

D-38 シード表現を用いたマイクロブログに存在する観光情報の極性判定

前田 裕十 井上 俊右 † 嶋田 和孝 † 遠藤 勉 †

†九州工業大学大学院情報工学府 †九州工業大学情報工学部

1 はじめに

近年、観光は様々な地域において基幹産業の1つになっている。観光の活性化は、地域の活性化に繋がると考えられる。そのような環境下で、Webに存在する多くの観光情報の分析は重要な役割を持っている。

本論文では、Twitterと呼ばれるマイクロブログ上に存在する観光情報を対象にして評判分析を行う。Twitterは、他のメディアと比べ、現地性やリアルタイム性、強い主観情報といった特徴を持っている。こうした特徴は、観光的なイベントや事象の分析に対して有用であると考えられる。

評判分析の方法としては、極性判定を行う。極性判定とは、その文が肯定的な意見(P)なのか、もしくは否定的な意見(N)なのかを分類することとする。自動的に獲得された訓練データに基づく機械学習を適用し、その有効性を検証する。

2 提案手法

テキスト極性判定については、様々なアプローチがある。その一つは機械学習を用いた手法である。Pangら[1]は映画のレビュー記事を対象とした極性判定において、いくつかの機械学習に基づく手法の有効性および比較を行っている。機械学習を用いれば、高い精度が見込めるが、一方で高い精度を得るためには十分な訓練データが必要となる。しかしながら一般に、訓練データとなるタグ付きコーパス作成のための作業は高コストである。

本論文では、人手によるタグ付きコーパスなしで、高精度な極性判定を行う枠組みを導入する。Turney[2]は、語やフレーズの極性値を自動的に算出するために、“excellent”や“poor”といった肯定および否定の極性をよく表す語との共起度を利用する手法を提案した。本手法もこれに倣い、シード語(種語)から自動的に訓練データを作成する。

まず、大量のtweet(Twitter上での投稿文)からシードを含むtweetを抽出する。シードと共起した単語群を肯定の表現もしくは否定の表現として仮定し、機械学習に適用する。なお、本論文では機械学習器としてSVMを用いる。

3 実験・考察

本節では提案手法の有効性を検証する。評価データとして、人手で116 tweetsをラベル付けた。そのうち、Pの数は64であり、Nの数は52であった。

本手法はシードを用いて訓練データを獲得する必要がある。そのためのタグなしデータとして、100万 tweetsを利用した。この100万 tweetsに対して、PまたはNのシードがマッチしたものがそれぞれPまたはNの訓練事例である。

シードはPやNを象徴する単語や顔文字、記号を主観的に選択した。その内いくつかのシードパターンにおける極性判定結果を表1に示す。なお、複数のシードの場合は、それらのうち、いずれかが含まれるtweetを訓練データとして用いるものとする。

表 1: シード毎の極性判定の精度

Pのシード	Nのシード	精度
	ToT	0.61
	残念	0.68
	orz	0.73
	orz, 残念	0.72
	orz, #), ない	0.75
良	orz, #), ない	0.64
, ^ ^)	orz, #), ない	0.78
, ^ ^),)	orz, #), ない	0.77

様々なシードパターンを試したところ、Pのシードに「, ^ ^)」, Nのシードに「orz, #), ない」を用いた場合に最も高い判定精度0.78を得た。精度を変化させるいくつかの原因を以下のように考察する。

「ToT」をシードに用いると極端に精度が低くなった。これは、「ToT」という顔文字が嬉し泣きや感涙といったPositive表現としても用いられることが比較的多いためであると考えられる。

「orz」は人間が跪く様を表したアスキーアートであり、落胆、失意、挫折の心理状態を表現する際に用いられる。「残念」も似たような意味を持つと考えられるが、「残念」と「orz」の精度に着目すると「orz」の方が有効に働いた。これは「残念」に比べ「orz」の方がより多様なNegativeな発言を網羅できるためだと考えられる。例えば「orz」を含むtweetでは、「食べ過ぎた orz」「起きれるかな...orz」のように、主観的疲労感や微弱な否定的感情を表現する際にも用いられる。それに比べると、「残念」を含むtweetが表現し得る範囲は限られている。

また、Twitterの文書には様々な文体が存在し、文体のドメイン毎に表現が異なる。たとえば、硬派な文体を使うユーザの間では顔文字はあまり用いられない。そのため、シードとして用いる語と、極性判定の対象となる文書の文体が一致していなければ期待する精度が得られない。

4 おわりに

本実験ではTwitter上の観光性のある文書に対して極性判定を行った。今後は感情や文体のドメイン毎に分割統制的に極性判定を行うことで精度の向上を目指す。

参考文献

- [1] B. Pang, L. Lee, and S. Vaithyanathan. Thumbs up? sentiment classification using machine learning techniques. In Proceedings of EMNLP, pp. 79-86, 2002.
- [2] P. D. Turney. Thumb up? or thumbs down? semantic orientation applied to unsupervised classification of reviews. In Proceedings of ACL, pp. 417-424, 2002.