

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5146979号  
(P5146979)

(45) 発行日 平成25年2月20日(2013.2.20)

(24) 登録日 平成24年12月7日(2012.12.7)

(51) Int. Cl. F 1  
G 0 6 F 17/27 (2006.01) G 0 6 F 17/27 M

請求項の数 6 (全 21 頁)

(21) 出願番号	特願2006-154497 (P2006-154497)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2
(22) 出願日	平成18年6月2日(2006.6.2)	(73) 特許権者	000208891 K D D I 株式会社 東京都新宿区西新宿二丁目3番2号
(65) 公開番号	特開2007-323475 (P2007-323475A)	(74) 代理人	100099933 弁理士 清水 敏
(43) 公開日	平成19年12月13日(2007.12.13)	(72) 発明者	隅田 英一郎 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成20年12月8日(2008.12.8)	(72) 発明者	菅谷 史昭 埼玉県ふじみ野市大原二丁目1番15号 株式会社K D D I 研究所内

最終頁に続く

(54) 【発明の名称】 自然言語における多義解消装置及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

自然言語文からなる入力文において、ある単語と、前記入力文において前記ある単語が置かれた文脈と、前記ある単語の意味を表す可能性のある複数の意味候補を含む意味候補の集合とが与えられると、当該意味候補の集合の中から、前記文脈において前記ある単語の意味として最も適切なものを選択する、自然言語における多義解消装置であって、

前記ある単語と、前記意味候補集合中の意味候補との組合せの各々について、所定のコーパスから、当該組合せを構成する語同土が共起する文書の集合を収集するための文書収集手段と、

前記文書収集手段によって前記組合せの各々について収集された前記文書の集合を学習データとして用い、前記ある単語と、その単語の文書中の文脈とが与えられると、前記ある単語の当該文脈中での意味として最適な意味候補を前記意味候補集合中から選択する分類器を自動的に作成するための分類器作成手段と、

前記入力文において、前記ある単語が置かれた文脈に基づいて、前記ある単語の意味として最適なものを、前記意味候補の集合の中から前記分類器を用いて選択するための分類実行手段とを含み、

前記文書収集手段は、前記ある単語と、前記意味候補集合中の意味候補との組合せの各々について、前記ある単語とその意味候補との両方を検索キーワードとして、インターネット上に存在する文書からなる仮想的コーパスから、当該組合せを構成する語同土が共起する文書の集合を検索し収集するための検索手段を含み、

10

20

前記分類器作成手段は、

前記文書収集手段によって前記組合せの各々について収集された前記文書の集合のうち、集合に含まれる文書の数が多いものを所定の基準に従って選択し、それら文書の集合に対応する意味候補のみを前記意味候補の集合の要素として選択する処理を行なうための意味候補選択手段と、

前記意味候補選択手段により選択された文書集合を学習データとして用い、前記ある単語と、その単語の文書中の文脈とが与えられると、前記ある単語の当該文脈中での意味として最適な意味候補を、前記意味候補集合中から選択する分類器を機械学習により自動的に作成するための機械学習手段とを含む、自然言語における多義解消装置。

【請求項 2】

前記機械学習手段は、

前記文書集合選択手段により選択された文書集合に含まれる文書の各々に対し、当該文書中における前記ある単語の位置の前後の所定範囲に存在する単語列から、当該文書中における前記ある単語の文脈の特徴量を表す、所定の構成の学習用の特徴量ベクトルを算出するための特徴量ベクトル算出手段と、

前記文書集合選択手段により選択された文書集合に含まれる文書の各々に対して前記特徴量ベクトル算出手段により算出された学習用の特徴量ベクトルと、当該文書の検索時に使用された意味候補とを組にして学習用データを作成し、当該学習用データを用いた機械学習により、前記学習用の特徴量ベクトルと同じ構成の分類用の特徴量ベクトルが与えられると、当該分類用の特徴量ベクトルに対応する文脈中における前記ある単語の意味として最適なものを、前記意味候補集合中から選択する所定の分類器を自動的に作成するための手段とを含む、請求項 1 に記載の自然言語における多義解消装置。

【請求項 3】

前記意味候補選択手段は、前記文書収集手段によって前記組合せの各々について収集された前記文書の集合のうち、集合に含まれる文書の数が多い所定の個数の集合を選択し、それら文書の集合に対応する意味候補のみを前記意味候補として選択する処理を行なうための手段を含む、請求項 1 又は請求項 2 に記載の、自然言語における多義解消装置。

【請求項 4】

前記意味候補選択手段は、前記文書収集手段によって前記組合せの各々について収集された前記文書の集合のうち、集合に含まれる文書の数が予め定められるしきい値より大きな集合を選択し、それら文書の集合に対応する意味候補のみを前記意味候補として選択する処理を行なうための手段を含む、請求項 1 又は請求項 2 に記載の、自然言語における多義解消装置。

【請求項 5】

前記収集するための手段は、前記ある単語と、前記意味候補集合中の意味候補との組合せの各々について、インターネット上に存在するウェブページからなる仮想的コーパスから、当該組合せを構成する語が共起するウェブページの集合を検索し、所定の定数を上限とした要素数の集合として収集するための手段を含む、請求項 1 に記載の自然言語における多義解消装置。

【請求項 6】

コンピュータにより実行されると、当該コンピュータを、請求項 1 ~ 請求項 5 のいずれかに記載の自然言語における多義解消装置として機能させる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は自然言語処理に関し、特に、単語の読み（日本語における仮名表記）、アクロニム（頭字語）のフルスペル、及び二つの言語の間での訳語の対応などに見られるあい

10

20

30

40

50

まい性を解消するための自然言語処理に関する。

【背景技術】

【0002】

自然言語には、あいまい性が常に付きまとう。例えば同形異音語という問題がある。同形異音語とは、読みが複数ある単語のことである。例えば英語の「bow」という単語には、「bow」（蝶型リボン）と「bow」（船首）という二つの読み方がある。日本語でもこうした例は多い。例えば「大平」という語は、「オオヒラ」とも、「タイハイ」とも、「オオダイラ」とも読める。

【0003】

こうしたあいまい性は、アクロニムにも存在する。例えば「ACL」というアクロニムは、「The Association for Computational Linguistics」、「Anterior Cruciate Ligament」、及び「Access Control List」のいずれとも解釈できる。同様のあいまい性は、翻訳の際の訳語の決め方等にも見出すことができる。

10

【0004】

人間の場合、こうしたあいまい性については、その語が生じた状況などに応じて適宜適切に判断をしたり、いずれかの手段を用いて候補をさがし、その中で状況に応じて最も適していると思われるものを選択したりすることで解決している。しかし、自然言語処理でそのような処理を実現することは困難である。

【0005】

このようなあいまい性は、自然言語処理において重大な問題となり得る。例えば、日本語テキストの読上げにおいて、同形異音語に遭遇した場合、適切な発音で読上げを行なうためには、その発音（かな表記）を決定する必要がある。さもないと、不適切な読上げが行なわれてしまう。

20

【0006】

こうした問題を解決するための提案が非特許文献1でなされている。非特許文献1では、予め、ある単語Wとその対応する意味 $S_i$ とを記述した学習データを人手で用意し、その学習データを用いて、ある単語Wが与えられたときに意味 $S_i$ のうちで適切なものを選択する分類器を作成する。

【非特許文献1】梅村祥之、清水司、「音声合成システムのための同形異音語の読み分け」、豊田中央研究所R&Dレビュー、2000年、第35巻第1号、67頁～74頁

30

【発明の開示】

【発明が解決しようとする課題】

【0007】

しかし、非特許文献1に開示された方法では、学習データを人手で用意する必要があり、時間と費用とがかさむという問題がある。また、限られた人の手によって学習データが作成されるので、学習データに偏りが生ずる可能性もあるため、信頼性が低いという問題もある。

【0008】

それ故に本発明の目的は、自然言語に伴うあいまい性を、容易に、かつ信頼性高く解決できる、自然言語における多義解消装置を提供することである。

40

【課題を解決するための手段】

【0009】

本発明の第1の局面によれば、自然言語における多義解消装置は、自然言語文からなる入力文において、ある単語と、入力文においてある単語が置かれた文脈と、ある単語の意味を表す可能性のある複数の意味候補を含む意味候補の集合とが与えられると、当該意味候補の集合の中から、文脈においてある単語の意味として最も適切なものを選択する、自然言語における多義解消装置であって、ある単語と、意味候補集合中の意味候補との組合せの各々について、所定のコーパスから、当該組合せを構成する語が共起する文書の集合を収集するための文書収集手段と、文書収集手段によって組合せの各々について収集された文書の集合を学習データとして用い、ある単語と、その単語の文書中の文脈とが与えら

50

れると、ある単語の当該文脈中での意味として最適な意味候補を意味候補集合中から選択する分類器を自動的に作成するための分類器作成手段と、入力文において、ある単語が置かれた文脈に基づいて、ある単語の意味として最適なものを、意味候補の集合の中から分類器を用いて選択するための分類実行手段とを含む。

【0010】

入力文中のある単語と、その単語が置かれた文脈と、その単語に意味を表す可能性のある複数の意味候補を含む意味候補の集合が与えられると、その単語と意味候補との組合わせの各々について、文書収集手段が所定のコーパスから当該組合せを構成する単語が共起する文書の集合を収集する。分類器作成手段は、組合せの各々について収集された文書の集合を学習データとして用いて分類器を作成する。この分類器は、ある単語と、その単語の文書中の文脈とが与えられると、ある単語の当該文脈中での意味として最適な意味候補を意味候補集合中から選択する機能を持つ。分類手段は、入力文中の単語と、その単語が置かれた文脈とを、このようにして作成された分類器に与え、その結果に基づいて、入力文中の単語の意味として最適なものを、意味候補の集合の中から選択する。

10

【0011】

すなわち、この装置では、ある単語と、その単語の文脈と、その単語に対応する可能性のある複数の意味候補とが与えられると、文脈から適切と思われる意味候補を自動的に選択できる。この作業には人手を介在させる必要はない。従って、容易に適切な意味候補を選択し、入力された単語の多義性を解消できる多義解消装置を提供できる。

【0012】

好ましくは、分類器作成手段は、文書収集手段によって組合せの各々について収集された文書の集合のうち、集合に含まれる文書の数が多いものを所定の基準に従って選択し、それら文書の集合に対応する意味候補のみを意味候補の集合の要素として選択する処理を行なうための意味候補選択手段と、意味候補選択手段により選択された文書集合を学習データとして用い、ある単語と、その単語の文書中の文脈とが与えられると、ある単語の当該文脈中での意味として最適な意味候補を、意味候補集合中から選択する分類器を機械学習により自動的に作成するための機械学習手段とを含む。

20

【0013】

収集された文書集合のうち、集合に含まれる文書の数が少ないものは意味候補選択手段により棄却される。集合に含まれる文書の数が少ないということは、その単語と、その集合に対応する意味候補とが共起する可能性が他と比較して少ないということである。従って、与えられた文脈におけるある単語の意味として不適切なものを排除できる。その結果、分類の信頼性を高めることができる。

30

【0014】

さらに好ましくは、機械学習手段は、文書集合選択手段により選択された文書集合に含まれる文書の各々に対し、当該文書中におけるある単語の位置の前後の所定範囲に存在する単語列から、当該文書中におけるある単語の文脈の特徴量を表す、所定の構成の学習用の特徴量ベクトルを算出するための特徴量ベクトル算出手段と、文書集合選択手段により選択された文書集合に含まれる文書の各々に対して特徴量ベクトル算出手段により算出された学習用の特徴量ベクトルと、当該文書の検索時に使用された意味候補とを組にして学習用データを作成し、当該学習用データを用いた機械学習により、学習用の特徴量ベクトルと同じ構成の分類用の特徴量ベクトルが与えられると、当該分類用の特徴量ベクトルに対応する文脈中におけるある単語の意味として最適なものを、意味候補集合中から選択する所定の分類器を自動的に作成するための手段とを含む。

40

【0015】

単語の文脈を、その単語の前後の所定範囲に存在する単語列から作成した学習用の特徴量ベクトルにより表す。こうした学習用の特徴量ベクトルを用いた機械学習により分類器を自動的に作成できる。その結果、適切な意味候補を手を介在させることなく自動的に選択し、入力された単語の多義性を解消できる多義解消装置を提供できる。

【0016】

50

より好ましくは、意味候補選択手段は、文書収集手段によって組合せの各々について収集された文書の集合のうち、集合に含まれる文書の数が多い所定の個数の集合を選択し、それら文書の集合に対応する意味候補のみを意味候補として選択する処理を行なうための手段を含む。

【0017】

入力された単語と、ある意味候補との組合せに対して収集された文書の集合に含まれる文書の数が多いということは、その組合せを構成する単語が共起する可能性が高いということである。従ってそうした意味候補は入力された単語に対する適切な意味候補である可能性が高い。また、この時点で意味候補の上限個数が設定されるので、以後の処理を安定した時間で完了できる。その結果、適切な意味候補を、人手を介在させることなく自動的に、信頼性高く、安定した時間で選択し、入力された単語の多義性を解消できる多義解消装置を提供できる。

10

【0018】

意味候補選択手段は、文書収集手段によって組合せの各々について収集された文書の集合のうち、集合に含まれる文書の数が予め定められるしきい値より大きな集合を選択し、それら文書の集合に対応する意味候補のみを意味候補として選択する処理を行なうための手段を含んでもよい。

【0019】

入力された単語と、ある意味候補との組合せに対して収集された文書の集合に含まれる文書の数があるしきい値より多いということは、その組合せを構成する単語が共起する可能性が高いということである。従ってそうした意味候補は入力された単語に対する適切な意味候補である可能性が高い。その結果、適切な意味候補を、人手を介在させることなく自動的に、かつ信頼性高く選択し、入力された単語の多義性を解消できる多義解消装置を提供できる。

20

【0020】

文書収集手段は、ある単語と、意味候補集合中の意味候補との組合せの各々について、インターネット上に存在するウェブページからなる仮想的コーパスから、当該組合せを構成する語が共起するウェブページの集合を検索し収集するための検索手段を含む。

【0021】

インターネット上のウェブページは、多数の人間により作成され維持されている。従ってそこで使用されている単語の用法は非常に数多い使用例をカバーしている。そのため、そうした文書をもとに分類器を作成すると、分類結果の偏りをなくし、信頼性を高めることができる。

30

【0022】

好ましくは、収集するための手段は、ある単語と、意味候補集合中の意味候補との組合せの各々について、インターネット上に存在するウェブページからなる仮想的コーパスから、当該組合せを構成する語が共起するウェブページの集合を検索し、所定の定数を上限とした要素数の集合として収集するための手段を含む。

【0023】

一つの集合について収集されるウェブページの個数に上限が設けられる。そのため、分類器の学習が過大な負荷となるおそれは小さい。その結果、適切な意味候補を、人手を介在させることなく自動的に、かつ信頼性高く安定して選択し、入力された単語の多義性を解消できる多義解消装置を提供できる。

40

【0024】

本発明の第2の局面に係るコンピュータプログラムは、コンピュータにより実行されると、当該コンピュータを、上記したいずれかの自然言語における多義解消装置として機能させるものである。

【発明を実施するための最良の形態】

【0025】

以下、本発明の実施の形態について図を参照して説明する。実施の形態は三つある。第

50

1の実施の形態は、日本語の入力文に対する音声合成において、複数の仮名表記（読み）を持つ語の仮名表記を決定する装置に関する。第2の実施の形態は、英語のアクセントに対し、英語の定義（フルスペル）を与える装置に関する。第3の実施の形態は、日本語から英語への翻訳において、日本語の単語に対し複数の英語の訳語が存在するときに、そのうちの一つを選択する装置に関する。すなわち、本発明において、ある単語の「意味」とは、日本語の場合に国語辞書にのっているような「意味」だけでなく、ある基準で見てその単語と等価であると評価できるような単語又は単語の集合又は文字列のことをいう。

【0026】

なお、以下の実施の形態の説明に用いる図面において、同一の部品には同一の参照符号を付してある。それらの名称及び機能も同一である。従って、それらについての詳細な説明は繰返さない。なお、後述するように、各実施の形態は、コンピュータハードウェアと、その上で実行されるコンピュータプログラムとにより実現可能である。従って、以下に示すブロック図中の機能ブロックの一部については、それを実現するためのコンピュータプログラムのフローチャート形式でその機能及び構成を示す。

10

【0027】

<第1の実施の形態>

[構成]

図1に、本発明の第1の実施の形態に係る音声合成システム30のブロック図を示す。図1を参照して、音声合成システム30は、音声合成の対象となる日本語の入力文を記憶するための入力文記憶部40と、入力文記憶部40から所定長を順次取出して記憶するための入力文バッファ42と、日本語の単語と、その仮名表記とを対応付けて記憶した複数の辞書からなる辞書群46と、入力文バッファ42に含まれる文を形態素解析して、漢字を含む単語があれば辞書群46を参照して仮名表記を検索し、仮名表記等の情報が付された形態素列を出力するための仮名変換部44とを含む。

20

【0028】

既に述べたように、漢字を含む単語の中には、複数の仮名表記を持つものがあり得る。音声合成のためには、それら複数の仮名表記の中で適切なものを選択する必要がある。音声合成システム30は、そのために、仮名変換部44及びいわゆるインターネット52に接続され、仮名変換部44がある単語Wについて複数の仮名表記候補 $R_k$  ( $k=1\sim K$ ;  $K$ は仮名表記候補の数)が存在することを検出したことに応答して、インターネット52上でその単語Wと仮名表記候補 $R_k$ とが共起するウェブページを、単語Wと仮名表記候補 $R_k$ の組合わせの各々について検索し、得られたウェブページのテキストを学習データとした機械学習による分類によって、単語Wにふさわしい仮名表記を決定して仮名変換部44に与えるための同形異音語解消処理部50とを含む。すなわち、このシステムでは、インターネット52上のウェブページの集合を、一つの仮想的なコーパスと見なして用例文書を収集している。

30

【0029】

音声合成システム30はさらに、仮名変換部44が出力する、仮名表記付入力文を記憶するための仮名表記入力文記憶部54と、音声合成のための、仮名表記に対応する音声を格納した音声データベース48と、仮名表記入力文記憶部54から仮名表記付入力文を読み出し、音声データベース48を参照して音声合成を行ない、アナログ音声信号を出力するための音声合成部56と、音声合成部56から出力されるアナログ音声信号を音声に変換するスピーカ58とを含む。

40

【0030】

本実施の形態では、同形異音語解消処理部50がインターネット52から検索するウェブページのテキストのうち、「スニペット」と呼ばれる部分を機械学習に用いる。「スニペット」とは、インターネットのいわゆる検索エンジンによる検索結果において、検索されたウェブページの内容を説明するための短文のことをいう。多くの場合、スニペットは、検索のキーワードとされた単語を含む部分のテキストからなる。なお、同形異音語解消処理部50によるウェブページの検索には、独自の検索プログラムを用いてもよいが、本

50

実施の形態では、既存の検索サービスサイトを利用し、単語Wと仮名表記候補R kとについてのAND検索をするクエリを検索サービスサイトに対して発行し、その結果を得ることで行なっている。なお、本実施の形態では、処理時間を安定させるため、検索件数の上限として、一回の検索について1000件という基準を設けている。

【0031】

音声合成部56による音声合成の部分は、本発明とは直接には関係しないため、その詳細についての説明はここでは行なわない。

【0032】

図2に、同形異音語解消処理部50の詳細なブロック図を示す。図2を参照して、同形異音語解消処理部50は、単語Wが与えられると、入力文バッファ42に記憶された入力文のうち、単語Wを中心とする所定長の窓に含まれる単語に基づいて行なう学習により、単語Wに関する所定の特徴ベクトルが与えられればその単語Wに対応する適切な仮名表記を出力するように学習可能な決定木82と、仮名変換部44から単語Wとその仮名表記候補R kとの組合せ(W, R k)を受け、それらが共起するウェブページのスニペットをインターネット52から収集し、その結果を用いて決定木82の学習を行なうための決定木作成部80と、仮名変換部44に接続され、仮名変換部44から、組合せ(W, R k)中の単語Wと、入力文中における単語Wを中心とする所定範囲の単語列85とが与えられると、それらから決定木82による分類に適合した分類用特徴ベクトルを作成し、出力するための分類用特徴ベクトル作成部84と、分類用特徴ベクトル作成部84から出力される分類用特徴ベクトルを決定木82に与え、その結果として決定木82から得られる、分類結果である仮名表記を仮名変換部44に与えるための分類実行部86とを含む。

【0033】

決定木作成部80は、単語Wとその仮名表記R kとの組合せ(W, R k)が与えられると、インターネット52上でそれらが共起するウェブページを検索するための検索部100と、検索部100により検索されたウェブページのスニペットの集合(以下単に「ウェブページの集合」と呼ぶ。)を組合せ(W, R k)ごとに記憶するための検索結果記憶部102と、組合せ(W, R k)のうちで、取得されたウェブページの件数の降順にウェブページの集合をソートし、件数が上位であるN件(Nは自然数)のみを選択することにより、決定木82のための学習データを作成するためのソート及び選択部104とを含む。本実施の形態では、このソート及び選択部104により選択された(W, R k)に含まれるN個の仮名表記候補R kが、仮名表記候補として残され、後の決定木の学習に用いられる。

【0034】

この処理では、他の文書集合は棄却され、それら文書集合の検索に用いられた仮名表記候補も棄却される。これは、単語Wと共起する頻度の低い仮名表記候補は候補として不適であると一般的に考えられるためである。もっとも、応用によってはそのように低頻度の仮名表記候補であっても棄却しない方がよい場合もあり得る。

【0035】

決定木作成部80はさらに、ソート及び選択部104により作成された学習データを記憶するための学習データ記憶部106と、検索対象となっている単語Wについて、学習データ記憶部106に記憶されている、その単語Wに関して検索された仮名表記候補R kのウェブページの各々から、所定の学習用特徴ベクトルを作成するための学習用特徴ベクトル作成部108と、学習用特徴ベクトル作成部108の作成した特徴ベクトルを記憶するための特徴ベクトル記憶部110と、特徴ベクトル記憶部110に記憶された特徴ベクトルを用いて決定木82を学習させるための決定木学習部112とを含む。

【0036】

図3は、図2に示す検索部100を実現するためのコンピュータプログラムのフローチャートである。図3を参照して、このプログラムは、ある単語Wについての仮名表記の候補R k(k = 1 ~ K)の各々について繰返されるステップ130 ~ 134の3つのステップを含む。

10

20

30

40

50

## 【 0 0 3 7 】

ステップ 1 3 0 では、クエリ「単語  $W$  and 単語  $R_k$ 」でウェブページを上限件数  $MAX = 1000$  件で検索する要求をインターネット上の検索エンジンに送信する。

## 【 0 0 3 8 】

ステップ 1 3 2 では、その検索結果として、単語  $W$  と仮名表記候補  $R_k$  ( $k = 1 \sim N$ ) とを含むスニペットの集合  $\{S_n(W, R_k)\}$  ( $n = 1 \sim L_k, k = 1 \sim K$ ) を取得する。ただしここで  $L_k$  は単語  $W$  と仮名表記候補  $R_k$  との組合せに対して得られた検索結果の数である。

## 【 0 0 3 9 】

ステップ 1 3 4 では、各集合  $S_n$  から仮名表記候補  $R_k$  を削除することで、検索結果のスニペットの集合  $\{(T_n(W), R_k) | n = 1 \sim L_k\}$  を作成する。

## 【 0 0 4 0 】

以上の 3 つのステップは、単語  $W$  に対する仮名表記候補  $R_k$  の全てに対して繰返される。

## 【 0 0 4 1 】

図 2 に示す検索部 1 0 0 の機能はこのようなプログラムで実現される。

## 【 0 0 4 2 】

なお、ソート及び選択部 1 0 4 によって、検索件数が上位  $N$  個のスニペットの集合  $\{(T_n(W), R_k) | n = 1 \sim L_k\}$  が抽出され、学習データ記憶部 1 0 6 に学習データとして記憶されるものとする。

## 【 0 0 4 3 】

図 4 に、図 2 に示す学習用特徴ベクトル作成部 1 0 8 の構成をブロック図形式で示す。図 4 を参照して、学習用特徴ベクトル作成部 1 0 8 は、学習データ記憶部 1 0 6 に記憶された学習データのスニペットの集合  $\{(T_n(W), R_k) | n = 1 \sim L_k\}$  に含まれる各スニペットから、そのスニペット中に存在する単語  $W$  をはさんで前後それぞれ  $M$  個 (合計  $2M$  個) の単語群 (これら合計  $2M$  個の単語群を「窓」と呼ぶ。) を抽出するための抽出部 1 5 0 と、学習データ記憶部 1 0 6 に記憶された学習データ 1 0 6 に出現する、単語  $W$  以外の語彙によって、決定木 8 2 (図 2 参照) の学習に用いる分類用特徴ベクトルの構成を決定するためのベクトル構成決定部 1 5 2 と、ベクトル構成決定部 1 5 2 により決定された特徴ベクトルの構成に従い、抽出部 1 5 0 によりスニペットごとに抽出された単語群に基づいて各スニペットの特徴ベクトルの要素を算出して、各スニペットの特徴ベクトルを作成し、特徴ベクトル記憶部 1 1 0 に記憶させるための要素算出部 1 5 4 とを含む。

## 【 0 0 4 4 】

図 5 に、単語  $W$  を中心とする「窓」の構成を模式的に示す。図 5 を参照して、学習用のスニペット 1 7 0 の単語列のうち、単語  $W$  を中心としてその前後に存在する単語列を、単語  $W$  を含めて、「 $W_{-m}, W_{-(m-1)}, W_{-(m-2)}, \dots, W_{-2}, W_{-1}, W, W_1, W_2, \dots, W_{m-2}, W_{m-1}, W_m$ 」と書くことができる。単語  $W$  を中心とし、その前の  $m$  個の単語からなる単語列 1 7 4 と、単語  $W$  より後の  $m$  個の単語からなる単語列 1 7 6 とを含む単語列により、窓長  $2m$  の窓 1 7 2 が構成される。本実施の形態では、窓長を  $2M$  とする。

## 【 0 0 4 5 】

ベクトル構成決定部 1 5 2 は、次のようにして特徴ベクトルの構成を決定する。すなわち、ベクトル構成決定部 1 5 2 は、学習データ記憶部 1 0 6 に存在する学習データ内に出現する単語の頻度を各単語について算出する。ベクトル構成決定部 1 5 2 はさらに、頻度が上位である  $H$  個の単語のみを選択する。ベクトル構成決定部 1 5 2 はさらに、特徴ベクトルの次元を  $H$  次元とし、1 番目  $\sim H$  番目の要素を、それぞれ頻度が 1 位  $\sim H$  位の単語に対応付ける。これにより特徴ベクトルの構成が決定される。この特徴ベクトルの要素数は  $H$  個である。各要素は 0 又は 1 の値をとる。各要素は、その要素に対応する単語がスニペット中の単語  $W$  を中心とする窓長  $2M$  の窓内に出現すると 1 の値となり、出現しないと 0

10

20

30

40

50



の値となる。

【0046】

従って、ある学習用のスニペット  $T_i$  について要素算出部 154 が行なう処理は次のような処理である。すなわち、要素算出部 154 は、このスニペット  $T_i$  に対応する  $H$  次元の特徴ベクトルの各要素について、対応する単語がスニペット  $T_i$  中の、単語  $W$  を中心とする窓長  $2M$  の窓の中に出現するか否かを調べる。その要素の値は、その単語が出現すれば 1、出現しなければ 0 となる。この処理を  $H$  個の要素の全てについて行なうことにより、スニペット  $T_i$  の特徴ベクトル  $V_i$  が算出される。この特徴ベクトル  $V_i$  と、その特徴ベクトルが得られた組合せ  $(W, R_k)$  の仮名表記候補  $R_k$  とを互いに関連付けて (特徴ベクトルに対する正解が仮名表記候補  $R_k$  であるとして) 決定木 82 の学習に用いる。

10

【0047】

図 2 に示す分類用特徴ベクトル作成部 84 が行なう分類用の特徴ベクトルの作成も、基本的にはこれと同様である。すなわち、分類用特徴ベクトル作成部 84 は、学習用特徴ベクトル作成部 108 のベクトル構成決定部 152 (図 4 参照) から、特徴ベクトルの各要素に対応する単語に関する情報を受け、処理対象となる単語  $W$  について、その単語  $W$  を中心とする窓長  $2M$  の窓内に所定の単語が出現するか否かによって、単語  $W$  に対する分類用の特徴ベクトルを作成する。すなわち、この特徴ベクトルは、学習用特徴ベクトル作成部 108 によって作成される特徴ベクトルと全く同じ構成となる。

【0048】

決定木学習部 112 は、機械学習によって決定木 82 の学習を行なう。この学習方式については機械学習の分野で慣用されている事項であるので、ここではその詳細な説明は行なわない。

20

【0049】

図 6 に、本実施の形態に係る要素算出部 154 により作成される決定木の一例である、「佐原」という単語に関する決定木 200 を示す。図 6 を参照して、この決定木は、4 つの中間のノード 210, 212, 214 及び 216 と、5 つの終端のノード 230, 232, 234, 236 及び 238 を含み、各ノード 210, 212, 214 及び 216 では、それぞれ窓内の単語が特定の条件を満たすか否かという質問がなされる。

【0050】

ノード 210 の質問は、単語「佐原」を中心とする窓長  $2M$  の窓内に、キーワード「千葉県」があるか、というものである。もしあればノード 230 に進み、「佐原」に対応する仮名表記として「さわら」が選択される。もしなければノード 212 に進む。なお、図 6 においては、「千葉県」のような具体的な単語について、窓内にあるか否かを聞いているが、実際の処理では、単語「佐原」の特徴ベクトル内において、単語「千葉県」に対応する要素 (ビット) の値が 1 か 0 かを調べることによってこの判定を行なっている。

30

【0051】

ノード 212 の質問は、キーワード「神奈川県」があるか、というものである。もしあればノード 232 に進み、「佐原」に対応する仮名表記として「さはら」が選択される。もしなければノード 214 に進む。

【0052】

ノード 214 の質問は、キーワード「成田」があるか、というものである。もしあればノード 234 に進み、「佐原」に対応する読みとして「さわら」234 が選択される。もしなければノード 216 に進む。

40

【0053】

ノード 216 の質問は、キーワード「横須賀」があるか、というものである。もしあればノード 236 に進み、「佐原」に対応する仮名表記として「さはら」が選択される。もしなければノード 238 に進み、「佐原」に対応する仮名表記として「さわら」が選択される。

【0054】

本実施の形態では、基本的に各単語に対し、決定木 200 が作成される。ある単語に対

50

応する特徴ベクトルが与えられると、その単語に対応する決定木を特徴ベクトルの各要素の値に従ってたどることにより、その単語の仮名表記が選択される。

【 0 0 5 5 】

[ 動作 ]

図 1 ~ 図 6 を参照して、上記した音声合成システム 3 0 は以下のように動作する。図 1 に示す入力文記憶部 4 0 には、音声合成の対象となる日本語の文が予め記憶される。そのうちの所定長部分が読出され、入力文バッファ 4 2 に記憶される。

【 0 0 5 6 】

仮名変換部 4 4 は、入力文バッファ 4 2 に記憶された文について辞書群 4 6 を参照して形態素解析を行なう。その結果、各単語の品詞、仮名表記（漢字の場合）、活用型、活用形などが決定される。もしも一つの単語について複数の仮名表記が得られた場合（すなわち同形異音語が存在する場合）、仮名変換部 4 4 は、その単語（単語 W とする。）と、仮名表記の組合せをそれぞれ同形異音語解消処理部 5 0 に与える。以下の説明では、構成のときに使用した表記を用いる。すなわち、ある単語 W に対して得られた K 個の仮名表記候補を仮名表記候補  $R_1 \sim R_K$  とする。

【 0 0 5 7 】

図 2 を参照して、検索部 1 0 0 は、単語 W と、仮名表記候補  $R_k$  ( $k = 1 \sim K$ ) との組合せ ( $W, R_k$ ) が与えられると、(単語 W and 単語  $R_k$ ) をクエリとしてインターネット 5 2 上の検索エンジンに検索件数上限 = 1 0 0 0 件という条件で検索要求を送信する（図 3 のステップ 1 3 0）。そして、この検索要求に応答して検索エンジンから得られたウェブページのスニペットの集合  $\{ S_n(W, R_k) \}$  ( $n = 1 \sim L_k$ ) を取得する（図 3 のステップ 1 3 2）。ここで  $L_k$  はクエリ (単語 W and 単語  $R_k$ ) に対して得られた検索結果（ウェブページ）の数である。このスニペットの集合の各々から単語  $R_k$  を削除して得られた検索結果のスニペットの集合が検索結果記憶部 1 0 2 に記憶される（図 3 のステップ 1 3 4）。これらスニペットの集合は、(単語 W, 仮名表記候補  $R_k$ ) の組合せごとに得られる。スニペットの集合の各々の要素の数  $L_k$  の上限 MAX は、本実施の形態では、上記したように 1 0 0 0 である。

【 0 0 5 8 】

検索部 1 0 0 は、単語 W と仮名表記候補  $R_k$  との組合せの各々に対し、上記した処理を実行する。すなわち、図 3 におけるステップ 1 3 0 ~ 1 3 4 の処理を各組合せに対し実行する。その結果、検索結果記憶部 1 0 2 には、これら組合せの各々について、検索結果のスニペットの集合  $\{ S_n(W, R_k) \}$  が記憶される。

【 0 0 5 9 】

ソート及び選択部 1 0 4 は、検索結果記憶部 1 0 2 に記憶されたスニペットの集合  $\{ S_n(W, R_k) \}$  を、その要素の数  $L_k$  をキーに降順にソートする。ソート及び選択部 1 0 4 はさらに、ソート結果のうち、上位 N 個のスニペットの集合  $\{ (T_n(W), R_k) \mid n = 1 \sim L_k \}$  を選択して、それらスニペットが得られた仮名表記要素  $R_k$  と関連付けて学習データ記憶部 1 0 6 に学習データとして記憶させる。すなわち学習データ記憶部 1 0 6 には、スニペットの集合のうち、検索結果の多かったものから順番に N 個が記憶される。

【 0 0 6 0 】

図 4 を参照して、学習用特徴ベクトル作成部 1 0 8 のベクトル構成決定部 1 5 2 は、学習データ記憶部 1 0 6 に学習データが記憶されると、これら学習データに出現する単語の頻度を各単語について算出する。ベクトル構成決定部 1 5 2 はさらに、出現頻度が上位 H 番目までの単語を選択する。特徴ベクトルの 1 番目 ~ H 番目の要素を出現頻度 1 位 ~ H 位の単語に対応付けることにより、特徴ベクトルの構成が決定される。ベクトル構成決定部 1 5 2 は、この特徴ベクトルの構成（すなわち特徴ベクトルの各要素に対応する単語に関する情報）を図 2 に示す分類用特徴ベクトル作成部 8 4 及び図 4 に示す要素算出部 1 5 4 に与える。

【 0 0 6 1 】

10

20

30

40

50

一方、抽出部 150 は、学習データ記憶部 106 に記憶されている各スニペットについて、単語 W を中心とする窓長 2 M の窓を抽出して要素算出部 154 に与える。

【0062】

要素算出部 154 は、ベクトル構成決定部 152 から与えられるベクトル構成に従い、抽出部 150 から与えられる窓に含まれる単語に基づいて、各スニペットの特徴ベクトルの各要素の値を算出する。その結果、各スニペットの特徴ベクトルが得られる。要素算出部 154 は、各スニペットを、そのスニペットが検索されたときの仮名表記候補 R k と関連付けて特徴ベクトル記憶部 110 に学習用データとして記憶させる。

【0063】

図 2 を参照して、決定木学習部 112 は、特徴ベクトル記憶部 110 に記憶された特徴ベクトルと、それら特徴ベクトルに関連付けられた仮名表記候補とを用いた機械学習により、決定木 82 の学習を行なう。

【0064】

以上の処理によって、決定木 82 は、ある単語 W を中心とする窓長 2 M の窓中の単語列、すなわち単語 W の文脈、を表す特徴ベクトルが与えられると、その文脈における単語 W の仮名表記として最適なものを出力するように機能するようになる。

【0065】

一方、仮名変換部 44 は、分類用特徴ベクトル作成部 84 に対し、同形異音語の解消を要求する単語 W と、入力文において単語 W を中心とする窓長 2 M の窓に含まれる単語列 85 とを与える。分類用特徴ベクトル作成部 84 は、単語 W について、仮名変換部 44 より与えられた、入力文中のその単語 W を中心とする窓長 2 M の窓に含まれる単語列 85 と、図 4 に示すベクトル構成決定部 152 から与えられたベクトル構成とによって、要素算出部 154 と同様の処理により単語 W の特徴ベクトルを作成し、分類実行部 86 に与える。

【0066】

分類実行部 86 は、この特徴ベクトルを決定木 82 に与える。決定木 82 は、単語 W を中心とする窓長 M から上記方法によって作成した特徴ベクトルが与えられると、単語 W の仮名表記として適切なものを出力するように学習済みである。分類実行部 86 は、この仮名表記を決定木 82 から得て、仮名変換部 44 に与える。

【0067】

仮名変換部 44 は、このようにして同形異音語解消処理部 50 から得られた仮名表記を、問題となった単語 W に形態素分析の結果と同様にして付加する。仮名変換部 44 はさらに、形態素解析が終わり、品詞、仮名表記（漢字の場合）、活用型、活用形などの情報が付された形態素列を音声合成部 56 に与える。この場合、同形異音語については既に同形異音語解消処理部 50 により解消されているため、一つの単語には一つの仮名表記しか付されていない。

【0068】

音声合成部 56 は、与えられた形態素列に基づき、形態素に付された仮名表記などを用いて音声データベース 48 から適切な音声波形を抽出し、波形接続処理によって合成音声波形データを作成し、さらにこの合成音声波形データをアナログ変換してスピーカ 58 に与える。スピーカ 58 はこの音声信号を音声に変換する。

【0069】

以上のように音声合成システム 30 によれば、入力文記憶部 40 に記憶された入力文に同形異音語が含まれていても、同形異音語解消処理部 50 によって同形異音語が解消され、一つの仮名表記のみがその単語に割当てられる。インターネット 52 上のウェブページをいわば仮想的なコーパスとして用い、自動的にこの同形異音語の解消のための決定木の学習が行なわれる。人手で学習データを作成する必要がなく、同形異音語の解消のための手間を従来と比較してはるかに少なくできる。さらに、インターネット 52 上で検索されるウェブページは多数の人により作成されたものであるため、少数の人が学習データを作成する場合と比較して、学習データの偏りが少なく、そのカバーする範囲も広がる。従って、同形異音語の解消の信頼性が従来より高くなるという効果がある。

10

20

30

40

50

## 【 0 0 7 0 】

## 〔コンピュータによる実現〕

上記した第 1 の実施の形態に係る音声合成システム 3 0 は、既に述べたようにコンピュータハードウェア及び当該コンピュータハードウェア上で実行されるコンピュータソフトウェアにより実現される。図 7 に音声合成システム 3 0 を実現するための一般的なコンピュータシステム 2 5 0 の外観を示し、図 8 にこのコンピュータシステム 2 5 0 の内部構成をブロック図形式で示す。

## 【 0 0 7 1 】

図 7 を参照して、コンピュータシステム 2 5 0 は、コンピュータ 2 6 0 と、いずれもコンピュータ 2 6 0 に接続されるモニタ 2 6 2、キーボード 2 6 6、マウス 2 6 8、マイクロホン 2 9 0 及び一対のスピーカ 5 8 とを含む。コンピュータ 2 6 0 には、DVD (Digital Versatile Disc) の再生及び記録が可能な DVD ドライブ 2 7 0 と、所定の規格に従った半導体メモリ記憶装置が装着可能なメモリポート 2 7 2 とが備えられている。コンピュータ 2 6 0 の内部構成については図 8 を参照して後述する。

10

## 【 0 0 7 2 】

図 8 を参照して、コンピュータ 2 6 0 は、図 7 に示す DVD ドライブ 2 7 0 及びメモリポート 2 7 2 に加え、CPU (中央演算処理装置) 2 7 6 と、CPU 2 7 6 に接続されたバス 2 8 6 と、いずれもバス 2 8 6 に接続された ROM (読出専用メモリ) 2 7 8、RAM (ランダムアクセスメモリ) 2 8 0、ハードディスク 2 7 4、ネットワークインタフェース 2 9 6、及びサウンドボード 2 8 8 を含む。

20

## 【 0 0 7 3 】

DVD ドライブ 2 7 0 には、DVD 2 8 2 が装着される。メモリポート 2 7 2 には半導体メモリ記憶装置 2 8 4 が装着される。CPU 2 7 6 は、バス 2 8 6 並びに DVD ドライブ 2 7 0 及びメモリポート 2 7 2 をそれぞれ介して、DVD 2 8 2 及びメモリ 2 8 4 をアクセスできる。

## 【 0 0 7 4 】

キーボード 2 6 6、マウス 2 6 8、モニタ 2 6 2 はいずれも図示しないインタフェースを介してコンピュータ 2 6 0 のバス 2 8 6 に接続される。スピーカ 5 8 及びマイクロホン 2 9 0 は、サウンドボード 2 8 8 に接続される。このコンピュータシステム 2 5 0 において、CPU 2 7 6 で実行される音声合成プログラムは、最終的にはデジタル形式の音声波形データを生成する。サウンドボード 2 8 8 はその音声波形データを CPU 2 7 6 から受取ると、アナログ信号に変換してスピーカ 5 8 を介して音声を発生させる処理をする。

30

## 【 0 0 7 5 】

上記実施の形態における入力文記憶部 4 0、辞書群 4 6、仮名表記入力文記憶部 5 4、音声データベース 4 8、検索結果記憶部 1 0 2、学習データ記憶部 1 0 6、特徴ベクトル記憶部 1 1 0 等は、RAM 2 8 0、ハードディスク 2 7 4、DVD ディスク 2 8 2、半導体メモリ記憶装置 2 8 4 のいずれでも実現できる。実際には、格納するデータの容量、読出し、書込みに要求される速度などによって、最も効率のよい記憶装置が各記憶部を実現するために選択される。

## 【 0 0 7 6 】

上記した第 1 の実施の形態に係る音声合成システム 3 0 を実現するためのコンピュータプログラムは、単一のプログラムでもよいし、複数のプログラムを組合せたものでもよい。特に、上記した各部の機能のうち、図 1 に示す仮名変換部 4 4 において行なわれる形態素解析処理、音声合成部 5 6 において行なわれる音声合成処理、図 2 に示す検索部 1 0 0 が実行するスニペットの検索処理、ソート及び選択部 1 0 4 が実行するソート及び選択処理、決定木学習部 1 1 2 が実行する決定木 8 2 の学習処理などについては、既に広く流布しているプログラムをそのまま使用できる。もちろん、これらプログラムは汎用に作成されているため、適切な調整を行なうことは要求されるが、それらはこの技術分野における通常の知識を持つ者にとっては、目的に照らして容易に実現できる範囲に留まる。

40

## 【 0 0 7 7 】

50

さらに、学習用特徴ベクトル作成部 108、分類用特徴ベクトル作成部 84 での処理についても、上記した説明に基づいて、当該技術分野の通常の知識を持つものであれば、仕様に応じて適宜実現することが可能である。

【0078】

これらプログラムは、例えば DVD ディスク 282 等のような記憶媒体に記憶され、又はインターネット 52 等のネットワークを通じて流通し、通常はハードディスク 274 等の不揮発外部記憶装置に記憶される。そして実行時にはハードディスク 274 から RAM 280 にコピーされ、CPU 276 内の図示しないプログラムカウンタにより指し示されるアドレスから読出された命令が CPU 276 により実行され、上記した所期の機能を実現する。コンピュータハードウェアそのものの動作形態については周知であるので、こ

10

【0079】

< 第 2 の実施形態 >

図 9 に、本発明の第 2 の実施の形態に係る、複数の定義を有する英語のアクロニムに対し、適切な定義を与える多義アクロニム解消システム 330 の構成をブロック図形式で示す。この多義アクロニム解消システム 330 は、アクロニムの近傍に、そのアクロニムの定義を与えている文書が多いこと、アクロニムの近傍に存在する単語は、その文書の分野によって何らかの傾向を持っていることを利用して、実施の形態 1 における同形異音語の解消と同じ原理によって、アクロニムに適切な定義を与えるものである。

【0080】

20

図 9 を参照して、この多義アクロニム解消システム 330 は、アクロニムを含む可能性のある入力文を記憶するための入力文記憶部 340 と、入力文記憶部 340 に記憶された入力文の所定部分を読込むための入力文バッファ 342 と、アクロニム及びその定義のリストよりなるデータからなる辞書群 346 と、入力文バッファ 342 に格納された入力文を形態素解析し、定義が付されていないアクロニムを見出すと、辞書群 346 によって当該アクロニムの定義を決定し、入力文中の当該アクロニムに当該定義を付して入力文を出力するためのアクロニム解釈部 344 とを含む。

【0081】

多義アクロニム解消システム 330 はさらに、アクロニム解釈部 344 から出力される、アクロニムに定義が付された入力文を記憶するためのアクロニム定義付入力文記憶部 354 と、アクロニム定義付入力文記憶部 354 に記憶された入力文の意味を理解するための文章理解装置 356 とを含む。

30

【0082】

既に述べたように、アクロニムの中には複数の定義を持つものもあり得る。そうした場合に、アクロニム解釈部 344 がアクロニムに複数の定義を付して出力することはできない。そうすると、文章理解装置 356 における文章理解の障害となるからである。従って、入力文中で定義されていないアクロニムに対し、複数の定義が辞書群 346 から見出された場合、何らかの手段によりそれらの中の適切な一つを自動的に選択できるようにする必要がある。

【0083】

40

こうした問題を解決するために、本実施の形態に係る多義アクロニム解消システム 330 は、アクロニム解釈部 344 及びインターネット 52 に接続され、アクロニム解釈部 344 から、アクロニムと、そのアクロニムに対して得られた複数の定義候補と、アクロニムの前後の所定の窓中に存在する単語列とが与えられると、インターネット 52 をコーパスとして用いた学習処理により、与えられた複数の定義候補のうち、与えられた単語列に対して最も適切と思われるものを選択し、アクロニム解釈部 344 に与える処理を行なうための多義アクロニム解消処理部 350 を含む。

【0084】

多義アクロニム解消処理部 350 の構成の詳細についてはここでは述べないが、多義アクロニム解消処理部 350 の構成及び動作は第 1 の実施の形態における同形異音語解消処

50

理部 50 と同様である。すなわち多義アクロニム解消処理部 350 は、以下の手順でアクロニムに対する適切な定義を決定する。

【0085】

(1) アクロニム A と定義候補  $D_k$  ( $k = 1 \sim K$  : K は定義候補の数) が与えられると、定義候補  $D_k$  の各々について、アクロニム A と定義候補  $D_k$  とが共起するウェブページのスニペットに対する検索要求をインターネット 52 上の検索エンジンに与える。

【0086】

(2) 検索結果として、アクロニム A と定義候補  $D_k$  とを含むスニペットの集合  $\{S_n(A, D_k)\}$  ( $n = 1 \sim L_k$ ) (ただし  $L_j$  ( $j = 1 \sim k$ ) はアクロニム A と定義候補  $D_k$  との組合せに対して検索されたスニペットの数を表す。) を取得する。

10

【0087】

(3) このスニペットの集合  $\{S_n(A, D_k)\}$  の各々から、定義候補  $D_k$  を削除することによって、検索結果のスニペットの集合  $\{(T_n(A), D_k) \mid n = 1 \sim L_k\}$  を作成する。

【0088】

(4) 上記した 3 つの処理を、全ての定義候補  $D_k$  に対して繰り返す。

【0089】

(5) 検索されたウェブページのスニペットの集合  $S_n$  を、それらに含まれるウェブページの数 (検索結果の数) の降順でソートし、さらにその内で上位 N 個のみを選択することで、N 個の学習用のスニペットの集合  $\{(T_n(A), D_k) \mid n = 1 \sim L_k\}$  が抽出され、学習データとして記憶される。

20

【0090】

(6) この学習データを用い、図 4 に示す学習用特徴ベクトル作成部 108 と全く同様にして学習用の複数個の特徴ベクトルが作成される。特徴ベクトルの作成の仕方も第 1 の実施の形態の場合と全く同様である。特徴ベクトルの作成時の窓長も第 1 の実施の形態と同様、 $2M$  と表すことにする。

【0091】

(7) これらの特徴ベクトルと、それら特徴ベクトルを与えたスニペットが検索されたときの検索に用いられた定義候補とを関連付けて学習用のデータとする。

【0092】

(8) この学習用のデータを用い、決定木の学習を行なう。この学習の結果、決定木は、入力文のうち、多義解消の対象となるアクロニム A を中心とする窓長  $2M$  に含まれる単語により作成される特徴ベクトルが与えられると、そのアクロニムに対する適切な定義を出力するようになる。

30

【0093】

(9) 入力文の中の、多義解消の対象となるアクロニム A を中心とし、窓長  $2M$  の窓から決定木のための特徴ベクトルを作成する。

【0094】

(10) この特徴ベクトルを決定木に与えることにより、決定木からはアクロニム A の定義を一つだけ選択する出力が得られる。この出力を多義アクロニム解消処理部 350 からアクロニム解釈部 344 に与えることにより、アクロニム解釈部 344 は当該アクロニムに対し、多義アクロニム解消処理部 350 から与えられたただ一つの定義を付して、アクロニム定義付入力文記憶部 354 に出力できる。

40

【0095】

< 第 3 の実施の形態 >

図 10 に、第 3 の実施の形態に係る日本語 - 英語の自動翻訳システム 430 のブロック図を示す。図 10 を参照して、この自動翻訳システム 430 は、日本語の入力文を記憶するための日本文記憶部 440 と、日本文記憶部 440 に記憶された日本文の所定量を記憶するための入力文バッファ 442 と、日本語から英語への 1 又は複数の辞書からなる辞書群 446 と、自動翻訳の前処理として、入力文バッファ 442 に記憶された日本文を形態

50

素解析し、各単語について辞書群 4 4 6 を参照して英語の訳語を割当て、出力するための訳語決定部 4 4 4 と、このように前処理された訳語付日本語を記憶するための訳語付日本語記憶部 4 5 4 と、訳語付日本語記憶部 4 5 4 に記憶された訳語付日本語を、その訳語を使用しながら英語に翻訳する自動翻訳装置 4 5 6 とを含む。

【 0 0 9 6 】

しかし、既に述べたとおり、入力される一つの日本語単語に複数の英語の訳語候補が存在する場合があります。そうしたときにそれら複数の英語の訳語候補を日本語単語にそのまま付して訳語決定部 4 4 4 から出力すると、自動翻訳装置 4 5 6 における翻訳に支障が生ずる。そのために、何らかの手段でこれら複数の訳語候補の中から適切なものを選択する必要があります。

10

【 0 0 9 7 】

そのために、本実施の形態に係る自動翻訳システム 4 3 0 は、訳語決定部 4 4 4 及びインターネット 5 2 に接続され、訳語決定部 4 4 4 から、日本語の単語と、その単語に対して得られた複数の訳語候補と、入力文における当該日本語の単語の前後の所定の窓中に存在する単語列とが与えられると、インターネット 5 2 をコーパスとして用いた学習処理により、与えられた複数の訳語候補のうち、与えられた単語列に対して最も適切と思われるものを選択し、訳語決定部 4 4 4 に与える処理を行なうための多義訳語解消処理部 4 5 0 を含む。

【 0 0 9 8 】

多義訳語解消処理部 4 5 0 の構成の詳細についてはここでは述べないが、多義訳語解消処理部 4 5 0 における処理が、第 1 の実施の形態の同形異音語解消処理部 5 0 における処理、及び第 2 の実施の形態の多義アクロニム解消処理部 3 5 0 における処理と同一であり、従ってその構成も同形異音語解消処理部 5 0 の構成と同一であることが理解されるであろう。

20

【 0 0 9 9 】

訳語決定部 4 4 4 は、入力文バッファ 4 4 2 中の文を読み出し、形態素解析して、各単語について辞書群 4 4 6 を参照することにより英語の訳語を割当て、訳語付日本語記憶部 4 5 4 に出力していく。複数の訳語候補が一つの日本語単語について出現した場合、訳語決定部 4 4 4 はその日本語単語と、複数の訳語候補とを多義訳語解消処理部 4 5 0 に引渡し、多義性の解消を依頼する。多義訳語解消処理部 4 5 0 は、第 1 の実施の形態における同形異音語解消処理部 5 0 と全く同じ動作によって決定木を作成し、入力文のうち、与えられた日本語単語の前後の窓内の単語列を用いて特徴ベクトルを作成し、決定木に与えることにより適切な訳語候補を得て、訳語決定部 4 4 4 に返す。訳語決定部 4 4 4 は問題となった日本語単語に、多義訳語解消処理部 4 5 0 から与えられたただ一つの訳語を付し、訳語付日本語記憶部 4 5 4 に出力する。従って、自動翻訳装置 4 5 6 における自動翻訳処理に支障が生ずることはない。

30

【 0 1 0 0 】

以上、第 1 ~ 第 3 の実施の形態の説明から明らかなように、本発明に係る多義性の解消、又はあいまい性の解消は、自然言語処理の分野の広い範囲にわたり、容易に適用できる。しかも、多義性の解消を行なう部分の仕組みは基本的に同一でよい。もちろん、解消処理の細部において様々な設計事項はあり得るが、ある分野で有効な方式は、基本的にそのままの形で他の分野についても適用可能である。

40

【 0 1 0 1 】

例えば日本語と英語との間の翻訳のみならず、あらゆる言語の間の単語の翻訳に、言語の相違にかかわらず本発明に係る多義性又はあいまい性の解消をする装置を適用できる。第 1 の実施の形態における同形異音語の解消を行なう機構も、言語にかかわらずほとんどそのまま適用できる。もちろん、言語に特有の調整が必要な場合もあり得るが（例えば日本語における形態素解析）、その部分は自然言語処理での前提として必ず前もって行なわれているとすれば、多義性又はあいまい性の部分の仕組みは言語に係らず同一でよい。

【 0 1 0 2 】

50

従って、自然言語処理の分野の広い領域において、本発明を適用することができ、しかもある領域から別の領域への移植も極めて簡単に実現できる。

【0103】

<可能な変形例>

上記した実施の形態では、適切な仮名表記、アクリニムの定義、及び訳語を決定するために、決定木を用いた。しかし本発明は決定木を用いるものには限定されず、インターネットから収集した学習データによって、対象となる単語又は単語列がおかれた文脈（環境）によって、目的物として複数のうちからどれを選択するかを機械学習により学習できるものであれば、どのような分類方法でも用いることができる。例えば、ナイーブベイズ、決定リスト、k - 最近隣法、オンラインアルゴリズム、最大エントロピー法、サポートベクトルマシン、ブースティングなどを利用できる。

10

【0104】

また、上記した実施の形態では学習データとしてウェブページのスニペットを収集したが、本発明がそのような実施の形態に限定されないことはもちろんである。例えばウェブページ全体を処理の対象としてもよい。また、例えば一つの単語Wと仮名表記Rkとの組み合わせに対して収集するウェブページの数の上限MAXを1000に限定しているが、この数が自由に変更できることはいうまでもない。また、このような限定を用いないことも可能である。

【0105】

さらに、上記した実施の形態では、問題となる単語と、その単語と対となるべきいくつかの候補が与えられると、その時点でインターネットにアクセスし、決定木を作成している。しかし本発明はそのような実施の形態には限定されない。例えば、予め何らかのテスト文に対し、上記したような処理をすることにより、テスト文中に含まれる、何らかのあいまい性を持ついくつかの単語について、そのあいまい性を解消するための分類装置を予め準備しておいてもよい。そうした分類装置を多数の単語に対して一つずつ予め準備しておけば、その単語が与えられてから分類装置の学習を行なったりする必要はなく、直ちに適切な答えを与えることができる。もしもそれら複数の分類装置ではあいまい性が解消できない単語であれば、そのときに上記実施の形態で示したように改めて一つの分類装置を作成して適切な答えを得るようにすればよい。

20

【0106】

また、上記した第1の実施の形態では、ソート及び選択部104により選択される仮名表記候補は、検索部100によりヒットしたウェブページの数の多い上位N件（Nは複数）であった。第2の実施の形態及び第3の実施の形態の場合も同様である。しかし本発明はそのような実施の形態には限定されない。例えば、ソート及び選択部104の処理でヒット数の多かった最上位の1件の仮名表記候補のみを単語Wの仮名表記として採用してもよい。この場合には、決定木は1:1の分類を行なうものとして機能する。もっとも、この方法では単語Wの文脈が全く考慮されないので、結果の信頼性は低く、あいまい性の解消とはいえない。

30

【0107】

また、第1の実施の形態のソート及び選択部104の処理で、ヒットしたウェブページの数の多い上位N件ではなく、所定のしきい値以上のウェブページがヒットしたものを全て仮名表記候補として選択してもよい。又は、全ヒット数のうち、上位から各候補の割合を積算し、所定割合を超えるまでのものを、その数にかかわらず全て仮名表記候補として採用してもよい。

40

【0108】

さらに、上記実施の形態では、一つの単語を単位としてその意味候補を決定している。しかし本発明はそのような実施の形態には限定されない。意味候補の集合を作成するための辞書の見出しとして、例えば複数の単語からなる句を設けておくことにより、その句の意味についても、複数の意味集合の中から適切なものを選択できるようになる。

【0109】

50



そして、そのようにして得られた分類装置を随時蓄積しておくことにより、直ちにあいまい性を解消できる単語が増加することになり、好ましい。

【 0 1 1 0 】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味及び範囲内でのすべての変更を含む。

【 図面の簡単な説明 】

【 0 1 1 1 】

【 図 1 】 本発明の第 1 の実施の形態に係る音声合成システム 3 0 のブロック図である。 10

【 図 2 】 図 1 に示す同形異音語解消処理部 5 0 のブロック図である。

【 図 3 】 図 2 の検索部 1 0 0 を実現するためのコンピュータプログラムの制御構造を示すフローチャートである。

【 図 4 】 図 2 の学習用特徴ベクトル作成部 1 0 8 のブロック図である。

【 図 5 】 「窓」の概念について説明するための図である。

【 図 6 】 決定木の一例を模式的に示す図である。

【 図 7 】 第 1 の実施の形態に係る音声合成システム 3 0 を実現するコンピュータシステム 2 5 0 の外観を示す図である。

【 図 8 】 図 7 に示すコンピュータシステム 2 5 0 の内部構成を示すブロック図である。

【 図 9 】 本発明の第 2 の実施の形態に係る多義アクロニム解消システム 3 3 0 のブロック図である。 20

【 図 1 0 】 本発明の第 3 の実施の形態に係る自動翻訳システム 4 3 0 のブロック図である。

【 符号の説明 】

【 0 1 1 2 】

3 0 音声合成システム

4 0 , 3 4 0 入力文記憶部

4 2 , 3 4 2 , 4 4 2 入力文バッファ

4 4 仮名変換部

4 6 , 3 4 6 , 4 4 6 辞書群 30

4 8 音声データベース

5 0 同形異音語解消処理部

5 2 インターネット

5 4 仮名表記入力文記憶部

5 6 音声合成部

5 8 スピーカ

8 0 決定木作成部

8 2 決定木

8 4 分類用特徴ベクトル作成部

8 6 分類実行部 40

1 0 0 検索部

1 0 2 検索結果記憶部

1 0 4 ソート及び選択部

1 0 6 学習データ記憶部

1 0 8 学習用特徴ベクトル作成部

1 1 0 特徴ベクトル記憶部

1 1 2 決定木学習部

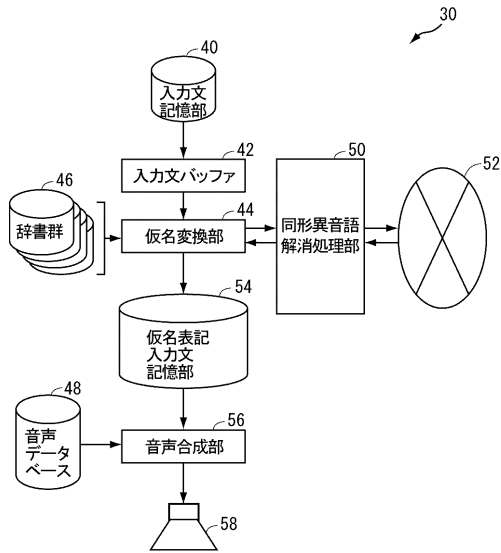
3 4 4 アクロニム解釈部

3 5 0 多義アクロニム解消処理部

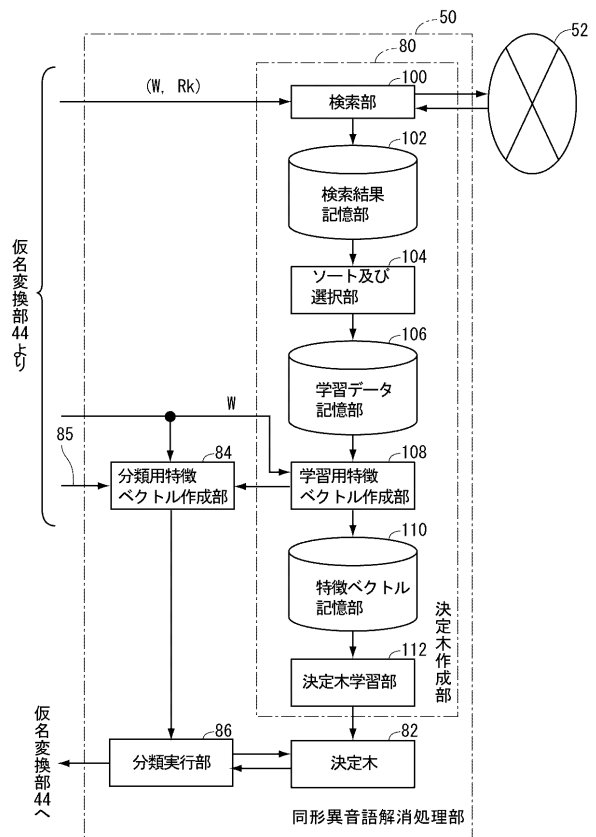
3 5 4 アクロニム定義付入力文記憶部 50

- 4 4 0 日本文記憶部
- 4 4 4 訳語決定部
- 4 5 0 多義訳語解消処理部
- 4 5 4 訳語付日本文記憶部
- 4 5 6 自動翻訳装置

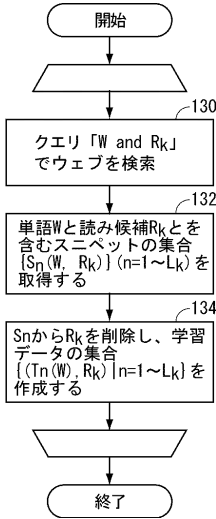
【図1】



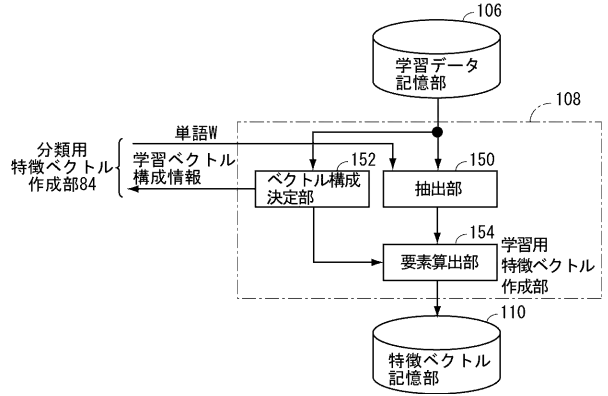
【図2】



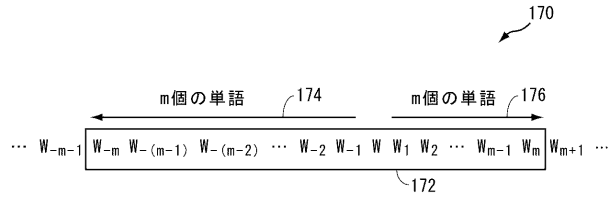
【図3】



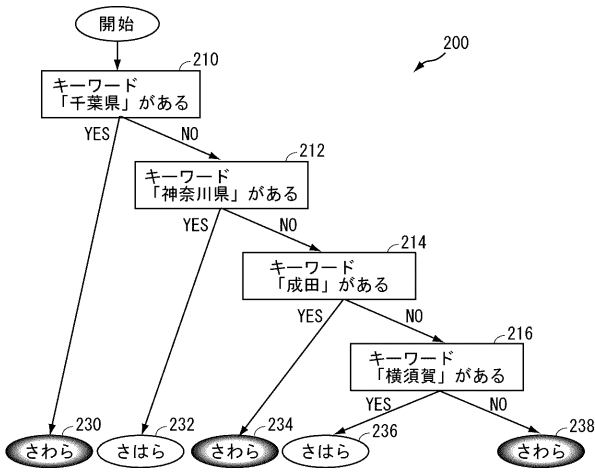
【図4】



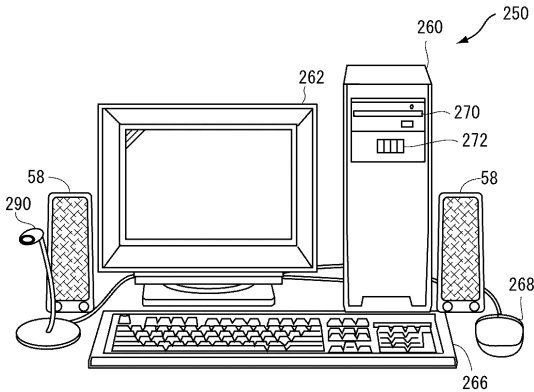
【図5】



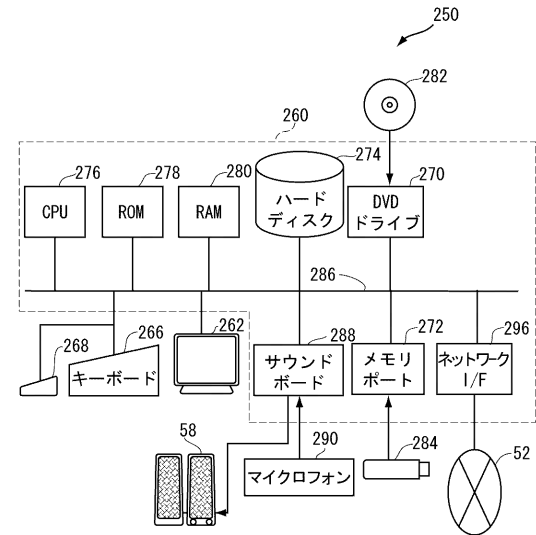
【図6】



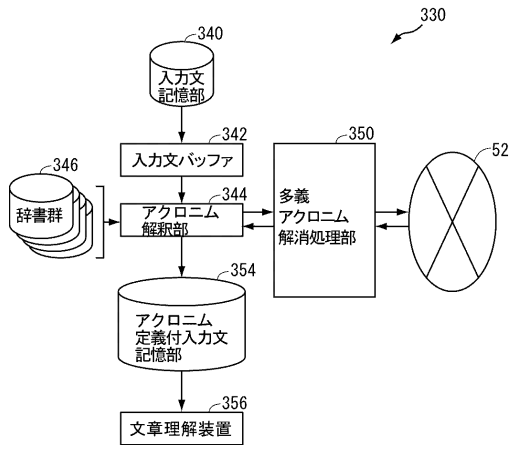
【図7】



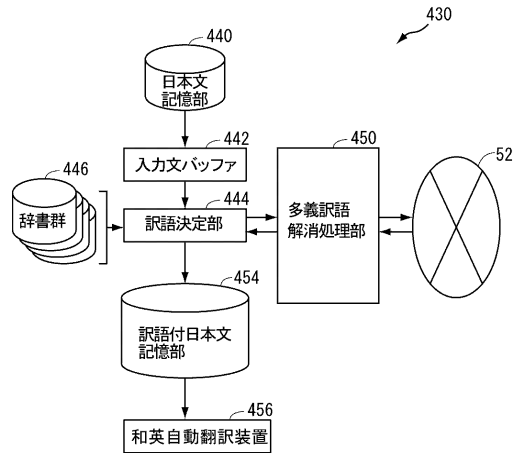
【図8】



【図9】



【図10】



---

フロントページの続き

審査官 成瀬 博之

- (56)参考文献 村田真樹 他4名,種々の機械学習手法を用いた多義解消実験,電子情報通信学会技術研究報告,日本,社団法人電子情報通信学会,2001年5月11日,第101巻第40号(NLC2001-1-8),7-14頁  
前田亮 他2名,言語横断情報検索におけるWeb文書群による訳語曖昧性解消,情報処理学会論文誌,日本,社団法人情報処理学会,2000年10月15日,第41巻No.SIG6(TOD7),12-21頁  
藤井敦,コーパスに基づく多義性解消,人工知能学会誌,日本,社団法人人工知能学会,1998年11月1日,第13巻第6号,904-911頁

(58)調査した分野(Int.Cl.,DB名)

G06F 17/20 - 17/28