

漢文読解問題の単語分散表現による解法

不破大智 (福井大学工学部)

小高知宏・黒岩丈介 (福井大学大学院工学研究科)

白井治彦 (福井大学工学部)・諏訪いずみ (仁愛女子短期大学)

1. はじめに

国立情報学研究所の「ロボットは東大に入れるか」プロジェクト[1]において、「東ロボくん」はセンター試験における英語や国語の偏差値を上げることに苦戦したという。そこで国語の中でも点数の低かった古典分野のうち漢文を対象を絞り、自然言語処理によって漢文の問題を解くためにはどのような方法があるか、どのような補助を行えば正答率が向上するかなどについて調べることとした。今回は「選択肢の中から本文中にある下線部の正しい解釈を選ぶ」という読解問題について、単語分散表現を用いて解くことが出来ないかを考える。

2. 文章の単語リスト化と単語の分散表現

単語リスト作成の流れを図1に示す。UD-Kanbun[2]は形態素解析の結果から各単語を対応する英単語に翻訳してくれるが、これをさらに日本語に訳すと二重翻訳になるうえ、漢文単語から英単語への翻訳は単語ごとに行われるため文章として翻訳することが出来ない。そこで選択肢側を英訳してから単語に分割することで、「下線部の単語リスト(英語)」と「選択肢の単語リスト(英語)」を用意する。なお選択肢の英訳は、作弄的な意識にならないよう機械翻訳(今回はGoogle翻訳)を使って行う。

完成した単語リストから fastText の学習済みモデル[3]を用いて各単語の分散表現を取得し、コサイン類似度を使って単語間の類似度を計算する。

3. コサイン類似度を用いた得点付け

図2のように下線部の単語リストと各選択肢の単語リストから総当たりに単語間のコサイン類似度を求め、その結果をもとに得点付けをする。なお、この得点付けは2つの方法で行う。1つ目の方法は、下線部の文を構成するそれぞれの単語について、その単語に対する類似度の最大値を得点に加算する方法である。2つ目の方法はその単語に対する類似度の平均値を得点に加算する方法である。得点の指標とすべく下線部の文と同一の単語で構成される単語リストFを追加して各選択肢の得点を計算したところ、次のような結果が得られた。

表1: 各選択肢の方法1, 方法2それぞれでの得点

	単語数	方法1での得点	方法2での得点
A	19	4.83	2.58
B	22	4.54	2.83
C	15	4.38	2.68
D	17	4.78	2.77
E	15	4.67	2.68
F	10	10.0	3.38

表1において青は正答である選択肢E, 赤は得点が最も高かった選択肢である。どちらの方法においても選択肢Eの得点は1位にならなかった。

4. 考察とまとめ

結果から、この方法で正答を目指すにはなんらかの改良が必要であるという知見が得られた。ほかの問題でも同様の結果になるのかの確認は必要だが、正答の得点が高くならなかったことについて、英語に変換するという手法をとっていることが原因の一つとして考えられる。より原文に近い状態のまま類似度を求められないか、コサイン類似度以外の類似度計算手法で解くことができないか検討することを今後の課題とする。

[1]国立情報学研究所, ロボットは東大に入れるか

<https://2lrobot.org/index.html> (参照 2021-04-13)

[2]漢字情報研究センター, 古典中国語のコーパスの研究
http://kanji.zinbun.kyoto-u.ac.jp/~yasuoka/kyodo_kenkyu/ (参照 2021-07-06)

[3]fastText, Word vectors for 157 languages

<https://fasttext.cc/docs/en/crawl-vectors.html>
(参照 2021-07-06)

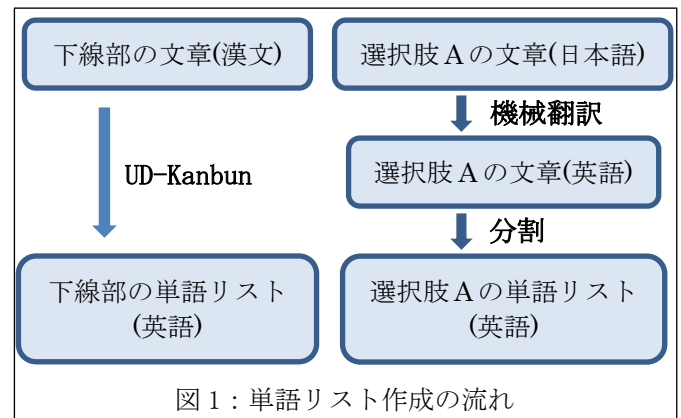


図1: 単語リスト作成の流れ

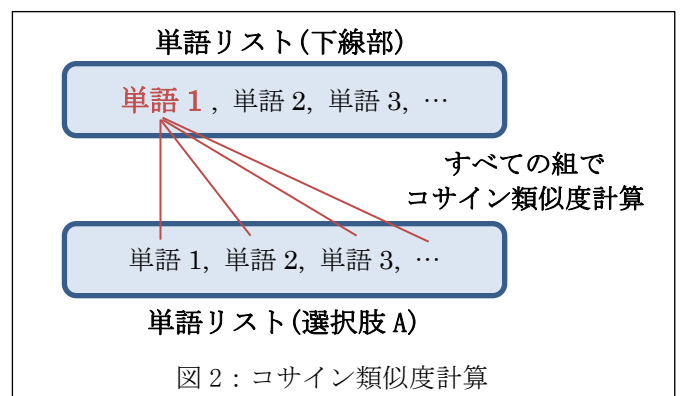


図2: コサイン類似度計算