

| | |
|--------------|---|
| Title | マイクロブログにおける皮肉表現を対象とした感情分析 |
| Author(s) | TUNGTHAMTHITI, PIYOROS |
| Citation | |
| Issue Date | 2016-09 |
| Type | Thesis or Dissertation |
| Text version | ETD |
| URL | http://hdl.handle.net/10119/13826 |
| Rights | |
| Description | Supervisor:白井 清昭, 情報科学研究科, 博士 |

| | | | |
|---------|--|----------------|-------------------|
| 氏名 | PIYOROS TUNGTHAMTHITI | | |
| 学位の種類 | 博士(情報科学) | | |
| 学位記番号 | 博情第 348 号 | | |
| 学位授与年月日 | 平成 28 年 9 月 23 日 | | |
| 論文題目 | Sentiment Analysis of Sarcasm on Microblogging | | |
| 論文審査委員 | 主査 | 白井 清昭 | 北陸先端科学技術大学院大学 准教授 |
| | | 飯田 弘之 | 北陸先端科学技術大学院大学 教授 |
| | | Nguyen Minh Le | 北陸先端科学技術大学院大学 准教授 |
| | | 長谷川 忍 | 北陸先端科学技術大学院大学 准教授 |
| | | 高村 大也 | 東京工業大学 准教授 |

論文の内容の要旨

Sentiment analysis of sarcasm in microblogging is important in a range of natural language processing (NLP) applications such as text mining and opinion mining. However, this is a challenging task, as the real meaning of a sarcastic sentence is the opposite of the literal meaning. Furthermore, microblogging messages are short and usually written in a free style that may include misspellings, grammatical errors, and complex sentence structures. This thesis proposes a novel method of sentiment analysis on microblogging that enables us to identify orientation and intensity of the sentiment expressed in the tweets, especially in the sarcastic tweets.

First, we introduce a novel method to identify sarcasm in tweets. It is an ensemble of two supervised classifiers: one is Support Vector Machine (SVM) with N-gram features, the other is SVM with our proposed features. Our features represent intensity of sentiment and contradiction of sentiment derived by a naive sentiment analysis of the tweet. In the sentiment contradiction feature, coherence among multiple sentences in the tweet is also considered, which is automatically identified by our proposed method based on unsupervised clustering algorithm. Furthermore, a way to expand concepts of unknown sentiment words is presented to compensate for insufficiency of a sentiment lexicon. Our method also considers punctuation and special symbols, which are frequently used in Twitter. Results of experiments using two datasets show that our proposed system outperforms baseline systems. The accuracy of sarcasm identification on two datasets is 83% or 76%.

Next, we propose a sentiment analysis system designed for handling sarcastic tweets. To train the model to guess the polarity and intensity of the sentiment in the sarcastic tweets, we used a rich set of features, that are our proposed features used for sarcasm recognition as well as the features grounded on several linguistic levels proposed by the previous work. A decision tree with these features is trained to classify the tweets into an 11-scale score in range of -5 to +5. The system is evaluated on the dataset released by the organizers of the SemEval 2015 task 11. The results show that our method largely outperforms the

systems proposed by the participants of the task on sarcastic and ironic tweets.

Finally, we propose a method for developing a sentiment analysis tool that can guess the fine-grained sentiment score for various types of the tweets. The system consists of two steps. At the first step, the given tweets are classified if they are sarcastic by our sophisticated sarcasm recognition method. At the second step, our sentiment analysis system designed for the sarcastic tweets is used to guess the sentiment scores of the tweets that are judged as sarcasm in the first step. On the other hand, for the tweets judged as non-sarcasm, the three existing sentiment analyzers are applied to guess the sentiment score. The results of the experiments show that our proposed two-steps sentiment analysis system outperforms any single sentiment analyzers on a data set consisting of both sarcastic and non-sarcastic tweets.

In addition, as for the application of the proposed method, our technique to recognize the sarcasm is integrated to an existing target-dependent sentiment analysis system. We also show that the integration can improve the performance via the experiments using a relatively small data set consisting of three targets.

Keywords: Sarcasm, Microblogging, Sentiment analysis, Coherence, Concept knowledge, Machine learning, Clustering

論文審査の結果の要旨

本論文は、皮肉を含むマイクロブログのテキスト(ツイート)に表明されているユーザの意見に対し、その極性スコアを推定する新しい手法を提案している。ここでの極性スコアとは、+5 から-5 の範囲の整数で、符号は肯定的か否定的かを、スコアの絶対値は意見の強さを表わす。複数の文からなるツイートに対し、それらの文の一貫性の有無と、文に出現する単語の極性の矛盾の有無を手がかりとして、皮肉を含むツイートの極性スコアを推測するモデルを構築する点に特長がある。

まず、ツイートが皮肉を含むか否かを判定する分類器を教師あり機械学習によって構築した。学習に用いる素性は以下の3つである。1つ目は、ツイートに含まれる単語の極性である。2つ目は、ツイート内の単語極性の矛盾の有無、すなわちツイート内に肯定的な単語と否定的な単語が含まれているか否かである。ただし、極性の矛盾は皮肉の存在を示唆するが、互いに無関係な文に肯定的な単語と否定的な単語が出現した場合には皮肉でない可能性が高い。そのため、提案手法では文の一貫性も学習素性に加える。文の一貫性とは、ここでは複数の文が言及しているトピックが同じことを指す。文の一貫性を判定するために、クラスタリングに基づく新しい手法を提案し、さらにクラスタリングの際に用いるツイートの特徴量ベクトルの重みを最適化する手法を考案した。3つ目は、句読点や顔文字など、マイクロブログ特有の表現の有無である。さらに、1つ目と2つ目の素性を得る際に、既存の感情語辞書に記載されていない単語の極性を推定するために、単語からその関連語を自動的に獲得する「概念拡張」と、関連語の中から文脈に合致したものだけを選別する「概念選択」の手法を提案した。

次に、皮肉を含むツイートの極性スコアを推定した。前述の皮肉の有無を判定するモデルで用いた素性と、先行研究で採用されている素性を採用し、これらの豊富な素性を用いて極性スコアを推定するモデルを機械学習した。さらに、皮肉を含むか否かに関わらず、任意のツイートの極性スコアを判定するモデルを構築した。まず、提案手法により皮肉であるかを判定し、皮肉である場合にはやはり本論文で提案する皮肉を含むツイートに特化した極性スコア推定モデルで、皮肉でない場合には既存の3つのツールで極性スコアを推測する。評価実験では、ツイートに対する極性スコアを推定するタスクにおいて、提案手法は過去の研究を上回る最高の成績を得た。

以上、本論文は、表層上の意味と真の意味が異なるために解析が難しい皮肉表現に焦点を合わせつつ、マイクロブログ上の意見の解析に取り組み、優れた成果を示したものであり、学術的に貢献するところが大きい。よって博士(情報科学)の学位論文として十分価値あるものと認めた。