

津軽方言音声認識のための Web からの方言資源獲得

武藤 由依[†] 山崎 善啓^{††}
[†] 東北大学工学部電気情報物理工学科

千葉 祐弥^{††} 伊藤 彰則^{††}
^{††} 東北大学大学院工学研究科

1. はじめに

近年、音声認識技術が普及し、医療用システムやコールセンターなどで利用されている。しかし一般的な音声認識システムは東京方言を対象としており、方言音声に対する認識性能はまだ十分とは言えない。

本研究では、方言音声の認識性能向上のため、特に言語モデルの改善を目指す。方言言語モデルを構築する上で問題となるのは、学習に使える言語資源が限られていることである。本稿では、文献 [1]と同様に方言言語資源を Web から獲得することを考える。文献 [1]に対して、本稿では対象方言を津軽方言とし、方言判定器に Multilayer Perceptron (MLP)を用いる方法を検討する。

2. Web からの津軽方言資源の獲得と言語モデル学習

図 1 に検討手法の概要を示す。検討手法ではまず、Web 検索エンジンを用いて津軽方言に関連する Web ページを獲得する。その後、取得した Web ページの各文に対して方言識別を行い、対象言語と判定されたものとあらかじめ保有しているコーパスを組み合わせることで言語モデルの学習に用いる。

ここで、方言判定モデルは文を入力とし、津軽方言か否かの 2 値分類を行う。識別器は隠れ層 3 層の MLP とした。あらかじめ仮名系列に変換された入力文に対して、文字レベルの One-hot 表現を平均したものを入力特徴量とした。

3. 実験条件

本研究では、Web 検索 API として Bing Web Search API を用いた。検索クエリは「津軽弁」、「津軽方言」、「青森弁」に設定し、2,137 件の Web ページを取得した。取得した Web ページは英字や記号等を除外した上で、方言判定モデルへ入力した。

方言判定モデルの学習には方言コーパス [2]を利用した。当該の方言コーパスには津軽方言談話音声の書き起こしと東京方言への対訳文が含まれる。ここでは、全 124 対話に含まれる発話を 3:1:1 の割合で学習セット、開発セット、テストセットに分割した。また、東京方言に対しては、対訳文に加えて新聞記事読み上げコーパス (JNAS) [4]の文書からランダムに選択した 519 文を加えた。MLP の隠れ層のノード数は事前の実験により 64 と設定した。ネットワークのパラメータは開発セットに対して Accuracy が最も高くなるように選択し、最終的な性能をテストセットによって評価した。

言語モデルの学習には、方言判定モデルの学習に用いた学習セットと Web から取得した方言文を組み合わせることで

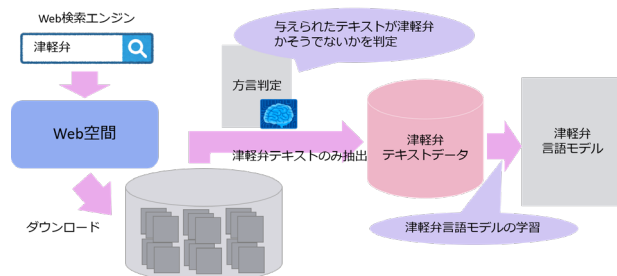


図 1 津軽弁言語モデル学習の流れ

いた。評価時には上記のテストセットに対して学習した言語モデルのパープレキシティを計算した。ここで、一般的な音声認識においては単語単位の言語モデルを用いるが、方言言語の形態素解析は容易ではないため、実験では文字単位での評価を行った。言語モデルの学習と評価には Palmkit を用いた。

4. 実験結果

まず、方言判定モデルの性能を評価した。学習したモデルのテストセットに対する方言の識別精度は 92.8%であり、高い精度が得られた。続いて、方言識別モデルによって津軽方言と判定された Web 上のテキストと方言識別モデルの学習セットを用いて言語モデルを構築し、テストセットに対するパープレキシティを計算した。結果として、55.7 のパープレキシティが得られたが、学習セットのみを用いて計算した数値である 28.3 と比較すると性能は低下した。これは、評価に用いた方言文が談話形式であるのに対して、Web 上に存在する方言文が書き言葉であることに起因していると考えられる。一般的な音声認識システムで入力される音声は談話形式と比較すると短く、Web 上に存在する発話文に近いものが多いため、今後は評価セットの選択方法に関しても検討を行う必要があると考えられる。

5. 結論

本稿では、津軽弁言語モデル作成のための津軽弁テキストデータの Web からの収集方法を検討した。今後は評価用のテキストセットの選択と系列モデルに基づく方言識別手法を検討する予定である。

参考文献

- [1] 廣田他, <https://doi.org/10.11517/pjsai.JSAI2013.0.2B13>, 2013
- [2] 国立国語研究所 (編), <https://doi.org/10.15084/00002241>, 2001-2008
- [3] K. Itou et al., <https://doi.org/10.1250/ast.20.199>, 1999