

多要素日本語複合語の構造解析

Structural Analysis of Multicomponent Japanese Compound Terms

小山 照夫
Teruo KOYAMA

何がわかる？

いろいろな分野の論文を探したり、論文の間の関係を発見する手がかりを与えるために、当該分野で用いられる用語の体系化を目指します。

まず、用語として重要な複合語の構造を判定することにより、用語の間の関係を明らかにすることを目的としています。

どんな研究？

複合語はいくつかの基本語の組み合わせとして用いられる用語ですが、3要素以上の語では、その構造がどのようになっているかが大切です。

用語を構成する要素の統計処理により、3要素の日本語複合語の構造を推定します。

研究の目的

大規模商業施設 は、「大規模商業 施設」？ それとも「大規模 商業施設」？

3形態素以上から構成される複合語では、語の構成される構造により、複合語の意味が異なってくる。このため、複合語の構造を正しく把握することは、用語の体系化にあたって重要な課題となる。

この研究では、構造の多義性が発生するもっとも簡単な、3つの形態素からなる日本語複合語について、その構造を正しく推定する方法を研究する。

研究の方法

テキストコーパスを形態素解析し、あらかじめ定められた、名詞系要素の接続からなる、複合語候補を抽出しておく。

これら候補のうち、3要素のもので、コーパス中に2回以上出現するものについて、どちらの2要素が先に結合していると考えられるかを、各形態素と、2要素接続の出現頻度を手がかりに、統計的手法により推定する。

手順

1. あらかじめ形態素解析をおこなったテキストコーパスから、一定の基準で、用語候補となる名詞関連要素の接続を抽出する。
2. このうち、3要素からなる用語候補で、2回以上出現しているものを選択する。ただし、特異性の高い接続を含むものをのぞく。
3. 用語ABCについて、接続ABおよびBCの接続強さを右の基準で推定し、強度の大きいほうが先に結合していると考ええる。

結果

用語候補として抽出されたもの 13,374

特異性の高い候補を除いた結果 10,952

用語としての精度 86% (ランダムサンプルした300データの評価結果)

結合順序が正しく評価できたもの

順序判定精度 82% (ランダムサンプルした300データの評価結果)

非用語でも結合順序が正しければOKとした)

結合強度

$$S = \frac{F(AB)}{\sqrt{F(A) \times F(B)}}$$

ただし

F(A) : 形態素Aの頻度

F(B) : 形態素Bの頻度

F(AB) : 接続ABの頻度