

State of the Art in Twitter Sentiment Analysis

Christian B.
Hotz-Behofsits
0929002
christian.hotz-
behofsits@tuwien.ac.at

Thomas Schmidleithner
1025525
tschmidleithner@auto.tuwien.ac.at

Dominik Pichler
1026045
dominik.pichler@aon.at

Matthias Reisinger
1025631
matthias.reisinger@web0.at

Florian Taus
0627918
florian.taus@hotmail.com

ABSTRACT

The extensive growth of user-generated content in social networks and the common usage of emoticons and hashtags has introduced new possibilities to classify these information. In this paper, an overview of the state-of-the-art regarding sentiment analysis of Twitter messages is provided. The usefulness of existing lexical resources and special expressions like smileys or hashtags is evaluated and an general introduction to sentiment analysis is given as part of the introduction.

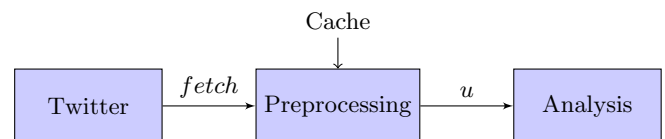
1. INTRODUCTION

Within the last years Twitter and similar social media platforms grew considerably in terms of user numbers and mainstream fame. Twitter has about 284 Mio. active users, generating approximately 500 Mio. posts¹ on a single day². These impressive numbers show the amount of information generated by the crowd. Furthermore Twitter is able to reach a vast number of potential customers especially for consumer businesses. As a result the named services gain more and more acceptance as opinion platforms and powerful tools for both, opinion leader as well as pollster. This trend influenced various marketing and sales strategies and created a new market for companies specialized on sentiment analysis of tweets, like tweetfeel³, Social Mention⁴ and Twitratr⁵.

The aim of sentiment analysis is to determine and try to measure positive and negative feelings, emotions and opinions written in a text. English as language allows to express

the same intent in different ways. The main challenge therefore consists in abstracting the intention of the writer from the grammatical and language specific rules.

In general sentiment analysis can be splitted into two steps: a preprocessing and a constitutive sentiment analysis phase. In Twitter sentiment analysis there is another step right bevor preprocessing: fetching the data from the application programming interface (API). This task is not trivial, because the official API is limited regarding the datasets fetched within certain intervals. This leads to the requirement of a caching mechanism to avoid fetching the same tweets multiple times and allow the usage of more datasets.



While sentiment analysis of conventional resources like news papers or articles is quite well investigated the special research area of sentiment analysis in terms of social network posts is relatively unexplored. Especially the possibility of prioritizing keywords⁶ and the common usage of emoticons lead to advanced possibilities to categorize feelings of posts.

2. RELATED WORK

Sentiment analysis a growing and well explored part of Natural Language Processing (NLP). There are several papers from Pang and Lee regarding sentiment analysis in general and related topics like the effects of various machine learning approaches [9][8]. Especially the last topic is a well studied field [7]. Some related research areas are document level classification, sentence-level classification and machine learning. Due the limited number of characters within a single twitter post, this topic is quite similar to sentence-level sentiment analysis. There are several recent papers regarding sentiment analysis in the context of twitter posts. Agarwal et al. published within "Sentiment analysis of Twitter data" general information and techniques which cover this area [1]. In "Twitter Sentiment Analysis : The Good the Bad and the OMG !" the utility of linguistic features for detect-

¹posts are also known as tweets

²<https://about.twitter.com/company>, Effective 12.11.2014

³www.tweetfeel.com

⁴www.socialmention.com

⁵www.twitratr.com

⁶known as hashtags

ing the sentiment of tweets are investigated. Saif, He, and Alani introduce a novel approach of adding semantics as additional features into the training set. For each extracted entity (e.g. Galaxy S) from Twitter Messages, they add a semantic concept (e.g. Samsung product) as an additional feature, and measure the correlation of the representative concept with negative/positive sentiment [10]. In *Twitter Sentiment Analysis* a algorithm is presented, which accurately classifies tweets as positive or negative, in respect to a query term[3].

3. DATA

Twitter is a social network and microblogging service that allows their users to compose and read so called *tweets*. The user registration is not limited to professional writers, authors or similar expert groups and there are no quality safeguards or writing guidelines.

Each tweet is a text message and can contain up to 140 characters including links, pictures⁷, and special tagged words. These special tagged words start with a well defined prefix and add semantic information to the prefixed string. Table 1 shows an overview of common prefixes supported by Twitter and their purpose. A hash (#) labels a word, also known as hashtag. It is used to mark important keywords within tweets. There is also a at-sign-prefix @, which allows to mention and link other Twitter users in a post. Tweet composers can insert URLs as part of their messages. Regarding the hard length limitation of a tweet, url-shortener are often used to insert longer internet addresses.

prefix	name of character	meaning
#	hashtag	keyword, index
@	at-sign	twitter username

Table 1: Common Twitter word prefixes

The combination of the 140 character restriction and the missing writing guidelines leads to another problem: Slang, emoticons and punctuation that are uncommon in longer text messages are heavily used.[2] In terms of sentiment analysis, tweet processing is more comparable to sentence level analysis as paragraph analysis.

4. PREPROCESSING

The data preprocessing step is an essential part of sentiment analysis. Its goal is to prepare data for the sentiment analysis and remove noisy, unnecessary, inconsistent or incomplete parts.[4] In a first step a tokenizer splits the tweet into conjugate parts, called words. The resulting list of words is filtered regarding the word class, the word itself and the intention of the sentiment analysis. In case of sentiment analysis of twitter posts it is common to remove the parts described in the following paragraphs.[4]

URLs. Uniform Resource Locators (URLs) refer to another site, which in fact can contain very important information, but the internet address on its own is not of interest, because

⁷pictures are replaced by a url targeting a public available image resource

in most cases it contains only auto-generated information, which can not be influenced by the poster.

Unnecesarry words. It is not an obvious task to classify if a word is unnecessary or not. But some word classes, like pronouns, articles or interrogatives do not contain any relevant information and can be skipped. (TODO:Referenz)

Retweets. Retweets reflect the position of another person and are not always created to express the same opinion (e.g. statement of sarcasm).

Repeated characters. Tweets contain often deformed words containing repeated characters like 'heeeey' or 'whuhuuu'. In this cases it is intended to reconstruct the correct english term and only if it is not possible or too sophisticated, the word is removed by the preprocessor.

4.1 Tokenizer

There are already several open source tokenizer implementations available. Some popular programming libraries for the english language include the Stanford NLP PTBTokenizer⁸ and the Apache OpenNLP Tokenizer⁹. These examples are quite versatile and can even be used for different tasks and languages. To overcome the special needs of sentiment analysis in terms of Twitter posts, highly specialized Tokenizer were created (e.g. Carnegie Mellon Twokenizer¹⁰). Some tokens like emoticons are difficult to handle. Although they are a good and popular resource for sentiment analysis. [5] It is a sophisticated task to recognize them correctly and some Tokenizers like the Stanford NLP PTBTokenizer even skips them per default. The Carnegie Mellon's Twokenizer recognizes the related parts and forwards them to the Tagger.

4.2 Tagger

Part-of-Speech (POS) tagging is a simple form of syntactic analysis and therefor very useful for Natural Language Processing. Tagging Tweets is very challenging, because of the twitter specific slang, emoticons and the lack of conventional orthography. spelling errors, bad structure, abbreviations, and slang There is a variety of freely downloadable POS taggers¹¹ e.g. Stanford POS tagger¹² but just a few are specialized on tagging twitter data like the TweetNLP tagger¹³. In comparison with the Stanford POS tagger the TweetNLP tagger reduces errors by 25%.

5. CLASSIFICATION

Sentiment analysis of twitter data, or more general, the idea of extracting sentiment and opinions from pieces of text is based on the more general principle of *classification* [9]. The broad aim of classification is to assign given textual units

⁸<http://nlp.stanford.edu/software/tokenizer.shtml>

⁹<http://opennlp.apache.org/documentation/manual/opennlp.html>

¹⁰<http://www.ark.cs.cmu.edu/TweetNLP/>

¹¹<http://www.nlp.stanford.edu/links/statnlp.html>

¹²<http://nlp.stanford.edu/software/tagger.shtml>

¹³<http://www.ark.cs.cmu.edu/TweetNLP/>

to a set of classes or categories or the apply some kind of regression or ranking. Sentiment analysis is a specialization of this approach which aims at assigning sentiment values to documents. One application might involve to classify an opinionated text by assigning one of two opposing sentiment polarities, i.e. classifying it as either *positive* or *negative*. This classification task is also referred to as *binary classification* or *sentiment polarity classification* [9]. But in general, the input to the sentiment classification process is not strictly opinionated which makes this task challenging. Therefore this kind of binary classification might not always be applicable. Different approaches that allow for a more fine grained classification might be appropriate, for example, based on a multi-point scale that allows for more than just two sentiment classes.

The characteristics of twitter messages introduce further challenges. Due to the structure of these messages, sentiment analysis of tweets is different from analysing conventional texts, such as review documents, in various ways. The length of at most 140 characters and the rather informal spelling style pose problems that have to be considered carefully when analysing the data. When preprocessing the texts these aspects already need to be handled so that the actual classification process is supplied with information that is considered useful for the analysis phase. After this step it is necessary to collect the relevant information as so called *features* and organize them into a *feature vector*.

5.1 Features

TODO

5.2 Methods

For the actual classification process extant literature distinguishes two general types of techniques, namely *machine learning* and *semantic orientation* [Ye20096527]. The machine learning approach is also referred to as *supervised learning*, since corresponding techniques generally require supervised training phases. Accordingly, semantic orientation is also called *unsupervised learning* because such a training phase is not necessary for the process to work. The most prominent machine learning methods comprise *support vector machines* and the *Naïve Bayes* classifier.

Naïve Bayes

TODO

$$P(c|d) = \frac{P(c)P(d|c)}{P(d)}$$

$$P_{NB}(c|d) := \frac{P(c)(\prod_{i=1}^m P(f_i|c)^{n_i(d)})}{P(d)}$$

Maximum Entropy

TODO

$$P_{ME}(c|d) := \frac{1}{Z(d)} \exp \left(\sum_i \lambda_{i,c} F_{i,c}(d, c) \right)$$

Support Vector Machines

TODO

Further techniques

TODO

6. CONCLUSIONS

APPENDIX

A. HEADINGS IN APPENDICES

References

- [1] Apoorv Agarwal et al. “Sentiment analysis of Twitter data”. In: *Association for Computational Linguistics* (2011), pp. 30–38.
- [2] Alexander Davies and Zoubin Ghahramani. “Language-independent Bayesian sentiment mining of Twitter”. In: (2011).
- [3] Alec Go, Lei Huang, and Richa Bhayani. *Twitter Sentiment Analysis*. Tech. rep. 2009, p. 17.
- [4] Govardhan Hemalatha Saradhi-Varma. “Preprocessing the Informal Text for efficient Sentiment Analysis”. In: *International Journal of Emerging Trends & Technology in Computer Science (IJETTCS)* 1 (2012), pp. 58–61.
- [5] Alexander Hogenboom et al. “Exploiting emoticons in sentiment analysis”. In: *Proceedings of the 28th Annual ACM Symposium on Applied Computing*. ACM. 2013, pp. 703–710.
- [6] Efthymios Kouloumpis, Theresa Wilson, and Johanna Moore. “Twitter Sentiment Analysis : The Good the Bad and the OMG !” In: *Artificial Intelligence* (2011), pp. 538–541.
- [7] Christopher D Manning and Hinrich Schütze. *Foundations of Natural Language Processing*. 2000, p. 678.
- [8] Bo Pang and Lillian Lee. *Opinion Mining and Sentiment Analysis*. 2008. arXiv: 0112017 [cs].
- [9] Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. “Thumbs up?: sentiment classification using machine learning techniques”. In: *Proceedings of the ACL-02 conference on Empirical methods in natural language processing - EMNLP '02*. Vol. 10. Association for Computational Linguistics, 2002, pp. 79–86.
- [10] Hassan Saif, Yulan He, and Harith Alani. “Semantic sentiment analysis of twitter”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Vol. 7649 LNCS. PART 1. 2012, pp. 508–524.