

(19)日本国特許庁 ( J P )

# (12) 公開特許公報 ( A )

(11)特許出願公開番号

## 特開2003 - 141110

( P 2 0 0 3 - 1 4 1 1 1 0 A )

(43)公開日 平成15年 5月16日 (2003.5.16)

(51)Int.Cl.<sup>7</sup>

識別記号

F I

テ-マコード (参考)

G06F 17/27

G06F 17/27

J 5B091

審査請求 有 請求項の数 3 O L (全12頁)

(21)出願番号 特願2001 - 331458( P 2001 - 331458)

(71)出願人 301022471

独立行政法人通信総合研究所

東京都小金井市貫井北町 4 - 2 - 1

(22)出願日 平成13年10月29日(2001.10.29)

(72)発明者 村田 真樹

東京都小金井市貫井北町 4 - 2 - 1 独立

行政法人通信総合研究所内

(72)発明者 井佐原 均

東京都小金井市貫井北町 4 - 2 - 1 独立

行政法人通信総合研究所内

(74)代理人 100090893

弁理士 渡邊 敏

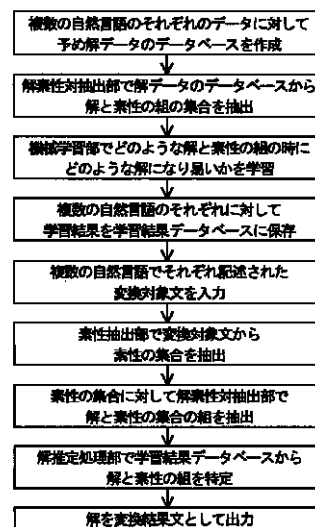
F タ-ム(参考) 5B091 AA01 AA15 BA03 CA01 EA01

(54)【発明の名称】複数言語入力での言語処理方法及び言語処理装置

(57)【要約】

【課題】 大量の変換規則を用意する必要がなく、複数の自然言語で記述された変換対象文から他の自然言語及び/又は同じ自然言語で記述された変換結果文への言語変換及び/又は複数の自然言語における言語解析を行うことができる複数言語入力での言語処理方法及び言語処理装置を提供する。

【解決手段】 複数の自然言語で記述された処理対象文から、他の自然言語及び/又は同じ自然言語で記述された処理結果文への処理を行う言語処理方法であって、前記処理を行う際に、前記他の自然言語及び/又は前記同じ自然言語への言語変換及び/又は前記同じ自然言語で記述されたどのような文又は文章になり易いかなどの言語解析を学習させる機械学習手法を用いたものである。



## 【特許請求の範囲】

【請求項 1】 複数の自然言語で記述された処理対象文に対して、他の自然言語及び／又は同じ自然言語で記述された処理結果文への言語変換及び／又は特定の言語現象を明らかにする言語解析を行う言語処理方法であつて、

前記言語変換を行う際に、前記他の自然言語及び／又は前記同じ自然言語で記述されたどのような文又は文章等になり易いかを学習及び／又は前記言語解析を行う際に、前記言語現象の解析においてどのような解析結果になり易いかを学習させる機械学習手法を用いたことを特徴とする複数言語入力での言語処理方法。

【請求項 2】 複数の自然言語で記述された処理対象文に対して、他の自然言語及び／又は同じ自然言語で記述された変換結果文への言語変換及び／又は特定の言語現象を明らかにする言語解析を行う言語処理装置であつて、

前記複数の自然言語で記述された処理対象文を入力する入力手段と、

前記入力手段により入力された前記複数の自然言語のデータに、前記他の自然言語及び／又は前記同じ自然言語への言語変換結果である解の情報が付与された解データ及び／又は特定の言語現象を明らかにする言語解析結果である解の情報が付与された解データを保存する解データベース部と、

前記解データから、前記解と解析に用いる細かい情報の 1 単位である素性の組を抽出する解素性対抽出部と、前記言語変換の際に、前記解と素性の組から、どのような解になり易いかを学習及び／又は前記言語解析の際に、前記解と素性の組から、どのような解析結果になり易いかを学習する機械学習部と、前記機械学習部で学習した結果を保存する学習結果データベース部と、

入力された複数の自然言語で記述された処理対象文から、素性を取り出す素性抽出部と、前記素性抽出部から取り出された素性の集合に対して、前記学習結果データベース部に保存された前記学習した結果から解を特定する解推定処理部と、を備えたことを特徴とする複数言語入力での言語処理装置。

【請求項 3】 前記解素性対抽出部は、解と素性の組及び候補と素性の組を抽出することを特徴とする請求項 2 記載の複数言語入力での言語処理装置。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】この発明は、ある自然言語で記述された変換対象文を、他の自然言語及び／又は同じ自然言語で記述された処理結果文に言語変換及び／又は特定の言語現象を明らかにする言語解析を行う言語処理に関し、特に、複数の自然言語で記述された処理対象文

を他の自然言語及び／又は同じ自然言語で記述された処理結果文に言語変換をする際及び／又は特定の言語現象を明らかにする言語解析をする際に、機械学習手法を用いる複数言語入力での言語処理方法及び言語処理装置に関するものである。

## 【0002】

【従来の技術】言語処理には、各言語の形態素解析、構文解析、格解析等を行う言語解析と他の言語への翻訳を行う言語変換とがある。ある自然言語から他の自然言語言語変換する従来の言語処理技術として、機械翻訳がある。機械翻訳では、ある自然言語で記述された文又は文章等を他の自然言語で記述された文又は文章等に言語変換する。また、同一の言語間における文又は文章の言語変換、例えば、要約文を自動生成あるいは文章を推敲する言語処理技術も用いられるようになってきている。

【0003】前記機械翻訳では、例えば、図 3 に示すように、CPU (中央演算処理装置)、メモリ、データ保存装置等からなるコンピュータ本体と周辺機器とから構成される言語処理装置 3 0 において、先ず、対象とする問題の答えである解のデータベースを作成して解データベース部 3 1 に保存しておく。前記解のデータベースには、入力されたある自然言語のデータに前記解の情報が付与されている。

【0004】次に、前記解データベース部 3 1 から各事例毎に、解素性対抽出部 3 2 で解と素性の集合の組を抽出する。前記素性は、解析に用いる情報の細かい 1 単位を意味し、前後の単語自体及び品詞、解析する単語自体及び品詞、解の単語及び品詞などである。

【0005】次に、前記解と素性の集合の組から、どのような素性の時にどのような解になり易いかを機械学習部 3 3 で学習する。この学習結果を解学習結果データベース部 3 4 に保存する。

【0006】ここまでは、予め準備しておく。ここから、先ず、解を求めたい文又は文章である変換対象文 3 5 を入力する。入力された変換対象文 3 5 から、素性抽出部 3 6 において素性の集合を取り出し、それらを解推定処理部 3 7 に渡す。

【0007】解推定処理部 3 7 では、渡された素性の集合の場合、どのような解になり易いかを前記解学習結果から特定する。最後に、特定された解を変換結果文 3 8 として出力する。

【0008】以上で示したように、機械翻訳では、機械学習を用い、ある自然言語で記述された文又は文章等から他の自然言語で記述されるどのような文又は文章になり易いかを特定して変換する。

【0009】また、前記形態素解析、構文解析、格解析等を行う言語解析においても同様に、解析に用いる素性を設定し、機械学習を用いてどのような解析結果になり易いかを学習させることが行われている。

【0010】また、ここで示した機械学習に基づく方法

の他に人手でパターンマッチ規則を作成し、これを用いて機械翻訳、言語解析を行うこともある。現状の実運用されている機械翻訳、言語解析ではむしろ、この人手で作成した規則に基づいて行っているものの方が主である。

【0011】また、同一自然言語間での文の言語変換処理では、一般に、変換前の語・句・文などのパターンと変換後の語・句・文などのパターンとの対からなる変換規則を大量に用意し、いわゆるパターン・マッチングによって入力文中に現れる処理前のパターンを探し出し、該当するパターンがあれば、それを処理後の語・句・文などのパターンに置き換える処理を行っている。

#### 【0012】

【発明が解決しようとする課題】しかしながら、前記従来の機械学習を用いた言語変換は、ある自然言語で記述された文又は文章などから他の自然言語及び/又は同じ自然言語で記述された文又は文章などへの言語変換に用いられているだけで、複数の自然言語で記述された文又は文章などから他の自然言語及び/又は同じ自然言語で記述された文又は文章などへの言語変換に用いられていなかった。また、前記従来の機械学習を用いた言語解析は、ある一つの自然言語における言語解析にしか用いられていなかった。

【0013】従って、前記言語変換あるいは前記言語解析を行おうとすると、従来のパターン・マッチングを用いて行うしかなく、この場合、大量の変換規則を用意しなければならないという問題点があった。

【0014】本発明は、前記従来の問題点を解決するためになされたもので、大量の変換規則を用意する必要がなく、複数の自然言語で記述された処理対象文から他の自然言語及び/又は同じ自然言語で記述された処理結果文への言語変換及び/又は複数の自然言語における言語解析を行うことができる複数言語入力での言語処理方法及び言語処理装置を提供することである。

【0015】また、機械学習を用いる方法は、複数の情報を素性によって容易に利用できるために、本課題の複数言語入力のように複数の情報が複雑に入力される課題に対して、まさにうってつけである。

#### 【0016】

【課題を解決するための手段】本発明は、複数の自然言語で記述された処理対象文に対して、他の自然言語及び/又は同じ自然言語で記述された処理結果文への言語変換及び/又は特定の言語現象を明らかにする言語解析を行う言語処理方法であって、前記言語変換を行う際に、前記他の自然言語で記述されたどのような文又は文章等になり易いかを学習及び/又は前記言語解析を行う際に、前記言語現象の解析においてどのような解析結果になり易いかを学習させる機械学習手法を用いたものである。

【0017】また、本発明は、複数の自然言語で記述さ

れた処理対象文に対して、他の自然言語及び/又は同じ自然言語で記述された処理結果文への言語変換及び/又は特定の言語現象を明らかにする言語解析を行う言語処理装置であって、前記複数の自然言語で記述された処理対象文を入力する入力手段と、前記入力手段により入力された前記複数の自然言語のデータに、前記他の自然言語及び/又は前記同じ自然言語への言語変換結果である解の情報が付与された解データ及び/又は特定の言語現象を明らかにする言語解析結果である解の情報が付与された解データを保存する解データベース部と、前記解データから、前記解と解析に用いる細かい情報の1単位である素性の組を抽出する解素性対抽出部と、前記言語変換の際に、前記解と素性の組から、どのような解になり易いかを学習及び/又は前記言語解析の際に、前記解と素性の組から、どのような解析結果になり易いかを学習する機械学習部と、前記機械学習部で学習した結果を保存する学習結果データベース部と、入力された複数の自然言語で記述された処理対象文から、素性を取り出す素性抽出部と、前記素性抽出部から取り出された素性の集合に対して、前記学習結果データベース部に保存された前記学習した結果から解を特定する解推定処理部とを備えたものである。

【0018】また、前記解素性対抽出部は、解と素性の組及び候補と素性の組を抽出するものでもよい。

#### 【0019】

【発明の実施の形態】以下、図面を参照して、本発明の実施形態について説明する。図1は、本発明による複数言語入力での言語処理装置の実施形態を示すブロック図である。

【0020】図1に示すように、本発明の実施形態の複数言語入力での言語処理装置10は、CPU(中央演算処理部)、データを一時保存するメモリ、データを保存するデータ保存部、例えば、ハードディスク等を有する本体部及び表示装置であるCRT等、必要に応じた周辺機器を備えたコンピュータで構成されており、複数の自然言語で記述された処理対象文を入力とし、これに対して他の自然言語及び/又は同じ自然言語への言語変換及び/又は特定の言語現象を明らかにする言語解析を行った結果を、処理結果文として出力する。

【0021】言語処理装置10は、複数の自然言語、ここでは2つの自然言語でそれぞれ記述された文又は文章などを入力する入力手段(図示せず)、例えばキーボードを備え、予め前記キーボードにより入力された前記2つの自然言語のそれぞれのデータに、他の自然言語及び/又は同じ自然言語への翻訳である言語変換及び/又は同じ自然言語への形態素解析、構文解析、格解析等の言語解析を行った処理結果である解の情報が付与された解データを保存する解データベース部11を有する。解データベース部11には、言語1と言語2についての解データがそれぞれ保存されている。

【0022】また、前記解の情報における解は、対象とする問題の答えであり、前記言語変換の場合、変換先の言語表現である。従って、前記解の情報は、変換先の言語表現に関する情報である。また、前記言語解析の場合、形態素解析であるならば、前記解は品詞であり、前記解の情報は品詞に関する情報である。

【0023】そして、言語処理装置10は、解データベース部11に保存されている解データから、前記解と解析に用いる細かい情報の1単位である素性の組(図中では、解-素性対と表す)とを抽出する解素性対抽出部(図中では、解-素性対抽出部12と表す)と、前記解と素性の組から、前記変換の際に、どのような解になり易いかを学習する機械学習部13と、機械学習部13で学習した結果を保存する学習結果データベース部14とを備えている。

【0024】解-素性対抽出部12では、解データベース部11に保存されている解データを取り出し、各事例ごとに、解と素性の組を抽出すると共に、機械学習部13で学習した結果、新たに得られた保存すべき解データを解データベース部11に保存する。

【0025】機械学習部13では、言語1と言語2のそれぞれについて、処理対象文15から処理結果文16に処理する際に、それぞれの解と素性の組からどのような解になり易いか、即ち、どのような解と素性の組み合わせの時に解である確率が高いかを学習し、学習した結果を学習結果データベース部14に保存する。この学習は、言語1と言語2のそれぞれに対して行い、学習結果はそれぞれ別々に分類され保存される。

【0026】更に、言語処理装置10は、入力された2つ自然言語の処理対象文15である言語1と言語2から、素性を取り出す素性抽出部17と、素性抽出部17から取り出された素性の集合に対して、学習結果データベース部14に保存された前記学習した結果から解を特定する解推定処理部18とを備えている。

【0027】処理対象文15の言語1と言語2は、それぞれ素性抽出部17に入力されてそれぞれの素性が取り出され、取り出されたそれぞれの素性の集合に対して、解-素性対抽出部12から解と素性の集合の組を取り出し、それを解推定処理部18に渡す。

【0028】解推定処理部18では、渡された解と素性の集合の組から、学習結果データベース部14に保存された学習した結果に基づき解を特定する。特定された解は、出力されて処理結果文16として得られ、必要に応じて保存される。

【0029】尚、解-素性対抽出部12は、解と素性の組を抽出するばかりでなく、必要に応じて、解と素性の組を抽出すると共に、解の候補となる解候補と素性の組を抽出するものでもよい。ここで、前記解の候補は、前記解以外の解の候補を意味する。

【0030】解と素性の集合の組を正例、解候補と素性

の集合の組を負例とすると、解若しくは解候補と素性の集合の組から、どのような解若しくは解候補と素性の集合の時に、正例である確率が高いかあるいは負例である確率が高いかを機械学習部13で学習し、その結果を学習結果データベース部に保存する。

【0031】解推定処理部18では、解-素性対抽出部12から抽出された解候補と素性の集合の全ての組に対して、渡された解候補と素性の組について正例、負例である確率を求め、最も正例である確率が高い解候補を解として出力し、処理結果文16として保存する。

【0032】次に、本発明による複数言語入力での言語処理方法の実施形態について説明する。本発明の実施形態の複数言語入力での言語処理方法は、複数の自然言語で記述された文又は文章等の処理対象文に対して、言語変換および言語解析を行う際に、前記他の自然言語及び/又は前記同じ自然言語で記述されたどのような文又は文章等になり易いか及び/又は言語解析においてどのような解析結果になり易いかを学習させる機械学習手法を用いたものである。

【0033】即ち、図2に示すように、複数の自然言語のそれぞれのデータに、他の自然言語及び/又は同じ自然言語への翻訳である言語変換及び/又は同じ自然言語への形態素解析、構文解析、格解析等の言語解析を行った処理結果である解の情報が付与された解データのそれぞれのデータベースを作成する。

【0034】次に、解素性対抽出部で、それぞれのデータベースから解データを取り出し、各事例ごとに、解と素性の組の集合を抽出する。

【0035】次に、機械学習部で、どのような解と素性の組の時にどのような解になり易いかを学習する。

【0036】次に、複数の自然言語のそれぞれについて学習した結果を学習結果データベースに保存する。

【0037】次に、複数の自然言語で記述された処理対象文を、それぞれ素性抽出部に入力し、それぞれの処理対象文からそれぞれの素性の集合を抽出する。

【0038】次に、それぞれの素性の集合に対して、解素性対抽出部から解と素性の集合の組を抽出し、それを解推定処理部に渡す。

【0039】次に、解推定処理部で、渡された解と素性の集合の組から、学習結果データベースから学習した結果に基づき解を特定する。

【0040】最後に、特定された解を処理結果文として出力する。

【0041】尚、解素性対抽出部は、解と素性の組及び解候補と素性の組を抽出するものでもよい。ここで、前記解の候補は、前記解以外の解の候補を意味する。

【0042】その場合には、解と素性の集合の組を正例、解候補と素性の集合の組を負例とすると、解若しくは解候補と素性の集合の組から、どのような解若しくは解候補と素性の集合の時に、正例である確率が高いかあ

るいは負例である確率が高いかを機械学習部で学習し、その結果を学習結果データベース部に保存する。

【0043】解推定処理部では、解素性対抽出部から抽出された解候補と素性の集合の全ての組に対して、渡された解候補と素性の組について正例、負例である確率を求め、最も正例である確率が高い解候補を解として出力する。

【0044】以上示したように、本発明の実施形態の複数言語入力での言語処理方法によれば、機械学習手法を用い、機械学習する際に、素性を複数の自然言語、例えば、2つの自然言語から取るだけで済むので、大量の変換規則を用意する従来の言語処理方法よりも処理が容易である。

【0045】

【実施例】（実施例1）日本語の言語解析である形態素解析の場合、図1に示す解・素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0046】日本語と英語の2つの自然言語のデータ（以下、対訳データという）の場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。対訳データは、「<きょうだい>で待つ。」と「I wait in Kyoto university.」であり、解は、「京大 名詞」である。

【0047】解析に用いる情報である素性として、次のものを用いる。1. 前の単語自体、2. 前の単語の品詞、3. 後の単語自体、4. 後の単語の品詞、5. 解析する単語自体、6. 解析する単語のとりうる品詞、7. 解の単語、8. 解の品詞、9. 日本語文と英語文の単語の一致数、10. 英語単語列、11. 解析する単語のとりうる品詞に解の品詞が含まれるか。

【0048】前記解データのデータベースから抽出される素性は、情報が無い時を<none>で記述すると、1.<none>、2.<none>、3.で（次の単語が「で」だけであることは既存の形態素解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい）、4.助詞（前記と同じ）、5.きょうだい、6.名詞（単語辞書を調べてとりうる品詞を探す。複数の品詞をとりうる場合もある）、7.京大、8.名詞、9.3（「待つ - wait」「で - in」「京大 - Kyoto university」の3つが一致する。この単語の一致の算出は現在解析対象としている部分（きょうだい）も含めて行う。また、現在解析対象としている部分については、解に記述している単語（京大）を用いて行う。日英で単語が一致するかどうかは対訳辞書を用いて行う）、10.「I」「wai

t」「in」「Kyoto」「university」（各単語が素性となる）、11.含まれる。

【0049】前記では解を用いているので、解と素性の集合の組は正例となる。ここで、「きょうだい」を「兄弟 名詞」とする解候補を用いた解候補と素性の集合の組である負例を考える。この場合、素性は次のようになる。1.<none>、2.<none>、3.で（次の単語が「で」だけであることは既存の形態素解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい）、4.助詞（前記と同じ）、5.きょうだい、6.名詞（単語辞書を調べてとりうる品詞を探す。複数の品詞をとりうる場合もある）、7.兄弟、8.名詞、9.2（「待つ - wait」「で - in」の2つが一致する。この単語の一致の算出は現在解析対象としている部分（きょうだい）も含めて行う。また、負例の場合は、現在解析対象としている部分については、解に記述していない他の候補の単語（兄弟）を用いて行う）、10.「I」「wait」「in」「Kyoto」「university」（各単語が素性となる）、11.含まれる。

【0050】前記情報を用いて、どのような場合に正例、負例になりやすいかを求めると、「日本語文と英語文の単語の一致数が多い」場合に正例になりやすく、「Kyotoという素性と京大という素性が共起する」場合に正例になりやすいということを学習する。

【0051】このような英語の情報をうまく利用したものから、「次の単語が助詞の『で』で、解の品詞が『名詞』である」場合に正例になりやすいという日本語ももとの性質を利用した学習も同時に行える。

【0052】「Kyotoという素性と京大という素性が共起する」場合に正例になりやすいといったことは、素性の共起を自動で考慮することができる機械学習システム、例えば、サポートベクトルマシン、を利用することで可能になる。また、素性の共起を自動で考慮しない学習アルゴリズムの場合は、人手で予めそのような共起を意味する素性を設定しておけばよい。

【0053】次に、入力する対訳データとして、「<とうだい>にいく。」と「I goto Tokyo university.」が与えられ、そのうち「とうだい」の部分形態素解析するように指示されたとする。この時、単語辞書等を調べて解の候補として、「灯台 名詞」と「東大 名詞」があがったとする。

【0054】先ず、「灯台 名詞」を対象とする。素性は、前記と同様な記述の仕方、1.<none>、2.<none>、3.に（次の単語が「に」だけであることは既存の形態素解析システムで特定してもよいし、もとの入力として与えられる対訳データにその情報があることにしてもよい）、4.助詞（前記同様）、5.とうだい、6.名詞（単語辞書を調べてとりうる品詞を探す。複数の品詞をとりうる場合もある）、7.灯

台、8. 名詞、9. 2 (「行く - go」、「に - to」の2つが一致する。この単語の一致の算出は現在解析対象としている部分(とうだい)も含めて行う。また、現在解析対象としている部分については、解の候補の単語(灯台)を用いて行う)、10. 「I」「go」「to」「university」(各単語が素性となる)、11. 含まれる。

【0055】次に、「東大 名詞」を対象とする。素性は、1. <none>、2. <none>、3. に(次の単語が「に」だけであることは既存の形態素解析システムで特定してもよいし、もとの入力として与えられる対訳データにその情報があることにしてもよい)、4. 助詞(前記同様)、5. とうだい、6. 名詞(単語辞書を調べてとりうる品詞を探す。複数の品詞をとりうる場合もある)、7. 東大、8. 名詞、9. 2 (「行く - go」「に - to」「東大 - Tokyo university」の3つが一致する。この単語の一致の算出は現在解析対象としている部分(とうだい)も含めて行う。また、現在解析対象としている部分については、解の候補の単語(東大)を用いて行う)、10. 「I」「go」「to」「university」(各単語が素性となる)、11. 含まれる。

【0056】ここで学習した結果と前記素性を用いて「灯台 名詞」と「東大 名詞」の正例である確率を求める。「日本語文と英語文の単語の一致数が多い」場合に正例になりやすいといった学習結果により、「東大 名詞」が正例である確率は、「灯台 名詞」のものよりも大きくなり、確率が大きい方の「東大 名詞」が解として出力される。ここでは、一単語を解析の対象とする場合のものを示した。

【0057】一文全体を解析の対象とする場合は、各単語ごとの解析を組み合わせることで実現できる。この場合、対訳データである「<とうだい>にいく。」と「I go to Tokyo university.」の解は、「東大 名詞、に助詞、行く 動詞」である。この時、解の候補を複数作り、その中から予め設定した評価値のよいものを選べばよい。解の候補は、「東大 名詞、に 助詞、行く 動詞」「灯台 名詞、に 助詞、行く 動詞」「と 接続詞、宇内 名詞、に 助詞、行く 動詞」である。予め設定する評価値は、各単語ごとに前記方法で前記正例である確率を求め、その一文全体での積とするとよい。この方法を高速に実現する手段として、ビタビアルゴリズム、ビームサーチが知られている。

【0058】また、前方から解析する方法をとる場合、前記素性の3に(次の単語が「に」だけであることは既存の形態素解析システムで特定してもよいし、もとの入力として与えられる対訳データにその情報があることにしてもよい)などの表現のうち、「もとの入力として与えられる対訳データにその情報があることにしてもよ

い」としていたが、前方から解析する方法をとる場合は、前方の解析結果があるため、「入力として与えられる対訳データにその情報」があることになる。また、後方のもも未だ解析していなかったとしても、候補を複数あげ、その一つ一つを解と仮定して解いていく場合は、後方の情報も「入力として与えられる対訳データにその情報」があることになる。

【0059】尚、解を「単語 品詞」としてその組み合わせによって一文全体の解を得る方法の他に、一文全体の解そのものを解として扱ってもよい。

【0060】(実施例2)英語の構文解析の場合、図1に示す解 - 素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0061】英語と日本語の2つの対訳データの場合を考える。前記解データは、英語と日本語の対訳データに解の情報が付与されたものである。以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。対訳データは、「She met a boy <with a picture> .」と「彼女は絵を持っている少年とあった。」であり、解は、「係り先 『a boy』」である。尚、対訳データ中で対象となるフレーズは“<”、“>”の記号で囲んでおく。

【0062】解析に用いる素性は、1. 解析するフレーズの意味的主辞の単語、2. 解析するフレーズの意味的主辞の単語の意味カテゴリ、3. 解析するフレーズの構文的主辞の単語、4. 解析するフレーズの構文的主辞の単語の意味カテゴリ、5. 解の係り先の単語、6. 解の係り先の単語の意味カテゴリ、7. 解の係り先の単語の品詞、8. 日本語単語列、9. 日本語文と英語文の対応する2つのフレーズのかかり受けの一致数、である。

【0063】前記から抽出される素性は、1. picture (「with a picture」の意味主辞は名詞句の主辞の「picture」となる。フレーズのどの部分が意味主辞で、どの部分が構文的主辞になるかは、予め文法を用いて定めておく)とよい)、2. 製品(どういう単語がどういう意味カテゴリになるかは、単語意味辞書を用いることで特定できる)、3. with、4. 前置詞、5. boy、6. 人、7. 名詞、8. 「彼女」「は」「絵」「を」「もっている」「少年」「と」「あった」(各単語が素性となる。また、これは既存の形態素解析システムで分割してもよいし、もとの対訳データで分割されたものが与えられていたとしてもよいし、もとの対訳データで分割されたものが与えられていたとしてもよい)、9. 4 (「she - met」と「彼女は - あった」、「met - a boy」と「少年と - あった」、「with - a picture」と「絵を - もっている」、「a boy - with」と「もっ

ている - 少年」の4つ。解析対象の「with」の係り先「a boy」との関係以外のフレーズの係り受けの情報は既存の構文解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい)。これらは、解の部分を用いているので、正例となる。

【0064】次に、負例として係り先を「met」としたものを考える。この場合、素性は、1. picture (「with a picture」の意味主辞は名詞句の主辞の「picture」となる。フレーズのどの部分が意味主辞で、どの部分が構文的な主辞になるかは、予め文法を用いて定めておく(とよい))、2. 製品(どういう単語がどういう意味カテゴリになるかは、単語意味辞書を用いることで特定できる)、3. with、4. 前置詞、5. meet、6. 知覚動詞、7. 動詞、8. 「彼女」「は」「絵」「を」「もっている」「少年」と「あった」(各単語が素性となる。また、これは既存の形態素解析システムで分割してもよいし、もとの対訳データで分割されたものが与えられていたとしてもよいし、もとの対訳データで分割されたものが与えられていたとしてもよい)、9. 3 (「she - met」と「彼女は - あった」、「met - a boy」と「少年と - あった」、「with - a picture」と「絵を - もっている」の3つ。解析対象のwithの係り先「a boy」との関係以外のフレーズの係り受けの情報は、既存の構文解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい)、となる。

【0065】前記情報を用いて、どういう場合に正例あるいは負例になりやすいかを求めると、「日本語文と英語文の構文リンクの一致数が多い」場合に正例になりやすいとか、「『もっている』という日本語単語素性と解析するフレーズの構文的な主辞の単語素性『with』と、解の係り先の単語の品詞素性『名詞』が共起する」場合に正例になりやすいとかを学習できる。このような日本語の情報をうまく利用するものの他に、英語の素性を用いるため、英語ももとの性質を利用した学習も同時に行える。

【0066】例えば、「製品を意味主辞にもつ『with』のフレーズは、『meet』よりも名詞にかかりやすい」など。尚、「『もっている』という日本語単語素性と解析するフレーズの構文的な主辞の単語素性『with』と、解の係り先の単語の品詞素性『名詞』が共起する」は、「with」を「もっている」と和訳する場合は、そのフレーズは名詞にかかりやすいということの意味する。

【0067】次に、入力する英語と日本語の対訳データとして、「She looked at a boy < with a telescope > .」「彼女は望遠鏡を持っている少年を見た。」が与えられ、そのうち「with a telescope」の部分の係り先

を求めるように指示されたとする。尚、「She looked at a boy with a telescope .」は構文的に曖昧な表現であり、「with a telescope」は「looked」にも「at a boy」にもかかりうる。「looked」にかかる場合は、「彼女は望遠鏡で少年を見た。」の意味になる。

【0068】ここまでのフレーズのまとめあげの結果などから、解の候補として「looked」と「a boy」があがったとして、先ず「looked」を対象とする。素性は、1. telescope (「with a telescope」の意味主辞は名詞句の主辞の「telescope」となる。フレーズのどの部分が意味主辞で、どの部分が構文的な主辞になるかは予め文法を用いて定めておく(とよい))、2. 製品(どういう単語がどういう意味カテゴリになるかは単語意味辞書を用いることで特定できる)、3. with、4. 前置詞、5. look、6. 知覚動詞、7. 動詞、8. 「彼女」「は」「望遠鏡」「を」「もっている」「少年」「見た」(各単語が素性となる。またこれは既存の形態素解析システムで分割してもよいし、もとの対訳データで分割されたものが与えられていたとしてもよい)、9. 3 (「she - looked」と「彼女は - 見た」、「looked - at a boy」と「少年を - 見た」、「with - a telescope」と「望遠鏡を - もった」の3つ。解析対象の「with」の係り先「a boy」との関係以外のフレーズの係り受けの情報は既存の構文解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい)、である。

【0069】次に、「a boy」を対象とする。素性は、1. telescope (「with a telescope」の意味主辞は名詞句の主辞の「telescope」となる。フレーズのどの部分が意味主辞で、どの部分が構文的な主辞になるかは予め文法を用いて定めておく(とよい))、2. 製品(どういう単語がどういう意味カテゴリになるかは、単語意味辞書を用いることで特定できる)、3. with、4. 前置詞、5. boy、6. 人、7. 名詞、8. 「彼女」「は」「望遠鏡」「を」「もっている」「少年」「見た」(各単語が素性となる。またこれは既存の形態素解析システムで分割してもよいし、もとの対訳データで分割されたものが与えられていたとしてもよい)、9. 4 (「she - looked」と「彼女は - 見た」、「looked - at a boy」と「少年を - 見た」、「with - a telescope」と「望遠鏡を - もった」、「a boy - with」と「もった - 少年」の4つ。解析対象の「with」の係り先「a boy」との関係以外のフレーズの係り受けの情報は既存の構文解析システムで特定してもよいし、もとの対訳データにその情報があることにしてもよい)、である。

【0070】ここで学習した結果と前記素性を用いて「looked」と「a boy」の正例である確率を求める。「日本語文と英語文の構文リンクの一致数が多い」場合に正例になりやすいとか、「『もっている』という日本語単語素性と解析するフレーズの構文的主辞の単語素性『with』と解の係り先の単語の品詞素性『名詞』が共起する」場合に正例になりやすいといった学習結果により、「a boy」が正例である確率は、「looked」のものよりも大きくなり、確率が大きい方の「a boy」が解として出力される。

【0071】ここでは、例えば、構文リンクの一致数を求める前記素性9が思うように動かなかったとしても（素性9は既存の構文解析システムなどを前提にするため、場合によっては動かない可能性がある）、「『もっている』という日本語単語素性と解析するフレーズの構文的主辞の単語素性「with」と解の係り先の単語の品詞素性『名詞』が共起する」場合に正例になりやすいという方の性質の方をうまく使うことで、「a boy」を正しく解として出力する。

【0072】機械学習手法は情報が一部不足した場合も他の情報をうまく利用することができる。ここでは、1つのフレーズの係り先の特定をするものを示した。構文解析でも形態解析と同様、一文全体を解析の対象とする場合は、各フレーズでの解析を組み合わせることで実現できる。

【0073】（実施例3）日本語の格解析の場合、図1に示す解-素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0074】日本語と英語の2つの対訳データの場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。対訳データは、「みかん<も>食べた。」と「We ate oranges, too.」であり、解は、「格を」である。

【0075】日本語の格解析とは、「は」「も」などの副助詞で表現されたり、連体節で表現されて（例、「食べたみかん」）、ガ格、ヲ格などの格助詞が消えている場合に、その消えた格を推定することを意味する。また、ここで求める格を意味関係にもつ体言と用言を単に体言と用言と書く。

【0076】解析に用いる素性は、1. 体言の単語自体、2. 体言の単語の意味カテゴリ、3. 用言の単語自体、4. 用言の単語の意味カテゴリ、5. 英語の単語2-gram列（2-gramとは2連続表現を意味し、単語2-gramは単語が2つ連続する表現を意味する）、6. 体言-用言に対応する英語表現の構文パターン、である。

【0077】前記から抽出される素性は、1. みかん、2. 食べ物、3. 食べる、4. 飲食関係の動詞、5. 「We ate」「ate oranges」「oranges, too」、6. VP NP（英語文を既存の構文解析システムなどでフレーズパターンを出力できるようにする。また、日本語に対する表現の特定は、日英翻訳辞書で単語逐語訳をして行う。もしくは、前記情報をまとめて、もとの対訳データで与えられるとしてもよい）、である。

10 【0078】前記素性の情報と解の情報「を」を利用して、機械学習手法により、どういう場合に「を」になりやすく、どういう場合に「が」になりやすいかなどを学習する。具体的には、英語表現の構文パターンが「VP NP」のときに「を」になりやすいとか、「NP VP」のときに「が」になりやすい、などを学習する。また、もとの日本語だけの情報に基づく体言「食べ物」、用言「飲食関係の動詞」のときに「を」になりやすいということも同時に学習する。

20 【0079】次に、入力する対訳データとして、「本<は>読んだ。」と「We read the book.」が与えられ、そのうち「は」の部分を格解析するように指示されたとする。この時、素性は、1. 本、2. 製品、3. 読む、4. 文書関係の動詞、5. 「We read」「read the」「the book」、6. VP NP、である。

30 【0080】前記素性でどの格になりやすいかを推定する。「VP NP」の素性の存在で「を」になりやすいと判定し、それが解として出力される。ところで、「V PNP」ならばいつでも「を」とは限らない。例えば、「I like apples.」だと、「りんごが好き」で「が」である。このような例外的現象も機械学習手法であると簡単に学習できる。例えば、前記素性3の用言が「好き」の場合は、「VP NP」でも、「が」と判定するように学習することになる。

【0081】（実施例4）単文の時制及びモダリティ表現の推定の場合、図1に示す解-素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

40 【0082】日本語と英語の2つの対訳データの場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。対訳データは、

50 「京大で待つ。」と「I wait in Kyoto university.」であり、解は、「現在」であり、また対訳データは、「京大に行く。」と「I go to Kyoto university.」であり、解は、「未来」である。ここでは、現在と未来しかあげていないが、過去、完了、要望、可能など、種類の



分類が考えられる。これらの分類は、文法書などを参考に予め決めておく。

【0083】前記と同様にして、素性は、1. 日本語文末文字列、2. 英語主節の動詞句表現の単語列、3. 日本語単語列、4. 英語単語列、である。

【0084】対訳データが、「京大に行く。」と「I will go to Kyotouniversity.」であり、解が、「未来」である場合、抽出される素性は、1. 「く」「行く」など、2. 「will go」「go」「will」、3. 「京大」「で」「待つ」、4. 「I」「will」「go」「to」「Kyoto」「university」である。

【0085】日本語だけでは「未来」などの時制を特定するのは難しいが、英語の前記素性2「will」などがあると「未来」などの時制を特定するのは容易である。また、英語主節の動詞句表現の特定はなんらかの構文解析システムが必要になる。場合によっては、その解析結果が間違ふ可能性もある。そのような場合は、日本語の文末表現の情報も使うことで、場合によっては、英語側の情報が誤っても日本語の方の情報でうまく行く場合がある。尚、実際の解析は省略する。

【0086】(実施例5)名詞句の指示性の推定の場合、図1に示す解-素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。名詞句指示性には、総称名詞句、定名詞句、不定名詞句があり、またこれを特定することで冠詞の生成などに役立つ。

【0087】日本語と英語の2つの対訳データの場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるので、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。対訳データは、「<犬>がいる。」と「There is <a dog> .」であり、解は、「不定名詞」であり、また対訳データは、「<その犬>は役に立つ。」と「The dog is useful」であり、解は、「定名詞」であり、また対訳データは、「<犬>は役に立つ。」と「The dog is useful.」であり、解は、「総称名詞」である。

【0088】前記と同様にして、素性は、1. 日本語周辺表層表現、2. 英語周辺表層表現、である。

【0089】対訳データが、「<その犬>は役に立つ。」と「The dog is useful.」であり、解が、「総称名詞」である場合、抽出される素性は、1. 「その」「役に立つ」など、2. 「The」「is」「useful」、である。

【0090】英語があると、冠詞が、定冠詞か否か、不定冠詞かで、「定名詞」「不定名詞」のどちらの可能性もないことが分かる。また、日本語で「その犬」のよう

に「その」が存在していると「総称名詞」の可能性はなくなる。そのような学習は、予め素性を適切に決めておくことと機械学習で自動で行うことができる。尚、実際の解析は省略する。

【0091】(実施例6)単文の時制及びモダリティ表現の中国語への翻訳の場合、図1に示す解-素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0092】日本語と英語の2つの対訳データの場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。

【0093】対訳データは、「あなたたちはあの映画を見ましたか?」と「Have you seen that film?」であり、解は、「看了」(「看」が見る、「了」が「~した」を意味する)であり、また対訳データは、「ここで写真をとってもよいですか?」と「May I have a picture here?」であり、解は、「可以照」(「照」が「写真をとる」、「可以」が「~してよい」を意味する)である。

【0094】これは、実施例4の分類カテゴリを単純に中国語の動詞表現にただけである。実施例4と同様に機械学習を用いると、日本語と英語の情報をうまく組み合わせる用いることができる。

【0095】また、ここでは動詞句表現を分類としたが、動詞句表現で使われる助動詞だけをとりあえず推定し、それを後の処理と組み合わせて翻訳に利用することも可能である。例えば、対訳データが、「あなたたちはあの映画をみましたか?」と「Have you seen that film?」であり、解が、「了」であり、対訳データが、「ここで写真をとってもよいですか?」と「May I have a picture here?」であり、解が、「可以」であるなど。尚、実際の解析は省略する。

【0096】(実施例7)英語単語の中国語単語への翻訳の場合、図1に示す解-素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0097】英語と日本語の2つの対訳データの場合を考える。前記解データは、英語と日本語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。

【0098】対訳データは、「May I have <a picture> here?」と「ここで写真をとってもよいですか?」であり、解は、「相」(「相」は「写真」を意味する)である。

【0099】英単語「a picture」の意味には写真の他にも絵の意味があり、絵を意味する中国語単語「画儿」との訳し訳の必要がある。日本語の「写真」という語の存在のおかげで正しく「相」と翻訳できる。動詞句が単語になっただけで実施例6とほぼ同様に扱える。尚、実際の解析は省略する。

【0100】(実施例8)中国語への翻訳における生成される中国語での構文構造推定の場合、図1に示す解 - 素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0101】日本語、英語及び中国語の3つの対訳データの場合を考える。前記解データは、日本語、英語及び中国語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。

【0102】対訳データは、「私は炒飯を食べます。」、「I have fried rice.」及び(中国語単語逐語訳列)「我<炒飯>吃」であり、解は、「炒飯」の係り先「吃」である。対訳データ中で係り先を求める中国語単語は“<”“>”の記号で囲んでおく。

【0103】解析に用いる情報である素性は、1. 解析するフレーズの主辞の単語、2. 解析するフレーズの主辞の単語の意味カテゴリ、3. 解析するフレーズの主辞の単語の品詞、4. 解の係り先の単語、5. 解の係り先の単語の意味カテゴリ、6. 解の係り先の単語の品詞、7. 日本語単語列、8. 英語単語列、9. 日本語文での係り受けと対応するか、10. 英語文での係り受けと対応するか、である。

【0104】前記解データのデータベースから抽出される「吃」を解とする場合の素性は、前記同様に表すと、1. 炒飯、2. 食べ物、3. 名詞、4. 吃、5. 飲食関係の動詞、6. 動詞、7. 「私」「は」「炒飯」「を」「食べます」、8. 「I」「have」「fried rice」、9. 対応する(日本語では「炒飯を」「食べます」にかかっている。また、この種の情報は既存のシステムで求めてもようし、入力で与えられているとしてもよい)、10. 対応する(英語では「fried rice」が「have」にかかっている)、である。これは正例となる。

【0105】また、「我」を解とする場合の素性は、1. 炒飯、2. 食べ物、3. 名詞、4. 我、5. 人、6. 名詞、7. 「私」「は」「炒飯」「を」「食べます」、8. 「I」「have」「fried rice」、9. 対応しない、10. 対応しない、である。これは負例となる。

【0106】以上の情報で学習すると、日本語文若しくは英語文での係り受けと対応すると正例である確率が高くなるように学習することになる。また、構文構造が対

応とれないようにしか解析できない場合もある。また、日本語、英語の構文構造の解析を失敗する場合もある。そのような場合は、前記1~8などの他の素性が役に立つことになる。尚、実際の解析は省略する。

【0107】(実施例9)中国語への一文翻訳の場合、図1に示す解 - 素性対抽出部12において、解と素性の集合及び解候補と素性の集合の組を抽出する方法を利用して解く。

【0108】日本語と英語の2つの対訳データの場合を考える。前記解データは、日本語と英語の対訳データに解の情報が付与されたものであるため、以下の対訳データの場合に、解データのデータベースの構成は、次のようになる。尚、対訳データ中で対象となる単語を“<”、“>”の記号で囲っておく。

【0109】対訳データは、「私は炒飯を食べます。」と「I have fried rice.」であり、解は、「我吃炒飯」である。解の部分には翻訳結果がはいっている。

【0110】前記素性の組を抽出する方法では、解の候補を作成する必要がある。簡単な方法としては、あらゆる単語逐語訳と、またその語順をあらゆる場合で並べ替えたものを全て解の候補とすればよい。そしてその中から、正例の確率が最も大きいものを選ぶとよい。

【0111】また、この方法で、解の候補の数が発散する場合には、問題を部分部分に分割し、各部分で正例の確率を算出し、その積が最大になるように部分部分を統合するようにしておけばよい。これは、一文全体の形態素解析を行うのと同様である。

【0112】(実施例10)機械学習手法の場合、教師信号が同じ形をしているものは併用して学習できる。例えば、実施例1の形態素解析では、対訳データは、「<きょうだい>で待つ。」と「I wait in kyoto university.」であり、解は、「京大 名詞」といった形のデータを教師信号として用いるが、「<きょうだい>で待つ。」「京大 名詞」のような対訳データでないものも教師信号と扱える。この場合、英語に関係する素性情報に欠けるが、その部分は無かったとしても日本語に関係する素性情報が残るので、その情報を使って学習することになる。

【0113】ところで、対訳データが、「<きょうだい>で待つ。」と「I wait in kyoto university.」であり、解が、「京大 名詞」であるものに、更に形態素情報もふったコーパスはあまりみないが、「<きょうだい>で待つ。」と「京大 名詞」のような形態素情報もふったコーパスは多く存在する。このデータも使って学習できると、学習データが多いため精度が向上する。

【0114】また、前記併用型の場合、大規模に使える単言語の学習データと、情報量の多い2言語対訳の学習データを併用するので、非常に強力である。また、既存

のシステムで学習を用いるものは、単言語の学習データを用いているので、この併用型は少なくとも既存の学習システムと同等程度の能力を確保した上で、さらに2言語対訳の学習データを併用するというものになっている。

【0115】

【発明の効果】本発明によれば、大量の変換規則を用意する必要がなく、複数の自然言語で記述された処理対象文から他の自然言語及び/又は同じ自然言語で記述された処理結果文への言語変換及び/又は複数の言語における言語解析を行うことができる複数言語入力での言語処理方法及び言語処理装置が得られる。

【図面の簡単な説明】

【図1】本発明による複数言語入力での言語処理装置の実施形態を示すブロック図である。

【図2】本発明による複数言語入力での言語処理方法の実施形態を示すフローチャートである。

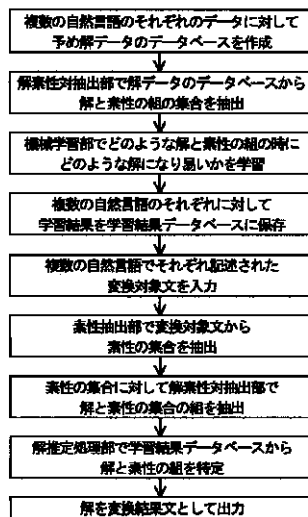
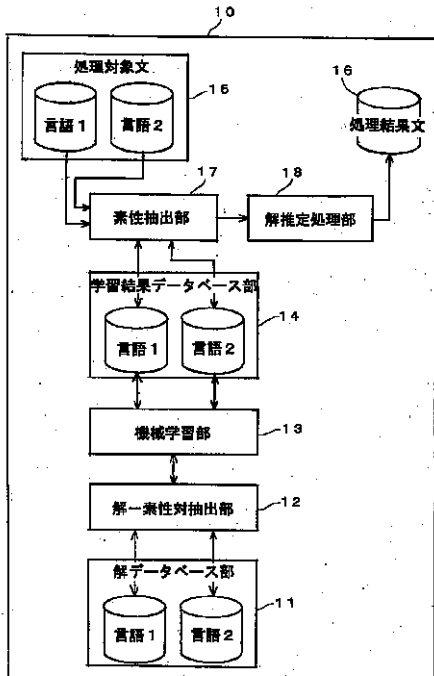
【図3】従来の機械学習手法を用いた単言語入力での言語処理装置のブロック図である。

【符号の説明】

- 10 言語処理装置
- 11 解データベース部
- 12 解 - 素性対抽出部
- 13 機械学習部
- 14 学習結果データベース部
- 15 処理対象文
- 16 処理結果文
- 17 素性抽出部
- 18 解推定処理部
- 30 言語処理装置
- 31 解データベース部
- 32 解 - 素性対抽出部
- 33 機械学習部
- 34 学習結果データベース部
- 35 変換対象文
- 36 素性抽出部
- 37 解推定処理部
- 38 変換結果文

【図1】

【図2】



【図 3】

