

学術論文の実験情報分類の評価

上田 和也[†] 新妻 弘崇^{††} 太田 学^{††}

[†] 岡山大学工学部情報系学科

^{††} 岡山大学大学院自然科学研究科

1. はじめに

近年では学術論文データベースが充実し、多くの論文が手軽に手に入るようになった。研究者にとって論文の実験の内容や結果などをまとめることは、研究の整理や比較のために必要であり、これらを自動でまとめることができれば便利である。本研究では、実験情報の書かれた文を5種類の機械学習により5クラスに分け、その分類精度を調べた。

2. 機械学習による実験情報の分類

本研究では、論文本文中の実験情報の記述を文ごとに実験結果、実験内容、データセット、考察、その他の5クラスに分類する。具体的には、実験情報の文をBag of Wordsによってベクトル化し、SVM、ランダムフォレスト、k近傍法、ロジスティック回帰、ナイーブベイズ分類器の5種類の機械学習を利用してこれらの5クラスに分類する。

3. 評価実験

実験データは、NTCIR-9[1]の論文40件を利用した。この論文40件中に実験情報の文が1,157文あり、その文の属する実験情報のクラスは表1の通りである。

表1の実験データから実験情報の文のベクトルを生成する。その後pythonのscikit-learnライブラリ[2]を利用し、先に述べた5種類の分類方法で分類する。実験では5分割交差検定による評価実験を10回行い、正解率の平均を求めた。各方法でパラメータはデフォルトの値を用いた。また、冠詞をストップワードとして除いた場合についても同様の実験を行った。

分類実験の分類正解率は表2のようになった。冠詞あり、なしのいずれにおいてもナイーブベイズ分類器の正解率が最も高く、ロジスティック回帰、SVMの正解率もこれに近い。冠詞を除くとナイーブベイズ分類器の正解率はわずかに向上したが、いずれの分類方法においても結果の差は小さかった。

いずれの分類器でも正しく分類できなかった文の例に“While the other runs are OPEN runs and our runs (including the baseline) are ORACLE runs, ours use summarization techniques only.”がある。この文のクラスは実験

表1 実験データの実験情報のクラス

実験情報のクラス	文数	割合(%)
実験結果	360	31.11
実験内容	414	35.78
データセット	142	12.27
考察	160	13.83
その他	81	7.00

表2 機械学習による分類結果

分類方法	正解率(%)	正解率(%)
	冠詞あり	冠詞なし
SVM	59.85	59.72
ランダムフォレスト	54.04	54.94
k近傍法	40.68	39.64
ロジスティック回帰	60.75	60.61
ナイーブベイズ分類器	62.45	63.14

内容であるが、分類器はすべて実験結果に分類した。これは単語 runs が実験結果クラスの文に多く出現したためだと考えられる。また、k近傍法のみで正しく分類できた例として“The LM adaptation improved the BLEU score by 0.9 points (from 39.14 to 40.04).”がある。この文のクラスは実験結果であるが、k近傍法以外では実験内容に分類された。文が短く単語数が少ないため、正しく分類できなかったものと考えられる。

また、冠詞なしで生成した文ベクトルの分類におけるscikit-learnのパラメータの影響について調べた結果、最も良い正解率がSVMで62.66%、ロジスティック回帰で62.84%、ナイーブベイズ分類器で63.87%となり、それぞれわずかに向上した。

4. まとめ

本研究では、機械学習を用いて学術論文の実験情報について書かれた文を5クラスに分類した。今後の課題としては、文のベクトル化の方法を工夫して正解率を向上させることや、分類した実験情報の活用法の検討などが挙げられる。

参考文献

- [1] NTCIR-9, <http://research.nii.ac.jp/ntcir/ntcir-9/index-ja.html>
- [2] scikit-learn, <http://scikit-learn.org/stable/>