

定型表現を利用した効率的な形態素解析の実現

2H-4

安藤 一秋 藤澤 貴之 獅々堀 正幹 青江 順一

徳島大学工学部知能情報工学科

1. はじめに

定型表現は、慣用表現と並び文章中に頻繁に出現する表現である。また、この表現は、それ自体が独自の意味を持つものと、持たないものがあるが、一つのまとまった単位で処理することで、処理効率や精度を上げることができる^[1]。

定型表現や慣用表現の利用は有効であるが、これらは、表現形式の客観的な定義が不明瞭なため、収集作業には時間と労力を要する^[2]。しかし、コーパスなどから自動収集する研究も多数報告されている^{[1][3]}。

このような定型表現を、OCRや音声認識システムなどの自然言語処理システムに導入すると、以下のような長所を生じる。1) 短い単語で検索された候補より、長い単語で検索された候補の方が情報量が多いため、各種処理システムの認識率や誤り訂正率を向上できる。2) 単語レベルだけでなく、構文・意味解析を交えた文脈レベルの誤り訂正が実現できる。

本稿では、有効とみられる付属語的定型表現の抽出結果、及び、それらの定型表現を利用した形態素解析手法について説明する。また、本手法の有効性を確認するため、約7万文のコーパス^[4]を対象にした実験結果を示す。更に、OCRの誤り訂正への応用も示す。

2. 定型表現

本稿では、付属語的定型表現を対象とする。付属語的定型表現には、「に関して」、「に基づく」などの文節と文節の関係を表し、助詞的な働きをする定型表現と、「かもしれない」、「なければならない」などの助述的な働きをする定型表現などがある。これらの表現は、一般に構成語順が不変であり、各構成語間の結合も強く、他の語の挿入も起こらないため、これらの表現を一語として扱うのが自然である。また、処理効率の面からみても一語として取り扱う効果が高く、実用的な場合も多い。

本稿で用いる定型表現は、首藤^[2]が抽出、分類

した付属語的な表現を基にしている。そして、コーパスから頻度を取り、頻度の高い表現と重要であると思われる表現を抽出した。以下にその例を示す。

【助詞的働きをする定型表現】

にもかかわらず、だけでなく、にわたって、について、によって、にとつて、のように、に対する、に対して、のために、とともに、を中心に、ではなく、とはいえ、ことから、に関する、としても、を通じて、に向けて、にしても、を使って、に関して、に基づく、として、による、のため、などの、ために、に対し、をして、とする、のでは、もので、.....

【助述的働きをする定型表現】

られている、になっている、のではない、てしまった、なければならない、かもしれない、わけではない、ないだろう、ものである、ねばならない、てはならない、からである、なのである、に違いない、ようになる、だけではない、なくてはならない、てはいけない、ように見える、そうである、に過ぎない、できている、よくなっている、.....

以下、これらの定型表現を利用した形態素解析について述べる。

3. 定型表現を利用した形態素解析

3. 1. 辞書

単語辞書は、自立語の表記と品詞コードを登録した自立語辞書と、助詞、助動詞などの付属語の表記と品詞コードを登録した付属語辞書を用いる。また、本実験で利用する定型表現は、付属語的表現を扱うので付属語辞書に登録する。

定型表現の辞書登録形式は、一般の単語と違い、表記と最前部形態素の品詞コード、最後部形態素の品詞コードを登録する。ここで、最前部形態素とは、定型表現を形態素に区切った結果の最初の形態素を指し、最前部形態素は、最後の形態素を指す。

例えば、「に関して」を辞書登録するとき、表記の他に最前部形態素の“に”の品詞コード（格助詞）と最後部形態素の“て”の品詞コード（接続助詞）を登録する。

3. 2. 形態素解析

本手法は、最長一致法に接続表を利用した解析手法^[5]を用いる。接続表の構成を以下に示す。

隣接した2つの形態素をmorpheme₁ (前形態素), morpheme₂ (後形態素)とする。接続表はmorpheme₁の品詞コードとmorpheme₂の品詞コードの対で構成され、2次元配列 connect で表す。morpheme₁の品詞コードをi, morpheme₂の品詞コードをjとすると、接続可否はconnect [i, j] の値αにより決定される。

- α = 1 : morpheme₁ と morpheme₂ は接続可能
 - α = 0 : morpheme₁ と morpheme₂ は接続不可能
- α = 1 で接続可能な場合、morpheme₂ を前形態素として解析を続ける。

但し、morpheme₁ または morpheme₂ が定型表現である場合は、以下の操作が必要となる。

- ・ morpheme₁ が定型表現である場合、morpheme₁ の最後部形態素の品詞コードを i とする。
- ・ morpheme₂ が定型表現である場合、morpheme₂ の最前部形態素の品詞コードを j とする。

“形態素に関する資料”の接続に関する解析例(図1)を以下に示す。

morpheme₁ を“形態素”とし、morpheme₂ を“に関する”とする。ここでは、morpheme₂ が定型表現であるので“形態素”の品詞コードと、morpheme₂ の最前部形態素である“に”の品詞コードで接続検定を行う。この場合、名詞と助詞であるから接続可能となる。次に、morpheme₁ を“に関する”とし、morpheme₂ を“資料”とする。ここで、morpheme₁ が定型表現であるので、morpheme₁ の最後部形態素である“関する”の品詞コードと“資料”の品詞コードで接続検定を行う。この場合は、動詞(連用形)と名詞であるから接続可能となる。

4. 実験結果

4.1. コーパスに対する実験

本手法の有効性を確認するため、約7万文コーパスに対して実験を行った。単語辞書には、約97000語の自立語辞書と約1800語の付属語辞書を用いた。

定型表現数、最長定型表現、最長定型表現文字

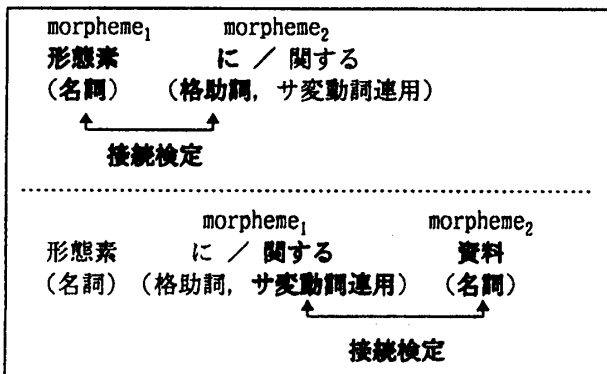


図1：“形態素に関する資料”の接続検定

表1:実験結果

定型表現数	1 5 5 3
最長定型表現文字数	1 1
平均定型表現文字数	5. 0
最大定型表現形態素数	6
平均定型表現形態素数	2. 8
定型表現含有率	1 1. 5 %
計算コスト減少率	6. 6 %

数、平均定型表現文字数、最大定型表現形態素数、平均定型表現形態素数、定型表現含有率、計算コスト減少率をまとめたのが表1である。ここで、定型表現含有率は、全付属語形態素中に、定型表現が含まれる割合を示したものである。計算コスト減少率は、定型表現を一つにまとめて処理することにより、接続に関する計算コストが減少する割合である。

実験結果より、定型表現を形態素解析に利用することで、計算コストの減少を約6.6%確認することができた。

4.2. OCRの誤り訂正への応用

認識率が約88.6%のサンプル文(2951文字)に対して、定型表現の有効性の確認を行った。このサンプル文中には、付属語的なものの認識誤りが約19.3%存在した。この誤りの中から、定型表現に関する誤りの例を示すと、「と・うより」(というより)、「らオ1ている」(られている)、「かあるのか」(があるのか)、「とみ扮1る」(とみられる)などがあつた。

これらの定型表現を辞書に登録すれば、誤り回復は容易に行え、その結果として、付属語の認識誤りが約12.0%に減少し、認識率は、約90.0%に向上する。

5. まとめ

本稿では、定型表現を利用した形態素解析の実験結果について述べた。定型表現を利用することで接続検定の計算コストの減少が確認できた。

今後の課題として、OCRや音声認識システムに定型表現を利用し、構文・意味処理を交えた文脈レベルの誤り訂正を行う予定である。

【参考文献】

[1]北研二,小倉健太郎,森元暹,矢野米雄：“仕事量基準を用いたコーパスからの定型表現の自動抽出”,情報処理学会論文,Vol. 34, No. 9(1993).
 [2]首藤公昭：“文節構造モデルによる日本語の機械処理”,福岡大学研究所報, No. 44(1979).
 [3]新納浩幸,井佐原均：“コーパスからの関係表現の自動抽出”,情報処理学会論文誌,Vol. 35, No. 11(1994).
 [4]EDR電子化辞書,(株)日本電子化辞書研究所.
 [5]田中穂積,“自然言語解析の基礎”,産業図書.