



Disambiguating context-dependent polarity of words: An information retrieval approach



Olga Vechtomova

Department of Management Sciences, University of Waterloo, 200 University Ave. W., Waterloo, ON, N2L 3G1, Canada

ARTICLE INFO

Article history:

Received 25 October 2016
Revised 16 March 2017
Accepted 27 March 2017

Keywords:

Sentiment analysis
Polarity disambiguation
Word polarity
Context-dependent polarity of words

ABSTRACT

The paper introduces *PolaritySim* – a novel approach to disambiguating context-dependent sentiment polarity of words. The task of resolving the polarity of a given word instance as positive or negative is addressed as an information retrieval problem. At the pre-processing stage, a vector of context features is built for each word w based on all its occurrences in the positive polarity corpus (consumer reviews with high ratings) and another vector – on its contexts in the negative polarity corpus (reviews with low ratings). Lexico-syntactic context features are automatically generated from dependency parse graphs of the sentences containing the word. These two vectors are treated as “documents”, one with positive and one with negative polarity. To resolve the contextual polarity of a specific instance of the word w in a given sentence, its context feature vector is built in the same way, and is treated as the “query”. An information retrieval (IR) model is then applied to calculate the similarity of the “query” to each of the two “documents”, with the polarity of the best matching “document” attributed to the “query”. The method uses no prior polarity sentiment lexicons or purposefully annotated training datasets. The only external resource used is a readily available corpus of user-rated reviews. Evaluation on different domains shows more effective performance compared to state-of-the-art baselines, Support Vector Machines (SVM) and Multinomial Naive Bayes (MNB) classifiers, on three out of four datasets. *PolaritySim*, SVM and MNB were also evaluated with an out-of-domain training corpus. The results indicate that *PolaritySim* is more effective and robust when used with an out-of-domain corpus compared to SVM and MNB. We conclude that an IR based approach can be an effective and robust alternative to machine learning approaches for disambiguating word-level polarity using either within-domain, or out-of-domain training corpora.

© 2017 Elsevier Ltd. All rights reserved.

1. Introduction

The popularity of online review sites has led to an abundance of content written by consumers. For example, a recently released Amazon corpus (McAuley, Targett, Shi, & van den Hengel, 2015) contains 142.8 million reviews across a wide range of categories covering the period from 1996 to 2014. Most consumer reviews have overall ratings, representing the reviewer’s satisfaction with the product or service. Rated reviews are readily available sources of rich contextual information representing how words are used in positive and negative contexts. We propose *PolaritySim* – an extensible method for identifying the context-dependent polarity of words expressing an opinion about another word or phrase (opinion target). The only external resource required is a review corpus with user-assigned numerical ratings. The method determines sentiment valence

E-mail address: ovechtom@uwaterloo.ca

<http://dx.doi.org/10.1016/j.ipm.2017.03.007>

0306-4573/© 2017 Elsevier Ltd. All rights reserved.

of words with ambiguous (e.g. “small”) or unambiguous (e.g. “beautiful”) sentiment, as well as words that do not carry sentiment valence on their own, but acquire it through context. For example, it correctly determines the negative polarity of “eat” in “This camcorder eats up tape”. The task of disambiguating the polarity of a given word instance as positive or negative is addressed as an information retrieval (IR) problem. At the pre-processing stage, we build one vector of all contexts of the word w in the positive set (i.e. reviews with high ratings) and another vector – of its contexts in the negative set (reviews with low ratings). The lexico-syntactic context features are automatically generated from the dependency parse graphs of all the sentences containing the word w in the positive or the negative corpus. The resulting positive and negative vectors are treated as “documents”. At run time, to determine the polarity of a specific instance of w in an unlabeled review, a context vector is built, which is treated as the “query”. The context features for this vector are derived only from the current sentence containing this instance of w . An information retrieval model is then applied to calculate the similarity of the “query” to each of the two “documents”.

The *PolaritySim* method is extensible in a number of ways. For example, the words in the context features could be expanded with related words or the feature set can be expanded with co-occurring patterns from adjacent sentences. Section 4.3 describes one such extension, whereby words in the context features are expanded with related words generated using a Word2Vec model.

The rest of the paper is organized as follows: Section 2 outlines the motivations and contributions of this work, Section 3 discusses related work, Section 4 presents the method, Section 5 describes the datasets and evaluation experiments, Section 6 contains the analysis of results, and Section 7 concludes the paper and suggests future research directions.

2. Motivation and contributions of the work

Most research efforts in the sentiment analysis field have been directed at identifying sentiment and its polarity at the sentence or document level. Two major sentiment analysis approaches to date have been: (a) lexicon-based and (b) machine learning based. In the first approach, the polarity of individual words is first determined by using a prior polarity lexicon, then possible polarity shifters are identified, usually by applying hand-crafted rules. Sentence or document level polarities are then calculated by using word counting methods. In the second approach, the machine learning models have to be trained on the training datasets manually labeled at the same level of granularity (phrase, sentence or document) as the test dataset. The main limitation of these approaches is their reliance on external resources, such as lexicons in the lexicon-based approaches and purposefully-built training datasets in the machine-learning based methods, which typically require substantial human effort to construct.

The main objective of this research is to develop a method for disambiguating contextual sentiment polarity at the lowest level of granularity – words, without relying on any purposefully-built training datasets and lexicons. Polarities of individual words are highly dependent on their context. Among the factors that can affect word polarity are: the target of opinion, e.g. “long rebooting time” (negative) vs. “long battery life” (positive), whether the word is used ironically or sarcastically, presence of phrases intensifying, reversing or diminishing the polarity of the word (e.g. “never”, “too”, “barely”, “hardly”, “even”). Instead of relying on prior lexicons, manually labeled datasets or a large number of handcrafted rules to capture different kinds of polarity shifters, the proposed method determines contextual polarities by using an IR approach and a large body of readily available user-rated reviews. It, therefore, eliminates the need for the manual effort required to build lexicons or datasets.

The method can be readily applied to different categories of user reviews, due to the availability of large datasets of user-rated reviews. It can also be applied to the categories and domains for which no user-rated reviews exist by using out-of-domain reference corpora.

The main theoretical contribution of this research is demonstrating that an IR approach to determining word-level contextual polarity can achieve performance that is comparable to or better than the performance of the state-of-the-art machine learning approaches. The paper also shows that the proposed approach is more robust with out-of-domain training corpora than the state-of-the-art machine learning approaches.

A number of practical applications can benefit from knowing word-level contextual polarity, such as generating text with custom recommendations for users based on existing reviews, extraction of specific positive and negative expressions referring to an entity, multi-document summarization of reviews, as well as question answering and information retrieval for complex information needs. For example, if a user has the information need: “Find a camera that works well in poor lighting conditions”, it would be useful for the system to know the contextual polarities of “sharp” and “low” in the sentence “The picture was sharp, even in low light.”, so that it can determine whether it is relevant to the user’s information need.

The specific contributions of this work are summarized below:

- The proposed method uses reference corpora with document-level positive and negative polarity labels to disambiguate context-dependent polarity of individual words in an unlabeled document;
- Readily available user reviews with overall numerical ratings are demonstrated to be effective positive and negative reference corpora for determining word-level contextual polarity;
- The method is compared to state-of-the-art baselines and proves to be more effective on three out of four datasets, achieving accuracy in the 83%–91% range in different subject domains using within-domain reference corpora;

Download English Version:

<https://daneshyari.com/en/article/4966396>

Download Persian Version:

<https://daneshyari.com/article/4966396>

[Daneshyari.com](https://daneshyari.com)