

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5708495号  
(P5708495)

(45) 発行日 平成27年4月30日(2015.4.30)

(24) 登録日 平成27年3月13日(2015.3.13)

(51) Int.Cl.		F I			
<b>G06F 17/30</b>	<b>(2006.01)</b>	G06F	17/30	320D	
<b>G06F 17/27</b>	<b>(2006.01)</b>	G06F	17/30	210D	
		G06F	17/27	635	

請求項の数 11 (全 23 頁)

(21) 出願番号	特願2011-545194 (P2011-545194)	(73) 特許権者	000004237
(86) (22) 出願日	平成22年12月3日 (2010.12.3)		日本電気株式会社
(86) 国際出願番号	PCT/JP2010/071696		東京都港区芝五丁目7番1号
(87) 国際公開番号	W02011/070980	(74) 代理人	100095407
(87) 国際公開日	平成23年6月16日 (2011.6.16)		弁理士 木村 満
審査請求日	平成25年11月8日 (2013.11.8)	(72) 発明者	水口 弘紀
(31) 優先権主張番号	特願2009-282304 (P2009-282304)		日本国東京都港区芝五丁目7番1号 日本電気株式会社内
(32) 優先日	平成21年12月11日 (2009.12.11)	(72) 発明者	久寿居 大
(33) 優先権主張国	日本国 (JP)		日本国東京都港区芝五丁目7番1号 日本電気株式会社内
		(72) 発明者	楠村 幸貴
			日本国東京都港区芝五丁目7番1号 日本電気株式会社内

最終頁に続く

(54) 【発明の名称】 辞書作成装置、単語収集方法、及び、プログラム

(57) 【特許請求の範囲】

【請求項1】

単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集する辞書増殖処理における、入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録手段と、

前記入出力過程記録手段に記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類手段と、

前記入出力過程記録手段に記録された情報を参照し、前記クラスタ分類手段が分類したクラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別手段と、

前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、  
を関連付けて出力する収集単語出力手段と、

を備えることを特徴とする辞書作成装置。

【請求項2】

単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集する辞書増殖手段をさらに備える、

ことを特徴とする請求項 1 に記載の辞書作成装置。

【請求項 3】

前記入出力過程記録手段は、複数回の入出力を繰り返した、入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する、

ことを特徴とする請求項 1 又は 2 に記載の辞書作成装置。

【請求項 4】

前記クラスタ分類手段は、前記入出力過程記録手段に記録されている情報から、前記辞書増殖処理で収集した単語のうち共通の単語を入力にする単語同士、又は共通の単語を出力する単語同士ほどその値が大きくなる値を示す単語間の結束度を算出し、算出した結束度に基づいて、単語をクラスタに分類する、

ことを特徴とする請求項 1 乃至 3 の何れか 1 項に記載の辞書作成装置。

【請求項 5】

前記同種判別手段は、前記入出力過程記録手段に記録されている情報に基づいて、クラスタ毎に、最初に入力を受け付けた入力単語から当該クラスタ内の各単語を出力するまでに要したターン数、及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数を算出し、算出したターン数の平均値を用いて、当該クラスタ内の単語が最初に入力を受け付けた入力単語と同種であるか異種であるかの判別をする

ことを特徴とする請求項 1 乃至 4 の何れか 1 項に記載の辞書作成装置。

【請求項 6】

前記辞書増殖処理で収集された単語を種類毎に、複数の単語グループに分類して記憶する、単語グループ記憶手段と、

所定の条件を満たす一の単語グループのなかから所定数の単語を選択する単語選択手段と、をさらに備え、

前記単語選択手段が選択した単語を入力単語とした前記辞書増殖処理を実行し、

前記同種判別手段は、前記入出力過程記録手段に記録された情報に基づいて、前記クラスタ分類手段が分類したクラスタ毎に、該クラスタ内の単語が前記単語選択手段が選択した入力単語と同じ種類の単語であるか否かを判別する、

ことを特徴とする請求項 1 乃至 5 の何れか 1 項に記載の辞書作成装置。

【請求項 7】

前記同種判別手段が判別した結果に基づいて、前記辞書増殖処理で収集された単語を前記単語グループ記憶手段に登録し、登録した単語グループのうち所定の条件を満たす単語グループがある場合に、前記単語選択手段に単語の選択を指示する再実行手段をさらに備え、

前記再実行手段は、収集単語を前記単語グループ記憶手段に登録する際、収集単語の属するクラスタが前記単語選択手段が選択した単語と同種の単語である場合には当該選択した単語と同じ単語グループに当該収集単語を登録し、異種であり且つ既に前記単語グループ記憶手段に記憶されている単語である場合には該記憶されている単語と同じ単語グループに収集単語を登録し、異種であり且つ未だ前記単語グループ記憶手段が記憶していない単語である場合には収集単語を新規の単語グループに登録する、

ことを特徴とする請求項 6 に記載の辞書作成装置。

【請求項 8】

前記入出力過程記録手段に記録されている情報から算出された、前記辞書増殖処理で収集した単語のうち共通の単語を入力にする単語同士、又は共通の単語を出力する単語同士ほどその値が大きくなる値を示す単語間の結束度を記憶する結束度記憶手段をさらに備え、

10

20

30

40

50

前記単語選択手段は、前記一の単語グループ内の単語間の結束度に基づいて、所定数の単語を選択する、

ことを特徴とする請求項 6 又は 7 に記載の辞書作成装置。

【請求項 9】

前記単語選択手段は、結束度の大きい順に単語を選択する割合、又は、結束度の小さい順に単語を選択する割合、が少なくとも予め設定されている条件情報に基づいて、所定数の単語を選択する、

ことを特徴とする請求項 8 に記載の辞書作成装置。

【請求項 10】

コンピュータが、単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集した辞書増殖処理における入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録ステップと、

コンピュータが、前記入出力過程記録ステップに記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類ステップと、

コンピュータが、前記入出力過程記録ステップに記録された情報を参照し、前記クラスタ分類ステップが分類したクラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別ステップと、

コンピュータが、前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、を関連付けて出力する収集単語出力ステップと、

を備えることを特徴とする単語収集方法。

【請求項 11】

コンピュータを、

単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集する辞書増殖処理における、入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録手段、

前記入出力過程記録手段に記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類手段、

前記入出力過程記録手段に記録された情報を参照し、前記クラスタ分類手段が分類したクラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別手段、

前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、を関連付けて出力する収集単語出力手段、

として機能させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、辞書作成装置、単語収集方法、及び、プログラムに関する。

10

20

30

40

50

## 【背景技術】

## 【0002】

少数の同種の単語を入力に、文献データやWebページ等から多数の同種の単語を収集した辞書を作成する辞書作成の手法が知られている。なお、ここでいう辞書とは、共通の上位概念を持つ同種の単語の集合のことである。

## 【0003】

上述した辞書作成の手法の一例が、非特許文献1に記載されている。この辞書作成の手法の概略を以下に示す。

## 【0004】

まず、収集の元となる少量の単語を入力する。以下、はじめに入力した単語をシード単語と呼ぶ。次に、Web検索エンジンを利用し、シード単語を含むWebページを収集する。次に、収集したWebページからシード単語とそれ以外の語を区切るパターンを作成する。そして、このパターンを使ってWebページから単語を抽出し、シード単語に追加する。なお、シード単語を入力してから単語が抽出されるまでをターンと呼ぶ。そして、単語が追加されたシード単語を用いて、さらにWebページを収集する。これを数ターン繰り返した後、抽出された単語をシード単語と同種の単語の集合(辞書)として出力する。

10

## 【0005】

このような辞書作成の手法では、新たにシード単語に追加される単語が、シード単語と異なる種類の単語である場合がある。例えば、レストラン名のシード単語を入力して、レストラン名の辞書を作成する際に、同じ文献に掲載されており、かつ、パターンが似ているラーメン店名やうどん店名などの単語が、新たにシード単語に追加されてしまう等の場合である。

20

このような場合、その異なる種類の単語から、さらに異なる種類の単語が次々にシード単語に追加されてしまい、シード単語と異なる種類の単語が多く収集されてしまい、辞書の精度が悪化することが知られている。

## 【0006】

このような事態を回避するために、各ターンで抽出される単語の信頼度を求め、特定の信頼度以上の単語のみをシード単語に追加して、次のターンで採用することが行われている。なお、この信頼度は、例えば、パターンの出現回数に基づく統計量や、パターンから検出された単語数に基づく統計量等が用いられる。非特許文献1では、信頼度として、単語のパターンによって抽出できたWebページの数を採用しており、抽出できたWebページの数が所定数以下の単語はシード単語に追加しないことで、異なる種類の単語が収集されることを防止している。

30

## 【先行技術文献】

## 【非特許文献】

## 【0007】

【非特許文献1】水口弘紀、河合英紀、土田正明、久寿居大、Web知識を利用したブートストラップによる辞書増殖手法、DEWS2007、2007

## 【発明の概要】

40

## 【発明が解決しようとする課題】

## 【0008】

上述した信頼度を用いて辞書を作成をした場合、信頼度が低い異なる種類の単語(異種単語)は、収集対象から除外されてシードに追加されない。従って、ユーザは、どのような異種単語がシード単語から収集されているのかを全く知ることができず、異種単語を再利用してさらに異なるグループの単語を収集するような活用ができない。

## 【0009】

本発明は、上記実情に鑑みてなされたものであり、どのような異種単語が収集されているのかをユーザに好適に出力することを可能にした辞書作成装置、単語収集方法、及び、プログラムを提供することを目的とする。

50

## 【課題を解決するための手段】

【0010】

上記目的を達成するため、本発明の第1の観点に係る辞書作成装置は、

単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集する辞書増殖処理における、入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録手段と、

前記入出力過程記録手段に記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類手段と、

10

前記入出力過程記録手段に記録された情報を参照し、前記クラスタ分類手段が分類したクラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別手段と、

前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、を関連付けて出力する収集単語出力手段と、

を備えることを特徴とする。

20

また、本発明の第2の観点に係る単語収集方法は、

コンピュータが、単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集した辞書増殖処理における入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録ステップと、

コンピュータが、前記入出力過程記録ステップに記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類ステップと、

コンピュータが、前記入出力過程記録ステップに記録された情報を参照し、前記クラスタ分類ステップが分類したクラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別ステップと、

30

コンピュータが、前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、を関連付けて出力する収集単語出力ステップと、

を備えることを特徴とする。

40

また、本発明の第3の観点に係るプログラムは、

コンピュータを、

単語の入力を受け付け、入力された入力単語に関連する単語を文書データから出力し、以降は所定の条件に達するまで出力した単語を前記入力単語に追加し、該入力単語に関連する単語を文書データから出力することを繰り返していくことで単語を収集する辞書増殖処理における、入力単語と該入力単語によって出力された出力単語との入出力の過程を示す情報を記録する入出力過程記録手段、

前記入出力過程記録手段に記録された情報に基づいて、前記入出力の過程における前記入力単語及び前記出力単語の類似度を用いて、前記辞書増殖処理で収集された単語を複数のクラスタに分類するクラスタ分類手段、

前記入出力過程記録手段に記録された情報を参照し、前記クラスタ分類手段が分類した

50

クラスタ毎に、該クラスタ内の単語が最初に入力を受け付けた入力単語からクラスタ内の各単語を出力するまでに要したターン数及び当該クラスタ内の各単語が最初に入力を受け付けた入力単語を出力するまでに要したターン数に基づいて、クラスタ内の単語が入力単語と同じ種類の単語であるか否かを判別する同種判別手段、

前記辞書増殖処理で収集された単語と、該単語が属するクラスタと、該クラスタを構成する単語が最初に入力を受け付けた入力単語と同じ種類の単語であるか否かを示す情報と、を関連付けて出力する収集単語出力手段、

として機能させるプログラムである。

【発明の効果】

【0011】

本発明によれば、辞書構築において収集された単語をクラスタリングし、クラスタ毎に、最初に入力した単語と同じ種類の単語であるか否かが判別される。従って、どのような異種単語が収集されているのかをユーザに好適に出力することができる。

【図面の簡単な説明】

【0012】

【図1】本発明の第1実施形態に係る辞書作成装置の構成を示す図である。

【図2】収集過程記憶部に記憶される情報の構成例を示す図である。

【図3】収集単語記憶部に記憶される情報の構成例を示す図である。

【図4】辞書作成処理の動作を説明するためのフローチャートである。

【図5】辞書増殖処理の動作を説明するためのフローチャートである。

【図6】クラスタリング処理の動作を説明するためのフローチャートである。

【図7】単語間の入出力の関係を示したグラフである。

【図8】同種判別処理の動作を説明するためのフローチャートである。

【図9】本発明の第2実施形態に係る辞書作成装置の構成を示す図である。

【図10】図10(A)及び図10(B)は、単語グループ記憶部に記憶される情報の構成例を示す図である。

【図11】辞書作成処理の動作を説明するためのフローチャートである。

【図12】単語グループ更新処理の動作を説明するためのフローチャートである。

【図13】本発明の第3実施形態に係る辞書作成装置の構成を示す図である。

【図14】収集単語記憶部に記憶される情報の構成例を示す図である。

【図15】各実施形態に係る辞書作成装置をコンピュータに実装する場合の、物理的な構成の一例を示すブロック図である。

【発明を実施するための形態】

【0013】

以下、本発明の実施形態について、図面を参照しながら詳細に説明する。なお、本発明は下記の実施形態及び図面によって限定されるものではない。本発明の要旨を変更しない範囲で下記の実施形態及び図面に変更を加えることが出来るのはもちろんである。また、図中同一または相当部分には同一符号を付す。

また、本発明で辞書とは、共通の上位概念を持つ同種の単語の集合のことである。

【0014】

(第1実施形態)

本発明の第1実施形態に係る辞書作成装置100について説明する。辞書作成装置100は、図1に示すように、入力部101と、辞書増殖部102と、クラスタリング部103と、種別判別部104と、出力部105と、文書記憶部106と、収集過程記憶部107と、収集単語記憶部108とを備える。

【0015】

入力部101は、キーボードやマウスなどから構成される。ユーザは、入力部101を介して、辞書(同種の単語の集合)を作成するためのサンプルとなる単語(シード単語)を入力する。

【0016】

10

20

30

40

50

辞書増殖部102は、非特許文献1に記載されているような従来の手法を用いて、シード単語と同種の単語を文書記憶部106に記憶されている文書内から収集する辞書増殖処理を行う。また、辞書増殖部102は、この辞書増殖処理において、どのような過程を経て単語が収集されたのかを示す情報を収集過程記憶部107に記録する。辞書増殖部102の行う辞書増殖処理の詳細については後述する。

【0017】

クラスタリング部103は、収集過程記憶部107に記憶されている情報に基づいて、辞書増殖部102が収集した単語を複数のクラスタに分類(クラスタリング)する。クラスタリング部103の行う処理の詳細については後述する。

【0018】

種別判別部104は、クラスタとそのクラスタに含まれる単語とを入力に、収集過程記憶部107に記憶されている情報を参照し、クラスタを構成する単語が、シード単語と同じ種類の単語であるか否かを判別する。種別判別部104の行う処理の詳細については後述する。

【0019】

出力部105は、種々の情報を出力する。例えば、出力部105は、辞書増殖処理によって収集された単語を、分類されたクラスタ毎に、シード単語と異種か同種かを示す情報を付して出力(表示)する。

【0020】

文書記憶部106は、辞書増殖部102による単語収集の対象となる各文書を定義するデータが記憶される。なお、各文書のデータにはID(文書ID)が付されている。

【0021】

収集過程記憶部107には、辞書増殖処理において、どのような入出力の過程を経て単語が収集されたのかを示す情報が記録される。具体的には、図2に示すように、収集過程記憶部107には、辞書増殖処理におけるターン毎に、当該ターンのターン数と、当該ターンで入力された入力単語と、該入力単語から生成されたパターンによって出力された出力単語とが対応付けられて記録される。

例えば、図2の先頭のエントリから、辞書増殖処理の1ターン目に、「レストランS」から作成されたパターンにより「レストランX」が抽出されたことがわかる。

【0022】

図1に戻り、収集単語記憶部108には、図3に示すように、収集された各単語と、各単語がどのクラスタに分類されているかを示すクラスタIDとが対応付けられて記憶される。また、各クラスタには、クラスタを構成する単語が、シード単語と同じ種類の単語であるのか(シード単語自体が当該クラスタに含まれる場合も同じ種類とする)、又は、異なる種類の単語であるのか、を示す情報が付与される。

例えば、図3から、「レストランA」と「レストランB」はクラスタ1に分類され、また、クラスタ1はシード単語と同じ種類の単語から構成されていることが分かる。同様に、「うどんC」と「うどんD」はクラスタ2に分類され、また、クラスタ2はシード単語と異なる種類の単語から構成されていることが分かる。

【0023】

続いて、辞書作成装置100で実施される処理の動作について説明する。

まず、ユーザは、入力部101を操作して、辞書(同種の単語の集合)を作成するためのサンプルとなる1乃至複数の単語(シード単語)を入力する。そして、入力したシード単語を元に、辞書を作成することを指示する。この指示操作に応じて、辞書作成装置100は、図4に示す辞書作成処理を行う。

【0024】

辞書作成処理が開始されると、まず、辞書増殖部102は、従来の手法で辞書増殖処理を行い、入力されたシード単語に関連する単語を収集する(ステップS100)。

【0025】

辞書増殖処理(ステップS100)の詳細について、図5のフローチャートを参照して説明

10

20

30

40

50

する。辞書増殖処理が開始されると、まず、辞書増殖部 102 は、ユーザによって入力されたシード単語を収集単語記憶部 108 に登録する（ステップS101）。そして、辞書増殖部 102 は、ターン数を示すカウンタ  $i$ （初期値 0）を 1 インクリメントする（ステップ S102）。

【 0 0 2 6 】

続いて、辞書増殖部 102 は、収集単語記憶部 108 に記憶されている単語のなかから所定数の単語をランダムに選択する（ステップS103）。そして、辞書増殖部 102 は、文書記憶部 106 に記憶されている文書のなかから、選択したシード単語が含まれている文書を検出する（ステップS104）。なお、ここでは、選択したシード単語を全て含む文書のみを検出してよいし、選択したシード単語のうち所定数のシード単語を含む文書を検出してよい。

10

【 0 0 2 7 】

続いて、辞書増殖部 102 は、検出した文書内における、ステップS103で選択したシード単語が出現する位置を特定し、シード単語とそれ以外の部分とを区切るパターンを作成する（ステップS105）。例えば、文書内でシード単語が出現する部分の前後の所定数の文字列を、パターンとして採用すればよい。

【 0 0 2 8 】

続いて、辞書増殖部 102 は、作成したパターンに合致する単語を、文書記憶部 106 に記憶されている文書から抽出する（ステップS106）。そして、辞書増殖部 102 は、抽出した単語を収集単語記憶部 108 に追加する（ステップS107）。

20

【 0 0 2 9 】

続いて、辞書増殖部 102 は、今回のターン数を示す情報（即ち、カウンタ  $i$  の値）と、ステップS103で選択した各単語（入力単語）と、入力単語から作成したパターンによりステップS106で抽出した単語（出力単語）とを対応付けて、収集過程記憶部 107 に記憶する（ステップS108）。

【 0 0 3 0 】

続いて、辞書増殖部 102 は、辞書増殖を終了させるための所定の終了条件を満たしているか否かを判別する（ステップS109）。終了条件としては、例えば、収集単語記憶部 108 に記憶した単語の数が所定数に達したか、又は、ターン数が所定数に達したか等の任意の条件を採用することが可能である。なお、後述するクラスタリング処理で収集した単語を適切にクラスタリングできるようにするために、ここでは、少なくとも 2 ターン以上は単語の収集を繰り返し実行するような終了条件を採用することが望ましい。

30

【 0 0 3 1 】

終了条件を満たしていないと判別した場合（ステップS109；No）、辞書増殖部 102 は、ステップS102～ステップS108を繰り返し、新たに単語が追加されたシード単語から単語を収集する処理を引き続き行う。

終了条件を満たしていると判別した場合（ステップS109；Yes）、辞書増殖部 102 は、辞書増殖処理を終了し処理をクラスタリング部 103 に移す。

【 0 0 3 2 】

図4に戻り、続いて、クラスタリング部 103 は、辞書増殖処理によって収集された単語をクラスタに分類するクラスタリング処理を行う（ステップS200）。

40

【 0 0 3 3 】

図6は、クラスタリング処理（ステップS200）の詳細を示すフローチャートである。クラスタリング処理が開始されると、まず、クラスタリング部 103 は、収集単語記憶部 108 から、未だ単語間の結束度を算出していない 2 つの単語を選択する（ステップS201）。

【 0 0 3 4 】

続いて、クラスタリング部 103 は、選択した 2 つの単語間の結束度を、収集過程記憶部 107 に記憶されている情報に基づいて算出する（ステップS202）。

【 0 0 3 5 】

50



なお、単語間の結束度とは、上述した辞書増殖処理において、共通の単語を入力にする単語同士、又は共通の単語を出力する単語同士ほど、その値が大きくなる指標のことである。例えば、2つの単語それぞれに入力される単語のうち共通の単語から2つの単語に入力される単語の割合と、2つの単語それぞれが出力する単語のうち2つの単語が共通の単語を出力する単語の割合と、の和を2つの単語間の結束度として算出することができる。

【0036】

より具体的には、2つの単語a,b間の結束度を $Sim(a,b)$ とすると、以下の式により、結束度を算出することができる。

$$Sim(a,b)=Sim\_in(a,b)+sim\_out(a,b)$$

【0037】

上式において、 $Sim\_in(a,b)$ は、単語a,bそれぞれに入力される単語のうち共通の単語から入力される単語の割合を示す値である。 $Sim\_in(a,b)=(\text{単語aと単語bの両方に入力される共通の単語の数}) / ((\text{単語aに入力される単語の数}) + (\text{単語bに入力される単語の数}))$ と求めることができる。

また、 $Sim\_out(a,b)$ は、2つの単語a,bそれぞれが出力する単語のうち共通の単語を出力する単語の割合を示す値である。 $Sim\_out(a,b)=(\text{単語aと単語bの両方から出力された共通の単語の数}) / ((\text{単語aが出力した単語の数}) + (\text{単語bが出力した単語の数}))$ と求めることができる。

【0038】

続いて、クラスタリング部103は、収集単語記憶部108に記憶されているシード単語の全ての組で、結束度を算出したか否かを判別する(ステップS203)。

【0039】

シード単語の全ての組で結束度を算出していない場合(ステップS203; No)、クラスタリング部103は、結束度未算出の2つのシード単語を選択して結束度を算出する処理(ステップS201、ステップS202)を繰り返す。

【0040】

シード単語の全ての組で結束度を算出した場合(ステップS203; Yes)、クラスタリング部103は、算出した結束度を類似度として、最短距離法、最長距離法、および、群平均法などの公知のクラスタリング手法を用いてクラスタリングを行い、収集単語記憶部108に記憶されているシード単語を複数のクラスタに分類する(ステップS204)。

そして、クラスタリング部103は、クラスタリングした結果を記録する(ステップS205)。具体的には、クラスタリング部103は、収集単語記憶部108に記憶されている単語に、クラスタに分類した結果が反映されるようにクラスタIDを付与する。以上でクラスタリング処理は終了する。

【0041】

このように、クラスタリング処理により、収集された単語間の結束度が算出され、算出された結束度に基づいて、収集単語が複数のクラスタに分類される。

【0042】

ここで、上述したクラスタリング処理について、具体例を挙げて説明する。図7は、図2に示すような情報が収集過程記憶部107に記憶されている場合の、辞書増殖処理のターン1からターン3の単語間の入出力の関係をグラフで示した図である。この図において、各単語はノードで表され、入力単語から出力単語の方向にアーク(矢印)で結ばれる。例えば、図7より、単語「レストランA」は、ターン2に「レストランX」と「レストランS」から作成されたパターンにより抽出されたことがわかる。また、ターン3では、単語「レストランA」から作成されたパターンにより「レストランE」と「レストランT」とが抽出されたことがわかる。

【0043】

ここで、「レストランA」と「レストランB」との間の結束度 $Sim(A,B)$ を算出する場合を考える。

「レストランA」に入力される単語は「レストランX」と「レストランS」であり、「

10

20

30

40

50

レストランB」に入力される単語は「レストランS」である。そして、このうち、「レストランS」が、「レストランA」と「レストランB」の両方に入力される。したがって、 $Sim_{in}(A,B)$ は、 $1/3$ となる。また、「レストランA」が出力する単語は「レストランE」と「レストランT」であり、「レストランB」が出力する単語は「レストランT」である。そして、このうち、「レストランT」が、「レストランA」と「レストランB」の両方から出力される。したがって、 $Sim_{out}(A,B)$ は、 $1/3$ となる。したがって、結束度 $Sim(A,B)=Sim_{in}(A,B)+Sim_{out}(A,B)=1/3+1/3=2/3$ と算出される。

【0044】

同様に、他の単語間の結束度についても、以下のように算出される。

レストランAとうどんCとの間の結束度： $Sim(A,C)=Sim_{in}(A,C)+Sim_{out}(A,C)=0+0=0$

レストランAとうどんDとの間の結束度： $Sim(A,D)=Sim_{in}(A,D)+Sim_{out}(A,D)=0+0=0$

レストランBとうどんCとの間の結束度： $Sim(B,C)=Sim_{in}(B,C)+Sim_{out}(B,C)=0+0=0$

レストランBとうどんDとの間の結束度： $Sim(B,D)=Sim_{in}(B,D)+Sim_{out}(B,D)=0+1/3=1/3$

うどんCとうどんDとの間の結束度： $Sim(C,D)=Sim_{in}(C,D)+Sim_{out}(C,D)=2/4+1/4=3/4$

【0045】

そして、これらの単語間の結束度を類似度として、公知のクラスタリングの手法を用いたクラスタリングがなされる。例えば、この結束度から、クラスタ1 { レストランA, レストランB }、クラスタ2 { うどんC, うどんD } の2つのクラスタが形成され、図3に示すように、収集単語記憶部108に記憶されている各単語に、クラスタIDが付与される。

【0046】

図4に戻り、続いて、種別判別部104は、クラスタリング処理で分類したクラスタが、最初に入力された単語(シード単語)と同種の単語から構成されるか否かを判別する同種判別処理を行う(ステップS300)。

【0047】

図8は、同種判別処理(ステップS300)の詳細を示すフローチャートである。同種判別処理が開始されると、まず、種別判別部104は、収集単語記憶部108から、同種判別を未だ行っていない1つのクラスタ、及び、当該クラスタに含まれる単語を選択する(ステップS301)。

【0048】

続いて、種別判別部104は、収集過程記憶部107を参照して、選択したクラスタ内の単語が、最初に入力された単語(シード単語)と同種の単語であるか否かを判別する(ステップS302)。なお、この判別は、クラスタ内の各単語のシード単語までの近さに基づいて行えばよい。

具体的には、種別判別部104は、シード単語からクラスタ内の各単語を出力するまでに要したターン数や、クラスタ内の各単語がシード単語を出力するまでに要したターン数を算出し、算出したターン数に基づいて、同種か異種かの判別をすればよい。

【0049】

続いて、種別判別部104は、判別結果を収集単語記憶部108に記録する(ステップS303)。

【0050】

続いて、種別判別部104は、収集単語記憶部108に記憶されているクラスタ全てで、上述の同種判別を実施したか否かを判別する(ステップS304)。

【0051】

同種判別未実施のクラスタがある場合(ステップS304; No)、種別判別部104は、そのクラスタを選択して同種判別をする処理(ステップS301~ステップS303)を繰り返す。

【0052】

同種判別を未実施のクラスタがない場合(ステップS304; Yes)、同種判別処理は終了

10

20

30

40

50

する。

【 0 0 5 3 】

このように、同種判別処理が実施されることにより、クラスタ毎に、クラスタを構成する単語がシード単語と同じ種類の単語であるか異なる種類の単語であるかが判別される。

【 0 0 5 4 】

続いて、上述した同種判別処理について、具体例を挙げて説明する。

前提として、図7に示すような入出力関係が、図2に示す収集過程記憶部 1 0 7 に記憶されている情報から得られているものとする。また、「レストラン A」と「レストラン B」がクラスタ 1、「うどん C」と「うどん D」がクラスタ 2 に分類されているものとする。また、同種判定に用いる閾値の値は0.6とする。なお、図7では、シード単語である「レストラン S」と「レストラン T」は、網掛けで示している。

【 0 0 5 5 】

まず、クラスタ 1 の同種判別について説明する。

クラスタ 1 内の単語「レストラン A」は、「レストラン S レストラン A」のルートにより、最短 1 ターンでシード単語「レストラン S」から出力される。若しくは、「レストラン A」は、「レストラン A レストラン T」のルートにより、最短 1 ターンでシード単語「レストラン T」を出力する。そのため、その最短のターン数 1 の逆数 1 を、「レストラン A」のシード単語までの近さを表す値とする。

同様に、クラスタ 1 内の単語「レストラン B」は、「レストラン S レストラン B」のルートにより、最短 1 ターンでシード単語「レストラン S」から出力される。若しくは、「レストラン B」は、「レストラン B レストラン T」のルートにより、最短 1 ターンでシード単語「レストラン T」を出力する。そのため、その最短のターン数 1 の逆数 1 を、「レストラン B」のシード単語までの近さを表す値とする。

したがって、クラスタ 1 全体でのシード単語までの近さは、「レストラン A」と「レストラン B」の近さの平均を取り 1 となる。この値は、閾値0.6以上であるため、クラスタ 1 は同種と判別され、その結果が収集単語記憶部 1 0 8 に記憶される。

【 0 0 5 6 】

続いて、クラスタ 2 の同種判別について説明する。

クラスタ 2 内の単語「うどん C」は、「レストラン S レストラン Z うどん C」又は「レストラン T レストラン W うどん C」等のルートにより、最短 2 ターンでシード単語「レストラン S」又は「レストラン T」から出力される。そのため、その最短のターン数 2 の逆数0.5を、「うどん C」のシード単語までの近さを表す値とする。

同様に、クラスタ 2 内の単語「うどん D」は、「レストラン S レストラン Z うどん D」又は「レストラン T レストラン W うどん D」等のルートにより、最短 2 ターンでシード単語「レストラン S」又は「レストラン T」から出力される。そのため、その最短のターン数 2 の逆数0.5を、「うどん D」のシード単語までの近さを表す値とする。

したがって、クラスタ 2 全体でのシード単語までの近さは、うどん Cとうどん Dの近さの平均を取り0.5となる。この値は、閾値0.6以下であるため、クラスタ 2 は異種と判別され、その結果が収集単語記憶部 1 0 8 に記憶される。

【 0 0 5 7 】

図4に戻り、続いて、出力部 1 0 5 は、収集単語記憶部 1 0 8 を参照して、収集され、クラスタに分類され、シード単語と同種か異種かを判別された単語を、それらの情報を関連付けて出力（表示）する（ステップS400）。例えば、出力部 1 0 5 は、「クラスタ 1 { レストラン A、レストラン B } : 同種、クラスタ 2 { うどん C、うどん D } : 異種」等と出力する。以上で、辞書作成処理は終了する。

【 0 0 5 8 】

このように、本実施形態では、辞書増殖処理によって収集された各単語は、クラスタに分類される。そして、各クラスタ毎に、シード単語と同じ種類の単語から構成されるか否かが判別されて出力される。従って、どのような異種の単語が収集されているのかをユーザに好適に出力することができる。

10

20

30

40

50

## 【 0 0 5 9 】

(第2実施形態)

第2実施形態に係る辞書作成装置200は、図9に示すように、第1実施形態の辞書作成装置100に、単語選択部201、再実行部202、および、単語グループ記憶部203が追加された構成である。なお、下記及び図面では、第1実施形態と同様のものについては、同一の符号を付す。また、第1実施形態と同様の構成要素の詳細な説明は、上記第1実施形態の説明に準じ、詳細な説明を省略する。

## 【 0 0 6 0 】

単語グループ記憶部203には、図10(A)、図10(B)に示すように、収集した単語と、該単語が属するグループの識別情報であるグループ名とが対応付けられて記憶される。

10

## 【 0 0 6 1 】

単語選択部201は、単語グループ記憶部203を参照して、未収集のグループを1つ選択し、選択したグループから所定数の単語を選択する。そして、単語選択部201は、選択した単語をシード単語とした辞書増殖処理の実行を辞書増殖部102に指示する。

## 【 0 0 6 2 】

再実行部202は、収集され、クラスタに分類され、シード単語と同種か異種かを判別された単語にグループ名を付して単語グループ記憶部203に追加する。そして、再実行部202は、未だ収集を行っていないグループがある場合には、そのグループから単語を選択することを単語選択部201に指示をする。

20

## 【 0 0 6 3 】

なお、その他の各部(入力部101、辞書増殖部102、クラスタリング部103、種別判別部104、出力部105、文書記憶部106、収集過程記憶部107、収集単語記憶部108)は、第1実施形態と同様の処理を行うため、ここでは説明を省略する。但し、辞書増殖部102が単語収集の起点とするシード単語は、単語選択部201が選択した単語である。

## 【 0 0 6 4 】

続いて、辞書作成装置200で実施される処理の動作について説明する。なお、予め、単語グループ記憶部203には、複数の単語が、グループ1として登録されている。また、このグループ1は、後述する収集未完グループであるとする。また、グループ1以外のグループは現時点では登録されていないものとする。

30

## 【 0 0 6 5 】

まず、ユーザは、入力部101を操作して、辞書を作成することを指示する。この指示操作に応じて、辞書作成装置200は、図11に示す辞書作成処理を行う。

## 【 0 0 6 6 】

辞書作成処理が開始されると、単語選択部201は、単語グループ記憶部203を参照して、未収集のグループ(即ちグループ1)に含まれる単語のなかから、予め設定されている数の単語をシード単語として選択する(ステップS50)。

## 【 0 0 6 7 】

続いて、辞書増殖部102は、第1実施形態と同様に辞書増殖処理を行い、シード単語と同種の単語を収集する(ステップS100)。但し、ここでは、ステップS50で選択された単語をシード単語とする。

40

## 【 0 0 6 8 】

続いて、クラスタリング部103は、第1実施形態と同様にクラスタリング処理を行い、辞書増殖処理によって収集された単語をクラスタに分類する(ステップS200)。

## 【 0 0 6 9 】

続いて、種別判別部104は、第1実施形態と同様に同種判別処理を行い、クラスタが、シード単語と同種の単語から構成されるか否かを判別する(ステップS300)。

## 【 0 0 7 0 】

続いて、再実行部202は、シード単語と同種か異種かを判別されたクラスタ毎に、該

50

クラスタを構成する単語を単語グループ記憶部 203 に登録して、グルーピングする単語グループ更新処理を行う（ステップS330）。

【0071】

図12に、単語グループ更新処理の詳細を示す。単語グループ更新処理が開始されると、まず、再実行部 202 は、上述のステップS200でクラスタリングしたクラスタのなかから未処理のクラスタを1つ選択する（ステップS331）。

【0072】

続いて、再実行部 202 は、ステップS300の同種判別処理の結果を参照して、選択したクラスタがシード単語と同種の単語から構成されているか否かを判別する（ステップS332）。

10

【0073】

シード単語と同種の場合（ステップS332；Yes）、再実行部 202 は、シード単語と同じグループ名を付して、選択したクラスタ内の単語を単語グループ記憶部 203 に登録する（ステップS333）。そして、ステップS337に処理を移す。

【0074】

シード単語と異種の場合（ステップS332；No）、再実行部 202 は、単語グループ記憶部 203 を参照して、選択したクラスタ内の単語のなかに、既に単語グループ記憶部 203 に記憶されている単語（既存単語）があるか否かを判別する（ステップS334）。

【0075】

既存単語があると判別された場合（ステップS334；Yes）、再実行部 202 は、その既存単語に付されているグループ名と同じグループ名を付して、選択したクラスタ内の単語を単語グループ記憶部 203 に登録する（ステップS335）。そして、ステップS337に処理を移す。

20

【0076】

既存単語がないと判別された場合（ステップS334；No）、再実行部 202 は、新たに発行したグループ名を付して、選択したクラスタ内の単語を単語グループ記憶部 203 に登録する（ステップS336）。そして、ステップS337に処理を移す。

【0077】

ステップS337では、再実行部 202 は、クラスタリングした全てのクラスタで、クラスタ内の単語を単語グループ記憶部 203 に登録する処理を行ったか否かを判別する。

30

【0078】

未だ単語グループ記憶部 203 に登録する処理を行っていないクラスタがある場合（ステップS337；No）、再実行部 202 は、未処理のクラスタを選択して、クラスタ内の単語を単語グループ記憶部 203 に登録する一連の処理（ステップS331～ステップS336）を繰り返す。

【0079】

全てのクラスタで、単語を単語グループ記憶部 203 に登録する処理を行った場合（ステップS337；Yes）、単語グループ更新処理は終了する。

【0080】

図11に戻り、続いて、再実行部 202 は、単語収集が未だ完了していないグループ（以下、収集未完グループという）があるか否かを判別する（ステップS360）。

40

例えば、以下に示す a)～d) の何れかの条件を満たすグループを収集未完グループと判断すればよい。

a) グループ内の単語数が一定数以上に達していないグループ。

b) グループ内の単語をシード単語とした辞書増殖処理を所定回数以上行っていないグループ。

c) グループに新たに追加された単語が一定数以上あるグループ。

d) a)～c) を所定の重み付けを付した割合で組み合わせた条件に合致するグループ。

【0081】

50

収集未完グループが有る場合（ステップS360；Yes）、再実行部202は、収集未完グループの1つからシード単語を選択することを単語選択部201に指示する。そして、シード単語から単語を収集して、クラスタリングし、シード単語と同種か異種かの判定を行い、グルーピングする処理が繰り返される（ステップS50～ステップS330）。

【0082】

収集未完グループが無い場合（ステップS360；No）、出力部105は、収集した単語を出力する。但し、単語の属するクラスタ、および、そのクラスタがシード単語を同種であるか否かを示す情報に加えて、単語が属するグループ名を単語グループ記憶部203から取得する。そして、これらの情報を、収集した単語と関連付けて出力（表示）するものとする。以上で、辞書作成処理は終了する。

10

【0083】

続いて、上述した辞書作成処理について、具体例を挙げて説明する。なお、前提として、図10（A）に示すように、収集未完グループであるグループ1のみが、単語グループ記憶部203には記憶されているものとする。

【0084】

従って、この状態で辞書作成処理が開始されると、まず、グループ1内の単語「レストランS」と「レストランT」が選択される（ステップS50）。続いて、この「レストランS」と「レストランT」とをシード単語とした辞書増殖処理が実行されて、単語が収集される（ステップS100）。そして、収集された単語は、その結束度に基づいてクラスタリングされ（ステップS200）、クラスタ毎に、シード単語「レストランS」「レストランT」と同種であるか否かが判別される（ステップS300）。ここでは、以下に示すようなクラスタ1～5が作成されたこととする。

20

- ・クラスタ1（同種）：「レストランA」「レストランB」
- ・クラスタ2（異種）：「うどんC」「うどんD」
- ・クラスタ3（同種）：「レストランX」「レストランZ」「レストランW」
- ・クラスタ4（同種）：「レストランS」「レストランT」
- ・クラスタ5（異種）：「うどんG」「うどんH」

【0085】

続いて、これらのクラスタ毎に、クラスタ内の単語をグループ化して単語グループ記憶部203に登録する単語グループ更新処理が実施される（ステップS330）。この場合、クラスタ1と、クラスタ3と、クラスタ4は、シード単語と同種と判定されているため、これらのクラスタ内の単語は、シード単語と同じグループ1の単語として単語グループ記憶部203に登録される（ステップS333）。

30

【0086】

また、クラスタ2とクラスタ5は、シード単語と異種の単語であり、また、これらのクラスタ内の単語は未だ単語グループ記憶部203に記憶されていない。従って、クラスタ2とクラスタ5内の単語は、それぞれ、グループ2、グループ3の新規のグループ名を付されて、単語グループ記憶部203に登録される（ステップS336）。

【0087】

そして、最終的には、図10（B）に示すように、クラスタ1～5内の単語がグループ名を付されて単語グループ記憶部203に登録される。

40

【0088】

続いて、収集未完のグループがある場合には、そのグループ（即ち、グループ2かグループ3）のうちの1つを選択して、選択したグループ内の単語を新たにシード単語とした単語収集を行う一連の処理が繰り返される。

【0089】

このように、第2実施形態では、異種単語がどの程度含まれているかだけでなく、同じような異種単語を新たなグループとして登録する。そして、そのグループ内の単語をシード単語として、さらに単語を収集することができる。これにより、初期に与えたシード単語と似ている単語も別グループとした単語収集を行うことができる。

50

## 【0090】

(第3実施形態)

第2実施形態では、グループ内の単語から、ランダムに選択した所定数の単語をシード単語として辞書増殖を行った。そのため、少ない収集回数で多くの単語を取得したい場合、収集回数が増えても収集される単語がシード単語と類似する精度を高くしたい場合、などといった種々の場面に応じた適切な単語の収集ができない。本実施形態では、種々の場面に応じた適切な単語の収集を可能とすることを特徴とする。

## 【0091】

第3実施形態に係る辞書作成装置300は、図13に示すように、第2実施形態の辞書作成装置200の単語選択部201が第二単語選択部301に置き換えられている。また、単語間結束度記憶部302が新たに追加されている。なお、下記及び図面では、第1実施形態、および、第2実施形態と同様のものについては、同一の符号を付す。また、第1実施形態、および、第2実施形態と同様の構成要素の詳細な説明は、上記第1実施形態、第2実施形態の説明に準じ、詳細な説明を省略する。

10

## 【0092】

第二単語選択部301は、単語グループ記憶部203を参照して、未収集のグループを1つ選択し、選択したグループに含まれる単語から複数の単語を選択する。この際、第二単語選択部301は、単語間結束度記憶部302を参照して、結束度が所定の条件を満たす単語を優先的に選択する。

## 【0093】

ここで、上記の所定の条件とは、例えば、「グループ内の単語のうち結束度の高い順に75%、残りの25%は結束度が低いものから順に選択する」などの条件である。結束度の高い単語のみを選択すると、頻繁に出現する単語のみが収集されるため、シード単語と類似の単語が収集される精度は高くなるが、収集される単語の数は少なくなり収集効率は悪化する。したがって、収集精度よりも収集効率を重視した単語収集を行いたい場合には、上記のような条件を採用することが望ましい。

20

また、収集効率よりも収集精度を重視した単語収集を行いたい場合には、「グループ内の単語のうち結束度の高い順に選択する」などの条件を採用することが望ましい。

なお、このような単語選択の条件を定義する条件情報が、予め、辞書作成装置300の記憶部に記憶されているものとする。

30

## 【0094】

単語間結束度記憶部302は、クラスタリング部103によって算出された単語間の結束度を記憶する。具体的には、図14に示すように、単語間結束度記憶部302には、2つの単語と、その2つの単語間の結束度が対応付けられて記憶される。例えば、図14の先頭のエントリから、「レストランS」と「レストランT」との間の結束度は0.9とわかる。

## 【0095】

なお、その他の各部(入力部101、辞書増殖部102、クラスタリング部103、種別判別部104、出力部105、文書記憶部106、収集過程記憶部107、収集単語記憶部108、再実行部202、単語グループ記憶部203)は、第2実施形態と同様の処理を行うため、ここでは説明を省略する。

40

## 【0096】

続いて、辞書作成装置300で実施される処理の動作について説明する。

なお、予め、収集の際に採用する結束度に関するグループから単語を選択するための条件が設定されているものとする。また、グループからは4つの単語を選択するものとする。

## 【0097】

ユーザは、入力部101を操作して、辞書を作成することを指示する。この指示操作に応じて、辞書作成装置300は、第2実施形態と同様の図11に示す辞書作成処理を行う。

## 【0098】

50

まず、第二単語選択部 301 は、単語グループ記憶部 203 を参照して、未収集のグループを 1 つ選択し、単語間結束度記憶部 302 を参照して、所定の条件に基づいて、選択したグループ内の単語のうちから所定数 (4 つ) の単語をシード単語として選択する (ステップ S50)。

【0099】

例えば、「グループ内の単語のうち結束度の高い順に 75%、残りの 25% は結束度が低いものから順に選択する」条件が設定されている場合を考える。即ち、結束度の高い単語を 3 つ、結束度の低い単語を 1 つ選択することとなる。

この場合、第二単語選択部 301 は、まず、グループ内の単語のうち、単語間の結束度が最も高い 2 単語を選択する。次に、第二単語選択部 301 は、その 2 つの単語それぞれと結束度が最も高い単語を 1 つ選択する。そして、第二単語選択部 301 は、これら 3 つの単語それぞれと、結束度の低い単語を 1 つ選択する。

【0100】

以降の処理は、第 2 実施形態と同様である。

即ち、辞書増殖部 102 は、第二単語選択部 301 によって選択された 4 つの単語をシード単語として、同種の単語を収集する辞書増殖処理を行う (ステップ S100)。続いて、クラスタリング部 103 が、収集された単語をクラスタリングする (ステップ S200)。なお、この際、クラスタリング部 103 は、クラスタリングするために算出した単語とその単語間の結束度とを、単語間結束度記憶部 302 に記録する。そして、種別判別部 104 が、クラスタ毎に、クラスタがシード単語と同種の単語から構成されるか否かを判別する (ステップ S300)。そして、再実行部 202 が、収集した単語をグルーピングする (ステップ S330)。そして、未収集のグループがある場合は (ステップ S360; Yes)、未収集のグループからシード単語を選択して単語を収集する処理を繰り返し、未収集のグループがない場合は (ステップ S360; No)、処理は終了する。

【0101】

このように、本実施形態では、グループ内の単語をランダムに選択するのではなく、単語間の結束度を考慮して単語を選択する。従って、種々の場面に対応した単語収集が可能となる。

【0102】

なお、本各実施形態は種々の変形、および、応用が可能である。

例えば、上記各実施形態では、文書記憶部 106 に記憶されている文書から単語を抽出したが、これに限らず、例えば、インターネット検索エンジンを用いて、インターネット上の Web ページから、単語を抽出してもよい。

【0103】

図 15 は、本発明の各実施形態に係る辞書作成装置 100、200、300 をコンピュータに実装する場合の、物理的な構成の一例を示すブロック図である。本発明の各実施形態に係る辞書作成装置 100、200、300 は、一般的なコンピュータ装置と同様のハードウェア構成によって実現することができる。辞書作成装置 100、200、300 は、制御部 21、主記憶部 22、外部記憶部 23、操作部 24、表示部 25 および入出力部 26 を備える。主記憶部 22、外部記憶部 23、操作部 24、表示部 25 および入出力部 26 はいずれも内部バス 20 を介して制御部 21 に接続されている。

【0104】

制御部 21 は CPU (Central Processing Unit) 等から構成され、外部記憶部 23 に記憶されている制御プログラム 30 に従って、前述した各実施形態における辞書作成処理を実行する。

【0105】

主記憶部 22 は RAM (Random-Access Memory) 等から構成され、外部記憶部 23 に記憶されている制御プログラム 30 をロードし、制御部 21 の作業領域として用いられる。

【0106】

外部記憶部 23 は、フラッシュメモリ、ハードディスク、DVD-RAM (Digital Ve

10

20

30

40

50



rsatile Disc Random-Access Memory)、DVD-RW(Digital Versatile Disc ReWritable)等の不揮発性メモリから構成され、上述の処理を制御部21に行わせるための制御プログラム30を予め記憶する。また、外部記憶部23は、制御部21の指示に従って、この制御プログラム30が記憶するデータを制御部21に供給し、制御部21から供給されたデータを記憶する。また、外部記憶部23は、上述した各実施形態における文書記憶部106、収集過程記憶部107、収集単語記憶部108、単語グループ記憶部203、および、単語間結束度記憶部302を物理的に実現する。

#### 【0107】

操作部24はキーボードおよびマウスなどのポインティングデバイス等と、キーボードおよびポインティングデバイス等を内部バス20に接続するインターフェース装置等から構成されている。操作部24を介して、シード単語や辞書作成処理の開始の指示が制御部21に供給される。

10

#### 【0108】

表示部25は、CRT(Cathode Ray Tube)またはLCD(Liquid Crystal Display)などから構成され、種々の情報を表示する。例えば、表示部25は、収集された各単語を、クラスタ毎に、シード単語と同種であるか異種であるかの情報を付して表示する。

#### 【0109】

入出力部26は、無線送受信機、無線モデムまたは網終端装置、およびそれらと接続するシリアルインタフェースまたはLAN(Local Area Network)インタフェース等から構成されている。例えば、入出力部26を介して、インターネット上のWebページから単語を収集してもよい。

20

#### 【0110】

図1、図9、および図13に示す辞書作成装置100、200、300の辞書増殖部102、クラスタリング部103、種別判別部104、出力部105、単語選択部201、再実行部202、および、第二単語選択部301の処理は、制御プログラム30が、制御部21、主記憶部22、外部記憶部23、操作部24、表示部25および入出力部26などを資源として用いて処理することによって実行する。

#### 【0111】

なお、前記のハードウェア構成やフローチャートは一例であり、任意に変更および修正が可能である。

30

#### 【0112】

また、制御部21、主記憶部22、外部記憶部23、操作部24、入出力部26および内部バス20などから構成される辞書作成装置100、200、300の処理を行う中心となる部分は、専用のシステムによらず、通常のコンピュータシステムを用いて実現可能である。たとえば、前記の動作を実行するためのコンピュータプログラムを、コンピュータが読み取り可能な記録媒体(フレキシブルディスク、CD-ROM、DVD-ROM等に格納して配布し、当該コンピュータプログラムをコンピュータにインストールすることにより、前記の処理を実行する辞書作成装置100、200、300を構成してもよい。また、インターネット等の通信ネットワーク上のサーバ装置が有する記憶装置に当該コンピュータプログラムを格納しておき、通常のコンピュータシステムがダウンロード等することで辞書作成装置100、200、300を構成してもよい。

40

#### 【0113】

また、辞書作成装置100、200、300の機能を、OS(オペレーティングシステム)とアプリケーションプログラムの分担、またはOSとアプリケーションプログラムとの協働により実現する場合などには、アプリケーションプログラム部分のみを記録媒体や記憶装置に格納してもよい。

#### 【0114】

また、搬送波にコンピュータプログラムを重畳し、通信ネットワークを介して配信することも可能である。たとえば、通信ネットワーク上の掲示板(BBS, Bulletin Board System)に前記コンピュータプログラムを掲示し、ネットワークを介して前記コンピュータプロ

50

グラムを配信してもよい。そして、このコンピュータプログラムを起動し、OSの制御下で、他のアプリケーションプログラムと同様に実行することにより、前記の処理を実行できるように構成してもよい。

【0115】

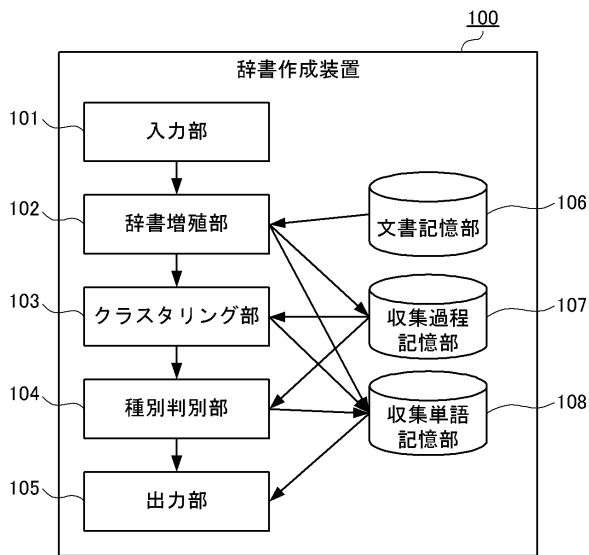
本発明は2009年12月11日に出願された日本国特許出願2009-282304号に基づく。本明細書中に日本国特許出願2009-282304号の明細書、特許請求の範囲、図面全体を参照として取り込むものとする。

【符号の説明】

【0116】

- 100 辞書作成装置
- 101 入力部
- 102 辞書増殖部
- 103 クラスタリング部
- 104 種別判別部
- 105 出力部
- 106 文書記憶部
- 107 収集過程記憶部
- 108 収集単語記憶部

【図1】



【図2】

収集過程記憶部

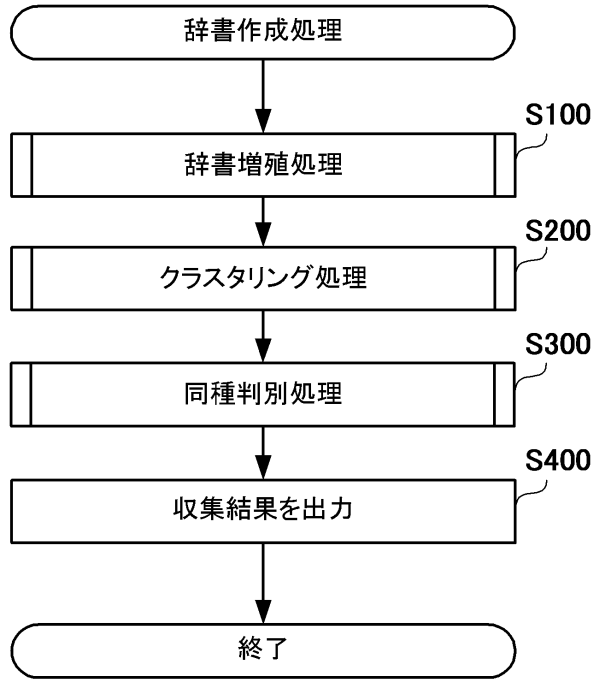
ターン数	入力単語	出力単語
1	レストランS	レストランX
1	レストランS	レストランZ
1	レストランS	レストランW
1	レストランT	レストランX
1	レストランT	レストランZ
1	レストランT	レストランW
2	レストランX	レストランA
2	レストランS	レストランA
2	レストランS	レストランB
2	レストランZ	うどんC
2	レストランW	うどんC
2	レストランZ	うどんD
2	レストランW	うどんD
3	レストランA	レストランE
3	レストランA	レストランT
3	レストランB	レストランT
3	うどんC	うどんG
3	うどんC	うどんH
3	うどんD	レストランT
3	うどんD	うどんH
...	...	...

【図3】

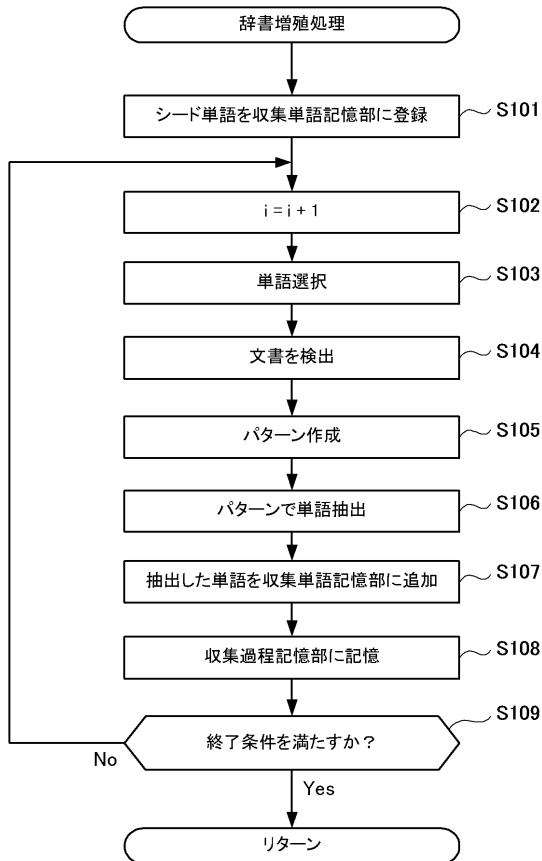
収集単語記憶部

収集単語	クラスターID	同種 or 異種?
レストランA	クラスター1	同種
レストランB		
うどんC	クラスター2	異種
うどんD		

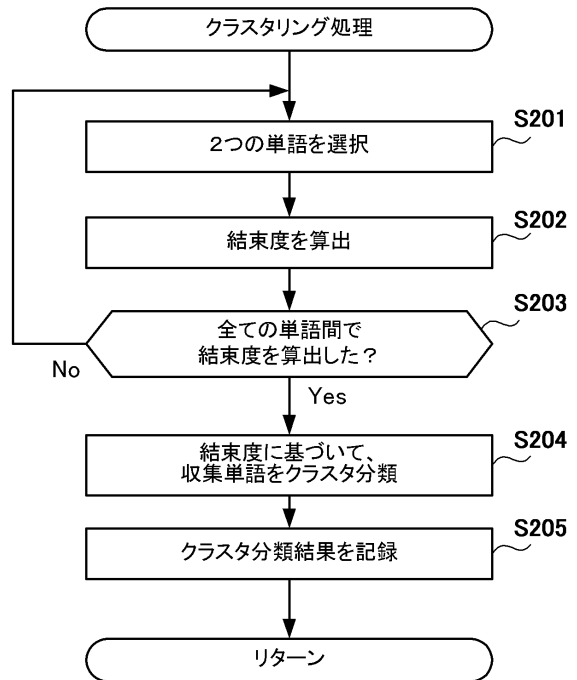
【図4】



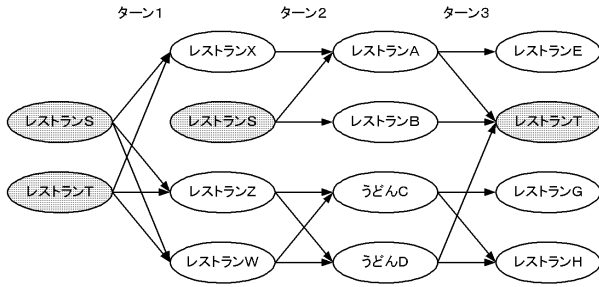
【図5】



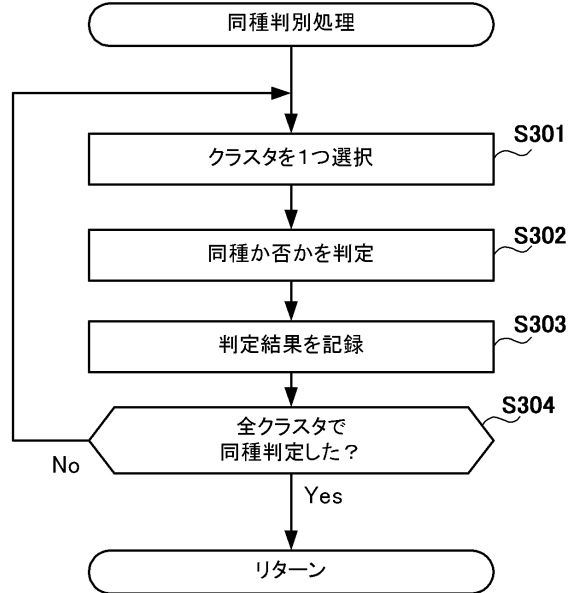
【図6】



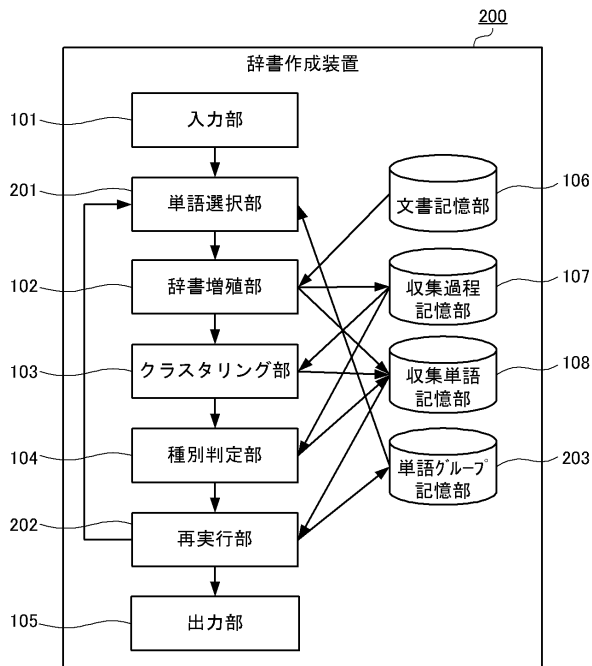
【図7】



【図8】



【図9】



【図10】

単語グループ記憶部

単語	グループ名
レストランS	グループ1
レストランT	

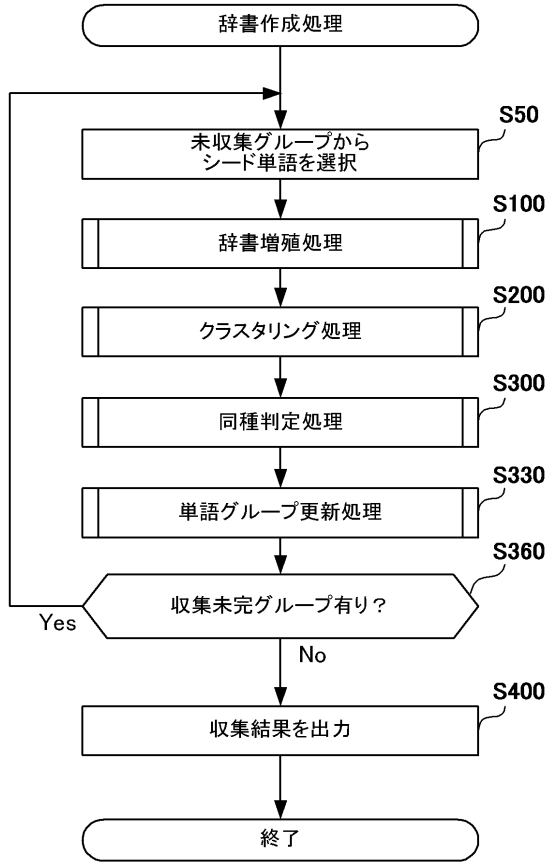
(A)

単語グループ記憶部

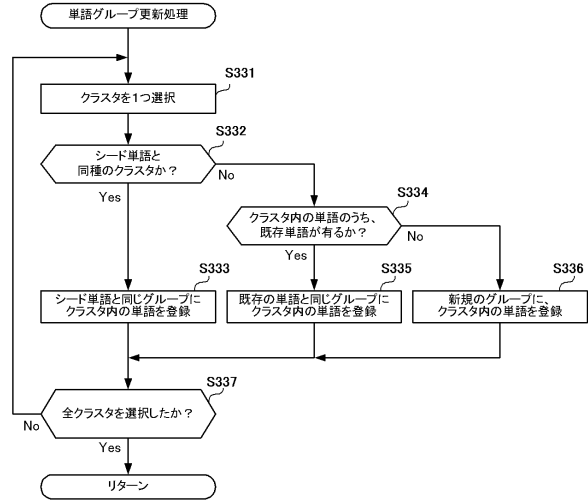
単語	グループ名
レストランS	グループ1
レストランT	
レストランX	
レストランW	
レストランZ	
レストランA	
うどんC	グループ2
うどんD	
うどんG	グループ3
うどんH	

(B)

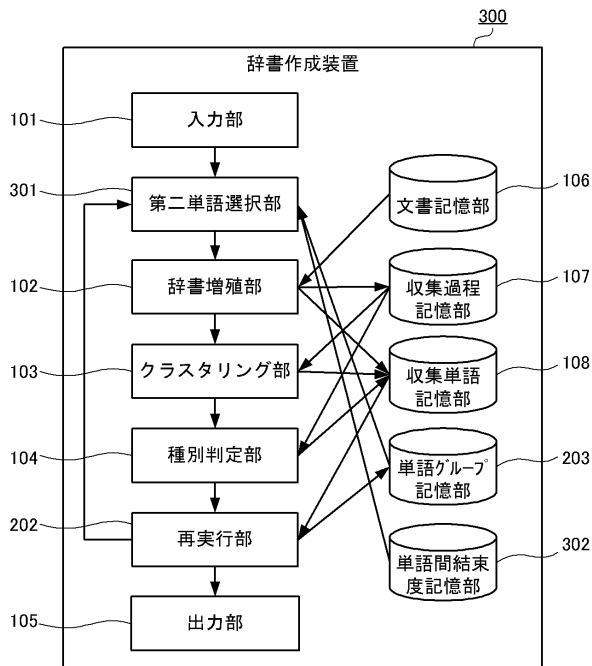
【図11】



【図12】



【図13】

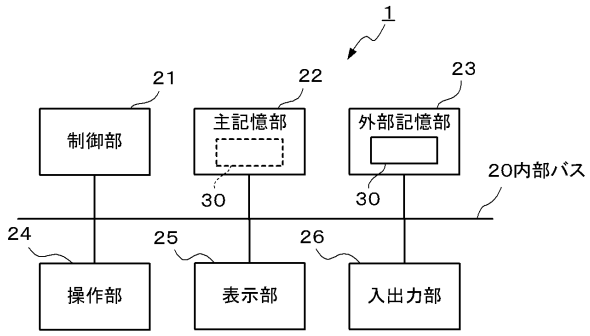


【図14】

収集単語記憶部

単語 1	単語 2	結束度
レストランS	レストランT	0.9
レストランA	レストランB	0.8
レストランX	レストランY	0.8
レストランW	レストランX	0.9
うどんC	うどんD	0.75
うどんD	うどんH	0.8
レストランS	レストランX	0.5
...	...	...

【図15】



## フロントページの続き

審査官 松田 直也

## (56)参考文献 特開2007-207218(JP,A)

水口 弘紀、河合 英紀、土田 正明、久寿居 大、Web知識を利用したブートストラップによる辞書増殖手法、電子情報通信学会 第18回データ工学ワークショップ論文集、[online]、日本、電子情報通信学会、2007年 6月 1日、[検索日 2014.6.3]、インターネット、URL、<http://www.ieice.org/~de/DEWS/DEWS2007/pdf/e8-5.pdf>

大島 裕明、田中 克己、正解語ペア漸増による関連語取得のための両方向構文パターン発見、第1回データ工学と情報マネジメントに関するフォーラム - DEIMフォーラム - 論文集、[online]、日本、電子情報通信学会、2009年 5月 9日、[検索日 2014.6.3]、インターネット、URL、<http://db-event.jp.org/deim2009/proceedings/files/B9-1.pdf>

## (58)調査した分野(Int.Cl.、DB名)

G06F 17/30

G06F 17/27