

4. 形態素解析辞書の改変

4.1 擬形態素の辞書への登録方法

擬形態素の左文脈 ID, 右文脈 ID, 単語コスト値を決定するために, まずは擬形態素を構成する形態素の情報を得る必要がある。

外国人向け動詞活用形の場合, 活用形を構成する形態素には規則性があるため, その規則に基づき形態素を辞書より取得する。表2にテ形の例を示す。

表2 活用形を構成する形態素の規則の例

動詞の活用	テ形の規則
五段活用	連用タ接続 + 助詞・接続助詞・テ
上一段活用	連用形 + 助詞・接続助詞・テ
下一段活用	連用形 + 助詞・接続助詞・テ
カ行変格活用	連用形 + 助詞・接続助詞・テ
サ行変格活用	連用形 + 助詞・接続助詞・テ

表2のように擬形態素を構成する形態素を取得し, その各形態素の情報によって, 擬形態素が作成できる。

擬形態素の文脈 ID は以下のように設定する。

左文脈 ID: 前方に位置する形態素の左文脈 ID

右文脈 ID: 後方に位置する形態素の右文脈 ID

擬形態素の単語コストは, 単語コストを C , 擬形態素を構成する形態素の数を n , i 番目の形態素の単語コストを a_i , i 番目と $i+1$ 番目の形態素間の接続コストを b_j とすると, 以下のように求められる。

$$C = \sum_{i=1}^n a_i + \sum_{j=1}^{n-1} b_j$$

※接続コストは辞書の接続コスト表から得られる。

例として, 五段活用の動詞のテ形「書いて」の擬形態素を作成するならば, 「書いて」は形態素「書い」と形態素「て」により構成されており, 左文脈 ID は「書い」の左文脈 ID である 687, 右文脈 ID は「て」の右文脈 ID である 307 を設定する。

単語コスト C は, 「書いて」を構成する形態素の数は「書い」と「て」の2つであるから,

$$\begin{aligned} C &= \text{「書い」の単語コスト} \\ &+ \text{「て」の単語コスト} \\ &+ \text{「書い」と「て」間の接続コスト} \\ &= 7883 + 5170 + (-7392) = 5651 \end{aligned}$$

※接続コストは形態素ペアの頻出度を表し, 多くの場合負の値をとる。

このようにして, 「書いて」の擬形態素情報が作成できる。

4.2 外国人向け活用形のコスト値の調整

活用形を擬形態素として登録するとき, 4.1 節の方法だけでは, 以下のように期待通りに形態素解析が行われない場合がある。

「彼は漢字を書けた」の形態素解析

期待する結果: 書けた (タ形)

実際の結果: 書け (マス形) + た

マス形は～マスに続く形の活用形であり, 「書け」のように連用形語幹のみで構成される。タ形である「書けた」はマス形「書け」と同じく連用形語幹を含むため, 解析の際に競合が起こる。

この競合の解消のため, 優先させたい活用形擬形態素に対しコスト値の調整を行う。コスト値は低いほど優先されるため, 優先したい活用形擬形態素の単語コスト値を, 任意の定数値で下げることで調整できる。

以下は外国人向け活用形の分類と競合の関係である。

- ① 語幹のみで構成される活用形(4種)
- ② 語幹+助詞で構成される活用形(3種, ①と競合)
- ③ 語幹+助動詞で構成される活用形(11種, ①と競合)
- ④ 語幹+助動詞+助詞で構成される活用形(6種, ①・③と競合)
- ⑤ 語幹+助動詞+助動詞で構成される活用形(6種, ①・③と競合)
- ⑥ 語幹+助動詞+助動詞+助詞で構成される活用形(2種, ①・③・⑤と競合)
- ⑦ 語幹+助動詞+助動詞+助動詞で構成される活用形(1種, ①・③・⑤と競合)

5. 考察

4.1 節において, 本稿では擬形態素を構成する形態素の取得に, 活用形を構成する形態素の規則性を利用したが, 規則の利用なしに単に形態素解析を行うことによっても形態素の取得は可能である。しかし, 擬形態素にしたい形態素のかたまりを形態素解析にかけると, 多くは前後の文脈に欠ける場合に, 期待する解析結果が得られないことがある。本稿のように大量の擬形態素を作成する場合には, すべての解析結果が正しいかどうかを検証するのは現実的でないため, 擬形態素の作成に形態素解析を利用する場合には, たとえば MeCab の制約付き解析のような, 補助的な手段が必要と考えられる。

6. おわりに

本稿では, 既存の辞書の情報を流用して擬形態素を辞書登録することで, 形態素解析器に外国人日本語学習者向けの解析を行わせることが可能であることを示した。

本稿にて作成した擬形態素に, 適宜必要な素性情報を付与することによって, さらに利用性が高まるものと考えられる。

謝辞

形態素解析器 MeCab, IPA 辞書の開発に携わった方々, 研究協力していただいた杉野勝也氏に感謝いたします。

参考文献

- [1] 田中よね, 牧野昭子, 重川明美, 御子神慶子, 古賀千世子, 石井千尋, “みんなの日本語 初級 I”, スリーエーネットワーク, 1998.
- [2] 田中よね, 牧野昭子, 重川明美, 御子神慶子, 古賀千世子, 沢田幸子, 新矢麻紀子, “みんなの日本語 初級 II”, スリーエーネットワーク, 1998.
- [3] 形態素解析器 MeCab, <http://taku910.github.io/mecab/>.
- [4] IPA 辞書, <https://osdn.jp/projects/ipadic/>.