

Sentiment Analysis of Soldiers' Tweets - Comparison with civilians (TBC)

Sumit Mukhija, Rachit Rastogi,
Chao Chen, Chen Wang, Chetan Prasad
School of Computer Science and Statistics, Trinity College Dublin
{mukhijas, rrastogi, chenc1, wangc5, cprasad}@tcd.ie

April 13, 2020

Abstract

The concern to veterans' mental health should be made. Existing works show that mental health changes caused by wars can be reflected in linguistic features of the social media texts. In order to detect and compare those changes we collected data from 20 soldiers' tweets and examined them with a list of positive and negative adjectives to identify the polarity and do a comparison with normal users' tweets. The total counts of tweets vary from 57 to 39,000. We identified the difference between normal users and soldiers and we did a close look to the result with discussion.

Keywords:

1 Introduction

Social media platforms and microblogging websites are some of the most popular online stages for people to express their views. Twitter, undeniably is one of the leading applications in this assortment. People use Twitter to post their real-time opinions in the form of tweets. These tweets can be analyzed and certain inferences can be extracted. These inferences can subsequently be used for academic and business purposes.

One of the primary reasons that make Twitter a feasible choice is the diverse nature of the users. In this research, we intend to analyze and compare the tweets of the war-veterans and the general public. We believe wars have an impact on soldiers' psychological and emotional states. We try to prove this hypothesis by comparing their tweets to the tweets posted by the civilians.

We collect public data using Twitter API and then process and count the words with a list of positive and negative adjectives to predict the polarity of the tweets. Then we examine a randomly collected dataset to compare the difference between tweets by veterans/soldiers and civilians.

(TBC due to the experiment implementation)

The remainder of the paper is organized as follows. We examine on the literature related to the topic, with papers related to previous works on the mental health of veterans, available databases on sentiment analysis and previous works done on sentiment analysis on social media in Section 2. In Section 3, we introduce our dataset and the experiment done on the dataset, with the results we have. In Section 4, we have a deep look into the result and bring the discussion. In Section 5 and 6, we conclude and bring up future works needed for the topic.

2 Background

2.1 Previous Work on Mental Health of Veterans

In order to make a medical diagnosis for patients, psychologists often use the linguistic content and expression of patients to judge their emotional changes and mental state according to previous research

in psychology and linguistic. The clinical diagnosis efficiency has been greatly improved because of the progress of science and technology, especially in computational linguistics. In addition, the widespread of social media such as Facebook, Twitter and Instagram, has provided mental researchers with a large scale of data. Therefore, they could easily use the collected dataset and machine learning techniques for sentiment analysis. Linguistic contents which users posted on social media have been proved to be the basis for evaluating a person's mental state (Weerasinghe et al., 2019) (Guntuku et al., 2017). However, the majority of research targets are normal people.

In this paper, veterans will be regarded as research targets. Westgate in (Leonard Westgate et al., 2015) has come up with a method about the evaluation of veterans' suicide risk. This paper will concentrate on analyzing the impact of the war on veterans' mental state through the Twitters posted by themselves before and after the war instead of focusing on the prediction of the suicide risk of veterans. In addition, the comparison with the twitters released by ordinary users will be presented. Finally, comprehensive sentiment analysis of veterans will be summarized.

2.2 Available Database on Sentiment Analysis

Sentiment analysis is a very useful method for analysing the sentiment of an article, tweets or reviews. The advances in Natural Language processing and linguistic research have led to the development of different methods of sentiment analysis. The sentiment analysis is an integral part of our work and it is the fulcrum on which our hypothesis hangs on. Determining not just the sentiment of a text but even the topic (Lin and He, 2009) on which it is written on is one of the interesting works in this field. Another work makes use of the lexicons in the text and word dictionaries to extract the sentiment behind it (Maas et al., 2011). There is another work that trains a model to learn not just the words but to also capture the sentiment behind it (Taboada et al., 2011).

In this work, we do sentiment analysis on tweets and we compare the tweets of two types of users. We utilise a simple model which basically just classifies the tweets into Happy and sad/depressing sentiments. This is sufficient considering our work is to identify the difference in the mental states between the soldiers who have served in wars and common people.

2.3 Sentiment and Emotion Analysis on Social Media

Emotion and sentiment are treated as different concepts in psychology. The definition of "emotion" is a complex psychological state, which plays an important role in operating motivators. For "sentiment", it is created based on emotion to refer to a mental attitude. A survey providing more information can be referred to by (Yue et al., 2018).

In sentiment analysis, the typical task is finding the polarities of the given texts. The tests are probably positive, negative or neutral. The approach is often counting the word using and produce scores due to the lexicons.

There are commonly two approaches - analysing users' social activities and calculate linguistic features of user-generated texts. The sentiment analysis mainly focuses on short texts(tweets) generated from Twitter accounts, since most of the data is public by default and easy to obtain online. Also, the tweets are short(limited to 280 characters) and often appears with spelling mistakes and slangs. A tweet often comes with other features like spreading tweets(retweet) from other accounts. These methods mentioned above make the analysis on tweets a paradigm to explore.

SentiWordNet makes use of Opinion Mining, which is understanding the opinion of text more than the topic (Esuli and Sebastiani, 2006). Sentiwordnet is a lexical resource which scores a text on three premises object, positivity and negativity. Synsets are the building blocks of the Sentiwordnet, they form a wordnet and the wordnet is associated with the three scores. The three scores help determine how objective, positive and negative the text is. Sentiwordnet is an open-source software which is free to use and helps in extracting the sentiment of the text. The latest version of this is SentiWordNet 3.0 (Baccianella, Esuli, and Sebastiani, Baccianella et al.) which is being used in this project. The latest version of Sentiwordnet uses an updated Wordnet compared to the older version. The algorithm used

is updated to include random walk step to refine the scores. There is also considerable improvement in the accuracy of Sentiwordnet 3.0. It is used in numerous projects for the analysis of reviews and other related matters to understand whether the text is subjective or objective.

Montejo-Ráez et al. (2012) has defined a work that uses SentiWordNet on Twitter data to identify the polarity of sentiment of the users. They extract weighted vector and use it in the SentiWordNet to determine the polarity making it an unsupervised solution. We are going to be using SentiWordNet on tweets in order to understand the differences between the tweets of a soldier and that of a normal user.

3 Experiment and Results

3.1 Experiment Setup

3.1.1 Data Collection

3.1.2 TBC

3.1.3 Sentiment & Emotion Analysis

We use lexicons to score the words in our corpus. SentiWordNet is used for sentiment polarity analysis and NRC Word-Emotion Association Lexicon (EmoLex) (Mohammad and Turney, 2013) is for emotion analysis. We filtered the tweets originated from only users themselves, which means retweets will be removed. The corpus is first preprocessed to remove elements mentioned in Tab.1:

Table 1: Elements to be removed when preprocessing

Element	Examples	Element	Examples
URLs	http://foo.bar	Blank spaces	
Mentions to other users	@Bot	Single letter words	a b c
Hashtags	#botRise	Numbers	1994 233
Twitter reserved words	RT via		

When we remove numbers we try to remain the years (from 1900 to 2100). We try not to remove punctuations and stopwords because we need to do Part-of-Speech (POS) tagging after tokenizing. Both tokenizing and POS tagging is done by NLTK (Bird et al., 2009).

Once the POS tags of words are generated. We use the words in SentiWordNet to determine the scoring for positiveness, negativeness and objectiveness, while using EmoLex to perform emotion analysis with 10 emotions. Scores of one tweet are generated calculating the means of the scores of all the words after preprocessing.

The result data is shown in Tab.???. And we counted words with top ?? frequency in soldiers and civilians corpora, the lists are shown in Tab.??.

3.2 Results

4 Discussion

5 Conclusion

6 Future Works

References

- Baccianella, S., A. Esuli, and F. Sebastiani. Sentiwordnet 3.0: an enhanced lexical resource for sentiment analysis and opinion mining.
- Bird, S., E. Klein, and E. Loper (2009). Natural language processing with python.
- Esuli, A. and F. Sebastiani (2006). Sentiwordnet: A publicly available lexical resource for opinion mining. In *LREC*, Volume 6, pp. 417–422. Citeseer.
- Guntuku, S. C., D. B. Yaden, M. L. Kern, L. H. Ungar, and J. C. Eichstaedt (2017). Detecting depression and mental illness on social media: an integrative review. *Current Opinion in Behavioral Sciences* 18, 43–49. Big data in the behavioural sciences.
- Leonard Westgate, C., B. Shiner, P. Thompson, and B. V. Watts (2015). Evaluation of veterans’ suicide risk with the use of linguistic detection methods. *Psychiatric Services* 66(10), 1051–1056. PMID: 26073409.
- Lin, C. and Y. He (2009). Joint sentiment/topic model for sentiment analysis. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM’ 09*, New York, NY, USA, pp. 375–384. Association for Computing Machinery.
- Maas, A. L., R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts (2011). Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1, HLT’ 11*, USA, pp. 142–150. Association for Computational Linguistics.
- Mohammad, S. M. and P. D. Turney (2013). Crowdsourcing a word-emotion association lexicon. 29(3), 436–465.
- Montejo-Ráez, A., E. Martínez-Cámara, M. T. Martín-Valdivia, and L. A. U. López (2012). Random walk weighting over sentiwordnet for sentiment polarity detection on twitter. In *Proceedings of the 3rd Workshop in Computational Approaches to Subjectivity and Sentiment Analysis*, pp. 3–10.
- Taboada, M., J. Brooke, M. Tofiloski, K. Voll, and M. Stede (2011). Lexicon-based methods for sentiment analysis. *Computational linguistics* 37(2), 267–307.
- Weerasinghe, J., K. Morales, and R. Greenstadt (2019). “Because... I was told... so much”: Linguistic indicators of mental health status on twitter. *Proceedings on Privacy Enhancing Technologies* 2019(4), 152–171.
- Yue, L., W. Chen, X. Li, W. Zuo, and M. Yin (2018). A survey of sentiment analysis in social media. *Knowledge and Information Systems*, 1–47.