

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号  
特開2022-51113  
(P2022-51113A)

(43)公開日 令和4年3月31日(2022.3.31)

(51)国際特許分類	F I	テーマコード(参考)
G 0 6 F 40/216(2020.01)	G 0 6 F 40/216	5 B 0 9 1
G 0 6 F 40/56(2020.01)	G 0 6 F 40/56	

審査請求 未請求 請求項の数 10 O L (全14頁)

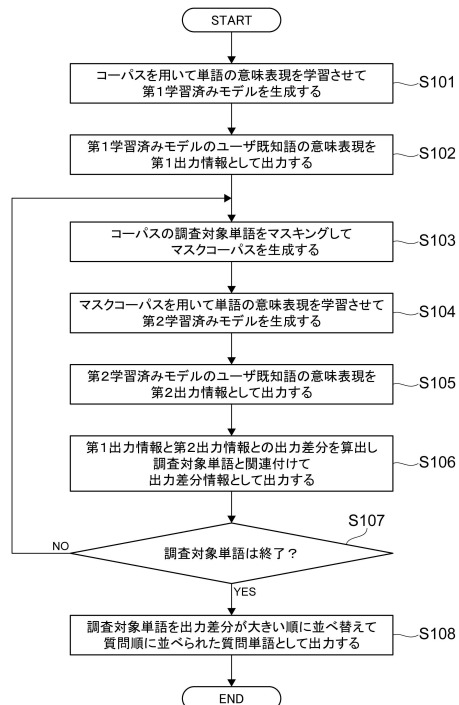
(21)出願番号	特願2020-157394(P2020-157394)	(71)出願人	000005496
(22)出願日	令和2年9月18日(2020.9.18)		富士フィルムビジネスイノベーション株式会社
			東京都港区赤坂九丁目7番3号
		(74)代理人	100104880
			弁理士 古部 次郎
		(74)代理人	100118108
			弁理士 久保 洋之
		(72)発明者	稲木 誓哉
			神奈川県横浜市西区みなとみらい六丁目1番 富士ゼロックス株式会社内
		Fターム(参考)	5B091 EA01

(54)【発明の名称】 情報出力装置、質問生成装置、及びプログラム

(57)【要約】 (修正有)【課題】アンケート等でユーザ既知語を調べる場合に比較して、ユーザ既知語を効率的に調べることを可能とする情報出力装置、質問生成装置及びプログラムを提供する。

【解決手段】質問生成装置は、プロセッサを備える。プロセッサは、特定の例文集合を用いて単語の意味表現を学習させた第1のモデルから得られた特定のユーザ既知語の意味表現と、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第2のモデルから得られた特定のユーザ既知語の意味表現との差分を算出し、差分に基づいて、対象単語がユーザ既知語である可能性に関する情報を出力する。

【選択図】図7



## 【特許請求の範囲】

## 【請求項 1】

プロセッサを備え、  
前記プロセッサは、  
特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出し、  
前記差分に基づいて、前記対象単語がユーザ既知語である可能性に関する情報を出力することを特徴とする情報出力装置。

10

## 【請求項 2】

前記プロセッサは、  
複数の対象単語の各対象単語について前記差分を算出することにより、複数の差分を算出し、  
前記対象単語がユーザ既知語である可能性に関する情報として、前記複数の対象単語の各対象単語についての前記差分に基づく順序で並べられた当該複数の対象単語を出力することを特徴とする請求項 1 に記載の情報出力装置。

## 【請求項 3】

前記差分に基づく順序は、当該差分が大きい順序であることを特徴とする請求項 2 に記載の情報出力装置。

20

## 【請求項 4】

前記第 2 のモデルは、前記特定の例文集合の前記対象単語を除く部分を用いて単語の意味表現を未学習モデルに新たに学習させたモデルであることを特徴とする請求項 1 に記載の情報出力装置。

## 【請求項 5】

前記特定の例文集合の前記対象単語を除く部分は、当該特定の例文集合の前記特定のユーザ既知語及び当該対象単語の少なくとも何れか一方を含む構成要素の当該対象単語を除く部分であることを特徴とする請求項 4 に記載の情報出力装置。

## 【請求項 6】

前記第 2 のモデルは、前記特定の例文集合の前記対象単語を除く部分を用いて単語の意味表現を学習済みモデルに更に学習させたモデルであることを特徴とする請求項 1 に記載の情報出力装置。

30

## 【請求項 7】

前記特定の例文集合の前記対象単語を除く部分は、当該特定の例文集合の当該対象単語を含む構成要素の当該対象単語を除く部分であることを特徴とする請求項 6 に記載の情報出力装置。

## 【請求項 8】

プロセッサを備え、  
前記プロセッサは、  
複数の対象単語の各対象単語について、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の当該各対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出することにより、複数の差分を算出し、  
前記複数の差分に基づいて、前記複数の対象単語を用いた質問を生成することを特徴とする質問生成装置。

40

## 【請求項 9】

前記プロセッサは、  
前記特定のユーザ既知語に代えて、前記質問に対するユーザの回答から把握される他のユーザ既知語を用いて、前記複数の差分を算出し、

50

前記複数の差分に基づいて、前記複数の対象単語を用いた質問を再生成することを特徴とする請求項 8 に記載の質問生成装置。

【請求項 10】

コンピュータに、

特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出する機能と、

前記差分に基づいて、前記対象単語がユーザ既知語である可能性に関する情報を出力する機能と

10

を実現させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報出力装置、質問生成装置、及びプログラムに関する。

【背景技術】

【0002】

あるパターンが他のパターンを含意しているような 2 つの言語パターンのペアを生成する技術は、知られている（例えば、特許文献 1 参照）。

【先行技術文献】

20

【特許文献】

【0003】

【特許文献 1】特許第 6551968 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

例えば、ユーザに単語を用いた質問を行い、その質問に対するユーザの回答を利用して何らかの処理を行うことがある。その際、ユーザがその単語の意味を知らないと、回答の質又は量が低下するので、その単語はユーザが意味を知っている単語（以下、「ユーザ既知語」という）であることが望ましい。ここで、ユーザ既知語を調べるには、アンケート等

30

【0005】

本発明の目的は、アンケート等でユーザ既知語を調べる場合に比較して、ユーザ既知語を効率的に調べることを可能とすることにある。

【課題を解決するための手段】

【0006】

請求項 1 に記載の発明は、プロセッサを備え、前記プロセッサは、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出し、前記差分に

40

基づいて、前記対象単語がユーザ既知語である可能性に関する情報を出力することを特徴とする情報出力装置である。

請求項 2 に記載の発明は、前記プロセッサは、複数の対象単語の各対象単語について前記差分を算出することにより、複数の差分を算出し、前記対象単語がユーザ既知語である可能性に関する情報として、前記複数の対象単語の各対象単語についての前記差分に基づく順序で並べられた当該複数の対象単語を出力することを特徴とする請求項 1 に記載の情報出力装置である。

請求項 3 に記載の発明は、前記差分に基づく順序は、当該差分が大きい順序であることを特徴とする請求項 2 に記載の情報出力装置である。

請求項 4 に記載の発明は、前記第 2 のモデルは、前記特定の例文集合の前記対象単語を除

50

く部分を用いて単語の意味表現を未学習モデルに新たに学習させたモデルであることを特徴とする請求項 1 に記載の情報出力装置である。

請求項 5 に記載の発明は、前記特定の例文集合の前記対象単語を除く部分は、当該特定の例文集合の前記特定のユーザ既知語及び当該対象単語の少なくとも何れか一方を含む構成要素の当該対象単語を除く部分であることを特徴とする請求項 4 に記載の情報出力装置である。

請求項 6 に記載の発明は、前記第 2 のモデルは、前記特定の例文集合の前記対象単語を除く部分を用いて単語の意味表現を学習済みモデルに更に学習させたモデルであることを特徴とする請求項 1 に記載の情報出力装置である。

請求項 7 に記載の発明は、前記特定の例文集合の前記対象単語を除く部分は、当該特定の例文集合の当該対象単語を含む構成要素の当該対象単語を除く部分であることを特徴とする請求項 6 に記載の情報出力装置である。

請求項 8 に記載の発明は、プロセッサを備え、前記プロセッサは、複数の対象単語の各対象単語について、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の当該各対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出することにより、複数の差分を算出し、前記複数の差分に基づいて、前記複数の対象単語を用いた質問を生成することを特徴とする質問生成装置である。

請求項 9 に記載の発明は、前記プロセッサは、前記特定のユーザ既知語に代えて、前記質問に対するユーザの回答から把握される他のユーザ既知語を用いて、前記複数の差分を算出し、前記複数の差分に基づいて、前記複数の対象単語を用いた質問を再生成することを特徴とする請求項 8 に記載の質問生成装置である。

請求項 10 に記載の発明は、コンピュータに、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、当該特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた当該特定のユーザ既知語の意味表現との差分を算出する機能と、前記差分に基づいて、前記対象単語がユーザ既知語である可能性に関する情報を出力する機能とを実現させるためのプログラムである。

#### 【発明の効果】

#### 【0007】

請求項 1 の発明によれば、アンケート等でユーザ既知語を調べる場合に比較して、ユーザ既知語を効率的に調べることが可能となる。

請求項 2 の発明によれば、複数の対象単語についてのユーザ既知語である可能性に基づく順序を知ることができる。

請求項 3 の発明によれば、複数の対象単語についてのユーザ既知語である可能性が高い順序を知ることができる。

請求項 4 の発明によれば、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させるモデルとして学習済みモデルを用意しなくても、ユーザ既知語を調べることが可能となる。

請求項 5 の発明によれば、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させるモデルとして未学習モデルを用意した場合において、特定の例文集合を削減することができる。

請求項 6 の発明によれば、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させるモデルとして未学習モデルを用意する場合に比較して、単語の意味表現を学習させたモデルの精度を向上することができる。

請求項 7 の発明によれば、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させるモデルとして学習済みモデルを用意した場合において、特定の例文集合を削減することができる。

請求項 8 の発明によれば、アンケート等でユーザ既知語を調べてユーザ既知語を用いた質

問を生成する場合に比較して、ユーザ既知語を用いた質問を効率的に生成することが可能となる。

請求項 9 の発明によれば、特定のユーザ既知語のみを用いて質問を生成する場合に比較して、ユーザ既知語を用いた質問が生成される可能性が高まる。

請求項 10 の発明によれば、アンケート等でユーザ既知語を調べる場合に比較して、ユーザ既知語を効率的に調べることが可能となる。

【図面の簡単な説明】

【0008】

【図 1】本発明の実施の形態における質問生成装置のハードウェア構成例を示した図である。

10

【図 2】本発明の実施の形態における質問生成装置の機能構成例を示したブロック図である。

【図 3】(a), (b) は本発明の実施の形態における質問生成装置で記憶されるコーパスの具体例を示した図である。

【図 4】(a), (b) は本発明の実施の形態における質問生成装置で記憶される学習済みモデルの具体例を示した図である。

【図 5】(a), (b) は本発明の実施の形態における質問生成装置で記憶される出力情報の具体例を示した図である。

【図 6】本発明の実施の形態における質問生成装置で記憶される出力差分情報の具体例を示した図である。

20

【図 7】本発明の実施の形態における質問生成装置の動作例を示したフローチャートである。

【発明を実施するための形態】

【0009】

以下、添付図面を参照して、本発明の実施の形態について詳細に説明する。

【0010】

[本実施の形態の概要]

本実施の形態は、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた特定のユーザ既知語の意味表現との差分を算出し、その差分に基づいて、対象単語がユーザ既知語である可能性に関する情報を出力する情報出力装置である。

30

【0011】

ここで、情報出力装置は、1つの対象単語について差分を算出し、この差分が閾値以上であれば、対象単語がユーザ既知語である可能性に関する情報として、対象単語がユーザ既知語と判断される旨の情報を出力するものでもよい。

【0012】

或いは、情報出力装置は、複数の対象単語の各対象単語について差分を算出することにより、複数の差分を算出し、対象単語がユーザ既知語である可能性に関する情報として、複数の対象単語の各対象単語についての差分に基づく順序で並べられた複数の対象単語を出力するものでもよい。

40

【0013】

情報出力装置は、これらの何れであってもよいが、以下では、後者であるものとして説明する。そして、単に複数の対象単語を出力するのではなく、複数の対象単語を用いた質問を生成するものとする。

【0014】

その場合、本実施の形態は、複数の対象単語の各対象単語について、特定の例文集合を用いて単語の意味表現を学習させた第 1 のモデルから得られた特定のユーザ既知語の意味表現と、特定の例文集合の各対象単語を除く部分を用いて単語の意味表現を学習させた第 2 のモデルから得られた特定のユーザ既知語の意味表現との差分を算出することにより、複

50

数の差分を算出し、複数の差分に基づいて、複数の対象単語を用いた質問を生成する質問生成装置となる。

【 0 0 1 5 】

従って、以下では、本実施の形態が質問生成装置である場合を例にとって説明する。

【 0 0 1 6 】

ここで、質問生成装置とは、ユーザに与える質問を生成する装置である。この装置は、例えば、質問に対するユーザの回答を利用して、目的のタスクを解くシステムにおいて、質問を生成する装置であってよい。タスクとしては、例えば、単語分類、単語間関連度予測がある。

【 0 0 1 7 】

システムが質問を与える方法としては、次のような方法が考えられる。

【 0 0 1 8 】

目的のタスクが単語分類タスクである場合は、システムが、単語と分類項目とを提示して、その単語に最も関連しそうな分類項目をユーザに質問する、という方法である。

【 0 0 1 9 】

目的のタスクが単語間関連度予測タスクである場合は、システムが、2つの単語を提示して、それらがどのくらい関連しているかをユーザに質問する、という方法である。

【 0 0 2 0 】

また、例文集合とは、何らかの例文を集めたものをいう。例文は、記事や書物等の一般に「文書」と呼ばれ得る比較的長い文であってもよいし、会話の文等の比較的短い文であってもよい。また、例文は、テキストデータとして記録された文だけでなく、例えば、音声データとして記録された文を含んでもよい。更に、例文は、自然言語処理の研究目的に限らず、如何なる目的で集められたものでもよい。以下では、例文集合としてコーパスを例にとって説明する。

【 0 0 2 1 】

更に、特定の例文集合の対象単語を除く部分とは、特定の例文集合に何らかの処理を行って対象単語が含まないようにされた部分のことをいう。この処理は、例えば、対象単語をマスクする処理でもよいし、対象単語を一時的に削除する処理でもよいが、以下では、前者の処理を例にとって説明する。

【 0 0 2 2 】

更に、単語の意味表現とは、単語の意味をベクトル化して表現したものをいう。但し、本実施の形態では、単語の意味表現により単語の意味の近さを計算できればよいので、単語の意味表現は、少なくとも単語の意味の近さを計算できる他の方法で表現したものであってもよい。

【 0 0 2 3 】

更にまた、差分に基づく順序とは、差分を用いて定められる順序をいう。差分に基づく順序は、例えば、差分が大きい順序でもよいし、差分が大きい順序を基本としつつ他の要素を加味した順序でもよい。ここで、他の要素は、他の複数のユーザ既知語を用いた場合の差分であってよい。例えば、特定のユーザ既知語のみを用いた場合の差分は小さいが、他の複数のユーザ既知語を用いた場合の差分の平均が大きい場合や分散が小さい場合に、順序を上げることが考えられる。或いは、他の要素は、対象単語の文法上の属性等であってもよい。以下では、差分に基づく順序として差分が大きい順序を用いた場合を例にとって説明する。

【 0 0 2 4 】

[ 質問生成装置のハードウェア構成 ]

図 1 は、本実施の形態における質問生成装置 1 0 のハードウェア構成例を示した図である。図示するように、質問生成装置 1 0 は、演算手段であるプロセッサ 1 1 と、記憶手段であるメインメモリ 1 2 及び H D D ( Hard Disk Drive ) 1 3 とを備える。ここで、プロセッサ 1 1 は、O S ( Operating System ) やアプリケーション等の各種ソフトウェアを実行し、後述する各機能を実現する。また、メインメモリ 1 2 は、各種ソフトウェアや

10

20

30

40

50

その実行に用いるデータ等を記憶する記憶領域であり、HDD13は、各種ソフトウェアに対する入力データや各種ソフトウェアからの出力データ等を記憶する記憶領域である。更に、質問生成装置10は、外部との通信を行うための通信I/F（以下、「I/F」と表記する）14と、ディスプレイ等の表示デバイス15と、キーボードやマウス等の入力デバイス16とを備える。

#### 【0025】

[質問生成装置の機能構成]

図2は、本実施の形態における質問生成装置10の機能構成例を示したブロック図である。図示するように、質問生成装置10は、コーパス記憶部21と、第1学習部22と、第1学習済みモデル記憶部23と、第1出力部24と、第1出力情報記憶部25とを備えている。また、質問生成装置10は、マスキング処理部31を備えている。更に、質問生成装置10は、マスクコーパス記憶部41と、第2学習部42と、第2学習済みモデル記憶部43と、第2出力部44と、第2出力情報記憶部45とを備えている。更にまた、質問生成装置10は、出力差分算出部51と、出力差分情報記憶部52と、ランキング処理部53と、質問単語記憶部54とを備えている。

#### 【0026】

コーパス記憶部21は、コーパスを記憶する。コーパスは、例えば、質問を行う分野における特定のコーパスである。コーパス記憶部21に記憶されたコーパスの具体例については後述する。

#### 【0027】

第1学習部22は、コーパス記憶部21に記憶されたコーパスを用いて単語の意味表現をモデルに学習させることにより第1学習済みモデルを生成する。本実施の形態では、特定の例文集合を用いて単語の意味表現を学習させた第1のモデルの一例として、第1学習済みモデルを用いている。ここで、第1学習部22は、コーパス記憶部21に記憶されたコーパスを用いて、全く学習していないモデルを学習させることにより、第1学習済みモデルを生成してもよい。或いは、第1学習部22は、コーパス記憶部21に記憶されたコーパスを用いて、既に学習したモデルを更新することにより、第1学習済みモデルを生成してもよい。

#### 【0028】

第1学習済みモデル記憶部23は、第1学習部22が生成した第1学習済みモデルを記憶する。第1学習済みモデル記憶部23に記憶された第1学習済みモデルの具体例については後述する。

#### 【0029】

第1出力部24は、第1学習済みモデル記憶部23に記憶された第1学習済みモデルから得られた特定のユーザ既知語の意味表現を第1出力情報として出力する。本実施の形態では、第1のモデルから得られた特定のユーザ既知語の意味表現の一例として、第1出力情報を用いている。

#### 【0030】

第1出力情報記憶部25は、第1出力部24が出力した第1出力情報を記憶する。第1出力情報記憶部25に記憶された第1出力情報の具体例については後述する。

#### 【0031】

マスキング処理部31は、コーパス記憶部21に記憶されたコーパスに対し、特定のユーザ既知語に対する寄与を調べたい対象の単語（以下、「調査対象単語」という）をマスクするマスキング処理を行うことにより、マスクコーパスを作成する。本実施の形態では、対象単語の一例として、調査対象単語を用いており、特定の例文集合の対象単語を除く部分の一例として、マスクコーパスを用いている。

#### 【0032】

マスクコーパス記憶部41は、マスキング処理部31が作成したマスクコーパスを記憶する。マスクコーパス記憶部41に記憶されたマスクコーパスの具体例については後述する。

。

10

20

30

40

50

## 【0033】

第2学習部42は、マスクコーパス記憶部41に記憶されたマスクコーパスを用いて単語の意味表現をモデルに学習させることにより第2学習済みモデルを生成する。本実施の形態では、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第2のモデルの一例として、第2学習済みモデルを用いている。ここで、第2学習部42は、マスクコーパス記憶部41に記憶されたマスクコーパスを用いて、全く学習していないモデルを学習させることにより、第2学習済みモデルを生成してもよい。この場合、第2学習済みモデルは、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を未学習モデルに新たに学習させたモデルの一例である。或いは、第2学習部42は、マスクコーパス記憶部41に記憶されたマスクコーパスを用いて、既に学習したモデルを更新することにより、第2学習済みモデルを生成してもよい。この場合、第2学習済みモデルは、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習済みモデルに更に学習させたモデルの一例である。

10

## 【0034】

第2学習済みモデル記憶部43は、第2学習部42が取得した第2学習済みモデルを記憶する。第2学習済みモデル記憶部43に記憶された第2学習済みモデルの具体例については後述する。

## 【0035】

第2出力部44は、第2学習済みモデル記憶部43に記憶された第2学習済みモデルから得られた特定のユーザ既知語の意味表現を第2出力情報として出力する。本実施の形態では、第2のモデルから得られた特定のユーザ既知語の意味表現の一例として、第2出力情報を用いている。

20

## 【0036】

第2出力情報記憶部45は、第2出力部44が出力した第2出力情報を記憶する。第2出力情報記憶部45に記憶された第2出力情報の具体例については後述する。

## 【0037】

出力差分算出部51は、複数の調査対象単語のそれぞれについて、第1出力情報記憶部25に記憶された第1出力情報と、その調査対象単語を選択した場合に第2出力情報記憶部45に記憶された第2出力情報との差分である出力差分を算出する。本実施の形態では、第1のモデルから得られた特定のユーザ既知語の意味表現と、第2のモデルから得られた特定のユーザ既知語の意味表現との差分を算出する手段の一例として、出力差分算出部51を設けている。また、本実施の形態では、複数の対象単語の各対象単語について、第1のモデルから得られた特定のユーザ既知語の意味表現と、第2のモデルから得られた特定のユーザ既知語の意味表現との差分を算出することにより、複数の差分を算出する手段の一例としても、出力差分算出部51を設けている。

30

## 【0038】

出力差分情報記憶部52は、複数の調査対象単語のそれぞれについて、その調査対象単語と、その調査対象単語を選択した場合に出力差分算出部51が算出した出力差分とを関連付けた出力差分情報を記憶する。

## 【0039】

ランキング処理部53は、複数の調査対象単語を、ユーザに与える質問で用いる単語（以下、「質問単語」という）として、出力差分情報記憶部52に記憶された出力差分が大きい順、つまり、ユーザ既知語である可能性が高い順に並べて出力する。これは、コーパス内に調査対象単語がある場合とない場合とで特定のユーザ既知語の意味表現が大きくずれるのであれば、調査対象単語がないとその特定のユーザ既知語の意味表現が得られないと考えられるので、調査対象単語はユーザ既知語と判断できる、という考え方に基づくものである。本実施の形態では、差分に基づいて、対象単語がユーザ既知語である可能性に関する情報を出力する手段の一例として、ランキング処理部53を設けている。また、本実施の形態では、複数の差分に基づいて、複数の対象単語を用いた質問を生成する手段の一例としても、ランキング処理部53を設けている。

40

50

## 【 0 0 4 0 】

質問単語記憶部 5 4 は、ランキング処理部 5 3 が出力した質問単語を、ランキング処理部 5 3 が並べた順序で記憶する。そして、タスクを実行するシステムが、質問単語記憶部 5 4 に記憶された質問単語を、質問単語記憶部 5 4 に記憶された順序で取り出して、ユーザに与える質問で用いることになる。

## 【 0 0 4 1 】

尚、これらの機能部は、ソフトウェアとハードウェア資源とが協働することにより実現される。具体的には、これらの機能部は、プロセッサ 1 1 が、これらを実現するプログラムを例えば HDD 1 3 からメインメモリ 1 2 に読み込んで実行することにより実現される。

## 【 0 0 4 2 】

次に、本実施の形態における質問生成装置 1 0 で記憶されるコーパスの具体例について説明する。

## 【 0 0 4 3 】

図 3 ( a ) は、コーパス記憶部 2 1 に記憶されるコーパスの具体例を示した図である。図示するように、コーパス記憶部 2 1 に記憶されるコーパスは、文書 2 1 1 , 2 1 2 , 2 1 3 , ... を含んでいる。そして、文書 2 1 1 は、文 2 1 1 1 , 2 1 1 2 , 2 1 1 3 , ... を含み、文書 2 1 2 は、文 2 1 2 1 , 2 1 2 2 , 2 1 2 3 , ... を含み、文書 2 1 3 は、文 2 1 3 1 , 2 1 3 2 , 2 1 3 3 , ... を含んでいる。ここで、ユーザ既知語  $n_1$  ,  $n_2$  ,  $n_3$  は、それぞれ、文 2 1 1 1 , 2 1 1 3 , 2 1 3 2 に存在するものとする。

## 【 0 0 4 4 】

図 3 ( b ) は、マスクコーパス記憶部 4 1 に記憶されるマスクコーパスの具体例を示した図である。図示するように、マスクコーパス記憶部 4 1 に記憶されるマスクコーパスは、コーパス記憶部 2 1 に記憶されるコーパスにおいて調査対象単語がマスクされたものになっている。ここでは、調査対象単語  $m_1$  ,  $m_2$  ,  $m_3$  が、それぞれ、文 4 1 1 1 , 4 1 2 3 , 4 1 3 2 に存在し、これらがマスクされているものとする。

## 【 0 0 4 5 】

ところで、図 3 ( a ) , ( b ) では、マスクコーパス記憶部 4 1 に記憶されるデータの単位を文としたが、これには限らない。データの単位は、より一般化し、文書の構成要素としてよい。文書の構成要素には、文以外に、段落、章、節等が含まれる。

## 【 0 0 4 6 】

また、図 3 ( b ) では、ユーザ既知語しか含まない文や、ユーザ既知語及び調査対象単語の何れも含まない文も、マスクコーパス記憶部 4 1 に記憶したが、これには限らない。ユーザ既知語しか含まない文や、ユーザ既知語及び調査対象単語の何れも含まない文は、マスクコーパス記憶部 4 1 に記憶しないようにしてもよい。

## 【 0 0 4 7 】

具体的には、第 2 学習部 4 2 が、全く学習していないモデルを学習させる場合は、ユーザ既知語及び調査対象単語の何れかを含む文のみをフィルタリングして、マスクコーパス記憶部 4 1 に記憶するとよい。つまり、図 3 ( b ) の例で言えば、文 4 1 1 1 , 4 1 1 3 , 4 1 2 3 , 4 1 3 2 をマスクコーパス記憶部 4 1 に記憶するとよい。これは、特定の例文集合の対象単語を除く部分が、特定の例文集合の特定のユーザ既知語及び対象単語の少なくとも何れか一方を含む構成要素の対象単語を除く部分である場合の一例である。

## 【 0 0 4 8 】

一方、第 2 学習部 4 2 が、既に学習したモデルを更新する場合は、調査対象単語を含む文のみをフィルタリングして、マスクコーパス記憶部 4 1 に記憶するとよい。つまり、図 3 ( b ) の例で言えば、文 4 1 1 1 , 4 1 2 3 , 4 1 3 2 をマスクコーパス記憶部 4 1 に記憶するとよい。更新前の学習済みモデルにユーザ既知語が含まれていると仮定できるからである。これは、特定の例文集合の対象単語を除く部分が、特定の例文集合の対象単語を含む構成要素の対象単語を除く部分である場合の一例である。

## 【 0 0 4 9 】

次に、本実施の形態における質問生成装置 1 0 で記憶される学習済みモデルの具体例につ

10

20

30

40

50

いて説明する。尚、以下では、Word2Vecを構成する2種類のモデルのうちCBOW (Continuous Bag-Of-Words) モデルにより単語の意味表現を学習させる場合を例にとって説明する。

【0050】

図4(a)は、第1学習済みモデル記憶部23に記憶される第1学習済みモデルの具体例を示した図である。ここでは、コーパス $X$ を入力としたCBOWモデルの出力である第1学習済みモデルを $Y$ と表記する。第1学習済みモデル $Y$ は、単語の意味表現を各行に持つ $V \times W$ の行列である。 $V$ は単語の数であり、 $W$ は意味表現の次元数である。以下、第1学習済みモデル $Y$ の単語 $v$ の行における次元 $w$ の意味表現を $Y_v(w)$ と表すことにする。図において、第1学習済みモデル $Y$ の1行目は、単語 $v_1$ の次元1, 2, 3, ...の意味表現を表している。また、2行目は、単語 $v_2$ の次元1, 2, 3, ...の意味表現を表し、3行目は、単語 $v_3$ の次元1, 2, 3, ...における意味表現を表している。

10

【0051】

図4(b)は、第2学習済みモデル記憶部43に記憶される第2学習済みモデルの具体例を示した図である。ここでは、マスキング処理部31が調査対象単語 $m_j$ をマスキングしたコーパス $X$ をコーパス $X_{m_j}$ とし、このコーパス $X_{m_j}$ を入力としたCBOWモデルの出力である第2学習済みモデルを $Y_{m_j}$ と表記する。第2学習済みモデル $Y_{m_j}$ も、単語の意味表現を各行に持つ $V \times W$ の行列である。以下、第2学習済みモデル $Y_{m_j}$ の単語 $v$ の行における次元 $w$ の意味表現を $Y_{v_{m_j}}(w)$ と表すことにする。図において、第2学習済みモデル $Y_{m_j}$ の1行目は、単語 $v_1$ の次元1, 2, 3, ...の意味表現を表している。また、2行目は、単語 $v_2$ の次元1, 2, 3, ...の意味表現を表し、3行目は、単語 $v_3$ の次元1, 2, 3, ...の意味表現を表している。

20

【0052】

次に、本実施の形態における質問生成装置10で記憶される出力情報の具体例について説明する。

【0053】

図5(a)は、第1出力情報記憶部25に記憶される第1出力情報の具体例を示した図である。図示するように、第1出力情報は、第1学習済みモデル $Y$ からユーザ既知語 $n_i$ に対応する行を抜き出したものである。ここでは、この抜き出された行である第1出力情報を $Y_{n_i}$ と表記する。第1出力情報 $Y_{n_i}$ は、単語の意味表現を要素に持つ $W$ 次元のベクトルである。

30

【0054】

図5(b)は、第2出力情報記憶部45に記憶される第2出力情報の具体例を示した図である。図示するように、第2出力情報は、第2学習済みモデル $Y_{m_j}$ からユーザ既知語 $n_i$ に対応する行を抜き出したものである。ここでは、この抜き出された行である第2出力情報を $Y_{n_i m_j}$ と表記する。第2出力情報 $Y_{n_i m_j}$ は、単語の意味表現を要素に持つ $W$ 次元のベクトルである。

【0055】

次に、本実施の形態における質問生成装置10で記憶される出力差分情報の具体例について説明する。

40

【0056】

図6は、出力差分情報記憶部52に記憶される出力差分情報の具体例を示した図である。図示するように、出力差分情報は、調査対象単語と、出力差分とを対応付けたものである。調査対象単語は $m_j$ であり、出力差分は $(n_i, m_j)$ である( $j = 1, 2, 3, \dots$ )。ここで、出力差分 $(n_i, m_j)$ は、第1出力情報 $Y_{n_i}$ と、調査対象単語 $m_j$ をマスクした場合の第2出力情報 $Y_{m_j n_i}$ との二乗距離として定義される。

【0057】

尚、その後、ランキング処理部53が、調査対象単語 $m_j$ を、出力差分 $(n_i, m_j)$ の大きい順に並べ替えて、質問単語記憶部54に記憶することになる。

【0058】

50

[ 質問生成装置の動作 ]

図 7 は、本実施の形態における質問生成装置 10 の動作例を示したフローチャートである。

【 0059 】

図示するように、質問生成装置 10 では、まず、第 1 学習部 22 が、コーパス記憶部 21 に記憶されたコーパスを用いて単語の意味表現を学習させて第 1 学習済みモデルを生成する（ステップ 101）。この第 1 学習済みモデルは、第 1 学習済みモデル記憶部 23 に記憶される。

【 0060 】

次に、第 1 出力部 24 が、第 1 学習済みモデル記憶部 23 に記憶された第 1 学習済みモデルからユーザ既知語の意味表現を抜き出して第 1 出力情報として出力する（ステップ 102）。この第 1 出力情報は、第 1 出力情報記憶部 25 に記憶される。

【 0061 】

一方、質問生成装置 10 では、マスキング処理部 31 が、コーパス記憶部 21 に記憶されたコーパスに対して調査対象単語をマスクするマスキング処理を行ってマスクコーパスを生成する（ステップ 103）。このマスクコーパスは、マスクコーパス記憶部 41 に記憶される。

【 0062 】

次に、第 2 学習部 42 が、マスクコーパス記憶部 41 に記憶されたコーパスを用いて単語の意味表現を学習させて第 2 学習済みモデルを生成する（ステップ 104）。この第 2 学習済みモデルは、第 2 学習済みモデル記憶部 43 に記憶される。

【 0063 】

次に、第 2 学習部 42 が、第 2 学習済みモデル記憶部 43 に記憶された第 2 学習済みモデルからユーザ既知語の意味表現を抜き出して第 2 出力情報として出力する（ステップ 105）。この第 2 出力情報は、第 2 出力情報記憶部 45 に記憶される。

【 0064 】

次いで、質問生成装置 10 では、第 1 出力情報記憶部 25 に記憶された第 1 出力情報と第 2 出力情報記憶部 45 に記憶された第 2 出力情報との出力差分を算出し、調査対象単語と関連付けて、出力差分情報として出力する（ステップ 106）。この出力差分情報は、出力差分情報記憶部 52 に記憶される。

【 0065 】

その後、質問生成装置 10 は、調査対象単語が終了したかどうかを判定する（ステップ 107）。つまり、着目すべき調査対象単語がなくなったかどうかを判定する。

【 0066 】

その結果、調査対象単語が終了していないと判定すれば、質問生成装置 10 は、処理をステップ 103 へ戻す。そして、他の調査対象単語に着目し、ステップ 103 ~ 106 の処理を行う。

【 0067 】

一方、調査対象単語が終了したと判定すれば、質問生成装置 10 は、処理をステップ 108 へ進める。

【 0068 】

そして、ランキング処理部 53 が、調査対象単語を出力差分が大きい順に並べ替えて、質問順に並べられた質問単語として出力する（ステップ 108）。この質問単語は、質問単語記憶部 54 に記憶される。

【 0069 】

[ 変形例 ]

上記実施の形態では言及しなかったが、システムは、ユーザから質問に対する回答が得られた時点で、新たなユーザ既知語を特定し、コーパス記憶部 21 に記憶されたコーパスにこれを反映させてもよい。ここで、新たなユーザ既知語は、ユーザがタスク中でその単語の意味を知っているかを明示的にシステムに伝えることで、特定されるようにするとよい

10

20

30

40

50

。これにより、質問生成装置 10 では、出力差分算出部 51 が、この新たなユーザ既知語が反映されたコーパスを用いて新たに出力差分情報を生成することにより、ユーザ既知語を再度予測するようにしてよい。そして、ランキング処理部 53 が、質問に用いる単語の順序をリアルタイムに更新してよい。この場合、出力差分算出部 51 は、特定のユーザ既知語に代えて、質問に対するユーザの回答から把握される他のユーザ既知語を用いて、複数の差分を算出する手段の一例であり、ランキング処理部 53 は、複数の差分に基づいて、複数の対象単語を用いた質問を再生成する手段の一例である。

【0070】

[プロセッサ]

本実施の形態において、プロセッサとは広義的なプロセッサを指し、汎用的なプロセッサ（例えば CPU : Central Processing Unit 等）や、専用のプロセッサ（例えば GPU : Graphics Processing Unit、ASIC : Application Specific Integrated Circuit、FPGA : Field Programmable Gate Array、プログラマブル論理デバイス等）を含むものである。

10

【0071】

また、本実施の形態におけるプロセッサの動作は、1つのプロセッサによって成すのみでなく、物理的に離れた位置に存在する複数のプロセッサが協働して成すものであってもよい。また、プロセッサの各動作の順序は、本実施の形態において記載した順序のみに限定されるものではなく、変更してもよい。

【0072】

20

[プログラム]

本実施の形態における質問生成装置 10 が行う処理は、例えば、アプリケーションソフトウェア等のプログラムとして用意される。

【0073】

即ち、本実施の形態を実現するプログラムは、コンピュータに、特定の例文集合を用いて単語の意味表現を学習させた第1のモデルから得られた特定のユーザ既知語の意味表現と、特定の例文集合の対象単語を除く部分を用いて単語の意味表現を学習させた第2のモデルから得られた特定のユーザ既知語の意味表現との差分を算出する機能と、差分に基づいて、対象単語がユーザ既知語である可能性に関する情報を出力する機能とを実現させるためのプログラムとして捉えられる。

30

【0074】

尚、本実施の形態を実現するプログラムは、通信手段により提供することはもちろん、CD-ROM等の記録媒体に格納して提供することも可能である。

【符号の説明】

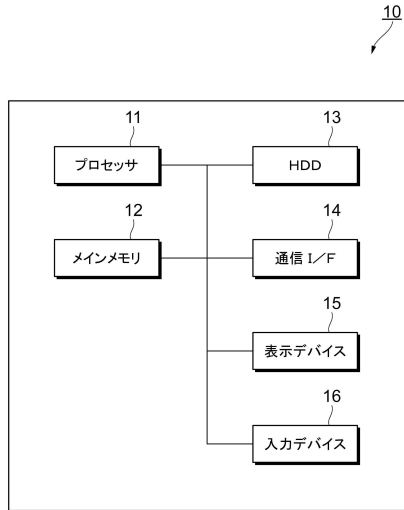
【0075】

10 ... 質問生成装置、21 ... コーパス記憶部、22 ... 第1学習部、23 ... 第1学習済みモデル記憶部、24 ... 第1出力部、25 ... 第1出力情報記憶部、31 ... マスキング処理部、41 ... マスクコーパス記憶部、42 ... 第2学習部、43 ... 第2学習済みモデル記憶部、44 ... 第2出力部、45 ... 第2出力情報記憶部、51 ... 出力差分算出部、52 ... 出力差分情報記憶部、53 ... ランキング処理部、54 ... 質問単語記憶部

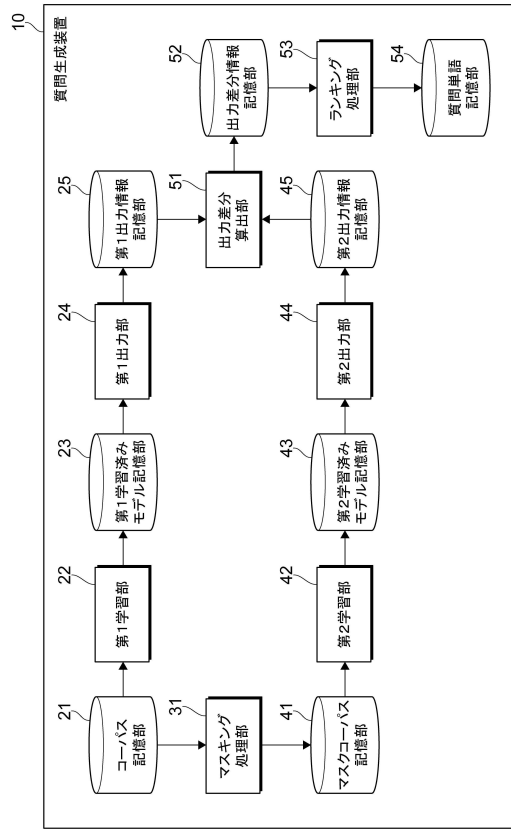
40

【 図 面 】

【 図 1 】



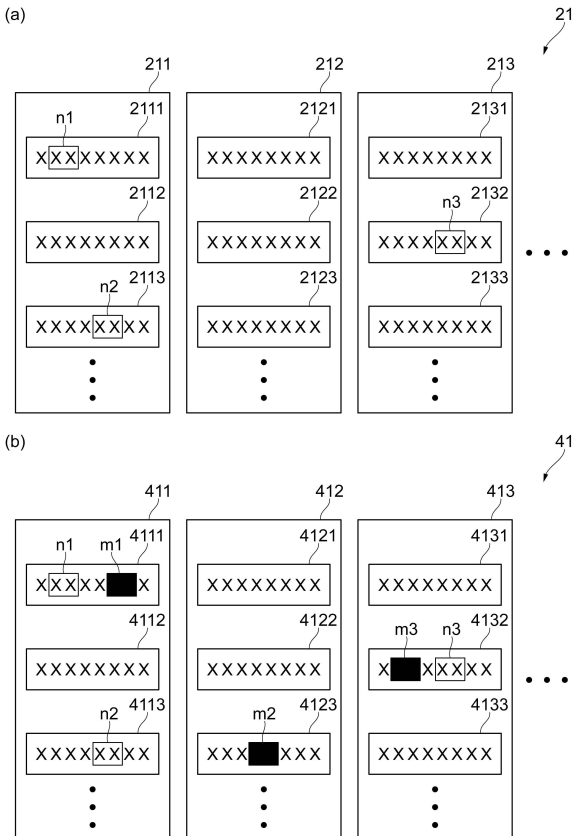
【 図 2 】



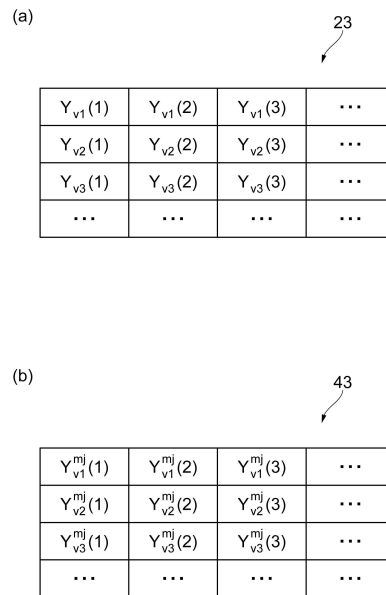
10

20

【 図 3 】



【 図 4 】

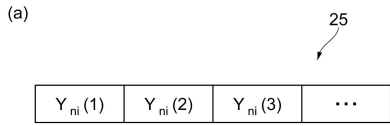


30

40

50

【 図 5 】

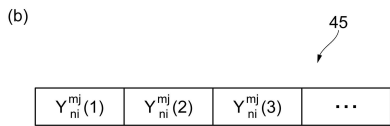


【 図 6 】

52

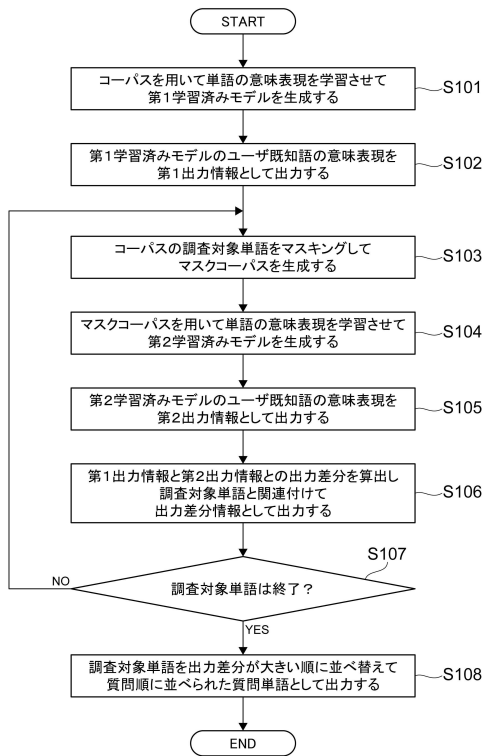
調査対象単語	出力差分
m1	$\delta(ni,m1)$
m2	$\delta(ni,m2)$
m3	$\delta(ni,m3)$
...	...

10



20

【 図 7 】



30

40

50