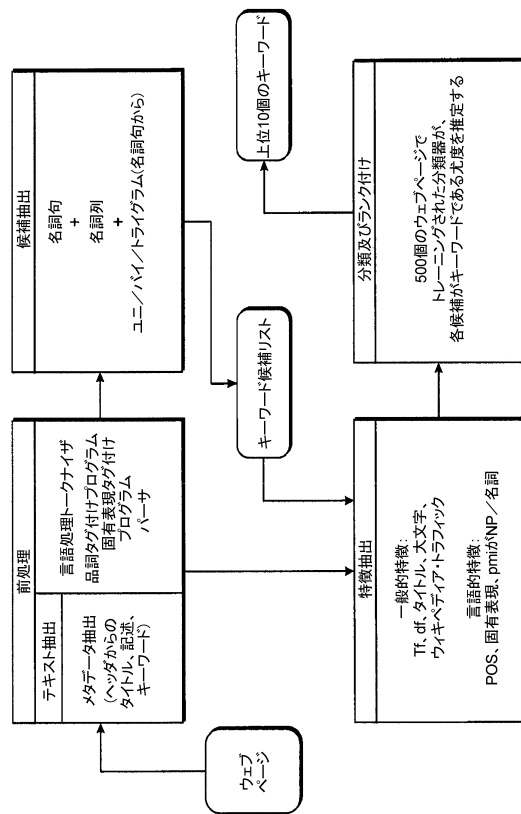
[16] Wu, Xiaoyuan及びBolivar, Alvaro. Keyword extraction for contextual advertisement. WWW '08: Proceeding of the 17th international conference on World Wide Web、ページ1195~1196、New York, NY, USA、2008。ACM。

[17] Yih, Wen-tau及びGoodman, Joshua及びCarvalho, Vitor R.. Finding advertising keywords on web pages. WWW '06: Proceedings of the 15th international conference on World Wide Web、ページ213~222、New York, NY, USA、2006。ACM。

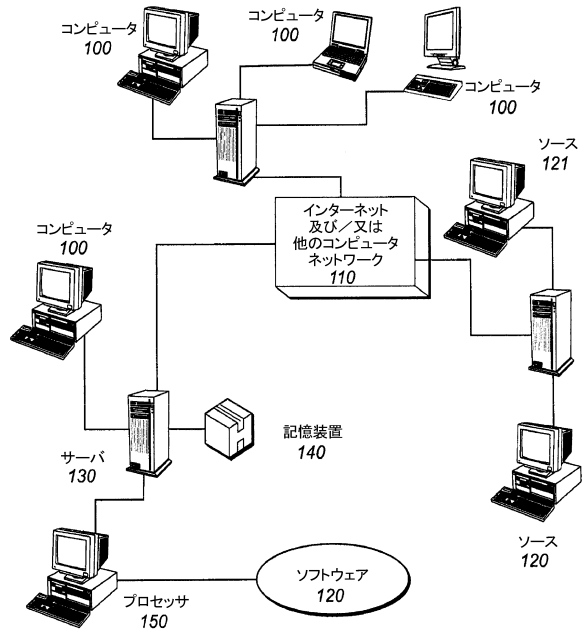
30

[18] OAK System, <http://nlp.cs.nyu.edu/oak/>。

【図 1】



【図 2】



フロントページの続き

(72)発明者 関根 聡

アメリカ合衆国， ニューヨーク州， スカースデール， ネルソン ロード 232

審査官 齊藤 貴孝

(56)参考文献 米国特許出願公開第2007/0112764 (US, A1)

特開2004-139553 (JP, A)

特開2008-065417 (JP, A)

特開2007-257390 (JP, A)

特開2005-085285 (JP, A)

特開平07-065018 (JP, A)

特開2009-271794 (JP, A)

特開2006-146705 (JP, A)

特開2010-204866 (JP, A)

上村 卓史、外2名，ウェブ閲覧における効率的なキーワード抽出とその利用，データベースとWeb情報システムに関するシンポジウム 情報処理学会シンポジウムシリーズ，日本，社団法人情報処理学会，2007年12月20日，第2007巻，第3号，p. 1 - 9

斎藤 一、外3名，キーワード地図構造モデリングによるグループ学習支援システムの構築，電子情報通信学会技術研究報告，日本，社団法人電子情報通信学会，2003年 7月17日，第103巻，第217号，p. 1 - 4

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

G06F 17/27