

(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2018-142188

(P2018-142188A)

(43) 公開日 平成30年9月13日(2018.9.13)

(51) Int.Cl. F I テーマコード (参考)
G06F 17/30 (2006.01) G06F 17/30 210A
 G06F 17/30 170A

審査請求 未請求 請求項の数 5 O L (全 16 頁)

(21) 出願番号 特願2017-36288 (P2017-36288)
 (22) 出願日 平成29年2月28日 (2017.2.28)

(71) 出願人 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番1号
 (74) 代理人 110002147
 特許業務法人酒井国際特許事務所
 (72) 発明者 片岡 正弘
 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
 (72) 発明者 尾上 聡
 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
 (72) 発明者 吉田 裕之
 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

(54) 【発明の名称】 解析プログラム、解析方法および解析装置

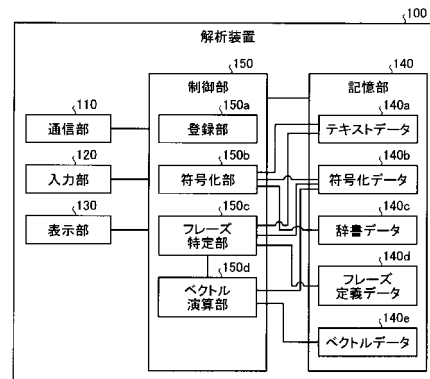
(57) 【要約】

【課題】 解析対象の文書に対する解析速度および解析精度の向上を図る。

【解決手段】 解析装置100は、解析対象の文書を単語単位で符号化した、複数の符号化単語を生成する。解析装置100は、符号化単語、または、複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、符号化単語、または、符号化フレーズに割り当てる。

【選択図】 図2

本実施例1に係る解析装置の構成を示す機能ブロック図



【特許請求の範囲】**【請求項 1】**

コンピュータに、
解析対象の文書を単語単位で符号化した、複数の符号化単語を生成し、
前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる
処理を実行させることを特徴とする解析プログラム。

【請求項 2】

第 1 符号化フレーズの前後に位置する前記符号化単語または他の第 2 符号化フレーズに対してSkip-gramによるモデル化を行い、前記第 1 符号化フレーズの前後に、前記符号化単語または前記第 2 符号化フレーズが出現する確率を特定することで、前記第 1 符号化フレーズのベクトル値を算出することを特徴とする請求項 1 に記載の解析プログラム。

10

【請求項 3】

フレーズを構成する第 1 単語と第 2 単語との組を定義したテーブルを基にして、前記解析対象の文書に含まれる第 1 単語を特定し、特定した第 1 単語と同じ文中に第 2 単語が含まれる場合には、特定した第 1 単語と同じ文中に含まれる第 2 単語を削除し、特定した第 1 単語の直後に前記第 2 単語を配置する処理を更にコンピュータに実行させることを特徴とする請求項 1 または 2 に記載の解析プログラム。

【請求項 4】

コンピュータが実行する解析方法であって、
解析対象の文書を単語単位で符号化した、複数の符号化単語を生成し、
前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる
処理を実行することを特徴とする解析方法。

20

【請求項 5】

解析対象の文書を単語単位で符号化した、複数の符号化単語を生成する符号化部と、
前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる演算部と
を有することを特徴とする解析装置。

30

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、解析プログラム等に関する。

【背景技術】**【0002】**

従来、文書を分散表現する手法として解析対象の文書を構成する形態素それぞれに基づいて、文書からベクトルを生成するWord2Vec技術が存在する。例えば、Word2Vec技術では、ある単語（形態素）と、ある単語に隣接する他の単語との関係に基づいて、各単語のベクトル値を算出する処理を行う。

40

【0003】

ここで、Word2Vec技術等により文書をベクトルを用いて分散表現する場合に、解析対象となる文書に含まれる「the」、「a」等の冠詞、「on」、「of」等の前置詞等の高頻度の単語の影響が過大となる。このため、Word2Vec技術では、高頻度の単語をストップワードとして文書から排除した後に、ベクトルによる分散表現を生成する。

【0004】

例えば、解析対象の文書「He takes care of his daughter」をWord2Vec技術では、ストップワードとなる「of」を除外した後、「He takes care his daughter」に含

50

まれる単語それぞれをベクトル化する。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2006-48685号公報

【特許文献2】特開2009-151757号公報

【非特許文献】

【0006】

【非特許文献1】Distributed Representations of Words and Phrases and their Compositionality, Tomas Mikolov et. al, pp. 3111-3119, Advances in Neural Information Processing Systems 26, 2013, Curran Associates, Inc.

10

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、上述した従来技術では、解析対象の文書に対する解析速度および解析精度が低いという問題がある。

【0008】

たとえば、Word2Vec技術により、ストップワードとして除外される「the」、「a」等の冠詞、「on」、「of」等の前置詞等は、特定の文字列において存在の有無により意味が変化する可能性がある。具体的には、「take care of」の「of」、「the Japanese」の「the」は、存在の有無により意味が変化するため、かかる「of」、「the」等を除外してベクトル化を行うと、本来の文書の意味が変わった状態のベクトル化がなされることから、生成されたベクトルを用いた解析の精度が低下する可能性がある。

20

【0009】

また、従来Word2Vec技術で用いられる解析において、ストップワードを含めて計算量が過大とならない解析手法は知られておらず、適切な計算時間により、目的の精度を得ることができない。

【0010】

1つの側面では、本発明は、解析対象の文書に対する解析速度および解析精度の向上を図ることができる解析プログラム、解析方法および解析装置を提供することを目的とする。

30

【課題を解決するための手段】

【0011】

第1の案では、コンピュータに下記の処理を実行させる。コンピュータは、解析対象の文書を単語単位で符号化した、複数の符号化単語を生成する。コンピュータは、符号化単語、または、複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、符号化単語、または、符号化フレーズに割り当てる。

【発明の効果】

【0012】

符号化したフレーズに対しても、値を割り当てたベクトルを生成することにより、解析対象の文書に対する解析速度および解析精度の向上を図ることができる。

40

【図面の簡単な説明】

【0013】

【図1】図1は、本実施例1に係る解析装置の処理の一例を説明するための図である。

【図2】図2は、本実施例1に係る解析装置の構成を示す機能ブロック図である。

【図3】図3は、フレーズ定義データのデータ構造の一例を示す図である。

【図4】図4は、本実施例1に係る解析装置の処理手順を示すフローチャートである。

【図5】図5は、本実施例2に係る解析装置の処理の一例を説明するための図である。

【図6】図6は、本実施例2に係る解析装置の構成を示す機能ブロック図である。

50

【図 7】図 7 は、離間フレーズテーブルのデータ構造の一例を示す図である。

【図 8】図 8 は、本実施例 2 に係る解析装置の処理手順を示すフローチャートである。

【図 9】図 9 は、オートマトンを用いたベクトル演算の一例を説明するための図である。

【図 10】図 10 は、解析装置と同様の機能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。

【発明を実施するための形態】

【0014】

以下に、本願の開示する解析プログラム、解析方法および解析装置の実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。

【実施例 1】

【0015】

図 1 は、本実施例 1 に係る解析装置の処理の一例を説明するための図である。ここでは、解析装置が、テキストデータ 10 a を符号化した符号化データ 10 b をベクトル化する場合について説明する。図 1 に示す例では、テキストデータ 10 a を「Every day we take care of our daughter」とする。

【0016】

解析装置は、所定のフレーズを定義したフレーズ定義テーブルと、テキストデータ 10 a とを比較して、テキストデータ 10 a に含まれる各単語のうち、所定のフレーズを構成する複数の単語の組を特定する。図 1 に示す例では、「take care of」が所定のフレーズとして特定される。

【0017】

解析装置は、単語とコードとを対応付けた辞書データと、テキストデータ 10 a に含まれる単語とを比較することで、テキストデータ 10 a の単語を単語単位で符号化することで、符号化データ 10 b を生成する。例えば、解析装置は、単語「Every」をコード A 1、単語「day」をコード A 2、単語「we」をコード A 3、単語「take」をコード A 4 に符号化する。解析装置は、単語「care」をコード A 5、単語「of」をコード A 6、単語「our」をコード A 7、単語「daughter」をコード A 8 に符号化する。

【0018】

解析装置は、複数のコードのうち、所定のフレーズを構成する単語のコードの組と、所定のフレーズを構成する単語に対応しないコードを特定する。以下の説明では、適宜、所定のフレーズを構成する単語のコードの組を、「符号化フレーズ」と表記する。所定のフレーズを構成する単語に対応しないコードを、「符号化単語」と表記する。図 1 に示す例では、コード A 4 ~ コード A 6 の組が、符号化フレーズ 15 となる。他のコード A 1 ~ A 3、A 7、A 8 は、それぞれ符号化単語となる。

【0019】

解析装置は、符号化単語および符号化フレーズの出現状況に応じて、符号化単語および符号化フレーズのベクトル値をそれぞれ算出することで、符号化データ 10 b をベクトル化する。

【0020】

解析装置が、符号化単語のベクトル値を算出する処理について説明する。ベクトル値の算出対象となる符号化単語を、対象単語と表記する。解析装置は、符号化データ 10 b 上において、対象単語の前方 2 つの符号化単語または符号化フレーズと、対象単語の後方 2 つの符号化単語に対して、Skip-gram によるモデル化を行い、対象単語の前後に符号化単語または符号化フレーズが出現する確率を特定することで、対象単語のベクトル値を算出する。

【0021】

例えば、解析装置は、コード A 3 のベクトル値を算出する場合には、コード A 1、コード A 2、符号化フレーズ 15、コード A 7 に対して、Skip-gram によるモデル化を行う。解析装置は、コード A 3 の前後に、コード A 1、コード A 2、符号化フレーズ 15、コード A 7 が出現する確率を特定することで、コード A 3 のベクトル値を算出する。解析装置

10

20

30

40

50

は、コード A 1、A 2、A 7、A 8 についても同様の処理を実行することで、各コードのベクトル値を算出する。

【0022】

解析装置が、符号化フレーズのベクトル値を算出する処理について説明する。ベクトル値の算出対象となる符号化フレーズを、対象フレーズと表記する。解析装置は、符号化データ 10 b 上において、対象フレーズの前方 2 つの符号化単語または符号化フレーズと、対象フレーズの後方 2 つの符号化単語に対して、Skip-gram によるモデル化を行い、対象フレーズの前後に符号化単語または符号化フレーズが出現する確率を特定することで、対象フレーズのベクトル値を算出する。

【0023】

例えば、解析装置は、符号化フレーズ 15 のベクトル値を算出する場合には、コード A 2、コード A 3、コード A 7、コード A 8 に対して、Skip-gram によるモデル化を行う。解析装置は、符号化フレーズ 15 の前後に、コード A 2、コード A 3、コード A 7、コード A 8 が出現する確率を特定することで、符号化フレーズ 15 のベクトル値を算出する。解析装置は、他の符号化フレーズについても同様の処理を実行することで、各コードのベクトル値を算出する。

【0024】

解析装置は、各符号単語、各符号化フレーズに対応するベクトル値を割り当てることで、符号化データ 10 b をベクトル化する。

【0025】

上記に記載したように、本実施例 1 に係る解析装置によれば、テキストデータ 10 a に含まれる単語からストップワードを除外すること無く、単語単位で符号化を行うことで、符号化データ 10 b を生成する。解析装置は、符号化単語のうち、所定のフレーズを構成する符号化単語を符号化フレーズとしてまとめ、符号化単語および符号化フレーズの出現状況に応じて、符号化データ 10 b をベクトル化する。このように、解析装置は、ストップワードを除外しないため、符号化データ 10 b をベクトル化の際の解析精度を向上させることができる。また、解析装置が利用する符号化では、符号化データ 10 b が単語単位で符号化されているため、Zip による符号化と異なり、符号化したままで各単語を区別できるため、復号化を行うことなく、文書に含まれる単語のベクトル値を算出でき、解析速度を向上できる。

【0026】

図 2 は、本実施例 1 に係る解析装置の構成を示す機能ブロック図である。図 2 に示すように、この解析装置 100 は、通信部 110 と、入力部 120 と、表示部 130 と、記憶部 140 と、制御部 150 とを有する。

【0027】

通信部 110 は、ネットワークを介して他の外部装置と通信を実行する処理部である。例えば、解析装置 100 は、後述するテキストデータ 140 a、辞書データ 140 c、フレーズ定義データ 140 d 等を、通信部 110 を介して、受信しても良い。

【0028】

入力部 120 は、解析装置 100 に各種の情報を入力する入力装置である。例えば、入力部 120 は、キーボードやマウス、タッチパネル等に対応する。

【0029】

表示部 130 は、制御部 150 から出力される各種の情報を表示する表示装置である。例えば、表示部 130 は、液晶ディスプレイやタッチパネル等に対応する。

【0030】

記憶部 140 は、テキストデータ 140 a と、符号化データ 140 b と、辞書データ 140 c と、フレーズ定義データ 140 d と、ベクトルデータ 140 e とを有する。記憶部 140 は、RAM (Random Access Memory)、ROM (Read Only Memory)、フラッシュメモリ (Flash Memory) などの半導体メモリ素子や、HDD (Hard Disk Drive) などの記憶装置に対応する。

10

20

30

40

50

【0031】

テキストデータ140aは、複数の単語を含む文字列データである。テキストデータ140aの一例は、図1に示したテキストデータ10aとなる。

【0032】

符号化データ140bは、テキストデータ140aに含まれる各単語を単語単位で符号化したデータである。符号化データ140bの一例は、図1に示した符号化データ10bとなる。

【0033】

辞書データ140cは、単語と、単語に対応するコードとを対応付けるデータである。

【0034】

フレーズ定義データ140dは、フレーズを構成する複数の単語の組み合わせを定義したデータである。図3は、フレーズ定義データのデータ構造の一例を示す図である。図3に示すように、フレーズ定義データ140dには、各種のフレーズが定義されている。図3に示すフレーズは一例であり、他のフレーズも含まれている。

【0035】

図3では一例として、フレーズを、符号化前の単語の組で定義しているがこれに限定されるものではない。例えば、フレーズ定義データ140dは、フレーズを構成する単語を、符号化後のコードによって定義しても良い。すなわち「in front of」であれば、「(inのコード) (frontのコード) (ofのコード)」によって、定義することができる。

【0036】

ベクトルデータ140eは、符号化データ140bに含まれる各符号化単語、各符号化フレーズに割り当てられたベクトル値を示す情報であり、符号化データ140bをベクトル化した情報である。

【0037】

制御部150は、登録部150aと、符号化部150bと、フレーズ特定部150cと、ベクトル演算部150dとを有する。制御部150は、CPU (Central Processing Unit) やMPU (Micro Processing Unit) などによって実現できる。また、制御部150は、ASIC (Application Specific Integrated Circuit) やFPGA (Field Programmable Gate Array) などのハードワイヤードロジックによっても実現できる。

【0038】

登録部150aは、通信部110または入力部120を介して、各種の情報を受け付けた場合に、受け付けた情報を記憶部140に登録する処理部である。例えば、登録部150aは、テキストデータ140a、辞書データ140c、フレーズ定義データ140dを受け付けた場合には、受け付けたテキストデータ140a、辞書データ140c、フレーズ定義データ140dを記憶部140に登録する。

【0039】

符号化部150bは、テキストデータ140aを符号化する処理部である。符号化部150bは、テキストデータ140aに含まれる各単語と、辞書データ140cとを比較して、辞書データ140cにヒットした単語を、単語単位で符号化する処理を繰り返し実行することで、符号化データ140aを生成する。

【0040】

フレーズ特定部150cは、フレーズ定義データ140dを基にして、符号化データ140bに含まれる各符号化単語のうち、符号化フレーズを構成する符号化単語の組を特定する処理部である。フレーズ特定部150cは、符号化データ140bのコードのうち、符号化フレーズを構成する符号化単語の組の情報を、ベクトル演算部150dに出力する。

【0041】

フレーズ定義データ140dに定義されたフレーズが、符号化前の単語の組で定義され

10

20

30

40

50

ている場合の処理について説明する。フレーズ特定部 150c は、テキストデータ 140a と、フレーズ定義データ 140d とを比較することで、テキストデータ 140a に含まれるフレーズを構成する単語の組を特定する。フレーズ特定部 150c は、特定したフレーズを構成する単語の組に対応する各コードを、辞書データ 140c を基にして特定する。フレーズ特定部 150c は、特定した各コードと、符号化データ 140b とを比較して、符号化データ 140b に存在する符号化フレーズを特定する。

【0042】

フレーズ定義データ 140d に定義されたフレーズが、符号化後のコードの組で定義されている場合の処理について説明する。フレーズ特定部 150c は、符号化データ 140b と、フレーズ定義データ 140d とを比較することで、符号化データ 140b に含まれる符号化フレーズを特定する。

10

【0043】

ベクトル演算部 150d は、符号化データ 140b に含まれる符号化単語および符号化フレーズの出現状況に応じて、符号化単語および符号化フレーズのベクトル値をそれぞれ算出し、割り当てることで、ベクトルデータ 140e を生成する処理部である。ベクトル演算部 150d は、演算部の一例である。

【0044】

ベクトル演算部 150d は、フレーズ特定部 150c から特定結果を受け付け、符号化データ 140b に含まれる各符号化単語のうち、符号化フレーズに対応する符号化単語と、符号化フレーズに対応しない符号化単語とを区別する。ベクトル演算部 150d は、符号化フレーズのベクトル値と、符号化フレーズに対応しない符号化単語のベクトル値とを算出する。

20

【0045】

ベクトル演算部 150d が、符号化単語のベクトル値を算出する処理について説明する。ベクトル値の算出対象となる符号化単語を、対象単語と表記する。ベクトル演算部 150d は、符号化データ 140b 上において、対象単語の前方 2 つの符号化単語または符号化フレーズと、対象単語の後方 2 つの符号化単語に対して、Skip-gram によるモデル化を行い、対象単語の前後に符号化単語または符号化フレーズが出現する確率を特定することで、対象単語のベクトル値を算出する。

【0046】

図 1 を用いて説明すると、ベクトル演算部 150d は、コード A 3 のベクトル値を算出する場合には、コード A 1、コード A 2、符号化フレーズ 15、コード A 7 に対して、Skip-gram によるモデル化を行う。ベクトル演算部 150d は、コード A 3 の前後に、コード A 1、コード A 2、符号化フレーズ 15、コード A 7 が出現する確率を特定することで、コード A 3 のベクトル値を算出する。解析装置は、コード A 1、A 2、A 7、A 8 についても同様の処理を実行することで、各コードのベクトル値を算出する。

30

【0047】

ベクトル演算部 150d が、符号化フレーズのベクトル値を算出する処理について説明する。ベクトル値の算出対象となる符号化フレーズを、対象フレーズと表記する。ベクトル演算部 150d は、符号化データ 140b 上において、対象フレーズの前方 2 つの符号化単語または符号化フレーズと、対象フレーズの後方 2 つの符号化単語に対して、Skip-gram によるモデル化を行い、対象フレーズの前後に符号化単語または符号化フレーズが出現する確率を特定することで、対象フレーズのベクトル値を算出する。

40

【0048】

図 1 を用いて説明すると、ベクトル演算部 150d は、符号化フレーズ 15 のベクトル値を算出する場合には、コード A 2、コード A 3、コード A 7、コード A 8 に対して、Skip-gram によるモデル化を行う。ベクトル演算部 150d は、符号化フレーズ 15 の前後に、コード A 2、コード A 3、コード A 7、コード A 8 が出現する確率を特定することで、符号化フレーズ 15 のベクトル値を算出する。ベクトル演算部 150d は、他の符号化フレーズについても同様の処理を実行することで、各コードのベクトル値を算出する。

50

【0049】

ここで、ベクトル演算部150dが、Skip-gramによるモデル化を行い、対象単語（対象フレーズ）のベクトル値を算出する処理は、例えば、文献（Tomas Mikolov他、「Distributed Representations of Words and Phrases and their Compositionality」）に記載した技術を利用する。

【0050】

図4は、本実施例1に係る解析装置の処理手順を示すフローチャートである。図4に示すように、解析装置100の符号化部150bは、テキストデータ140aを読み込む（ステップS101）。符号化部150bは、テキストデータ140aを、辞書データ140cを基にして、単語単位に符号化することで、符号化データ140bを生成する（ステップS102）。

10

【0051】

解析装置100のフレーズ特定部150cは、フレーズ定義データ140dを基にして、符号化データ140bに含まれる各コードのうち、符号化単語（符号化フレーズに含まれない符号化単語）と、符号化フレーズとを特定する（ステップS103）。

【0052】

解析装置100のベクトル演算部150dは、符号化フレーズに含まれない符号化単語について、符号化単語の出現状況に応じて、ベクトル値を算出する（ステップS104）。ベクトル演算部150dは、符号化フレーズについて、符号化フレーズの出現状況に応じて、ベクトル値を算出する（ステップS105）。ベクトル演算部150dは、符号化データ140bに対するベクトルデータ140eを生成する（ステップS106）。

20

【0053】

上記に記載したように、解析装置100によれば、テキストデータ140aに含まれる単語からストップワードを除外すること無く、単語単位で符号化を行うことで、符号化データ140bを生成する。解析装置100は、符号化単語のうち、所定のフレーズを構成する符号化単語を符号化フレーズとしてまとめ、符号化単語および符号化フレーズの出現状況に応じて、符号化データ140bをベクトル化する。このように、解析装置100は、ストップワードを除外しないため、符号化データ140bをベクトル化する際の解析精度を向上させることができる。また、解析装置100が利用する符号化では、符号化データ140bが単語単位で符号化されているため、Zipによる符号化と異なり、符号化したまま各単語を区別できるため、復号化を行うことなく、文書に含まれる単語のベクトル値を算出でき、解析速度を向上できる。

30

【実施例2】

【0054】

図5は、本実施例2に係る解析装置の処理の一例を説明するための図である。ここでは、解析装置が解析するテキストデータ20aを「We take lunch out」とする。テキストデータ20aに含まれる「take、out」はフレーズに対応する単語の組であるが、各単語が離れているため、実施例1で説明した解析装置100が利用するフレーズ定義データ140dにヒットせず、符号化フレーズのベクトル値を算出できない場合がある。以下の説明では、それぞれが離間した単語により構成されるフレーズを「離間フレーズ」と表記する。

40

【0055】

本実施例2に係る解析装置は、テキストデータ20aを走査して、離間フレーズを検出した場合には、離間フレーズを構成する各単語が連続するように、テキストデータ20aの単語を並び変えることで、テキストデータ21aを生成する。例えば、図5に示す例では、解析装置100は、テキストデータ20aの「out」を削除し、削除した「out」を「take」の直後に配置することで、テキストデータ21aを生成する。解析装置は、テキストデータ21a（テキストデータ21aを符号化した符号化データ）に基づいて、テキストデータ21aをベクトル化する。解析装置がテキストデータ21aに基づいて、テキストデータ21aをベクトル化する処理は、上述した実施例1の処理と同様である。

50

【0056】

上記処理を実行することで、本実施例2に係る解析装置によれば、離間フレーズがテキストデータ20aに存在する場合でも、離間フレーズを符号化フレーズとして特定することができる。このため、離間フレーズを構成する各単語のコード毎のベクトル値を算出することを抑止して、解析精度を向上させることができる。

【0057】

図6は、本実施例2に係る解析装置の構成を示す機能ブロック図である。図6に示すように、この解析装置200は、通信部210と、入力部220と、表示部230と、記憶部240と、制御部250とを有する。このうち、入力部220、表示部230に関する説明は、図2で説明した入力部120、表示部130に関する説明と同様であるため、説明を省略する。

10

【0058】

通信部210は、ネットワークを介して他の外部装置と通信を実行する処理部である。例えば、解析装置200は、後述するテキストデータ240a、離間フレーズテーブル240c、辞書データ240d、フレーズ定義データ240e等を、通信部210を介して、受信しても良い。

【0059】

記憶部240は、テキストデータ240aと、テキストデータ241aと、符号化データ240bと、離間フレーズテーブル240cと、辞書データ240dと、フレーズ定義データ240eと、ベクトルデータ240fとを有する。記憶部240は、RAM、ROM、フラッシュメモリなどの半導体メモリ素子や、HDDなどの記憶装置に対応する。

20

【0060】

テキストデータ240aは、複数の単語を含む文字列データである。テキストデータ240aの一例は、図5に示したテキストデータ20aとなる。テキストデータ241aは、図5で説明したように、離間フレーズの単語が連続するように単語が並び換えられテキストデータ21aに対応するデータである。

【0061】

符号化データ240bは、テキストデータ241aに含まれる各単語を単語単位で符号化したデータである。

【0062】

離間フレーズテーブル240cは、離間フレーズに関する情報を定義したテーブルである。図7は、離間フレーズテーブルのデータ構造の一例を示す図である。図7に示すように、この離間フレーズテーブル240cは、主単語、副単語、フレーズを対応付ける。主単語は、離間フレーズのうち、最初に現れる単語である。副単語は、離間フレーズのうち、主単語の次に現れる単語である。フレーズは、離間フレーズを構成する単語を連続して並べたものである。

30

【0063】

辞書データ240d、フレーズ定義データ240eに関する説明は、図2で説明した辞書データ140c、フレーズ定義データ140dに関する説明と同様である。

【0064】

ベクトルデータ240fは、符号化データ240bに含まれる各符号化単語、各符号化フレーズに割り当てられたベクトル値を示す情報であり、符号化データ240bをベクトル化した情報である。

40

【0065】

制御部250は、登録部250aと、離間フレーズ処理部250bと、符号化部250cと、フレーズ特定部250dと、ベクトル演算部250eとを有する。制御部250は、CPUやMPUなどによって実現できる。また、制御部250は、ASICやFPGAなどのハードワイヤードロジックによっても実現できる。

【0066】

登録部250aは、通信部210または入力部220を介して、各種の情報を受け付け

50

た場合に、受け付けた情報を記憶部 2 4 0 に登録する処理部である。例えば、登録部 2 5 0 a は、テキストデータ 2 4 0 a、離間フレーズテーブル 2 4 0 c、辞書データ 2 4 0 d、フレーズ定義データ 2 4 0 e を受け付けた場合には、受け付けた各データを記憶部 2 4 0 に登録する。

【 0 0 6 7 】

離間フレーズ処理部 2 5 0 b は、テキストデータ 2 4 0 a と、離間フレーズテーブル 2 4 0 c とを比較して、テキストデータ 2 4 0 a に含まれる離間フレーズを特定する。離間フレーズ処理部 2 5 0 b は、特定した離間フレーズの単語が連続するように並び変えを行うことで、テキストデータ 2 4 1 a を生成する。以下において、離間フレーズ処理部 2 5 0 b の処理の一例について説明する。

10

【 0 0 6 8 】

離間フレーズ処理部 2 5 0 b は、テキストデータ 2 4 0 a と、離間フレーズテーブル 2 4 0 c とを比較して、離間フレーズテーブル 2 4 0 c の主単語にヒットする単語を、テキストデータ 2 4 0 a から特定する。離間フレーズ処理部 2 5 0 b は、主単語にヒットする単語が存在する場合には、ヒットした単語と同じ文中で、ヒットした単語から後ろ方向に、所定語数未満の位置に、副単語（主単語に対応する副単語）にヒットするか否かを判定する。離間フレーズテーブル 2 4 0 c は、主単語および副単語にヒットした場合には、係る主単語、副単語を離間フレーズとして特定する。

【 0 0 6 9 】

例えば、離間フレーズ処理部 2 5 0 b は、テキストデータ 2 4 0 a と、離間フレーズテーブル 2 4 0 c とを比較し、主単語「take」がヒットしたものとする。離間フレーズ処理部 2 5 0 b は、ヒットした主単語「take」と同じ文中で、「take」から後ろ方向に所定語数未満の位置に、副単語「out」が存在する場合には、離間した「take」、「out」を離間フレーズであると特定する。

20

【 0 0 7 0 】

離間フレーズ処理部 2 5 0 b は、離間フレーズを特定すると、離間フレーズの副単語を削除し、主単語の直後に副単語を配置する処理を実行する。離間フレーズ処理部 2 5 0 b は、各離間フレーズについて、上記処理を繰り返し実行することで、テキストデータ 2 4 1 a を生成する。

【 0 0 7 1 】

なお、離間フレーズ処理部 2 5 0 b は、主単語をフレーズに置き換え、副単語を削除することで、主単語と副単語とが連続するように置き換えを行っても良い。

30

【 0 0 7 2 】

符号化部 2 5 0 c は、テキストデータ 2 4 1 a を符号化する処理部である。符号化部 2 5 0 c は、テキストデータ 2 4 1 a に含まれる各単語と、辞書データ 2 4 0 d とを比較して、辞書データ 2 4 0 d にヒットした単語を、単語単位で符号化する処理を繰り返し実行することで、符号化データ 2 4 0 b を生成する。

【 0 0 7 3 】

フレーズ特定部 2 5 0 d は、フレーズ定義データ 2 4 0 e を基にして、符号化データ 2 4 0 b に含まれる各符号化単語のうち、符号化フレーズを構成する符号化単語の組を特定する処理部である。フレーズ特定部 2 5 0 d は、符号化データ 2 4 0 b のコードのうち、符号化フレーズを構成する符号化単語の組の情報を、ベクトル演算部 2 5 0 e に出力する。フレーズ特定部 2 5 0 d に関するその他の処理は、図 2 に示したフレーズ特定部 1 5 0 c の処理と同様である。

40

【 0 0 7 4 】

ベクトル演算部 2 5 0 e は、符号化データ 2 4 0 b に含まれる符号化単語および符号化フレーズの出現状況に応じて、符号化単語および符号化フレーズのベクトル値をそれぞれ算出し、割り当てることで、ベクトルデータ 2 4 0 f を生成する処理部である。ベクトル演算部 2 5 0 e に関するその他の処理は、図 2 で説明したベクトル演算部 1 5 0 e に関する処理と同様である。

50

【0075】

図8は、本実施例2に係る解析装置の処理手順を示すフローチャートである。図8に示すように、解析装置200の離間フレーズ処理部250bは、テキストデータ240aを読み込む(ステップS201)。離間フレーズ処理部250bは、離間フレーズテーブル240cとテキストデータ240aとを比較して離間フレーズを特定する(ステップS202)。

【0076】

離間フレーズ処理部250bは、離間フレーズにヒットしない場合には(ステップS203, No)、ステップS205に移行する。離間フレーズ処理部250bは、離間フレーズにヒットした場合には(ステップS203, Yes)、ステップS204に移行する。離間フレーズ処理部250bは、離間フレーズに対応する副単語を移動し、主単語の直後に副単語を配置する(ステップS204)。

10

【0077】

解析装置200の符号化部250cは、テキストデータ241aを、辞書データ240dを基にして、単語単位に符号化することで、符号化データ240bを生成する(ステップS205)。

【0078】

解析装置200のフレーズ特定部250dは、フレーズ定義データ240eを基にして、符号化データ240bに含まれる各コードのうち、符号化単語(符号化フレーズに含まれない符号化単語)と、符号化フレーズとを特定する(ステップS206)。

20

【0079】

解析装置200のベクトル演算部250eは、符号化フレーズに含まれない符号化単語について、符号化単語の出現状況に応じて、ベクトル値を算出する(ステップS207)。ベクトル演算部250eは、符号化フレーズについて、符号化フレーズの出現状況に応じて、ベクトル値を算出する(ステップS208)。ベクトル演算部250eは、符号化データ240bに対するベクトルデータ240fを生成する(ステップS209)。

【0080】

上記に記載したように、解析装置200によれば、離間フレーズがテキストデータ20aに存在する場合でも、離間フレーズを符号化フレーズとして特定することができる。このため、離間フレーズを構成する各単語のコード毎のベクトル値を算出することを抑止して、解析精度を向上させることができる。

30

【0081】

ところで、本実施例で説明した解析装置100(200)は、各符号化単語とベクトル値との関係が既知の場合には、各符号化単語とベクトル値とを対応付けたオートマトンを用いて、テキストデータのベクトル化を行ってもよい。

【0082】

図9は、オートマトンを用いたベクトル演算の一例を説明するための図である。図9に示す例では、テキストデータを「Every day we take care of our daughter」とする。解析装置100は、図1と同様にして、テキストデータ10aに含まれる各単語を、単語単位で符号化することで、符号化データ10bを生成する。

40

【0083】

続いて、解析装置100は、符号化データ10aとオートマトン50とを比較して、各符号化単語をベクトル値に変換する。ここで、オートマトン50は、各符号化単語とベクトル値とを対応付けた情報である。なお、オートマトン50は、ストップワード等に対応する符号化単語に対応するベクトル値を「0」に設定しておく。

【0084】

図9に示す例では、オートマトン50による変換により、コードA1~A5、コードA6~A8が、ベクトル値V1~V7に変換される。なお、コードA6は、ストップワードに対応する符号化単語であるため、ベクトル値は0に設定される。

【0085】

50

上記のように、オートマトン50を利用して、符号化データ10bをベクトル化することで、ベクトル化する処理を高速化することが可能となる。

【0086】

次に、上記実施例に示した解析装置100、200と同様の機能を実現するコンピュータのハードウェア構成の一例について説明する。図10は、解析装置と同様の機能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。

【0087】

図10に示すように、コンピュータ300は、各種演算処理を実行するCPU301と、ユーザからのデータの入力を受け付ける入力装置302と、ディスプレイ303とを有する。また、コンピュータ300は、記憶媒体からプログラム等を読み取る読み取り装置304と、ネットワークを介して他のコンピュータとの間でデータの授受を行うインタフェース装置305とを有する。また、コンピュータ300は、各種情報を一時記憶するRAM306と、ハードディスク装置307とを有する。そして、各装置301~307は、バス308に接続される。

【0088】

ハードディスク装置307は、離間フレーズ処理プログラム307a、符号化プログラム307b、フレーズ特定プログラム307c、ベクトル演算プログラム307dを有する。CPU301は、離間フレーズ処理プログラム307a、符号化プログラム307b、フレーズ特定プログラム307c、ベクトル演算プログラム307dを読み出してRAM306に展開する。

【0089】

離間フレーズ処理プログラム307aは、離間フレーズ処理プロセス306aとして機能する。符号化プログラム307bは、符号化プロセス306bとして機能する。フレーズ特定プログラム307cは、フレーズ特定プロセス306cとして機能する。ベクトル演算プログラム307dは、ベクトル演算プロセス306dとして機能する。

【0090】

離間フレーズ処理プロセス306aの処理は、離間フレーズ処理部250bの処理に対応する。符号化プロセス306bの処理は、符号化部150b、250cの処理に対応する。フレーズ特定プロセス306cの処理は、フレーズ特定部150c、250dの処理に対応する。ベクトル演算プロセス306dの処理は、ベクトル演算部150d、250eの処理に対応する。

【0091】

なお、各プログラム307a~307dについては、必ずしも最初からハードディスク装置307に記憶させておかなくても良い。例えば、コンピュータ300に挿入されるフレキシブルディスク(FD)、CD-ROM、DVDディスク、光磁気ディスク、ICカードなどの「可搬用の物理媒体」に各プログラムを記憶させておく。そして、コンピュータ300が各プログラム307a~307dを読み出して実行するようにしても良い。

【0092】

以上の各実施例を含む実施形態に関し、さらに以下の付記を開示する。

【0093】

(付記1) コンピュータに、

解析対象の文書を単語単位で符号化した、複数の符号化単語を生成し、

前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる

処理を実行させることを特徴とする解析プログラム。

【0094】

(付記2) 第1符号化フレーズの前後に位置する前記符号化単語または他の第2符号化フレーズに対してSkip-gramによるモデル化を行い、前記第1符号化フレーズの前後に、前記符号化単語または前記第2符号化フレーズが出現する確率を特定することで、前記第1

10

20

30

40

50

符号化フレーズのベクトル値を算出することを特徴とする付記 1 に記載の解析プログラム。

【 0 0 9 5 】

(付記 3) フレーズを構成する第 1 単語と第 2 単語との組を定義したテーブルを基にして、前記解析対象の文書に含まれる第 1 単語を特定し、特定した第 1 単語と同じ文中に第 2 単語が含まれる場合には、特定した第 1 単語と同じ文中に含まれる第 2 単語を削除し、特定した第 1 単語の直後に前記第 2 単語を配置する処理を更にコンピュータに実行させることを特徴とする付記 1 または 2 に記載の解析プログラム。

【 0 0 9 6 】

(付記 4) コンピュータが実行する解析方法であって、
解析対象の文書を単語単位で符号化した、複数の符号化単語を生成し、
前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる
処理を実行することを特徴とする解析方法。

10

【 0 0 9 7 】

(付記 5) 第 1 符号化フレーズの前後に位置する前記符号化単語または他の第 2 符号化フレーズに対してSkip-gramによるモデル化を行い、前記第 1 符号化フレーズの前後に、前記符号化単語または前記第 2 符号化フレーズが出現する確率を特定することで、前記第 1 符号化フレーズのベクトル値を算出することを特徴とする付記 4 に記載の解析方法。

20

【 0 0 9 8 】

(付記 6) フレーズを構成する第 1 単語と第 2 単語との組を定義したテーブルを基にして、前記解析対象の文書に含まれる第 1 単語を特定し、特定した第 1 単語と同じ文中に第 2 単語が含まれる場合には、特定した第 1 単語と同じ文中に含まれる第 2 単語を削除し、特定した第 1 単語の直後に前記第 2 単語を配置する処理を更にコンピュータに実行させることを特徴とする付記 4 または 5 に記載の解析方法。

【 0 0 9 9 】

(付記 7) 解析対象の文書を単語単位で符号化した、複数の符号化単語を生成する符号化部と、
前記符号化単語、または、前記複数の符号化単語の組み合わせに割り当てられた符号化フレーズに関する、前記解析対象の文書の出現状況に応じてそれぞれ生成されたベクトル値を、前記符号化単語、または、前記符号化フレーズに割り当てる演算部と
を有することを特徴とする解析装置。

30

【 0 1 0 0 】

(付記 8) 演算部は、第 1 符号化フレーズの前後に位置する前記符号化単語または他の第 2 符号化フレーズに対してSkip-gramによるモデル化を行い、前記第 1 符号化フレーズの前後に、前記符号化単語または前記第 2 符号化フレーズが出現する確率を特定することで、前記第 1 符号化フレーズのベクトル値を算出することを特徴とする付記 7 に記載の解析装置。

40

【 0 1 0 1 】

(付記 9) フレーズを構成する第 1 単語と第 2 単語との組を定義したテーブルを基にして、前記解析対象の文書に含まれる第 1 単語を特定し、特定した第 1 単語と同じ文中に第 2 単語が含まれる場合には、特定した第 1 単語と同じ文中に含まれる第 2 単語を削除し、特定した第 1 単語の直後に前記第 2 単語を配置する離間フレーズ処理部を更に有することを特徴とする付記 7 または 8 に記載の解析装置。

【 符号の説明 】

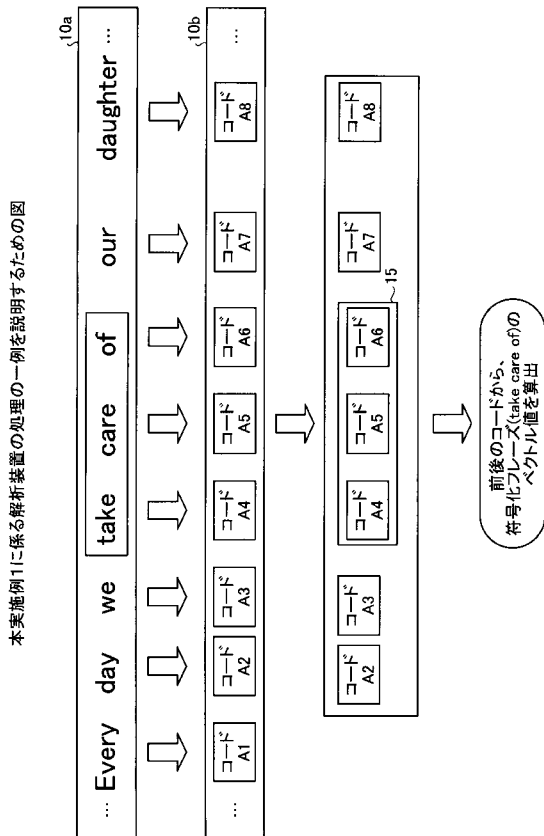
【 0 1 0 2 】

1 0 0、2 0 0 解析装置
1 1 0、2 1 0 通信部
1 2 0、2 2 0 入力部

50

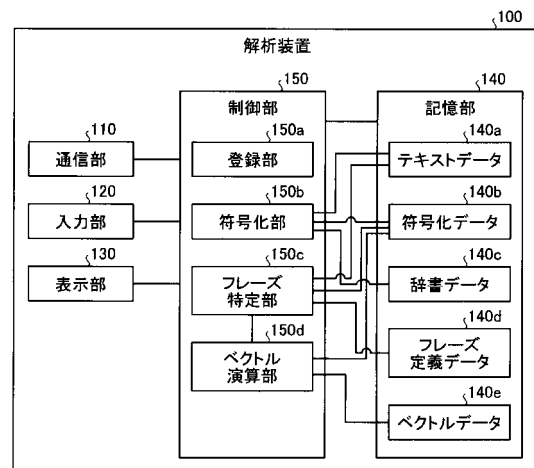
- 1 3 0、2 3 0 表示部
- 1 4 0、2 4 0 記憶部
- 1 5 0、2 5 0 制御部

【 図 1 】



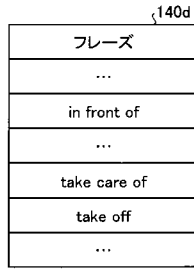
【 図 2 】

本実施例1に係る解析装置の構成を示す機能ブロック図



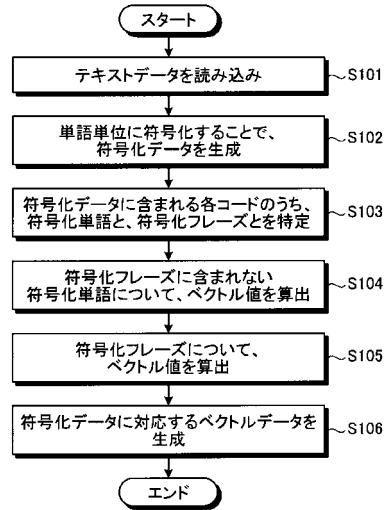
【 図 3 】

フレーズ定義データのデータ構造の一例を示す図



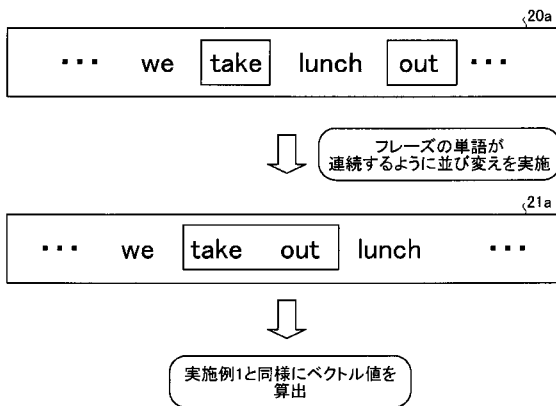
【 図 4 】

本実施例1に係る解析装置の処理手順を示すフローチャート



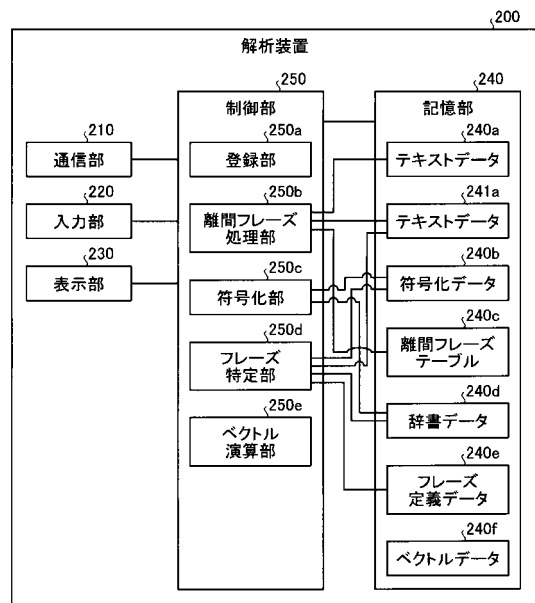
【 図 5 】

本実施例2に係る解析装置の処理の一例を説明するための図



【 図 6 】

本実施例2に係る解析装置の構成を示す機能ブロック図



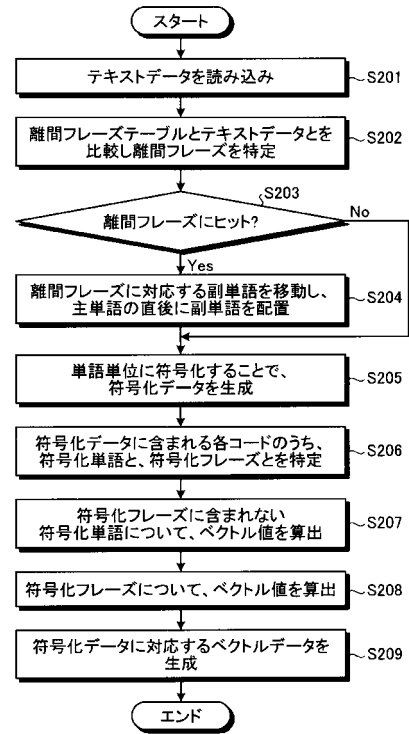
【 図 7 】

離間フレーズテーブルのデータ構造の一例を示す図

主単語	副単語	フレーズ
...
figure	out	figure out
take	out	take out
turn	off	turn off
put	out	put out
make	out	make out
...

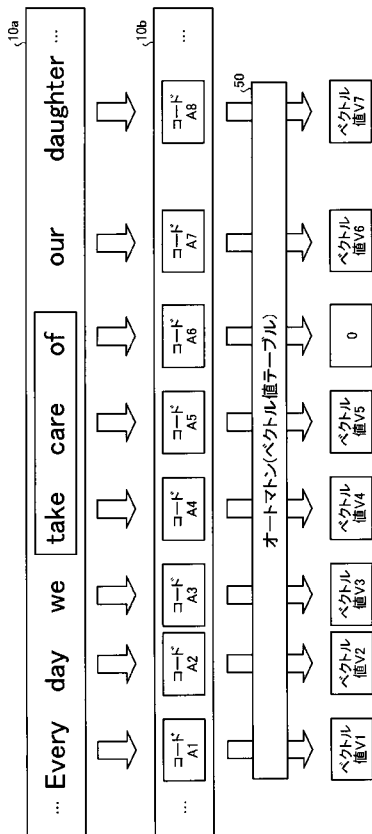
【 図 8 】

本実施例2に係る解析装置の処理手順を示すフローチャート



【 図 9 】

オートマトンを用いたベクトル演算の一例を説明するための図



【 図 10 】

解析装置と同様の機能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

