

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-219764
(P2019-219764A)

(43) 公開日 令和1年12月26日(2019.12.26)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 16/00 (2019.01)	G06F 17/30	320C
G06F 21/62 (2013.01)	G06F 17/30	340Z
	G06F 21/62	345

審査請求 未請求 請求項の数 16 O L (全 24 頁)

(21) 出願番号 特願2018-114944 (P2018-114944)
(22) 出願日 平成30年6月15日 (2018.6.15)

(出願人による申告)平成24年度、国立研究開発法人科学技術振興機構、戦略的創造研究推進事業・総括実施型研究(ERATO)「河原林巨大グラフプロジェクト」に係る委託業務、産業技術力強化法第19条の適用を受ける特許出願

(71) 出願人 504202472
大学共同利用機関法人情報・システム研究機構
東京都立川市緑町10番3号
(74) 代理人 100205084
弁理士 吉浦 洋一
(72) 発明者 河原林 健一
東京都千代田区一ツ橋二丁目1番2号 大学共同利用機関法人情報・システム研究機構 国立情報学研究所内
(72) 発明者 町出 智也
東京都千代田区一ツ橋二丁目1番2号 大学共同利用機関法人情報・システム研究機構 国立情報学研究所内

最終頁に続く

(54) 【発明の名称】 情報検索システム

(57) 【要約】

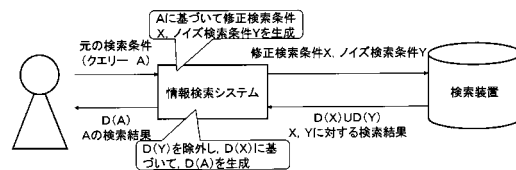
【課題】

情報を検索するための情報検索システムを提供することを目的とする。

【解決手段】

ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部と、第1の検索単語の意味解析を用いて、第1の検索単語を修正する第3の検索単語を特定する検索条件修正処理部と、第2の検索単語と第3の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部と、を有する情報検索システムである。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

情報を検索するための情報検索システムであって、
前記情報検索システムは、
ユーザの実際の検索条件である第 1 の検索単語の意味解析に基づいて、ノイズとなる第 2 の検索単語を特定するノイズ処理部と、
前記第 1 の検索単語の意味解析に基づいて、前記第 1 の検索単語を修正する第 3 の検索単語を特定する検索条件修正処理部と、
前記第 2 の検索単語と前記第 3 の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部と、
を有することを特徴とする情報検索システム。

10

【請求項 2】

情報を検索するための情報検索システムであって、
前記情報検索システムは、
ユーザの実際の検索条件である第 1 の検索単語の意味解析に基づいて、ノイズとなる第 2 の検索単語を特定するノイズ処理部と、
前記第 1 の検索単語と前記第 2 の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部と、
を有することを特徴とする情報検索システム。

【請求項 3】

前記ノイズ処理部は、
前記第 1 の検索単語と同じクラスに属する単語を用いて、クラスタリング耐性のある前記第 2 の検索単語を特定する、
ことを特徴とする請求項 1 または請求項 2 に記載の情報検索システム。

20

【請求項 4】

前記ノイズ処理部は、
前記第 1 の検索単語に基づいて、少なくとも二以上の手法により、クラスタリング耐性のある前記第 2 の検索単語を特定し、
各手法による前記第 2 の検索単語の数または割合が変動する、
ことを特徴とする請求項 1 から請求項 3 のいずれかに記載の情報検索システム。

30

【請求項 5】

前記ノイズ処理部は、
前記第 1 の検索単語と同じクラスに属する単語から複数の単語を特定することで単語群を構成し、
前記構成した単語群に対して、高密度クラスタから前記第 2 の検索単語を特定するクラスタ手法、前記単語群を分割することで前記第 2 の検索単語を特定する分割手法、前記単語群を構成する単語からランダムに前記第 2 の検索単語を特定するランダム手法、のいずれか一以上の手法を用いることで、ノイズ単語を特定する、
ことを特徴とする請求項 1 から請求項 4 のいずれかに記載の情報検索システム。

【請求項 6】

前記ノイズ処理部は、
前記クラスタ手法として、前記構成した単語群を用いて、前記第 1 の検索単語とは異なるクラスタを構成する複数の単語を特定することで、前記第 2 の検索単語を特定する、
ことを特徴とする請求項 5 に記載の情報検索システム。

40

【請求項 7】

前記ノイズ処理部は、
前記クラスタ手法として、前記構成した単語群を用いて、頻出頻度に基づく単語群を構成し、
前記頻出頻度に基づく単語群において、前記第 1 の検索単語からの距離と類似性に基づいて特定した単語を用いてクラスタを生成することで、前記第 2 の検索単語を特定する、

50

ことを特徴とする請求項 5 または請求項 6 に記載の情報検索システム。

【請求項 8】

前記ノイズ処理部は、

前記分割手法として、前記構成した単語群を用いて、前記第 1 の検索単語とは非類似であり、かつ類似する単語同士を、前記第 2 の検索単語として特定する、

ことを特徴とする請求項 5 から請求項 7 のいずれかに記載の情報検索システム。

【請求項 9】

前記ノイズ処理部は、

前記分割手法として、前記構成した単語群を複数に分割し、分割した単語群における単語と前記第 1 の検索単語との類似性を用いて、前記第 2 の検索単語を特定する、

ことを特徴とする請求項 5 から請求項 8 のいずれかに記載の情報検索システム。

【請求項 10】

前記検索条件修正処理部は、

ベクトル化した前記第 1 の検索単語とノイズベクトルとを用いて演算することで、前記第 3 の検索単語を特定する、

ことを特徴とする請求項 1、請求項 3 から請求項 9 のいずれかに記載の情報検索システム。

【請求項 11】

前記情報検索システムは、

前記検索装置から受け付けた前記第 2 の検索単語に対応する検索結果を除外し、前記検索装置から受け付けた前記第 1 の検索単語または前記第 2 の検索単語に対応する検索結果に基づいて、前記第 1 の検索単語に対する検索結果を出力する検索結果処理部、

を有することを特徴とする請求項 1 から請求項 10 のいずれかに記載の情報検索システム。

【請求項 12】

前記情報検索システムは、

前記第 2 の検索単語と前記第 3 の検索単語とを出力することで、前記第 1 の検索単語を推測させる処理部、

を有することを特徴とする請求項 1 から請求項 11 のいずれかに記載の情報検索システム。

【請求項 13】

情報を検索するための情報検索システムであって、

前記情報検索システムは、

ユーザの実際の検索条件であるオリジナル検索条件をベクトル化し、ベクトル化した前記オリジナル検索条件を用いて修正検索条件を特定する検索条件修正処理部と、

前記オリジナル検索条件に基づいて、ノイズとなるノイズ検索条件を特定するノイズ処理部と、

前記修正検索条件と前記ノイズ検索条件とを検索装置に送り、検索結果を受け付ける検索装置処理部、

を有することを特徴とする情報検索システム。

【請求項 14】

コンピュータを、

ユーザの実際の検索条件である第 1 の検索単語の意味解析に基づいて、ノイズとなる第 2 の検索単語を特定するノイズ処理部、

前記第 1 の検索単語の意味解析に基づいて、第 3 の検索単語を特定する検索条件修正処理部、

前記第 2 の検索単語と前記第 3 の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部、

として機能させることを特徴とする情報検索プログラム。

【請求項 15】

10

20

30

40

50

コンピュータを、
 ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部、
 前記第1の検索単語と前記第2の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部、
 として機能させることを特徴とする情報検索プログラム。

【請求項16】

コンピュータを、
 ユーザの実際の検索条件であるオリジナル検索条件をベクトル化し、ベクトル化した前記オリジナル検索条件を用いて修正検索条件を特定する検索条件修正処理部、
 前記オリジナル検索条件に基づいて、ノイズとなるノイズ検索条件を特定するノイズ処理部
 前記修正検索条件と前記ノイズ検索条件とを検索装置に送り、検索結果を受け付ける検索装置処理部、
 として機能させることを特徴とする情報検索プログラム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報を検索するための情報検索システムに関する。とくに、検索者が入力をしたキーワードなどの検索条件を、検索エンジンなどの検索装置に知られずに検索を行うことができる情報検索システムに関する。

20

【背景技術】

【0002】

インターネットやデータベースから、所望の情報を得るために、検索エンジンなどの検索装置が用いられている。とくにインターネットでの検索エンジンは、無数にあるウェブサイトから検索条件にヒットするウェブサイトを特定するために有益である。

【0003】

検索装置は有益な面があるものの、検索装置に入力された検索条件を蓄積して解析をすることで、当該検索者の関心や興味の対象、思想などの一定の傾向を把握することが可能となる。そのため検索装置に入力する検索条件を、極力、検索装置に把握されることを回避したい要望がある。しかし、検索装置は、検索条件に基づいて情報の検索を行うので、検索装置に適切な検索条件を入力しないと、所望の情報が記載された検索結果が得られないこととなる。

30

【0004】

そこで、検索装置に、ユーザの実際の検索条件を把握されにくくする一方、検索装置からは所望の検索結果を得ることができるためのシステムが検討されており、たとえば特許文献1、特許文献2がある。

【先行技術文献】

【特許文献】

【0005】

40

【特許文献1】特開平11-259512号公報

【特許文献2】特許第5306356号

【発明の概要】

【発明が解決しようとする課題】

【0006】

特許文献1のシステムは、入力された検索言語を、類似語、上位概念語に変換をすることで、データ検索サーバには、直接、入力された検索言語が把握されないようにするシステムである。

【0007】

特許文献2のシステムは、検索語を文字単位で分解することで抽出した検索語における

50

正当文字列に、ノイズ語のノイズ文字列を組み合わせてクエリを生成することで、検索エンジンに、検索語が把握されないようにするシステムである。

【0008】

特許文献1のシステムでは類似語や上位概念語に変換がされているので、クラスタリングをすることで、検索言語を類推することができてしまう課題がある。また、特許文献2のシステムでは、単語単位で文字列を分解していても、特許文献1と同様に、クラスタリングをすることで、元の単語を推定できてしまう課題がある。

【課題を解決するための手段】

【0009】

そこで本発明者らは、上記課題に鑑み、実際の検索条件の特定が困難である情報検索システムを発明した。特に、クライスタリング耐性の高いノイズを用いる場合には、クラスタリング耐性を高めることができる。

10

【0010】

第1の発明は、情報を検索するための情報検索システムであって、前記情報検索システムは、ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部と、前記第1の検索単語の意味解析に基づいて、前記第1の検索単語を修正する第3の検索単語を特定する検索条件修正処理部と、前記第2の検索単語と前記第3の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部と、を有する情報検索システムである。

【0011】

20

第2の発明は、情報を検索するための情報検索システムであって、前記情報検索システムは、ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部と、前記第1の検索単語と前記第2の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部と、を有する情報検索システムである。

【0012】

第1の発明、第2の発明を用いることで、実際の検索条件である第1の検索単語に対して意味解析をして特定したノイズとなる第2の検索単語を検索装置に送ることができる。これによって、実際の検索条件の特定を困難とすることができる。

【0013】

30

上述の発明において、前記ノイズ処理部は、前記第1の検索単語と同じクラスに属する単語を用いて、クラスタリング耐性のある前記第2の検索単語を特定する、情報検索システムのように構成することができる。

【0014】

ノイズとする第2の検索単語について、クラスタリング耐性となる単語を用いることで、検索装置側においてクラスタ解析を行ったとしても、実際の検索条件である第1の検索単語を特定することが困難となる。

【0015】

上述の発明において、前記ノイズ処理部は、前記第1の検索単語に基づいて、少なくとも二以上の手法により、クラスタリング耐性のある前記第2の検索単語を特定し、各手法による前記第2の検索単語の数または割合が変動する、情報検索システムのように構成することができる。

40

【0016】

複数の手法を用いて第2の検索単語を特定することで、検索装置側において、実際の検索条件である第1の検索単語を特定することがさらに困難となる。

【0017】

上述の発明において、前記ノイズ処理部は、前記第1の検索単語と同じクラスに属する単語から複数の単語を特定することで単語群を構成し、前記構成した単語群に対して、高密度クラスタから前記第2の検索単語を特定するクラスタ手法、前記単語群を分割することで前記第2の検索単語を特定する分割手法、前記単語群を構成する単語からランダムに

50

前記第2の検索単語を特定するランダム手法、のいずれか一以上の手法を用いることで、ノイズ単語を特定する、情報検索システムのように構成することができる。

【0018】

ノイズとする第2の検索単語を特定するためには、本発明のような方法を一または複数用いるとよい。

【0019】

上述の発明において、前記ノイズ処理部は、前記クラスタ手法として、前記構成した単語群を用いて、前記第1の検索単語とは異なるクラスタを構成する複数の単語を特定することで、前記第2の検索単語を特定する、情報検索システムのように構成することができる。

10

【0020】

上述の発明において、前記ノイズ処理部は、前記クラスタ手法として、前記構成した単語群を用いて、頻出頻度に基づく単語群を構成し、前記頻出頻度に基づく単語群において、前記第1の検索単語からの距離と類似性に基づいて特定した単語を用いてクラスタを生成することで、前記第2の検索単語を特定する、情報検索システムのように構成することができる。

【0021】

これらの発明の処理を実行することで、ノイズとする第2の検索単語について、第1の検索単語とは異なるクラスタに属する単語により構成することができる。そのため、検索装置側においてクラスタ解析を行ったとしても、実際の検索条件である第1の検索単語を特定することが困難となる。

20

【0022】

上述の発明において、前記ノイズ処理部は、前記分割手法として、前記構成した単語群を用いて、前記第1の検索単語とは非類似であり、かつ類似する単語同士を、前記第2の検索単語として特定する、情報検索システムのように構成することができる。

【0023】

上述の発明において、前記ノイズ処理部は、前記分割手法として、前記構成した単語群を複数に分割し、分割した単語群における単語と前記第1の検索単語との類似性を用いて、前記第2の検索単語を特定する、情報検索システムのように構成することができる。

【0024】

これらの発明の処理を実行することで、ノイズとする第2の検索単語について、第1の検索単語とは類似していない単語により構成することができる。そのため、検索装置側においてクラスタ解析を行ったとしても、実際の検索条件である第1の検索単語を特定することが困難となる。

30

【0025】

上述の発明において、前記検索条件修正処理部は、ベクトル化した前記第1の検索単語とノイズベクトルとを用いて演算することで、前記第3の検索単語を特定する、情報検索システムのように構成することができる。

【0026】

本発明のように構成することで、第1の検索単語そのものではないが、意味が近い単語を特定し、第3の検索単語を特定することができる。

40

【0027】

上述の発明において、前記情報検索システムは、前記検索装置から受け付けた前記第2の検索単語に対応する検索結果を除外し、前記検索装置から受け付けた前記第1の検索単語または前記第2の検索単語に対応する検索結果に基づいて、前記第1の検索単語に対する検索結果を出力する検索結果処理部、情報検索システムのように構成することができる。

【0028】

第2の検索単語はノイズであるので、その検索結果は不要である。したがって、第2の検索単語による検索結果を除外して、最終的な検索結果を出力すればよい。

50

【 0 0 2 9 】

上述の発明において、前記情報検索システムは、前記第2の検索単語と前記第3の検索単語とを出力することで、前記第1の検索単語を推測させる処理部、を有する情報検索システムのように構成することができる。

【 0 0 3 0 】

本発明の情報検索システムの効果は、そのまま認識しにくい。そこで、本発明のように構成することで、本発明の効果を認識させることができる。

【 0 0 3 1 】

第13の発明は、情報を検索するための情報検索システムであって、前記情報検索システムは、ユーザの実際の検索条件であるオリジナル検索条件をベクトル化し、ベクトル化した前記オリジナル検索条件を用いて修正検索条件を特定する検索条件修正処理部と、前記オリジナル検索条件に基づいて、ノイズとなるノイズ検索条件を特定するノイズ処理部と、前記修正検索条件と前記ノイズ検索条件とを検索装置に送り、検索結果を受け付ける検索装置処理部と、を有する情報検索システムである。

10

【 0 0 3 2 】

本発明を用いることで、検索装置において、実際の検索条件である第1の検索条件の特定を困難とすることができる。この場合、検索条件としてはベクトル表現できる情報であればよく、単語に限らず、画像情報、音情報であっても同様に実現することができる。

【 0 0 3 3 】

第1の発明は、本発明のプログラムをコンピュータに読み込ませて実行することで実現できる。すなわち、コンピュータを、ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部、前記第1の検索単語の意味解析に基づいて、第3の検索単語を特定する検索条件修正処理部、前記第2の検索単語と前記第3の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部、として機能させる情報検索プログラムのように構成することができる。

20

【 0 0 3 4 】

第2の発明は、本発明のプログラムをコンピュータに読み込ませて実行することで実現できる。すなわち、コンピュータを、ユーザの実際の検索条件である第1の検索単語の意味解析に基づいて、ノイズとなる第2の検索単語を特定するノイズ処理部、前記第1の検索単語と前記第2の検索単語とを検索装置に送り、検索結果を受け付ける検索装置処理部、として機能させる情報検索プログラムのように構成することができる。

30

【 0 0 3 5 】

第13の発明は、本発明のプログラムをコンピュータに読み込ませて実行することで実現できる。すなわち、コンピュータを、ユーザの実際の検索条件であるオリジナル検索条件をベクトル化し、ベクトル化した前記オリジナル検索条件を用いて修正検索条件を特定する検索条件修正処理部、前記オリジナル検索条件に基づいて、ノイズとなるノイズ検索条件を特定するノイズ処理部、前記修正検索条件と前記ノイズ検索条件とを検索装置に送り、検索結果を受け付ける検索装置処理部、として機能させる情報検索プログラムのように構成することができる。

40

【 発明の効果 】

【 0 0 3 6 】

本発明の情報検索システムを用いることによって、実際の検索条件の特定を困難とする情報検索システムを発明した。特に、クラスタリング耐性の高いノイズを用いる場合には、クラスタリング耐性を高めることができる。

【 図面の簡単な説明 】

【 0 0 3 7 】

【 図 1 】 本発明の情報検索システムの概念の一例を示す図である。

【 図 2 】 本発明の情報検索システムのシステム構成を示すブロック図の一例である。

【 図 3 】 本発明の情報検索システムを実現するコンピュータのハードウェア構成の一例を示す図である。

50

- 【図 4】本発明の情報検索システムの処理プロセスの一例を示すフローチャートである。
- 【図 5】検索条件修正処理の処理プロセスの一例を示すフローチャートである。
- 【図 6】ノイズ処理の全体の処理プロセスの一例を示すフローチャートである。
- 【図 7】クラスタ手法の処理プロセスの一例を示すフローチャートである。
- 【図 8】分割手法の処理プロセスの一例を示すフローチャートである。
- 【図 9】ランダム手法の処理プロセスの一例を示すフローチャートである。
- 【図 10】修正検索条件を特定する処理の一例を模式的に示す図である。
- 【図 11】クラスタ手法の処理の一例を模式的に示す図である。
- 【図 12】分割手法の処理の一例を模式的に示す図である。
- 【図 13】単語群 X のみを用いた場合の匿名性と再構築可能性の精度の関係を示す図である。 10
- 【図 14】単語群 X と単語群 Y とを用いた場合の匿名性と再構築可能性の精度の関係を示す図である。
- 【図 15】情報検索システムによる検索結果と、実際の検索条件「kyoto」を入力した場合の比較例を示す図である。
- 【図 16】情報検索システムによる検索結果と、実際の検索条件「kyoto」を入力した場合の比較例を示すほかの図である。
- 【図 17】情報検索システムによる検索結果と、実際の検索条件「kyoto」を入力した場合の比較例を示すほかの図である。
- 【図 18】情報検索システムによる検索結果と、実際の検索条件「kyoto」を入力した場合の比較例を示すほかの図である。 20
- 【図 19】情報検索システムによる検索結果と、実際の検索条件「nagasaki」を入力した場合の比較例を示す図である。
- 【図 20】情報検索システムによる検索結果と、実際の検索条件「nagasaki」を入力した場合の比較例を示す図である。
- 【図 21】実施例 2 において、単語群 X の単語、単語群 Y の単語をそれぞれ表示した状態の画面を示す図である。
- 【図 22】実施例 2 において、単語群 Y の単語を削除し、単語群 X の単語のみを表示した状態の画面を示す図である。
- 【図 23】実施例 2 において、正解を表示した状態を示す画面である。 30
- 【図 24】単語 A の検索結果 $D(A)$ と単語群 X の検索結果 $D(X_i)$ との関係を示す図である。
- 【発明を実施するための形態】
- 【0038】
- 本発明の情報検索システム 1 の全体の概念の一例を図 1 に示す。また、本発明の情報検索システム 1 のシステム構成のブロック図の一例を図 2 に示す。情報検索システム 1 では、情報の検索を行うユーザが利用するユーザ端末 4 と、情報の検索を行う検索サーバなどの検索装置 3 と、情報検索システム 1 の各処理を実行するための制御端末 2 とを用いる。情報検索システム 1 における制御端末 2 は、コンピュータによって実現される。コンピュータのハードウェア構成の一例を図 3 に示す。なお、制御端末 2 とユーザ端末 4、制御端末 2 と検索装置 3、制御端末 2 とユーザ端末 4 と検索装置 3 とが一体的に構成されていてもよい。 40
- 【0039】
- コンピュータはプログラムの演算処理を実行する CPU などの演算装置 70 と、情報を記憶する RAM やハードディスクなどの記憶装置 71 と、ディスプレイなどの表示装置 72 と、情報の入力を行う入力装置 73 と、演算装置 70 の処理結果や記憶装置 71 に記憶する情報などの各種情報を通信する通信装置 74 とを有している。なお、コンピュータがタッチパネルディスプレイを備えている場合には表示装置 72 と入力装置 73 とが一体的に構成されていてもよい。タッチパネルディスプレイは、携帯電話やスマートフォン、タブレット型コンピュータなどの可搬型通信端末などで利用されることが多いが、それに限 50

定するものではない。

【0040】

タッチパネルディスプレイは、そのディスプレイ上で、直接、所定の入力デバイス（タッチパネル用のペンなど）や指などによって入力を行える点で、表示装置72と入力装置73の機能が一体化した装置である。

【0041】

情報検索システム1の制御端末2は一台のコンピュータによって実現されていてもよいが、その機能が複数のコンピュータによって実現されていてもよい。この場合のコンピュータとして、たとえばクラウドサーバであってもよい。

【0042】

さらに、本発明の情報検索システム1における各処理部は、その機能が論理的に区別されているのみであって、物理上あるいは事実上は同一の領域を為していてもよい。

【0043】

検索装置3は、インターネットの情報を検索するための検索エンジンサーバや、各種の情報を記憶するデータベースサーバなど、情報を検索するための装置である。なお、検索装置3としては、検索エンジンサーバやデータベースサーバに限定するものではなく、情報を検索するための装置であればよい。

【0044】

情報検索システム1における制御端末2は、検索条件受付処理部21と検索条件修正処理部22とノイズ処理部23と検索装置処理部24と検索結果処理部25とを有する。

【0045】

検索条件受付処理部21は、ユーザ端末4から、ユーザが実際に検索をしたい検索条件を含むクエリの入力を受け付ける。検索条件としては、ベクトル表現できる情報であればいかなる情報であってもよい。本明細書では、検索条件として、キーワードなどの単語の場合を説明するが、画像情報、音情報などでも同様の処理を実行することで実現できる。たとえば単語が画像情報、音情報になっている場合には、画像情報をOCR認識してテキスト化した後に処理を実行し、音情報を音声認識技術に基づいてテキスト化した後に処理を実行してもよい。また、画像情報における各画素の色情報に基づいてベクトル化してもよいし、音情報における周波数情報に基づいてベクトル化し、以降の処理を実行してもよい。

【0046】

検索条件修正処理部22は、検索条件受付処理部21で受け付けた実際の検索条件に基づいて、検索装置3に入力するための修正した検索条件（修正検索条件）を特定する処理である。修正検索条件を特定する処理にはさまざまな方法を用いることができる。たとえば実際の検索条件が単語である場合、その単語そのものではないが、その単語に近い単語を修正検索条件として特定する。この場合、検索条件修正処理部22は、実際の検索条件の単語の意味解析に基づいて、修正検索条件となる単語を特定する。なお、意味解析とは、単語エンベディング（Word embedding）であって、自然言語解析における技術である。すなわち、ある単語とほかの単語の意味や概念などが類似しているか否かなど、単語の意味関係を自動的に解析するための技術である。

【0047】

検索条件受付処理部21で受け付けた実際の検索条件に基づいてコサイン類似度を用いて、修正検索条件を特定することができる。すなわち、検索条件修正処理部22は、検索条件受付処理部21で入力を受け付けた実際の検索条件をベクトル化し、それにノイズベクトルを演算、たとえば加算する。そして、演算したベクトルの点からコサイン類似度に基づき近傍検索（コサイン類似度が一定の範囲内にあるか）をすることで、修正検索条件を特定する。修正検索条件を特定する処理の一例を模式的に示すのが図10である。

【0048】

たとえば、以下のような処理を実行すればよい。検索条件受付処理部21で受け付けたクエリにおける実際の検索条件における単語をAとした場合、検索条件修正処理部22は

10

20

30

40

50

、単語 A に対応するベクトル v ($v = v_1, v_2, \dots, v_{300}$) を取得する。単語 A に対応するベクトル v の取得方法としては、GloVe (Global Vector for Word Representation), word2vec, fasttext などを用いる方法があるが、それに限定するものではない。また、本明細書では、単語のベクトル化に GloVe を用いるので、300次元のベクトルで説明するが、それに限定するものではない。

【0049】

そして検索条件修正処理部 22 では、単語 A に対応するベクトル v に対して、同次元のノイズベクトル n ($n = n_1, n_2, \dots, n_{300}$) を取得する。なお、ノイズベクトルにおける各実数 n_i は、たとえばガウス分布のノイズを用いることができるが、それに限定するものではない。

10

【0050】

以上のようにして検索条件修正処理部 22 で単語 A に対応するベクトル v とノイズベクトル n とを取得すると、それぞれを演算、たとえば加算することでベクトル v' を算出する。そして、ベクトル v' の点から距離、たとえばコサイン類似度 (コサイン距離) が近い (ベクトル v' の点からコサイン類似度が一定範囲内にある)、任意の m 個のベクトル x (x_1, x_2, \dots, x_m) を、上述の GloVe のデータセット (単語 A をベクトル化した際のモデルのデータセット) から特定する。そして、特定したベクトル x (x_1, x_2, \dots, x_m) に対応する単語 X_1, X_2, \dots, X_m の単語群 X を特定することで、検索条件である単語 A (第 1 の検索単語) に対応する修正検索条件である単語 (第 3 の検索単語) の単語群 X を特定することができる。単語群 X を構成する単語の数 m は、任意の数でよく、複数、たとえば 10 個から 20 個程度とすることができるが、それに限定するものではない。

20

【0051】

なお、検索条件修正処理部 22 は、上述の処理のほか、たとえば、単語とそれに類似、関連する単語、上位概念の単語をあらかじめ対応づけて記憶しておき、その対応関係に基づいて、検索条件受付処理部 21 で受け付けた検索条件における単語 A に対応する単語の単語群 X を特定してもよい。

【0052】

ノイズ処理部 23 は、検索条件受付処理部 21 で入力を受け付けた検索条件に対するノイズとなる検索条件を特定する処理である。ノイズ処理としては、無関係の検索条件を付加するほか、クラスタリングが困難となるノイズを付加することが好ましい。たとえば実際の検索条件が単語である場合、実際の検索条件の単語の意味解析に基づいて、ノイズとする単語を特定することができる。

30

【0053】

ノイズ処理部 23 におけるノイズ処理としては、高密度クラスタからノイズを選択するクラスタ手法、分割手法、ランダム手法などがあり、これらの手法のいずれか一以上によって得られた検索条件を特定するとよい。また、上記の 3 手法に限定するものではなく、上記の 3 手法以外、あるいは上記の 3 手法と組み合わせ、ほかの手法を用いることも可能である。

40

【0054】

上述と同様に、検索条件受付処理部 21 で受け付けたクエリにおける実際の検索条件における単語を A とした場合、ノイズ処理部 23 は、単語 A と同じクラスに属する単語から、ノイズの候補となる所定数、たとえば 1000 個の単語 w_1, \dots, w_{1000} をランダムに特定し、単語群 W を構成する。なお、単語 A と同じクラスに属する単語 w は、たとえば、ウィキペディアの Ontology クラスのデータセットを用いることなどで特定することができるが、それに限定するものではない。

【0055】

そして、ノイズ処理部 23 は、特定した単語群 W に対して、クラスタ手法、分割手法、ランダム手法のいずれかまたは複数の手法による処理を実行し、ノイズとなる単語 $Y_1,$

50

Y_2, \dots, Y_z による単語群 Y を取得する。単語群 Y を構成する単語の数 z は任意の数とすることができ、好ましくは複数、たとえば 10 個 ~ 20 個程度とすることができるが、それに限定するものではない。ノイズとなる単語群 Y は、クラスタリングに対する耐性が高い単語により構成されることが好ましい。

【0056】

ノイズ処理部 23 は、複数の手法を用いる場合、ノイズとなる単語群 Y の各単語について、各手法により得られる単語の単語数を任意の割合または数として設定することができる。たとえば単語群 Y の単語数を 10 個とする場合、クラスタ手法による単語を 5 個 (50%)、分割手法による単語を 3 個 (30%)、ランダム手法による単語を 2 個 (20%) のように設定することができる。各手法による単語の割合や数は、毎回、変更してもよいし、固定でもよい。

10

【0057】

ノイズ処理部 23 におけるクラスタ手法は、検索条件となる単語 A と同じクラスに属する単語群 W の単語 w に基づいて、単語 A とは異なるクラスタを構成可能な複数の単語を特定してノイズとする単語群 Y を構成することで、クラスタリングに対する耐性を高める。ノイズ処理部 23 におけるクラスタ手法は、以下のように実行する。クラスタ手法によるノイズとなる検索条件を特定する処理の一例を模式的に示すのが図 11 である。

【0058】

まず、単語群 W の単語 w のなかから、単語 A と頻出頻度が近い単語を特定することで、クラスタの中心の候補となる単語の単語群 $S (s_1, s_2, \dots, s_i)$ を構成する。そして、特定した単語群 S の単語 s のなかから、使用する単語群 $S' (s'_1, s'_2, \dots, s'_j)$ (ただし $j < i$) を特定する。この特定の際には、単語 A のベクトルから適度に離れており、意味が類似している単語 s' を優先して特定をすることが好ましい。すなわち、単語群 S' における単語 s' の特定は、単語 A のベクトルに対して、同次元のノイズベクトル (単語 A から適度に離れる値として設定するベクトル) を加算等の演算をして算出し、その加算したベクトルの点から、コサイン類似度 (コサイン距離) が一定の閾値以上である単語群 S における単語 s を、単語 s' として特定する。

20

【0059】

特定した単語群 S' において、それぞれの単語 s'_1, s'_2, \dots, s'_j に近い単語を所定数特定し、単語のクラスタ C を生成する。たとえば各クラスタの単語数は 3 ~ 10 個とするが、それに限定するものではない。そして、密集度の高いクラスタ C から順番に、クラスタにおける単語を特定し、あらかじめ設定した数になったら、それらをノイズの単語群 $Y (Y_1, Y_2, \dots, Y_z)$ として特定をする。

30

【0060】

このような処理を実行することで、単語 A とは相違する密集度の高いクラスタ C における単語を、ノイズの単語として特定できるので、仮に検索装置 3 側でクラスタリングをしたとしても、実際の単語 A の特定が困難となり、クラスタリングに対する耐性が高くなる。

【0061】

また、ノイズ処理部 23 における分割手法は、検索条件となる単語 A と同じクラスに属する多数の単語を分割、たとえば 2 分割 (ただし単語数は同数ではない) し、その分割によって構成される単語群における単語を用いて、検索条件となる単語 A とは似ていない単語同士の単語群を生成することを、所定条件を充足するまで繰り返し、条件充足後の単語群から、単語 A と似ている複数の単語を特定してノイズとする単語群 Y を構成することで、クラスタリングに対する耐性を高める。ノイズ処理部 23 における分割手法は、以下のように実行する。分割手法によるノイズとなる検索条件を特定する処理の一例を模式的に示すのが図 12 である。

40

【0062】

まず実際の検索条件である単語 A の点を取る超平面をランダムに特定し、その超平面において、単語群 $W (w_1, \dots, w_{1000})$ を、単語群 W_1 、単語群 W_2 の 2 つに分

50

割をする。ただし，単語群 W_1 の単語数は，単語群 W_2 の単語数より多いとする。そして，単語群 W_1 において，単語 A と類似していない単語を，単語群 W_1 から所定割合または所定数，たとえば 10% だけ消去し，消去した単語を新たに単語群 W として構成する。なお，単語 A と類似している単語か否かは，単語 A と，単語群 W_1 における比較対象となる単語とのコサイン類似度（コサイン距離）が一定の閾値以上であるかで特定可能である。

【0063】

以上の処理を所定条件，たとえば $|W| < 2y$ (y は任意の値) となるまで繰り返す。

【0064】

上記の所定条件を充足した場合，その単語群 W の単語のなかから，単語 A に類似している単語を z 個特定し，それらを単語群 $Y (Y_1, Y_2, \dots, Y_z)$ として特定をする。

10

【0065】

以上のような分割手法を用いることで，実際の検索条件における単語 A (第1の検索単語) とは非類似であって，かつ，また単語群 Y を構成する単語自体は意味が近い単語 (第2の検索単語) をノイズの単語として特定できるので，単語 A とは異なるクラスタとなりやすい単語をノイズの単語とすることができる。そのため，仮に検索装置 3 側でクラスタリングをしたとしても，実際の単語 A の特定が困難となり，クラスタリングに対する耐性が高くなる。

【0066】

さらに，ノイズ処理部 23 におけるランダム手法は，単語群 $W (w_1, \dots, w_{100})$ の中から，ランダムに z 個の単語を特定し，それらを単語群 $Y (Y_1, Y_2, \dots, Y_z)$ とする。

20

【0067】

ランダムにノイズとなる単語を特定することで，クラスタリングに対する耐性を高めることができる。

【0068】

以上のような処理をノイズ処理部 23 が実行することで，ノイズとする検索条件 Y を特定できる。とくに，一つの手法のみならず，複数の手法を組み合わせることで，クラスタリングに対する耐性は，一層，高くすることができる。

【0069】

検索装置処理部 24 は，検索条件修正処理部 22 で特定した単語群 X における単語と，ノイズ処理部 23 で特定した単語群 Y における単語とを，それぞれ検索装置 3 に送ることで検索処理を実行させる。この際には，単語群 X における単語，単語群 Y における単語をランダムな順番で検索装置 3 に送るとよい。そして，各単語に対する検索結果を受け付ける。なお，少なくとも，検索装置 3 に送った単語群 X における単語と，その検索結果とを対応づけて記憶しておく。

30

【0070】

検索結果処理部 25 は，検索装置 3 から受け付けた検索結果に基づいて，ユーザ端末 4 に送る検索結果を出力する。検索結果処理部 25 は，検索装置 3 から受け付けた単語群 Y の単語 Y_1, Y_2, \dots, Y_z に対する検索結果 $D(Y_i) (1 \leq i \leq z)$ をユーザ端末 4 に送る検索結果から除外し，単語群 X における単語 X_1, X_2, \dots, X_m に対する検索装置 3 での各検索結果 $D(X_i) (1 \leq i \leq m)$ に基づいて検索結果を生成する。たとえば各検索結果 $D(X_1), D(X_2), \dots, D(X_m)$ をソートすることで，検索結果を生成する。また検索結果を生成する際に，検索結果におけるページランクを用いてもよいし，ページランクの重み付けなどを用いてソートをしてよい。さらに検索結果処理部 25 は，検索装置 3 から受け付けた検索結果の群 $D(X_i)$ に対して，単語 A に基づいて検索を行うことで，検索結果を生成してもよい。なお，検索結果の生成は，公知の方法を用いることができる。

40

【0071】

検索結果 $D(X_i)$ は，単語群 X の単語 X_1, X_2, \dots, X_m に対する検索装置 3

50

での検索結果である。そして単語群 X における単語 X_1, X_2, \dots, X_m は、単語 A に対応する単語ベクトルに近いものを特定している。すなわち、単語 A と単語群 X の単語とは共起性が高い（同一の文に同時に現れやすい）。そのため、単語 A の検索結果 $D(A)$ は、共起性の高い単語群 X に対する検索結果 $D(X_i)$ に基づいて生成することができる。この関係を模式的に示すのが図 24 である。

【0072】

検索結果処理部 25 は、以上のように生成した検索結果をユーザ端末 4 に送る。

【0073】

以上のような処理を実行することで、ユーザが入力をした検索条件は検索装置 3 側に知られることなく、精度のよい検索結果を得ることができる。

【実施例 1】

【0074】

つぎに本発明の情報検索システム 1 を用いて情報の検索を行う場合の処理プロセスの一例を、図 4 乃至図 9 のフローチャートを用いて説明する。なお、本発明の処理は一例であって、その処理、とくに検索条件修正処理部 22、ノイズ処理部 23 の処理などの順序を適宜、変更することは可能である。

【0075】

ユーザが、自らが入力する検索条件としての単語を知られずに検索装置 3 で検索を行うことを所望する場合、ユーザ端末 4 において実際の検索条件としての単語 A を入力すると、単語 A を含むクエリがユーザ端末 4 から制御端末 2 に送られる。そして、制御端末 2 の検索条件受付処理部 21 で、単語 A を含むクエリを受け付け (S100)、検索条件修正処理部 22 において、検索条件としての単語 A を修正する、検索条件修正処理を実行する (S110)。

【0076】

すなわち、検索条件修正処理部 22 は、Glove などの公知のモデルを用いることで、単語 A の 300 次元の単語ベクトル v ($v = v_1, v_2, \dots, v_{300}$) を取得する (S200)。また、検索条件修正処理部 22 は、300 次元のノイズベクトル n ($n = n_1, n_2, \dots, n_{300}$) を取得する (S210)。

【0077】

このように取得した単語ベクトル v とノイズベクトル n とをそれぞれ加算することでベクトル v' ($v' = v'_1, v'_2, \dots, v'_{300}$) を算出し (S220)、ベクトル v' の点からコサイン類似度 (コサイン距離) が一定の範囲内にある、任意の m 個のベクトル x (x_1, x_2, \dots, x_m) を、Glove のデータセットを参照することで特定をする (S230)。そして検索条件修正処理部 22 は、特定した各ベクトル x (x_1, x_2, \dots, x_m) に対応する単語 X_1, X_2, \dots, X_m を特定し、それらを修正検索条件の単語群 X とする (S240)。たとえば、 m は 10 個とすることができるが、数を増減してもよい。

【0078】

以上のように修正検索条件の単語群 X の単語 X_1, X_2, \dots, X_m を特定する。

【0079】

また、ノイズ処理部 23 は、単語 A に基づいて、ノイズとする単語の単語群 Y を特定するノイズ処理を実行する (S120)。

【0080】

ノイズ処理部 23 は、まず、ウィキペディアの Ontology クラスのデータセットを参照し、単語 A と同じクラスに属する単語から、十分に大きな数、たとえば 1000 個程度以上の単語 w を特定する (S300)。これらの単語 w によって構成される単語群を、単語群 W とする。たとえば単語群 W は、単語 w_1, \dots, w_{1000} により構成される。

【0081】

そしてノイズ処理部 23 は、ノイズとする単語群 Y の単語数を 10 個とし、その比率を

10

20

30

40

50

、たとえばクラスタ手法による単語数が5個、分割手法が3個、ランダム手法が2個と決定をすると、各手法によって、単語群Yにおけるノイズとする単語 Y_1, Y_2, \dots, Y_{10} を特定する処理を実行する(S310, S320, S330)。

【0082】

まずクラスタ手法によりノイズとする単語 Y_1, Y_2, \dots, Y_5 を特定するには(S310)、ノイズ処理部23は、単語Aの頻出頻度と、単語 w_1, \dots, w_{1000} のそれぞれの頻出頻度とを比較することで、単語Aの頻出頻度から所定範囲内の頻出頻度にある単語wを特定し、その特定した単語により単語群S(s_1, s_2, \dots, s_i)を構成する(S400)。この単語群Sにおける単語 s_1, s_2, \dots, s_i は、クラスタの中心の候補となる単語である。

10

【0083】

そして、単語Aの単語ベクトルに、任意に設定する同次元のノイズベクトルを加算する。そして、その加算したベクトルの点から、コサイン類似度(コサイン距離)が一定の閾値以上である単語群Sの単語を特定し、特定した単語により単語群S'(s'_1, s'_2, \dots, s'_j)(ただし $j < i$)を構成する(S410)。

【0084】

以上のように特定した単語群S'を構成する各単語 s'_1, s'_2, \dots, s'_j のうち、これらの各単語の単語ベクトルの点からコサイン類似度(コサイン距離)が一定の範囲内にある任意の数(たとえば3~10個程度)のベクトルを、Gloveのデータセットを参照することで特定をする。そして特定した各ベクトルに対応する単語を特定することで、一つのクラスタCを構成する。そして、単語のクラスタCを一または複数構成する(S420)。このようにすることで、単語群S'を構成する各単語 s'_1, s'_2, \dots, s'_j に近い単語に基づいて単語のクラスタCを構成することができる。

20

【0085】

そして各クラスタCにおける単語の密集度が高いクラスタから順番に、そのクラスタにおける単語を特定し、あらかじめ設定した数、ここでは5個になったら、それらをノイズの単語群Y(Y_1, Y_2, \dots, Y_5)として特定をする(S430)。

【0086】

以上のような処理をノイズ処理部23が実行することで、クラスタ手法によるノイズとする単語 Y_1, Y_2, \dots, Y_5 を特定できる。

30

【0087】

つぎに、ノイズ処理部23が分割手法により、ノイズとする単語 Y_6, Y_7, Y_8 を特定するには(S320)、ノイズ処理部23は、まず、単語Aのベクトルの点を通る超平面をランダムに特定することで(S500)、S300で特定した単語群Wにおける単語wについて、単語群W1、単語群W2に分割をする(S510)。このとき、単語数が多い領域を単語群W1、少ない領域を単語群W2とする。

【0088】

そして、ノイズ処理部23は、単語Aと、単語群W1における各単語とのコサイン類似度(コサイン距離)を比較し、コサイン類似度に基づいてソートをする。そして、コサイン類似度が低い順に下から、たとえば10%程度の単語群W1における単語を、単語Aに類似していない単語として、単語群W1から消去する。そしてこの消去した各単語を、新たな単語群Wとして構成する(S520)。

40

【0089】

S520で特定した単語群Wの単語に基づいて、S500乃至S520の処理を、所定条件、たとえば $|W| < 2y$ (yは任意の値)となるまで繰り返す(S530)。

【0090】

そして所定条件を充足した場合、最終的な単語群Wにおける単語と、単語Aとのコサイン類似度(コサイン距離)を比較し、コサイン類似度に基づいてソートをする。そして、コサイン類似度が高い順に上から、分割手法によるノイズの単語数分(ここでは3個)の単語を特定することで、分割手法によるノイズとする単語 Y_6, Y_7, Y_8 を特定する(

50

S 5 4 0)。

【 0 0 9 1 】

さらに、ノイズ処理部 2 3 がランダム手法により、ノイズとする単語 Y_9, Y_{10} を特定するには (S 3 3 0)、S 3 0 0 で特定した単語群 $W (w_1, \dots, w_{1000})$ の中から、ランダムに 2 個の単語を特定し、それらをランダム手法によるノイズとする単語 Y_9, Y_{10} とする (S 6 0 0)。

【 0 0 9 2 】

ノイズ処理部 2 3 が以上のような処理を実行することで、ノイズとする単語群 Y を構成する単語 Y_1, Y_2, \dots, Y_{10} を特定することができる。

【 0 0 9 3 】

検索装置処理部 2 4 は、検索条件修正処理部 2 2 で特定した修正検索条件の単語群 X の単語 X_1, X_2, \dots, X_{10} 、ノイズ処理部 2 3 で特定したノイズとする単語群 Y の単語 Y_1, Y_2, \dots, Y_{10} を、たとえばランダムや所定の規則に基づいて検索装置 3 に送ることで、各単語に基づく検索処理を検索装置 3 に実行させる (S 1 3 0)。なおこの際に、検索装置処理部 2 4 は、単語群 X 、単語群 Y の各単語のほか、単語 A を検索装置 3 に送ってもよい。

【 0 0 9 4 】

そして、検索装置処理部 2 4 は、検索装置 3 に送った各単語に基づく検索結果を受け付け (S 1 4 0)、検索結果処理部 2 5 が、ユーザ端末 4 に送る検索結果の生成処理を行う (S 1 5 0)。すなわち、検索結果処理部 2 5 は、検索装置処理部 2 4 で受け付けた検索結果のうち、単語群 Y における単語 Y_1, Y_2, \dots, Y_{10} に対する検索結果 $D (Y_1), D (Y_2), \dots, D (Y_{10})$ を除外し、単語群 X における単語 X_1, X_2, \dots, X_{10} に対する検索装置 3 での各検索結果 $D (X_1), D (X_2), \dots, D (X_{10})$ に基づいてソートするなど公知の手法を用いることで、単語 A に対する検索結果 $D (A)$ を生成する。そして検索結果処理部 2 5 は、S 1 5 0 で生成した検索結果を、単語 A に対する検索結果 $D (A)$ として、ユーザ端末 4 に送る (S 1 6 0)。

【 0 0 9 5 】

ユーザ端末 4 でこの検索結果を受け付けることで、ユーザは、自らが入力した単語 A に対する検索結果 $D (A)$ を取得することができる。

【 0 0 9 6 】

検索条件の匿名化 (検索装置 3 に検索条件を知られないようにすること) と、検索結果の精度とはトレードオフの関係にある。本発明の情報検索システム 1 において、匿名性は、実際の検索条件である単語 A との間の平均コサイン類似度を用い、以下の数 1 で示される。

(数 1)

$$\alpha = 1 - \frac{1}{|Q(A)|} \sum_{i=1}^n \text{sim}(v(A), v(X_i))$$

ここで、 $v (A)$ は単語 A のベクトルであり、 $v (X_i)$ は修正検索条件である単語 X のベクトルであり、 $Q (A)$ は、単語 A に基づく修正検索条件の単語群 $X (X_1, X_2, \dots, X_n)$ である。

【 0 0 9 7 】

また、修正検索条件 X に基づく検索結果による、検索結果の再構築可能性の精度は、以下の数 2 で示される。

(数 2)

$$\rho = \frac{|D(A) \cap D'(A)|}{|D(A)|}$$

ここで $D (A)$ は、単語 A に基づく検索結果であり、 $D' (A)$ は、以下の数 3 で示される修正検索条件 X による検索結果を用いて再構成された検索結果である。

10

20

30

40

50

(数3)

$$D'(A) = \bigcup_{i=1}^n D(X_i)$$

【0098】

そして、匿名性と、検索結果の再構築可能性の精度は、以下の数4の関係性が成立する。

(数4)

$$\log p = \frac{ck}{2d}(c + 2(1 - \alpha) \|v(A)\|_2) - \log Z$$

10

【0099】

出願人による本発明の情報検索システム1における、単語群Xのみを用いた場合の匿名性と再構築可能性の精度の関係を図13に、単語群Xと単語群Yとを用いた場合の匿名性と再構築可能性の精度の関係を図14に示す。図13(a)および図14(a)は修正検索条件の単語Xを特定するにあたりノイズベクトルnを用いない場合であり、図13(b)および図14(b)はノイズベクトルが小さい場合であり、図13(c)および図14(c)はノイズベクトルが大きい場合である。

【0100】

図13と図14は、本発明が実験的に確認できていることを示している。つまり、匿名性と再構築性がトレードオフの関係にあり、匿名性を上げれば再構築性は小さくなり、匿名性を下げれば再構築性は大きくなる。匿名性はノイズの大小で制御できるため、ノイズの選び方により、匿名性が再構築性のどちらを重要視したいかを選択できる。図13と図14の比較から、単語群Yを使用した方(図14)が使用しない方(図13)より全般的に匿名性が向上することがわかる。

20

【0101】

また、図15乃至図20に、本発明の情報検索システム1による検索結果と、実際の検索条件を入力した場合の比較例を示す。なお、図15乃至図20で用いた検索装置3はwikipediaである。図15乃至図18は実際の検索条件として「kyoto」を用いており、図15および図16では強いノイズ(ノイズベクトルが大きい)を、図17および図18では弱いノイズ(ノイズベクトルが小さい)の場合を示している。また図15乃至図18では

30

【0102】

図15では単語群Xの単語として「tokyo,copenhagen,hokkaido,nagoya,osaka,japan,kansai,seoul,fukuoka,chiba」が、単語群Yの単語として「arkansas,pueblo,saitama,conway,john,rosario,owen sound,armenia,patti,lyons,laporte,knowle west,columbus,north berwick,surat,patterson,millbrook,san diego,gill,walnut」が特定されている。そして、実際の検索条件「kyoto」で検索した場合と比較して、10個中9個の検索結果が一致している(左側のコラムが実際の検索条件に基づく検索結果、右側のコラムが本発明の情報検索システム1に基づく検索結果であり、左側のコラムに表示される「E」が一致している検索結果である)。また、匿名性は0.778、検索結果の再構築可能性は0.421(ただし上位100の検索結果では0.71)である。

40

【0103】

また図16では単語群Xの単語として「vasteras,nagoya,seoul,cmom,waseda,osaka,ginza,joad,tokyo,yokohama」が、単語群Yの単語として「ina,valley,islampur,rudbar,qasemabad,habibabad,alexander,wollongong,first,mehrdasht,humboldt,price,lara,perth,hayden,dauphin,hath,kuhsar,jahanabad,nosratabad」が特定されている。そして、実際の検索条件「kyoto」で検索した場合と比較して、10個中7個の検索結果が一致している。また、匿名性は0.796、検索結果の再構築可能性は0.35(ただし上位100の検索結果では0.59)である。

【0104】

50

また図 17 では単語群 X の単語として「nagoya,osaka,japan,copenhagen,seoul,tokyo,oslo,unfccc,treaty,nara」が、単語群 Y の単語として「fernando,street,lugo,gray,walton,madhubani,stoney,mineral,english,nassau,sulphur,spring,durango,clay,rain,buenavista,gatineau,mari,lacey,foster」が特定されている。そして実際の検索条件「kyoto」で検索した場合と比較して、すべての検索結果が一致している。また、匿名性は 0.769、検索結果の再構築可能性は 0.442 (ただし上位 100 の検索結果では 0.76) である。

【0105】

また図 18 では単語群 X の単語として「japan,hiroshima,osaka,nagasaki,oslo,seoul,nagoya,tokyo,copenhagen,treaty」が、単語群 Y の単語として「columbus,saint-louis,henderson,sidney,murray,roy,wolf,fox,hunter,clarksville,fountain,madison,marsha,monroe,seneca,southside,belvedere,huntingdon,avondale,afonso」が特定されている。そして、実際の検索条件「kyoto」で検索した場合と比較して、10 個中 9 個の検索結果が一致している。また、匿名性は 0.787、検索結果の再構築可能性は 0.409 (ただし上位 100 の検索結果では 0.66) である。

【0106】

さらに、図 19 および図 20 は実際の検索条件として「nagasaki」を用いており、いずれも強いノイズ(ノイズベクトルが大きい)の場合を示している。また図 19 および図 20 ではノイズ処理部 23 が、クラスタ手法とランダム手法を用いて単語群 Y を特定している。

【0107】

図 19 では単語群 X の単語として、「iacono,niigata,bombing,bombed,hiroshima,bombs,a-bombing,osaka,sanfexce,hijrah」が、単語群 Y の単語として「anjar,clark,jennings,lakeland,alexander,marshall,apple,valley,james,belleair,jefferson,band,thompson,ripley,morrison,taft,minneapolis,brooklyn heights,franklin,anderson」が特定されている。そして、実際の検索条件「nagasaki」で検索した場合と比較して、10 個中 6 個の検索結果が一致している。また、匿名性は 0.833、検索結果の再構築可能性は 0.539 (ただし上位 100 の検索結果では 0.47) である。

【0108】

図 20 では単語群 X の単語として、「devastated,kiel,suburbs,niigata,prefecture,incinerated,bombings,inundated,bombed,hiroshima」が、単語群 Y の単語として「grants,lakeview,howard,on,woodland,horton,lakeside,rudbar,nosratabad,baker,melrose park,davis,valley,victor,logan,manor,haymana,va,wollongong,bloomfield」が特定されている。そして、実際の検索条件「nagasaki」で検索した場合と比較して、10 個中 7 個の検索結果が一致している。また、匿名性は 0.836、検索結果の再構築可能性は 0.539 (ただし上位 100 の検索結果では 0.51) である。

【0109】

以上のように、図 15 乃至図 20 の比較結果に基づけば、実際の検索条件を匿名化した上で、検索結果の再構築可能性も一定の精度を有している。とくにウェブサイトを検索する検索エンジンの場合には、検索結果としてせいぜい上位 10 位程度しか参照しないことも多い。そうすると、検索結果の再構築可能性も 6 割以上であるなど、十分に実用に耐えることができる。

【実施例 2】

【0110】

本発明の情報検索システム 1 を用いた、別の実施態様として、たとえば検索条件修正処理部 22 による単語群 X の単語と、ノイズ処理部 23 による単語群 Y の単語とを表示させ、実際の検索条件を推測させる処理を設けることも可能である。この場合、実際の検索条件は、制御端末 2 が任意に特定をすることで、それに基づいて検索条件修正処理部 22 で単語群 X の単語を、ノイズ処理部 23 で単語群 Y の単語を、それぞれ特定する。そして単語群 X、単語群 Y の各単語をユーザ端末 4 にランダムの順番で表示させることで、それら

10

20

30

40

50

の単語から，実際に入力された単語を推測させる，一種のゲーム感覚で，情報検索システム 1 における効果を体感することも可能である。

【0111】

図 2 1 は単語群 X の単語，単語群 Y の単語をそれぞれ表示した状態を示す画面である。そして，図 2 2 は，図 2 1 の状態で分からなかった場合（分からないことを示す操作を受け付けた，または回答として入力を受けた単語が誤っていた場合）に，単語群 Y の単語を削除し，単語群 X の単語のみを表示した状態を示す画面である。図 2 3 は正解を表示した状態を示す画面である。

【0112】

このように，本実施例の処理を実行することで，ゲーム感覚で本発明の情報検索システム 1 の効果を体感することもできる。

10

【産業上の利用可能性】

【0113】

本発明の情報検索システム 1 を用いることによって，実際の検索条件の特定を困難とする情報検索システム 1 を発明した。特に，クラスタリング耐性の高いノイズを用いる場合には，クラスタリング耐性を高めることができる。

【符号の説明】

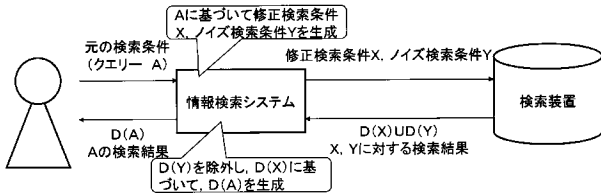
【0114】

- 1：情報検索システム
- 2：制御端末
- 3：検索装置
- 4：ユーザ端末
- 21：検索条件受付処理部
- 22：検索条件修正処理部
- 23：ノイズ処理部
- 24：検索装置処理部
- 25：検索結果処理部
- 70：演算装置
- 71：記憶装置
- 72：表示装置
- 73：入力装置
- 74：通信装置

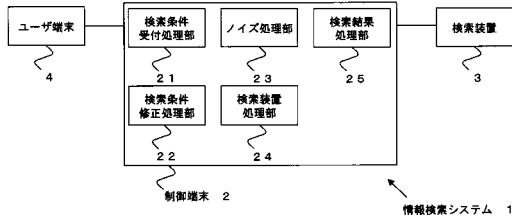
20

30

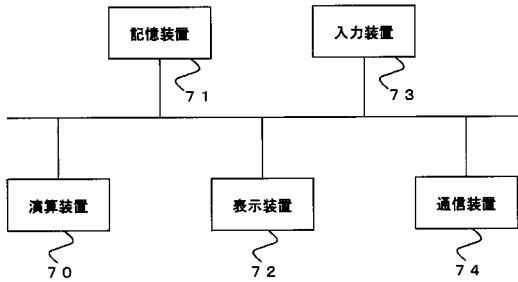
【 図 1 】



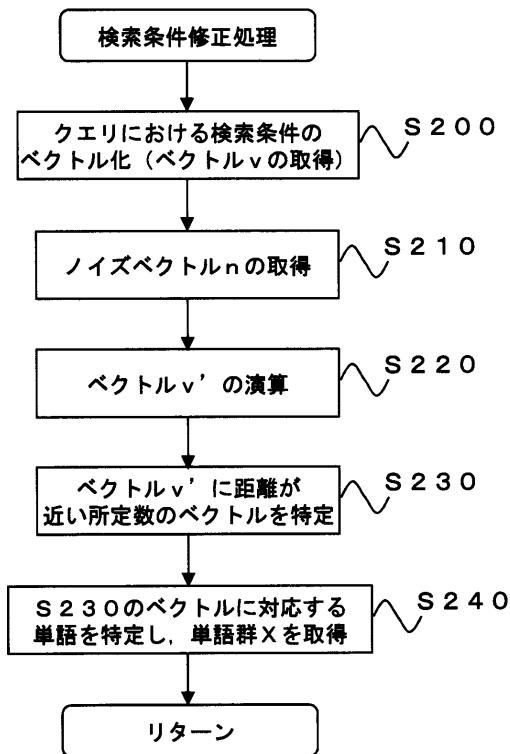
【 図 2 】



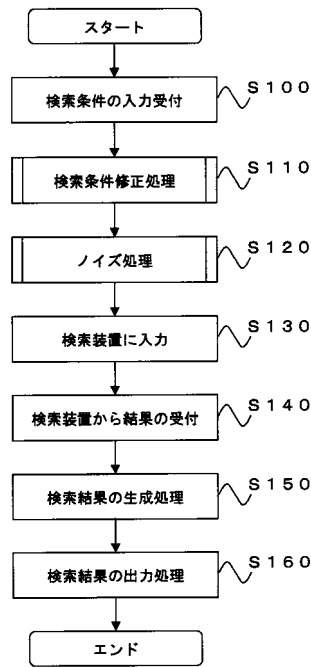
【 図 3 】



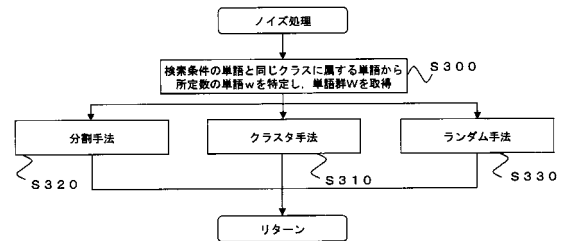
【 図 5 】



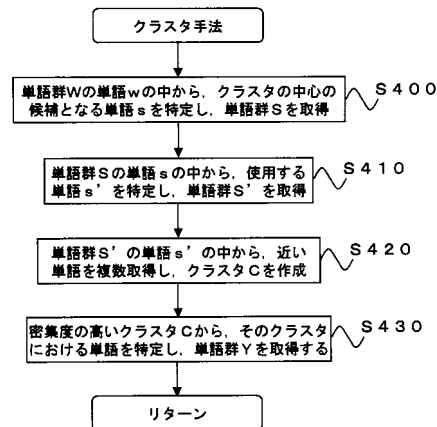
【 図 4 】



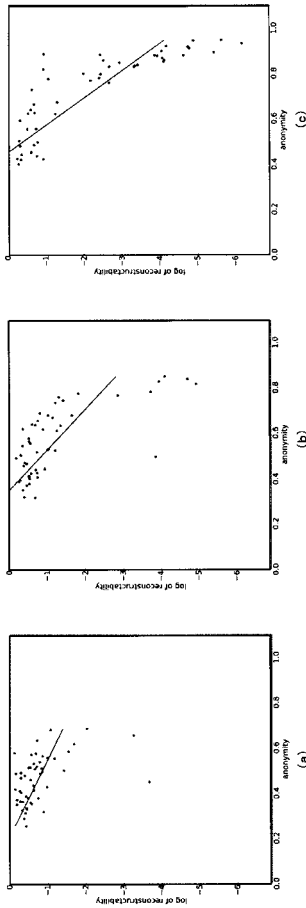
【 図 6 】



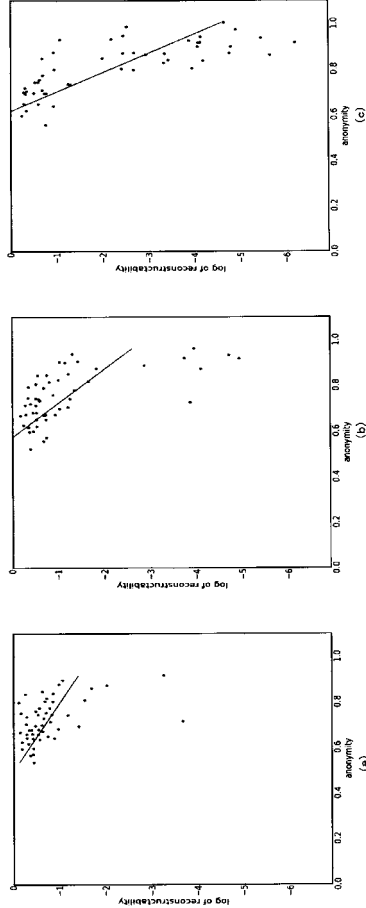
【 図 7 】



【 図 1 3 】



【 図 1 4 】



【 図 1 5 】

query: kyoto (up to three words available)
 noise level of semantically related terms: high-level
 number of unrelated distractor terms: 20 distractor method: 2-divison-random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Kyoto	1	Kyoto
2	Kyoto Protocol	2	Kyoto Prefecture
3	Kyoto Prefecture	3	Kyoto Protocol and government action
4	Kyoto Protocol and government action	4	Kyoto Station
5	Kyoto Station	5	Japanese garden
6	Japanese garden	6	Nara Line
7	Nara Line	7	Views on the Kyoto Protocol
8	Views on the Kyoto Protocol	8	Karasuma Line
9	Karasuma Line	9	List of National Treasures of Japan (crafts: others)
10	List of National Treasures of Japan (crafts: others)	10	Kansai dialect

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia
 ('E' means "exist in the right results of our system")

Anonymity: 0.778 (0.493), Reconstructability: 0.421 (0.71)

Note 1: Value (0.493) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.778.
 Note 2: Value (0.71) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.421.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

tokyo, arkansas, pueblo, copenhagen, saitama, conway, hokkaido, john, rosario, owen sound, armenia, patti, nagoya, osaka, lyons, laporte, knowle west, columbus, north berwick, surat, japan, patterson, millbrook, san diego, kansai, gill, seoul, walnut, fukuoka, chiba

【 図 1 6 】

query: kyoto (up to three words available)
 noise level of semantically related terms: high-level
 number of unrelated distractor terms: 20 distractor method: 2-divison-random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Kyoto	1	Kyoto
2	Kyoto Protocol	2	Kyoto Prefecture
3	Kyoto Prefecture	3	Kyoto Station
4	Kyoto Protocol and government action	4	Japanese garden
5	Kyoto Station	5	Nara Line
6	Japanese garden	6	Karasuma Line
7	Nara Line	7	List of National Treasures of Japan (crafts: others)
8	Views on the Kyoto Protocol	8	Kansai dialect
9	Karasuma Line	9	Kyoto University
10	List of National Treasures of Japan (crafts: others)	10	Japan National Route 1

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia
 ('E' means "exist in the right results of our system")

Anonymity: 0.796 (0.624), Reconstructability: 0.35 (0.59)

Note 1: Value (0.624) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.796.
 Note 2: Value (0.59) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.35.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

ina, valley, västerås, nagoya, islampur, seoul, rudbar, qasemabad, cmom, waseda, osaka, habibabad, alexander, wollongong, first, mcbrdashit, humboldt, price, lara, perth, ginza, joad, tokyo, hayden, dauphin, hat, kuhsar, jahanabad, yokohama, nosratabad

【 17 】

query: kyoto (up to three words available)
 noise level of semantically related terms: low-level
 number of unrelated distractor terms: 20 distractor method: 2-diverse+random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Kyoto	1	Kyoto
2	Kyoto Protocol	2	Kyoto Protocol
3	Kyoto Prefecture	3	Kyoto Prefecture
4	Kyoto Protocol and government action	4	Kyoto Protocol and government action
5	Kyoto Station	5	Kyoto Station
6	Japanese garden	6	Japanese garden
7	Nara Line	7	Nara Line
8	Views on the Kyoto Protocol	8	Views on the Kyoto Protocol
9	Karasuma Line	9	Karasuma Line
10	List of National Treasures of Japan (crafts: others)	10	List of National Treasures of Japan (crafts: others)

(E means "exist in the right results of our system")

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia

Anonymity : 0.769 (0.48), Reconstructability : 0.442 (0.76)

Note 1: Value (0.48) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.769.
 Note 2: Value (0.76) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.442.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

nagoya, fernando, street, lugo, gray, walton, madhubani, stoney, osaka, mineral, english, japan, copenhagen, nassau, sulphur, seoul, spring, durango, clay, rain, tokyo, buena vista, gatineau, mari, oslo, unfccc, treaty, lacey, foster, nara

【 18 】

query: kyoto (up to three words available)
 noise level of semantically related terms: low-level
 number of unrelated distractor terms: 20 distractor method: 2-diverse+random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Kyoto	1	Kyoto
2	Kyoto Protocol	2	Kyoto Prefecture
3	Kyoto Prefecture	3	Kyoto Protocol and government action
4	Kyoto Protocol and government action	4	Kyoto Station
5	Kyoto Station	5	Japanese garden
6	Japanese garden	6	Nara Line
7	Nara Line	7	Views on the Kyoto Protocol
8	Views on the Kyoto Protocol	8	Karasuma Line
9	Karasuma Line	9	List of National Treasures of Japan (crafts: others)
10	List of National Treasures of Japan (crafts: others)	10	Kansai district

(E means "exist in the right results of our system")

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia

Anonymity : 0.787 (0.478), Reconstructability : 0.409 (0.66)

Note 1: Value (0.478) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.787.
 Note 2: Value (0.66) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.409.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

japan, hiroshima, osaka, columbus, saint-louis, nagasaki, oslo, seoul, henderson, nagoya, sidney, murray, tokyo, roy, wolf, fox, hunter, clarksville, fountain, madison, marsa, copenhagen, treaty, monroe, seneca, southside, belvedere, huntingdon, avondale, afonso

【 19 】

query: nagasaki (up to three words available)
 noise level of semantically related terms: high-level
 number of unrelated distractor terms: 20 distractor method: cluster+random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Nagasaki	1	Nagasaki
2	Nagasaki Prefecture	2	First Inno Nagasaki
3	First Inno Nagasaki	3	Takasaki Nagai
4	Nagasaki Main Line	4	Nagasaki Station (Nagasaki)
5	Takasaki Nagai	5	Nagasaki Peace Park
6	Dejima	6	List of cities in Japan
7	Japan National Route 57	7	Hibokusha
8	Nagasaki Station (Nagasaki)	8	Slope car
9	Nagasaki Peace Park	9	Hiroshi Motoblima
10	List of cities in Japan	10	Jocho Itoh

(E means "exist in the right results of our system")

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia

Anonymity : 0.833 (0.592), Reconstructability : 0.539 (0.47)

Note 1: Value (0.592) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.833.
 Note 2: Value (0.47) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.539.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

anjer, clark, jennings, iacono, niigata, lakeand, bombing, bombed, alexander, marshall, apple, valley, james, hiroshima, belleair, jefferson, band, thompson, ripley, morrison, taft, bombs, minneapolis, a-bombing, brooklyn heights, osaka, sanfrecco, franklin, hijrah, anderson

【 20 】

query: nagasaki (up to three words available)
 noise level of semantically related terms: high-level
 number of unrelated distractor terms: 20 distractor method: cluster+random

Search Results in wikipedia titles

Results of input query		Results of our system	
1	Nagasaki	1	Nagasaki
2	Nagasaki Prefecture	2	First Inno Nagasaki
3	First Inno Nagasaki	3	Takasaki Nagai
4	Nagasaki Main Line	4	Japan National Route 57
5	Takasaki Nagai	5	Nagasaki Station (Nagasaki)
6	Dejima	6	Nagasaki Peace Park
7	Japan National Route 57	7	List of cities in Japan
8	Nagasaki Station (Nagasaki)	8	Hibokusha
9	Nagasaki Peace Park	9	Slope car
10	List of cities in Japan	10	Hiroshi Motoblima

(E means "exist in the right results of our system")

Note: Ranking is sorted by the frequency of the input query in the text of each wikipedia

Anonymity : 0.836 (0.615), Reconstructability : 0.539 (0.51)

Note 1: Value (0.615) in the left parentheses is the anonymity in case of semantically related terms: the actual value is 0.836.
 Note 2: Value (0.51) in the right parentheses is the reconstructability in case of top 100 search results: the actual value is 0.539.
 Note 3: Go to hyperlink pages for the definitions of the anonymity and reconstructability. (Click each word).

Semantically related terms (10 items), Unrelated distractor terms (20 items)

grants, lakeview, howard, on, woodland, devastated, kiel, horton, lakeside, suburbs, rudbar, nosratabad, niigata, baker, melrose park, prefecture, davis, valley, victor, incinerated, bombings, inundated, logan, bombed, manor, haymana, va, wollongong, bloomfield, hiroshima

【 ☒ 2 1 】

Game start

noise level of related terms: number of distractor terms: distractor method:

Semantically related terms + Unrelated distractor terms (10 items + 10 items)

huron, nagoya, bole, grammar, spoken, lyons, konni, speaking, linguistic, english, lombard, vocabulary, ruma, word, seneca, kaba, translation, learning, kodi, meaning

Enter your prediction query

Check the answer

Original query:

Cosine similarity:

Levenshtein distance:

Back

【 ☒ 2 2 】

Game start

noise level of related terms: number of distractor terms: distractor method:

Semantically related terms + Unrelated distractor terms (10 items + 10 items)

huron, nagoya, bole, grammar, spoken, lyons, konni, speaking, linguistic, english, lombard, vocabulary, ruma, word, seneca, kaba, translation, learning, kodi, meaning

Enter your prediction query

Check the answer

Your answer city is wrong. Please predict again from semantically related terms.

Original query: ??

Cosine similarity: 0.367

Levenshtein distance: 8

Back

【 ☒ 2 3 】

Game start

noise level of related terms: number of distractor terms: distractor method:

Semantically related terms + Unrelated distractor terms (10 items + 10 items)

huron, nagoya, bole, grammar, spoken, lyons, konni, speaking, linguistic, english, lombard, vocabulary, ruma, word, seneca, kaba, translation, learning, kodi, meaning

Enter your prediction query language

Check the answer

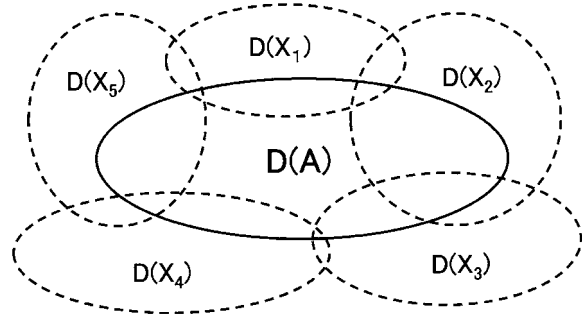
Original query: language

Cosine similarity: 1.0

Levenshtein distance: 0

Back

【 ☒ 2 4 】



フロントページの続き

(72)発明者 ボレガラ ダヌシカ

東京都千代田区一ツ橋二丁目1番2号 大学共同利用機関法人情報・システム研究機構 国立情報学研究所内