

(51)Int.Cl. ⁷ G06F 17/28	識別記号	F I G06F 17/28	テ-マコード [*] (参考) Z 5B091 C
--	------	-------------------	--

審査請求 有 請求項の数 2 O L (全14頁)

(21)出願番号 特願2001 - 201010 (P 2001 - 201010)	(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 7月 2日 (2001.7.2)	(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
特許法第30条第 1項適用申請有り 2001年 3月 9日 社 団法人電子情報通信学会発行の「電子情報通信学会技術 研究報告 信学技報 V o l . 100 , N o . 698」に発表	(72)発明者 馬 青 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 100119161 弁理士 重久 啓子 (外 2名)

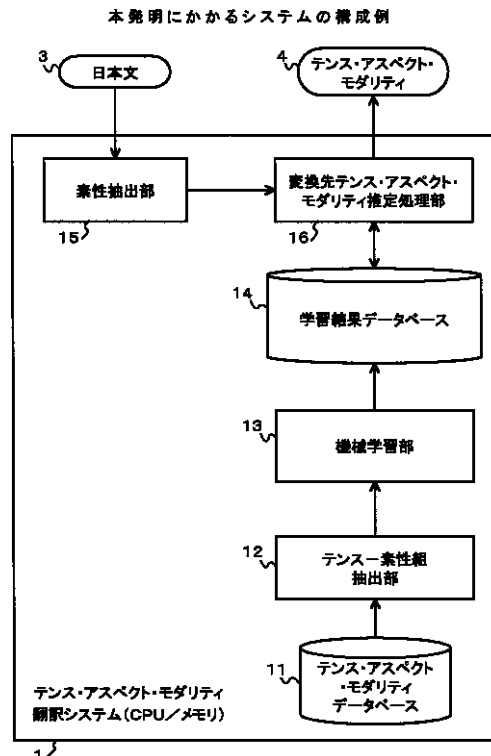
最終頁に続く

(54)【発明の名称】 テンス・アスペクト・モダリティ翻訳処理方法 , テンス・アスペクト・モダリティ翻訳システム

(57)【要約】

【課題】 機械翻訳の際に変換先のテンス・アスペクト・モダリティを精度よく翻訳する。

【解決手段】 テンス - 素性組抽出部12は、テンス・アスペクト・モダリティ・データベース11から、事例ごとにテンス・アスペクト・モダリティ(以下「テンス等」という)と素性の集合の組とを抽出する。機械学習部13は、抽出した組からどのような素性のときにどのようなテンス等となりやすいかを学習し、その結果をデータベース14に保存する。素性抽出部15は、訳出したい日本文3が入力されると素性の集合を抽出する。変換先テンス・アスペクト・モダリティ推定処理部16は、学習結果データベース14を参照し、その素性の集合から日本文3の素性の場合にどのようなテンス等になりやすいかを推定し、推定したテンス等4を出力する。



【特許請求の範囲】

【請求項 1】 コンピュータにより一の言語から他の言語へ翻訳処理をする際に、変換元言語から変換先言語のテンス・アスペクト・モダリティを翻訳する方法であって、予め備えられた変換元言語の事例と当該事例のテンス・アスペクト・モダリティとを記憶するデータベースから、テンス・アスペクト・モダリティと当該事例または当該事例に関連するデータから抽出した、複数の形式の素性の集合とからなるテンス - 素性組を事例ごとに抽出する過程と、前記テンス - 素性組のうちのすべてまたはいくつかの素性を用いて、機械学習法により変換先言語のテンス・アスペクト・モダリティを判定するための学習データを作成し保存する過程と、変換元言語の入力文から、当該入力文の素性の集合を抽出する過程と、前記入力文の素性の集合のうちのすべてまたはいくつかの素性をもとに、前記学習データを参照して前記入力文のテンス・アスペクト・モダリティを推定する過程とを備えることを特徴とするテンス・アスペクト・モダリティ翻訳処理方法。

【請求項 2】 請求項 1 に記載のテンス・アスペクト・モダリティ翻訳処理方法において、前記機械学習法として、決定リスト法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の手法を用いることを特徴とするテンス・アスペクト・モダリティ翻訳処理方法。

【請求項 3】 コンピュータにより一の言語から他の言語へ翻訳処理をする際に、変換元言語から変換先言語のテンス・アスペクト・モダリティを翻訳するシステムであって、予め備えられた変換元言語の事例と当該事例のテンス・アスペクト・モダリティとを記憶するデータベースから、テンス・アスペクト・モダリティと当該事例または当該事例に関連するデータから抽出した、複数の形式の素性の集合とからなるテンス - 素性組を事例ごとに抽出するテンス - 素性組抽出手段と、前記テンス - 素性組のうちのすべてまたはいくつかの素性を用いて、機械学習法により変換先言語のテンス・アスペクト・モダリティを判定するための学習データを作成し保存する機械学習手段と、変換元言語の入力文から、当該入力文の素性の集合を抽出する素性抽出手段と、前記入力文の素性の集合のうちのすべてまたはいくつかの素性をもとに、前記学習データを参照して前記入力文のテンス・アスペクト・モダリティを推定する変換先テンス・アスペクト・モダリティ推定処理手段とを備えることを特徴とするテンス・アスペクト・モダリティ翻訳システム。

【請求項 4】 請求項 3 に記載のテンス・アスペクト・モダリティ翻訳システムにおいて、前記機械学習部では、決定リスト法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の手法を用いて学習することを特徴とするテンス・アスペクト・モダリティ翻訳システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、コンピュータによる翻訳システムの分野において、特に、機械学習法を用いてテンス（時制）、アスペクト（相）、またはモダリティ（様相）を翻訳する翻訳処理方法および翻訳システムに関するものである。

【0002】

【従来の技術】近年、WWW (World Wide Web) などのインターネットの発展とともに機械翻訳の必要性が高まり続けている。この機械翻訳において、テンス・アスペクト・モダリティは翻訳が難しい問題として知られている。

【0003】テンス・アスペクト・モダリティとは、動詞部分の時制（テンス）、進行形や完了形などの相（アスペクト）、または英文における助動詞相当句である様相（モダリティ）についての情報である。

【0004】従来、テンス・アスペクト・モダリティの表現は、人手により作成されたルールによって扱われていた。しかしながら、近年では、用例ベース（k 近傍法）の方法などのコーパスデータにもとづくアプローチでも処理されるようになってきた。用例ベースの方法では、集めた事例ごとに、どの場合にどの時制などを使うかを記したデータを対応づけた用例データベースを用意しておき、入力された文と良く似た事例に基づいてテンス・アスペクト・モダリティを翻訳するといったことが行なわれていた。[参考文献 1] 村田真樹 馬青 内元清貴 井佐原均、用例ベースによるテンス・アスペクト・モダリティの日英翻訳、人工知能学会誌、Vol.16, No.1, 2001 参考文献 1 に記載されている研究では、日本文から英文への機械翻訳のテンス・アスペクト・モダリティの判定の際に、日本文のテンス・アスペクト・モダリティは文末に表されることに着目して、入力された日本文の文末の所定の長さの文字列と、予め用意したコーパスデータとの類似度を k 近傍法により判断してテンス・アスペクト・モダリティを決定する手法を用いている。k 近傍法とは、最もよく似た一つの事例の代わりに、最もよく似た k 個の事例を用い、この k 個の事例での多数決によって分類先を定める手法である。

【0005】

【発明が解決しようとする課題】しかし、人手でルールを記述し、このルールをもとにテンス・アスペクト・モダリティを分類する方法では、人的資源の問題や、人手による作業の精度などの問題がある。

【0006】また、入力文とよく似た事例を使う手法では、入力文と事例の類似度を定義する必要があり、例えば文末の文字列のように類似度を定義することができるような平易な情報しか扱うことができなかった。そのため、参考文献 1 に記載された研究の手法において、文末の文字列の情報のみによってテンス・アスペクト・モダ

リティの分類を判定することで、判定結果の精度が低くなる場合が生じる。

【0007】例えば、実例データ「もう行きました。」のテンス・アスペクト・モダリティが「過去完了」である場合に、「昨日行きました。」という文が入力されたとする。この入力文の正しいテンス・アスペクト・モダリティは「過去」であるにもかかわらず、文字列「ました\$(\$ =文末)」の表示の類似度から、実例データと同様に「過去完了」と判定されてしまう場合がある。

【0008】したがって、参考文献1の研究の手法のように文末の一致する文字列だけでなく、例えば、この場合の「昨日」のように、文末の文字列とは異なる形式の情報を合わせて用いることが有効であると考えられる。

【0009】しかし、テンス・アスペクト・モダリティを解析するための情報(素性)として、形態素情報(形態素素性)、意味解析情報(単語素性など)、構文解析情報(構文解析素性)などの異なる形式の素性を組み合わせて用いることが有効であるとしても、参考文献1の研究で用いたk近傍法のような類似度を定義する必要があり判定手法では、複数の形式の素性をを用いることができないという問題があった。

【0010】本発明は、上記問題点の解決を図り、変換元言語の事例やその事例に関連するデータから抽出した異なる形式の素性を取り扱うことができる機械学習手法を用いて、どのような素性の場合にどのようなテンス・アスペクト・モダリティになるかを学習し、その学習結果を用いて入力文の変換先のテンス・アスペクト・モダリティを精度よく翻訳できる手段を提供することを目的とする。

【0011】

【課題を解決するための手段】上記課題を解決するため、本発明に係る方法は、予め備えられた変換元言語の事例と当該事例のテンス・アスペクト・モダリティとを記憶するデータベースから、テンス・アスペクト・モダリティと当該事例または当該事例に関連するデータから抽出した、複数の形式の素性の集合とからなるテンス - 素性組を事例ごとに抽出する過程と、前記テンス - 素性組のうちのすべてまたはいくつかの素性を用いて、機械学習法により変換先言語のテンス・アスペクト・モダリティを判定するための学習データを作成し保存する過程と、変換元言語の入力文から、当該入力文の素性の集合を抽出する過程と、前記入力文の素性の集合のうちのすべてまたはいくつかの素性をもとに、前記学習データを参照して前記入力文のテンス・アスペクト・モダリティを推定する過程とを備えることを特徴としている。

【0012】本発明では、従来の手法のように、文末の文字列のように単一の種類の素性だけを用いてテンス・アスペクト・モダリティの解析を行うのではなく、文字列の他、一文全体の形態素素性、意味的素性、構文的素性、前文のテンス・アスペクト・モダリティ、または対

訳データの該当する構成部分データなど、二以上の異なる形式の素性を任意に用いて解析処理を行う点が、従来の手法と異なる。

【0013】また、本発明では、多くの形式の素性を自由に用いることができる、類似度を設定する必要のない種々の機械学習手法を用いて解析処理を行う点が、類似度の定義を必要とする従来のk近傍法のような手法による判定と異なる。

【0014】以上の本発明に係る処理方法またはシステムは、その処理過程や手段、構成、要素をコンピュータに実行させるプログラムによっても実現することができる。このプログラムは、コンピュータが読み取り可能な、可搬媒体メモリ、半導体メモリ、ハードディスクなどの適当な記録媒体に格納することができ、これらの記録媒体に記録して提供され、または通信インタフェースを介して種々の通信網を利用した送受信により提供される。

【0015】

【発明の実施の形態】以下に、本発明の実施の形態を図を用いて説明する。本実施の形態では、日本文から英文への翻訳に際しテンス・アスペクト・モダリティを翻訳する場合を例とする。

【0016】図1に、本発明のシステム構成例を示す。図1中、1は本発明に係るテンス・アスペクト・モダリティ翻訳システムを表す。テンス・アスペクト・モダリティ翻訳システム1はCPU、メモリなどで構成され、テンス・アスペクト・モダリティ・データベース11、テンス - 素性組抽出部12、機械学習部13、学習結果データベース14、素性抽出部15、変換先テンス・アスペクト・モダリティ推定処理部16を持つ。

【0017】テンス - 素性組抽出部12は、予め用意しておいたテンス・アスペクト・モダリティ用コーパスであるテンス・アスペクト・モダリティ・データベース11から、事例ごとに、テンス・アスペクト・モダリティと事例の素性の集合との組を抽出する手段である。

【0018】機械学習部13は、テンス - 素性組抽出部12で抽出されたテンス・アスペクト・モダリティと素性の集合との組から、どのような素性のときに、どのようなテンス・アスペクト・モダリティになりやすいかを機械学習法により学習し、その学習結果を学習結果データベース14に保存する手段である。

【0019】素性抽出部15は、入力された日本文3から素性の集合を抽出し、それらを変換先テンス・アスペクト・モダリティ推定処理部16へ渡す手段である。

【0020】変換先テンス・アスペクト・モダリティ推定処理部16は、学習結果データベース14を参照して、渡された素性の集合の場合に、変換先の言語においてどのようなテンス・アスペクト・モダリティになりやすいかを推定し、日本文3の変換先のテンス・アスペクト・モダリティ4を出力する手段である。

【0021】本発明の処理の流れの概略を説明する。図2は、図1に示すシステムの処理フローチャートである。

【0022】図2に示す処理を開始する前に、テンス・アスペクト・モダリティ・データベース11を予め用意しておく。テンス・アスペクト・モダリティ・データベース11は、機械翻訳用の日英の対訳コーパスであり、日本語と英語の対訳データにテンス・アスペクト・モダリティの情報が付与されている。

【0023】日英の対訳データに付与するテンス・アスペクト・モダリティの分類として、例えば以下のものを用いる。以下の分類は、対訳の英語文の動詞がどのような形になっているかによって定められる。

(1)各助動詞相当語句 (be able to, be going to, can, have to, had better, may, must, need, ought, shall, used to, will の12種類) がつくかどうかと、{現在形, 過去形}と{進行形, 進行形でない}と{完了, 完了でない}のすべての組み合わせ(助動詞相当語句が複数つく場合も許している。) : 2¹⁵種類

- (2)命令形(1種類)
- (3)名詞句(1種類)
- (4)分詞構文(1種類)
- (5)動詞省略(1種類)
- (6)間投詞, 挨拶文など(1種類)
- (7)日本語と英語で動詞の対応がとれない場合(1種類)

(8)作業不可(1種類)
ただし、上記の分類のうち、「(3)名詞句」から「(8)作業不可」までの6つの分類はテンス・アスペクト・モダリティの分類としては扱う必要がないか、もしくはテンス・アスペクト・モダリティの翻訳を行なう必要がないと思われるので、本形態では省略している。

【0024】また、これらの分類は、「英語の主節の動詞部分」と「日本語の主節の動詞に対応する英語の動詞部分」の二か所にふられる。しかし、日英翻訳において日本語のテンス・アスペクト・モダリティに対応するのは「英語の主節の動詞部分」であろうと考えられるので、本発明に係るテンス・アスペクト・モダリティ翻訳システム1では、日本文3を与えて変換先の「英語の主節の動詞部分」のテンス・アスペクト・モダリティの分類を推定し、推定結果であるテンス・アスペクト・モダリティ4を出力することとしている。

【0025】ステップS1:まず、テンス-素性組抽出部12により、用意されたテンス・アスペクト・モダリティ・データベース11から、各事例ごとに、テンス・

$$f_{max} = \operatorname{argmax}_{f_j \in F} \max_{a_i \in A} \tilde{p}(a_i | f_j) \quad (2)$$

【0035】また、
【0036】
【数2】

アスペクト・モダリティと事例の素性の集合との組を抽出する。

【0026】テンス-素性組抽出部12では、素性の集合として、文字列素性、形態素素性、単語素性、構文的素性、一前文のテンス・アスペクト・モダリティの情報、英文対訳データの動詞部分など、種々の形式の素性のうち、所定の素性を抽出することができる。

【0027】図3に抽出する素性の集合とテンス・アスペクト・モダリティの組を示す。図3に示すように、テンス-素性組抽出部12により、テンス・アスペクト・モダリティ・データベース11の事例「もう登録しました。」から、テンス・アスペクト・モダリティと、文字列素性「もう登録しました\$」, 「う登録しました\$」, …, 「た\$」, 単語素性「もう」, 「登録」「し」「まし」「た」などの素性の集合との組を抽出する。なお、ここでは、抽出した文字列素性には入力された文全体の形態素列と区別できるように末尾に\$をつけている。また、文末表現の正規化のため、句点などは消している。

【0028】ステップS2:続いて、機械学習部13により、抽出されたテンス・アスペクト・モダリティと素性の集合との組から、どのような素性のときにどのようなテンス・アスペクト・モダリティになりやすいかを機械学習し、その学習結果を学習結果データベース14に保存する。

【0029】機械学習では、例えば、所定の長さの文末の文字列素性、事例の全文の形態素素性、単語素性のうち、いくつかの素性を用いて処理を行ってもよい。

【0030】機械学習の手法は、種々の形式の素性の集合を扱うことができるような機械学習法であればよく、例えば、以下に示すような決定リスト法、最大エントロピー法、サポートベクトルマシン法などを用いる。

【0031】(1)決定リスト法
決定リスト法は、あらかじめ設定しておいた素性 f_j ($F, 1 \leq j \leq k$)のうち、いずれか一つの素性のみを文脈として各分類の確率値を求めて、その確率値が最も大きい分類を求める分類とする方法である。

【0032】ある文脈 b で分類 a を出力する確率は以下の式によって与えられる。

$$p(a | b) = p(a | f_{max}) \quad (1)$$

ただし、 f_{max} は以下の式によって与えられる。

【0034】
【数1】

【0037】は、素性 f_j を文脈に持つ場合の分類 a_i

の出現の割合である。

【0038】具体的には、各素性ごとに、どのようなテンス・アスペクト・モダリティの分類になるのかの確率を求めておき、入力文のすべての素性のうち最大確率の素性の分類を用いてテンス・アスペクト・モダリティの分類を推定する。

【0039】決定リスト法にもとづく分類は簡便ではあるが、ある一つの素性のみを文脈としてテンス・アスペクト・モダリティの分類の推定を行なうので、機械学習の手法としては少々貧弱なものとなっている。

$$\sum_{a \in A, b \in B} p(a, b) g_j(a, b) = \sum_{a \in A, b \in B} \tilde{p}(a, b) g_j(a, b) \quad (3)$$

for $\forall f_j (1 \leq j \leq k)$

$$H(p) = - \sum_{a \in A, b \in B} p(a, b) \log(p(a, b)) \quad (4)$$

【0042】ただし、A、Bは分類と文脈の集合を意味し、 $g_j(a, b)$ は文脈bに素性 f_j があつて、なおかつ分類がaの場合1となり、それ以外で0となる関数を意味する。また、

【0043】

【数4】

$$\tilde{p}(a, b)$$

【0044】は、既知データでの(a, b)の出現の割合を意味する。

【0045】式(3)は確率pと出力と素性の組の出現を意味する関数gをかけることで出力と素性の組の頻度の期待値を求めることになっており、右辺の既知データにおける期待値と、左辺の求める確率分布に基づいて計算される期待値が等しいことを制約として、エントロピー最大化(確率分布の平滑化)を行なつて、出力と文脈の確率分布を求めるものとなっている。最大エントロピー法の詳細については、以下の参考文献2および参考文献3を参照されたい。

[参考文献2] Eric Sven Ristad, Maximum Entropy Modeling for Natural Language, (ACL/EACL Tutorial Program, Madrid, 1997)

[参考文献3] Eric Sven Ristad, Maximum Entropy Modeling Toolkit, Release 1.6beta, (<http://www.mnemoniic.com/software/memt>, 1998)

本手法の利用は、もともと訳出対象となる入力データの素性と同一な素性が既知のコーパス中に多数存在することは稀であり、コーパスでの素性の出現確率をそのまま使用することはできないが、等価な状態は等価な確率を持つことを前提に確率分布を平滑化すれば、コーパスに

【0040】(2)最大エントロピー法

最大エントロピー法は、あらかじめ設定しておいた素性 $f_j (1 \leq j \leq k)$ の集合をFとするとき、以下の式(3)を満足しながらエントロピーを意味する式(4)を最大にするときの確率分布 $p(a, b)$ を求め、その確率分布にしたがって求まる各分類の確率のうち、もっとも大きい確率値を持つ分類を求める分類とする方法である。

【0041】

【数3】

$$\sum_{a \in A, b \in B} \tilde{p}(a, b) g_j(a, b) \quad (3)$$

for $\forall f_j (1 \leq j \leq k)$

$$H(p) = - \sum_{a \in A, b \in B} p(a, b) \log(p(a, b)) \quad (4)$$

おける素性出現の期待値を用いることができると考えられることによる。

【0046】(3)サポートベクトルマシン法

サポートベクトルマシン法は、空間を超平面上で分割することにより、2つの分類からなるデータを分類する手法である。図4にサポートベクトルマシン法のマージン最大化の概念を示す。図4において、白丸は正例、黒丸は負例を意味し、実線は空間を分割する超平面を意味し、破線はマージン領域の境界を表す面を意味する。図4

(A)は、正例と負例の間隔が狭い場合(スモールマージン)の概念図、図4(B)は、正例と負例の間隔が広い場合(ラージマージン)の概念図である。

30

【0047】このとき、2つの分類が正例と負例からなるものとするとき、学習データにおける正例と負例の間隔(マージン)が大きいものほどオープンデータで誤った分類をする可能性が低いと考えられ、図4(B)に示すように、このマージンを最大にする超平面を求めそれを用いて分類を行なう。

【0048】基本的には上記のとおりであるが、通常、学習データにおいてマージンの内部領域に少数の事例が含まれてもよいとする手法の拡張や、超平面の線形部分を非線型にする拡張(カーネル関数の導入)がなされたものが用いられる。この拡張された方法は、以下の識別関数を用いて分類することと等価であり、その識別関数の出力値が正か負かによって二つの分類を判別することができる。

40

【0049】

【数5】

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^l \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right) \quad (5)$$

$$b = \frac{\max_{i,y_i=-1} b_i + \min_{i,y_i=1} b_i}{2}$$

$$b_i = \sum_{j=1}^l \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_i)$$

【0050】ただし、 \mathbf{x} は識別したい事例の文脈 (素性の集合) を、 \mathbf{x}_i と y_i ($i = 1, \dots, l$, $y_i \in \{1, -1\}$) は学習データの文脈と分類先を意味し、関数 sgn は、

$$\text{sgn}(x) = 1 \quad (x \geq 0) \quad (6)$$

- 1 (otherwise) であり、また、各 α_i は式 (8) と式 (9) の制約のもと式 (7) を最大にする場合のものである。

【0051】

【数6】

$$L(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \quad (7)$$

$$0 \leq \alpha_i \leq C \quad (i = 1, \dots, l) \quad (8)$$

$$\sum_{i=1}^l \alpha_i y_i = 0 \quad (9)$$

【0052】また、関数 K はカーネル関数と呼ばれ、様々なものが用いられるが、本形態では以下の多項式のものを用いる。

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^d$$
 C , d は実験的に設定される定数である。後述する具体例では C はすべての処理を通して 1 に固定した。また、 d は、1 と 2 の二種類を試している。ここで、 $\alpha_i > 0$ となる \mathbf{x}_i は、サポートベクトルと呼ばれ、通常、式 (5) の和をとっている部分はこの事例のみを用いて計算される。つまり、実際の解析には学習データのうちサポートベクトルと呼ばれる事例のみしか用いられない。

【0054】なお、拡張されたサポートベクトルマシン法の詳細については、以下の参考文献 4 および参考文献 5 を参照されたい。

[参考文献 4] Nello Cristianini and John Shawe-Taylor, An Introduction to Support Vector Machines and other kernel-based learning methods, (Cambridge University Press, 2000)

[参考文献 5] Taku Kudoh, Tinysvm: Support Vector machines, (<http://cl.aist-nara.ac.jp/taku-ku//software/TinySVM/index.html>, 2000)

サポートベクトルマシン法は、分類の数が 2 個のデータを扱うもので、通常これにペアワイズ手法を組み合わせることで、分類の数が 3 個以上のデータを扱うことになる。

【0055】ペアワイズ手法とは、 N 個の分類を持つデータの場合、異なる二つの分類先のあらゆるペア (N

個) を用いる。

【0053】

$$(10)$$

($N - 1$) / 2 個) を作り、各ペアごとにどちらがよいかを 2 値分類器 (ここではサポートベクトルマシン法によるもの) で求め、最終的に $N(N - 1) / 2$ 個の 2 値分類器の分類先の多数決により、分類先を求める方法である。

【0056】本形態における 2 値分類器としてのサポートベクトルマシンは、サポートベクトルマシン法とペアワイズ手法を組み合わせることによって実現するものであり、以下の参考文献 6 により工藤氏が作成した Tiny SVM を利用している。

[参考文献 6] 工藤拓 松本裕治, Support vector machine を用いた chunk 同定, 自然言語処理研究会, 2000-NL-140, (2000)

具体的には、テンス・アスペクト・モダリティの各分類ごとに最大マージンの超平面を求めておく。そして、入力された未知の素性が、例えばテンスが現在 (正例) であるか過去 (負例) であるかなどについて、それぞれについて正例の領域と負例の領域のどちらに属するかを判定していき、その判定結果の多数決により、最終的にテンス・アスペクト・モダリティを推定する。

【0057】ステップ S3: テンス・アスペクト・モダリティを訳出したい日本語 3 が入力される。

【0058】ステップ S4: 素性抽出部 15 では、テン

ス - 素性組抽出部 1 2 での処理とほぼ同様に、入力された日本語 3 から素性の集合を取り出し、それらを変換先 テンス・アスペクト・モダリティ 推定処理部 1 6 へ渡す。図 5 に、抽出した素性の集合の例を示す。入力された日本語「もう行きました。」から、文字列素性「もう行きました \$」, 「う行きました \$」, …, 「た \$」および単語素性「もう」, 「行き」, 「まし」, 「た」などが抽出される。

【0059】ステップ S 5 : 変換先 テンス・アスペクト・モダリティ 推定処理部 1 6 では、学習結果データベース 1 4 をもとに、渡された素性の集合の場合にどのような テンス・アスペクト・モダリティ になりやすいかを特定し、特定した テンス・アスペクト・モダリティ 4 を出力する。例えば、「過去,完了」, 「過去,完了,進行形」, 「現在,shall 付」, 「過去;be able to付」などのデータを出力する。

【0060】入力された日本語「もう行きました。」の素性の集合について、学習結果データベース 1 4 に記憶された学習結果を使用すると「過去完了」でよいとわかるので、テンス・アスペクト・モダリティ 4 として「過去完了」を出力する。この場合に、従来のように単純に文末の文字列だけを用いて推定したときは、「ました \$」の表現の一致からテンスが「現在完了」と推定されてしまう。しかし、本発明では、文字列以外に全文の単語素性のうちのいくつかの素性を用いて学習した結果を参照することにより「過去完了」と正しく推定できる。

【0061】また、入力された日本語 3 が「昨日行きました」である場合には、同様に、従来の手法では「ました \$」の表現の一致からテンスが「現在完了」と推定されてしまう。しかし、学習結果データベース 1 4 に「昨日登録しました。」にテンスが「過去」とであるという学習結果が記憶されていれば、「昨日」という単語素性をもとに正しく「過去」と推定することができる。

【0062】以上では、主に素性の集合として形態素素性と単語素性を用いた場合を例に説明したが、テンス - 素性組抽出部 1 2 では、形態素素性や単語素性だけでなく、一前文(事例)のテンス・アスペクト・モダリティを素性として抽出してもよい。これは、テンス(時制)は継続しやすいという性質を利用するものである。すなわち、一前文に付与されたテンスが「現在」であれば、対象事例も「現在」で記述したほうがよいと学習するものである。特に論文の実験の記述で時制を統一する必要がある場合には有用である。

【0063】また、テンス - 素性組抽出部 1 2 では、事例の英文対訳データの該当する語句(動詞部分)などを素性として抽出してもよい。これは、訳出される英語文の構造がかわると用いるべきモダリティも変化する場合があることにもとづく。

【0064】以下の事例で説明する。
例文 1) 和文: 彼は質実な生活を { 送っている }

英訳: He { lives } a sober and simple life.

例文 2) 和文: 彼は情性的に怠惰な生活を { 送っている }

英訳: He { is leading } a lazy life out of habit.

例えば、例文 1 の「送っている」のモダリティは「現在形」であるが、例文 2 の「送っている」のモダリティは「進行形」である。これらはほとんど意味の同じ文であり同じモダリティを持っていると考えてもよいのだが、訳出に用いる動詞を「live」と「lead」とかえただけでこのような違いが出てくる。

【0065】なお、この場合に、入力される日本語 3 に仮訳された英文が付与されることになり、素性抽出部 1 5 では、日本語 3 に付与された英文の動詞部分が素性として抽出される。

【0066】高品質な処理を行ないたい場合、このように、素性の取り出しの際に、日本語側だけでなく英文対訳データのうち機械翻訳システムの構造解析部が想定している英語側の構造(あるいは動詞)を取り出すということが、テンス・アスペクト・モダリティの翻訳の向上に有効である。

【0067】以下、本発明を一実施例によりさらに詳細に説明する。

【0068】〔テンス・アスペクト・モダリティ・データベース〕図 6 に、実施例で用いるテンス・アスペクト・モダリティ・データベース 1 1 である対訳コーパスの一部を示す。この対訳コーパスは、例えば、以下の参考文献 7 にもとづいて作成する。

[参考文献 7] 村田真樹 内山将夫 内元清貴 馬青井佐原均, 機械学習を用いた機械翻訳用モダリティコーパスの修正, 言語処理学会第 7 回年次大会, (2001)
図 6 中、英語側の文には以下の二か所のタグが付与されている。

・英語の主節の動詞部分を < v >, < / v > のタグで囲む。

・日本語の主節の動詞に対応する英語の動詞部分を < v j >, < / v j > のタグで囲む。

【0069】また、日本語側の文の先頭に“ c ”や“ d ”といった記号がふられているが、これらはこの対訳データのテンス・アスペクト・モダリティを意味する。例えば、“ c ”は can を、“ d ”は過去形を意味する。

【0070】図 6 に示すコーパスの一つめのデータには“ , ”があるが、これは < v j > を用いるときに使われるもので、“ , ”の左に < v > で囲まれた動詞に対するテンス・アスペクト・モダリティが、右に < v j > で囲まれた動詞に対するテンス・アスペクト・モダリティが記述される。なお、このコーパスでは現在形の出現が多いのでその場合はタグをふらなかつた。このため、“ , ”の左右が空欄となつてこの部分には“ , ”だけが付与されている。

【0071】また、「日本語に対応する英語の動詞部分」と「英語の主節の動詞部分」が一致する場合は「英語の主節の動詞部分」のタグのみ付与した。また、「日本語に対応する英語の動詞部分」の方はそれほど綿密にタグ付与は行なっておらず、「日本語に対応する英語の動詞部分」と「英語の主節の動詞部分」が一致しない場合にもタグ付与をしなかった場合もある。

【0072】日英の対訳データに付与するテンス・アスペクト・モダリティの分類として、前述した分類のうち、以下の(1)および(2)を用いた。

(1) 各助動詞相当語句 (be able to, be going to, can, have to, had better, may, must, need, ought, shall, used to, will の12種類) がつくかどうかと、{現在形, 過去形}と{進行形, 進行形でない}と{完了, 完了でない}のすべての組み合わせ(助動詞相当語句が複数つく場合も許している。): 2^{15} 種類

(2) 命令形(1種類)

本発明に係るシステムで扱うテンス・アスペクト・モダリティの分類は英語の表層表現に基づいて定めたものであり、日本語文だけを与えてこの分類を推定できれば、モダリティ表現の日英翻訳ができあがる。このため、本例では、原則として、テンス・アスペクト・モダリティの分類を示すタグと日本語文のみを用いている。また、前述したように、これらの分類は、「英語の主節の動詞部分」と「日本語の主節の動詞に対応する英語の動詞部分」の二か所にふられるが、日本文3を与えて「英語の主節の動詞部分」のテンス・アスペクト・モダリティの分類を推定することを問題設定としている。

【0073】本例では、テンス・アスペクト・モダリティ・データベース11として、以下の二種類の対訳コーパスを用いた。

・K社和英辞典の例文(事例総数は39,660個, 分類の総数は46個)

・白書データ(事例総数は5,805個, 分類の総数は30個)

これらのコーパスは、人手により確認しながらタグづけを行ない、さらに参考文献7および以下の参考文献8に示すコーパス修正の方法を利用して作成しており、非常に高精度なものとなっている。

[参考文献8] 村田真樹 内山将夫 内元清貴 馬青 井佐原均, 決定リスト, 用例ベース手法を用いたコーパス誤り検出・誤り訂正, 自然言語処理研究会, 2000-NL-136(2000)

【0074】〔抽出する素性〕本例では、日本文3の入力を与えられたときにテンス・アスペクト・モダリティ4として分類を出力する。このため、素性は入力される日本文3から取り出すことになる。ここでは、素性集合として以下の三種類のものに対して処理を行った。

(1) 素性集合F1

日本語文末の1~10gramの文字列と入力された文全体

の形態素列を素性の集合とする。

例: 「ない\$」「しなかった\$」「今日」「は」「走る」

この場合に、素性の数は、K社データで230,134個, 白書データで25,958個となる。

(2) 素性集合F2

日本語文末の1~10gramの文字列を素性の集合とする。

例: 「ない\$」「しなかった\$」

10 この場合に、素性の数は、K社データで199,199個, 白書データで16,610個となる。

(3) 素性集合F3

入力された文全体の形態素列を素性の集合とする。

例: 「今日」「は」「走る」

この場合に、素性の数は、K社データで30,935個, 白書データで9,348個となる。

【0075】入力された文を形態素列に分解するには、JUMANを用いた。JUMANの詳細な説明については、以下の参考文献9に記載されている。

20 [参考文献9] 黒橋禎夫 長尾真, 日本語形態素解析システムJUMAN使用説明書, version 3.6 (京都大学大学院工学研究科, 1998)

素性集合F1は、素性集合F2と素性集合F3との組合わせである。素性集合F2は、上記の参考文献1の研究を参考にして作成したものであり、日本語文においてテンス・アスペクト・モダリティを示す表現は文末の動詞にあらわれることが多いことから、日本語文の文末の文字列を素性としている。素性集合F3は、「明日」「昨日」などの副詞もテンス・アスペクト・モダリティを示す表現であり、用いるべきだと考えて作成したもので、入力された文全体の形態素列とするものである。

〔機械学習によるテンス・アスペクト・モダリティの分類〕本例では、機械学習の手法として、決定リスト法、最大エントロピー法、サポートベクトルマシン法を用いた。さらに、本発明に係るシステムで用いる機械学習法と従来の手法との処理結果の比較のためにk近傍法を用いた処理も行った。

【0076】k近傍法は、素性集合だけでなく事例同士の類似度を定義する必要がある。しかし、本例では素性集合F1と素性集合F3は入力された文全体の形態素列をも素性の集合とするので、類似度の定義が困難である。そのため、k近傍法では素性集合F2だけを用いることにする。素性集合F2での類似度の定義としては、事例間で一致した文字列の最長がx-gramのとき、類似度をxとすることにした。

【0077】なお、他の機械学習の手法としては、他にC4.5などの決定木学習を利用する方法があるが、本例では、種々の問題で決定木学習手法が他の手法に比べて劣っていること、また、本例で扱う問題は属性の種類

と精度が落ちるであろうことの二つの理由により用いていない。

【0078】〔第1の例〕まず、K社和英辞典の例文のデータを用いた処理を行なった。その処理結果の精度を図7に示す。本例では、クローズとオープンとの二種類の処理を行なった。オープンの実験は10分割のクロスバリデーションで行なった。図7の括弧内の数字はクローズでの精度を意味する。

【0079】この処理結果から以下のことがわかる。

・決定リスト法は素性集合F2を用いるときに、k近傍法と同程度の精度を得ている。

・最大エントロピー法は、k近傍法または決定リスト法に比べて高い精度を得ている。

・サポートベクトルマシン法は、常に他の手法に比べて高い精度をあげている。

・素性集合の比較としては、最大エントロピー法および決定リスト法では、素性集合F2が最も精度が高く、素性集合F1のように形態素の情報を追加すると逆に精度が下がる結果となっている。これは、素性が増えても不要な素性も増えるために精度が低下したのと思われる。

・サポートベクトルマシン法での素性集合の比較では、素性集合F1で最も高い精度をあげている。これは、サポートベクトルマシン法では形態素の情報の追加が効果があったことを意味する。他の手法では形態素の情報の追加では逆に精度が下がったので、サポートベクトルマシン法では不要な素性を除去し有用な素性を選択する素性選択の能力も他の手法に比べて高いと推測される。

【0080】この結果に対し、手法の理論的な側面からは以下のような説明をつけることができる。

・決定リスト手法は、ある一つの素性のみから解を求める方法のため、不要な素性が多い場合その不要な素性のみを文脈として解を求めてしまいがちになり、不要な素性が多い場合精度が低下する。

・最大エントロピー法は、常にほとんどすべての素性を用いるので、不要な素性が多い場合には精度が低下する。

・これらに対し、サポートベクトルマシン法では、サポートベクトルとなる事例のみを用いそれ以外の事例を用いないといった事例を捨てる操作があるため、多くの不要な素性をこの事例とともに捨てることになり、不要な素性が多くてもそれほど精度低下を招かない傾向がある。

【0081】以上のように、全手法通じて最も精度が高かったのは、 $d = 1$ 、素性集合F1のときのサポートベクトルマシン法であった。

【0082】上記の結果のうち、サポートベクトルマシン法において、素性集合F1を用いる方が素性集合F2を用いるよりも良かった、すなわち、形態素の情報の追加が効果があった、という結果が有意なことなのかを調

べるために符合検定を行なった。これは、 $d = 1$ の方が精度がよかったので $d = 1$ で行なった。全事例39、660個のうち、素性集合F1で正しく素性集合F2で誤った事例は648個であり、素性集合F2で正しく素性集合F1で誤った事例は427個であったが、これを符合検定にかけると0.00000001%（計算では8桁で切っていたため、実際の値はこの値よりももっと小さい可能性がある。）以下の危険率で有意な差があると判定された。このことにより、サポートベクトルマシン法において、形態素の情報を追加する効果があったことは、ほぼ間違いないと考えてよい。

【0083】次に形態素の情報といっても、実際にどのような素性が有効に働いているかを調べることにした。これは、素性集合F1で正しく素性集合F2で誤った事例は、648個に偏って出現している素性を調べることによって行なうことにした。ここでは二項検定を利用して、全事例39、660個での出現確率よりも有意水準1%で大きいと判断されたものを偏って出現しているものとした。

【0084】この有効に働いたと思われる形態素素性の頻度の大きいもの上位20個を図8に示す。図8では、「もう」「最近」「だろう」「まだ」「なければ」「ましよう」「あす」など、テンス・アスペクト・モダリティの推定に役立ちそうな形態素素性が得られており、実際にこういった素性によって精度が向上したものと推測される。

【0085】〔第2の例〕次に白書データを用いて処理を行った。この場合には、精度の良かったサポートベクトルマシン法を用いて行なった。本例でも10分割のクロスバリデーションを行なうことでオープンでの精度を求めている。処理結果の精度を図9に示す。

【0086】この処理結果より以下のことがわかる。

・白書データの精度は、最大で64.67%であった。

・白書データでも文末文字列のみを用いる素性集合F2よりも一文全体の形態素情報も加えて用いる素性集合F1の方が高い精度を得ている。また、白書データでは素性集合F2よりも、一文全体の形態素情報を用いる素性集合F3の方が精度が高い。これらの結果はさらに、一文全体の形態素情報の素性としての有効性を確かめるものとなっている。

【0087】〔第3の例〕次に、K社データを学習データとして、白書データをテストデータとしたような場合、すなわち異分野のデータを教師データとした処理を行なった。本例により、異なる分野のデータを用いると精度がどのようにかわるのかを調べることができた。この処理では精度の良かった $d = 1$ および $d = 2$ のデータを対象として、素性集合F1のサポートベクトルマシン法を用いて行なった。この処理でも学習データとテストデータが重なる場合は重なった部分において10分割のクロスバリデーションを行なうことでオープンでの精度

を求めている。処理結果の精度を図10に示す。

【0088】この処理結果により以下のことがわかる。
・異なる分野のデータを用いると精度が非常に下がった。(白書データを学習データとしてK社データを解析したり、K社データを学習データとして白書データを解析したりすると、精度は10%~20%程度に落ちた。)このことから、入力されるデータと同分野の学習データを用いることが有効であることがわかる。

【0089】人手で書いた規則を用いる手法では異分野に適應したシステムを作るのが難しい。これに対して本発明のような機械学習を用いる方法であれば、学習データをかえて学習し直すことにより、分野ごとに適應したシステムを作るのが容易となる。

・K社と白書の両方を学習データとして用いた場合は、精度はほとんどかわらないか、もしくは少し下がる程度であった。このことから、学習データは多ければよいというのではなく、異分野のデータの場合は、混在させて学習データを用いてもそれほど効果がないことがわかる。

【0090】本例において、テンス・アスペクト・モダリティの翻訳の処理を、k近傍法も含めて様々な機械学習手法を用いて行なった。また検証のため、機械学習法のうちどの方法がもっともよいかを調べた。

【0091】従来の手法(参考文献1等)では、テンス・アスペクト・モダリティの翻訳の際に、素性として文末の文字列しか用いていなかった。本発明では文末の文字列以外にその一文中の形態素情報を追加して用いた。

【0092】その結果、従来では用いていなかった一文中の形態素情報の利用が、処理精度を向上させる効果があることを検定を用いて確認した。

【0093】また、機械学習手法として、決定リスト法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の手法を用いても、従来のk近傍法よりも高い精度でテンス・アスペクト・モダリティの翻訳を行なうことができた。

【0094】特に、サポートベクトルマシン法による方法が最も高い精度を得て、従来手法のk近傍法による手法よりも高い精度で、テンス・アスペクト・モダリティの翻訳を行なうことができた。

【0095】また、異なる分野(本例では、K社英和辞典データと白書データ)のコーパスを用いた処理の例を行なった。この処理では、異なる分野のデータを用いると精度が格段に落ちることを確認し、異なる分野ごとにテンス・アスペクト・モダリティの翻訳システムを構築する必要があることを確認した。このことは、異分野に適應するシステムを人手で作成することが困難であることを考えれば、機械学習手法を用いる本発明の有用性を示すことになる。

【0096】以上、本発明をその実施の態様により説明したが、本発明はその主旨の範囲において種々の変形が

可能である。例えば、本発明の実施の形態では、本発明に係るテンス・アスペクト・モダリティ翻訳システム1は独立して構成されるものとして説明してきたが、他の機械翻訳システムの一部として構成されることも可能である。

【0097】また、機械学習部13で用いる機械学習法は、決定リスト法、最大エントロピー法、サポートベクトルマシン法に限らず、異なる形式の素性を組み合わせることができる方法であればどのような手法であってもよく、また、テンス-素性組抽出部12または素性抽出部15で抽出する素性は、対象となる事例もしくはは入力文から抽出可能な素性であれば種類は限定されないことは当然である。

【0098】

【発明の効果】以上説明したように、本発明は以下のような格別の効果を奏する。

・本発明では、テンス・アスペクト・モダリティの翻訳の問題で、サポートベクトルマシン法に代表される、複数の形式の解析情報(素性)を取り扱うことができるような機械学習法を用いる。これにより、従来の類似度を用いるk近傍法よりも高い精度でテンス・アスペクト・モダリティの翻訳を行なうことができる、テンス・アスペクト・モダリティ翻訳処理方法および翻訳システムを提供することができる。

・本発明では、テンス・アスペクト・モダリティの翻訳の際に、文末の文字列以外にその一文中の形態素情報を新たに追加して用いている。これにより、素性として文末の文字列しか用いていなかった従来の手法に比べて精度の高い翻訳を行うことができる。

・本発明では、機械学習部13で用いるテンス・アスペクト・モダリティと素性の集合の組を抽出するテンス・アスペクト・モダリティ・データベースとして、種々の分野のコーパスを利用することができ、さらにそのコーパスにもとづいて人手によらずに学習結果を取得することができる。これにより、異分野ごとに適應するテンス・アスペクト・モダリティ翻訳システムを容易に実現することができる。

【図面の簡単な説明】

【図1】本発明にかかるシステムの構成例を示す図である。

【図2】本発明にかかるシステムの処理フローチャートである。

【図3】テンス・アスペクト・モダリティと素性の集合の組の例を示す図である。

【図4】サポートベクトルマシン法におけるマージン最大化を説明するための図である。

【図5】入力文からの素性の集合の抽出の例を示す図である。

【図6】実施例におけるテンス・アスペクト・モダリティ・データベースの一部の例を示す図である。

【図7】第1の例におけるテンス・アスペクト・モダリティの翻訳の精度を比較するための図である。

【図8】有効に働いたと思われる形態素素性の例を示す図である。

【図9】第2の例におけるテンス・アスペクト・モダリティの翻訳処理の精度を比較するための図である。

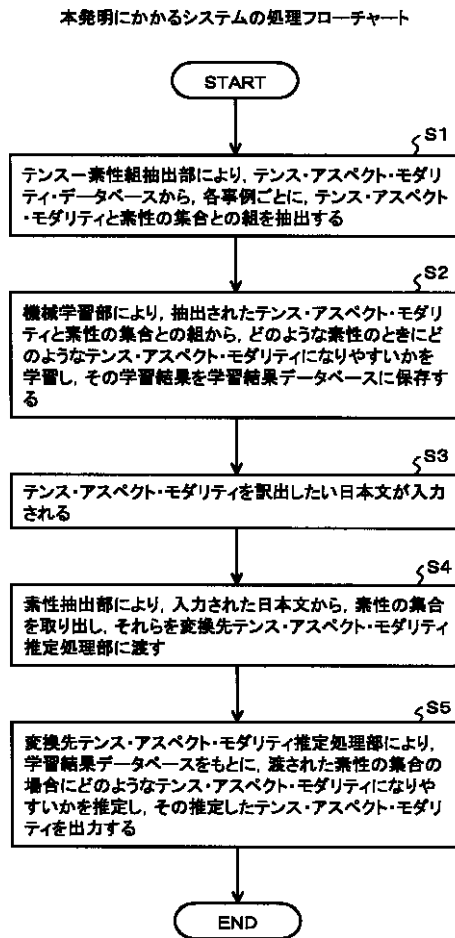
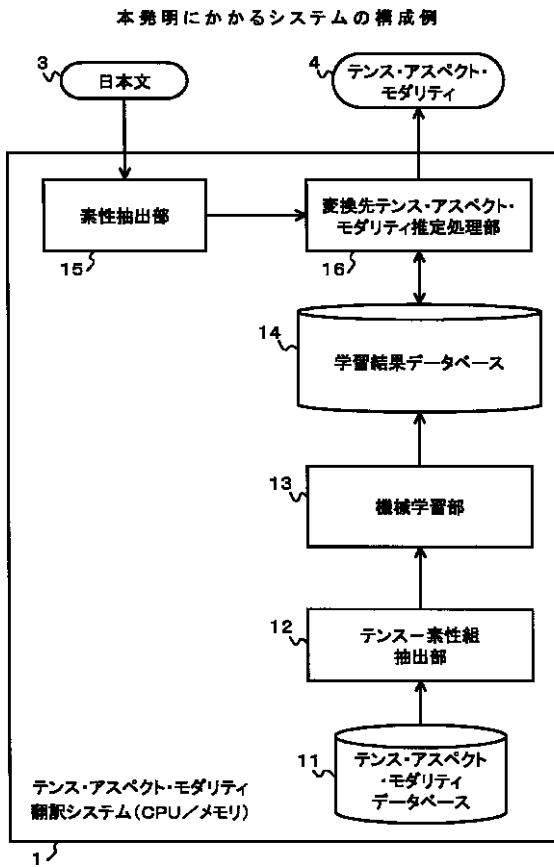
【図10】第3の例におけるテンス・アスペクト・モダリティの翻訳処理の精度を比較するための図である。

【符号の説明】

- 1 テンス・アスペクト・モダリティ翻訳システム
- 11 テンス・アスペクト・モダリティ・データベース
- 12 テンス-素性組抽出部
- 13 機械学習部
- 14 学習結果データベース
- 15 素性抽出部
- 16 変換先テンス・アスペクト・モダリティ推定処理部

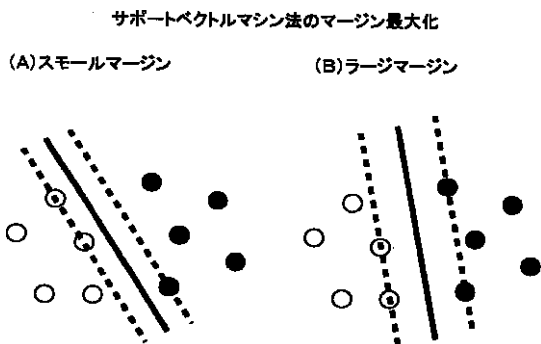
【図1】

【図2】



【図4】

【図9】



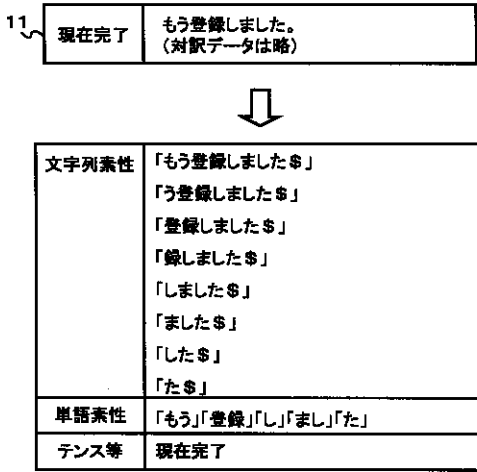
第2の例におけるテンス・アスペクト・モダリティの翻訳の精度

機械学習の手法	素性集合F1	素性集合F2	素性集合F3
サポートベクトル(d=1)	60.10% (99.81%)	56.61% (89.87%)	56.14% (96.67%)
サポートベクトル(d=2)	64.67% (99.81%)	58.74% (89.87%)	62.07% (99.83%)

<括弧()内の数字はクローズの場合>

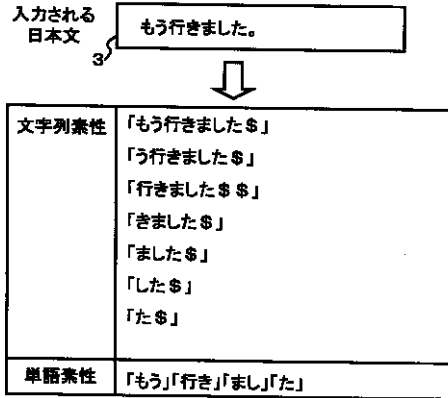
【 図 3 】

テンス・アスペクト・モダリティと素性の集合との組の例



【 図 5 】

素性の集合の抽出の例



【 図 7 】

第1の例におけるテンス・アスペクト・モダリティの翻訳の精度

機械学習の手法	素性集合F1		素性集合F2		素性集合F3	
K近傍法(k=1)	-	-	79.36%	(98.50%)	-	-
K近傍法(k=3)	-	-	80.35%	(83.94%)	-	-
K近傍法(k=5)	-	-	80.43%	(82.39%)	-	-
K近傍法(k=7)	-	-	80.39%	(81.71%)	-	-
K近傍法(k=9)	-	-	80.22%	(81.30%)	-	-
決定リスト	74.19%	(98.21%)	80.23%	(98.18%)	67.90%	(88.58%)
最大エントロピー	80.37%	(88.87%)	81.16%	(83.85%)	75.35%	(84.15%)
サポートベクトル(d=1)	82.48%	(98.70%)	81.93%	(98.50%)	78.68%	(96.68%)
サポートベクトル(d=2)	82.28%	(98.48%)	81.37%	(98.48%)	79.01%	(98.74%)

<括弧()内の数字はクローズの場合>

【 図 6 】

テンス・アスペクト・モダリティ・データベースの一部の例

'	この子供はああ言えばこう言うから小憎らしい This child always talks back to me, and this <v> is </v> why I <vj> hate </vj> him.
d	彼がああおくびょうだとは思わなかった I <v> did not think </v> he was so timid.
c	ああ忙しくては休む暇もないはずだ Such a busy man as he <v> cannot have </v> any spare time.

【 図 10 】

第3の例におけるテンス・アスペクト・モダリティの翻訳の精度

処理対象のデータ	学習に用いたデータ		
	K社と白書	K社だけ	白書だけ
K社データ(d=1)	82.44% (99.71%)	82.48% (98.70%)	65.31% -
K社データ(d=2)	82.31% (98.74%)	82.28% (98.48%)	51.92% -
白書データ(d=1)	60.02% (99.79%)	47.65% -	60.10% (99.81%)
白書データ(d=2)	64.01% (99.83%)	49.53% -	64.67% (99.81%)

<括弧()内の数字はクローズの場合>

【図8】

有効に働いたと思われる形態素素性の例

頻度	形態素素性
221	が
121	いる
89	ない
28	ます
23	なら
23	きた
22	もう
19	中
18	最近
16	だろう
16	ました
14	まだ
13	れる
12	なければ
11	ましよう
11	すっかり
10	られ
9	っ
8	あす
8	なんて

【手続補正書】

【提出日】平成14年7月29日(2002.7.29)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】特許請求の範囲

【補正方法】変更

【補正内容】

【特許請求の範囲】

【請求項1】 コンピュータにより一の言語から他の言語へ翻訳処理をする際に、変換元言語から変換先言語のテンス・アスペクト・モダリティを翻訳する方法であって、予め備えられた変換元言語の事例と当該事例の変換先言語におけるテンス・アスペクト・モダリティとの組を記憶するデータベースであるテンス・アスペクト・モダリティ・データベースにアクセスする過程と、前記テンス・アスペクト・モダリティ・データベースの各事例ごとに、テンス・アスペクト・モダリティと当該テンス・アスペクト・モダリティに対応する事例から抽出した単語素性および文字列素性を含む複数の形式の素性の集合とからなるテンス・素性組を抽出する過程と、前記テンス・素性組を教師データとして用いて複数の素性の出現パターンについて、それぞれのパターンのときになりやすいテンス・アスペクト・モダリティを、決定リスト

法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の機械学習法により学習する過程と、前記機械学習する過程における学習結果を、入力文の変換先言語のテンス・アスペクト・モダリティを判定するための学習データとして学習結果データベースに保存する過程と、変換元言語の入力文から、当該入力文の素性の集合を抽出する過程と、前記入力文の素性の集合をもとに、前記機械学習法により、前記学習データベースに保存された学習データを参照して前記素性の集合の素性の出現のパターンについて、なりやすいテンス・アスペクト・モダリティを特定し、前記入力文のテンス・アスペクト・モダリティの推定解として出力する過程とを備えることを特徴とするテンス・アスペクト・モダリティ翻訳処理方法。

【請求項2】 コンピュータにより一の言語から他の言語へ翻訳処理をする際に、変換元言語から変換先言語のテンス・アスペクト・モダリティを翻訳するシステムであって、予め備えられた変換元言語の事例と当該事例の変換先言語におけるテンス・アスペクト・モダリティとの組を記憶するテンス・アスペクト・モダリティ・データベースと、前記テンス・アスペクト・モダリティ・データベースの各事例ごとに、テンス・アスペクト・モダリティと当該テンス・アスペクト・モダリティに対応す

る事例から抽出した単語素性および文字列素性を含む複数の形式の素性の集合とからなるテンス・素性組を抽出するテンス・素性抽出手段と、前記テンス・素性組を教師データとして用いて複数の素性の出現パターンについて、それぞれのパターンのときになりやすいテンス・アスペクト・モダリティを、決定リスト法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の機械学習法により学習する機械学習手段と、前記機械学習手段における学習結果を、入力文の変換先言語のテンス・アスペクト・モダリティを判定するための学習データとして学習結果データベースに保存する学習結果データベースと、変換元言語の入力文から、当該入力文の素性の集合を抽出する素性抽出手段と、前記入力文の素性の集合をもとに、前記機械学習法により、前記学習データベースに保存された学習データを参照して前記素性の集合の素性の出現のパターンについて、なりやすいテンス・アスペクト・モダリティを特定し、前記入力文のテンス・アスペクト・モダリティの推定解として出力する変換先テンス・アスペクト・モダリティ推定処理手段とを備えることを特徴とするテンス・アスペクト・モダリティ翻訳システム。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0011

【補正方法】変更

【補正内容】

【0011】

【課題を解決するための手段】上記課題を解決するた

め、本発明に係る方法は、コンピュータにより一の言語から他の言語へ翻訳処理をする際に、変換元言語から変換先言語のテンス・アスペクト・モダリティを翻訳する方法であって、予め備えられた変換元言語の事例と当該事例の変換先言語におけるテンス・アスペクト・モダリティとの組を記憶するデータベースであるテンス・アスペクト・モダリティ・データベースにアクセスする過程と、前記テンス・アスペクト・モダリティ・データベースの各事例ごとに、テンス・アスペクト・モダリティと当該テンス・アスペクト・モダリティに対応する事例から抽出した単語素性および文字列素性を含む複数の形式の素性の集合とからなるテンス・素性組を抽出する過程と、前記テンス・素性組を教師データとして用いて複数の素性の出現パターンについて、それぞれのパターンのときになりやすいテンス・アスペクト・モダリティを、決定リスト法、最大エントロピー法、またはサポートベクトルマシン法のいずれか一の機械学習法により学習する過程と、前記機械学習する過程における学習結果を、入力文の変換先言語のテンス・アスペクト・モダリティを判定するための学習データとして学習結果データベースに保存する過程と、変換元言語の入力文から、当該入力文の素性の集合を抽出する過程と、前記入力文の素性の集合をもとに、前記機械学習法により、前記学習データベースに保存された学習データを参照して前記素性の集合の素性の出現のパターンについて、なりやすいテンス・アスペクト・モダリティを特定し、前記入力文のテンス・アスペクト・モダリティの推定解として出力する過程とを備えることを特徴としている。

フロントページの続き

(72)発明者 内元 清貴
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

F ターム(参考) 5B091 AA15 AB13 BA03 CC03 CC16

EA01