

Lexical and Non-Lexical Analysis of Soldiers' Tweets - With Sentiment and Emotion Metrics

Chao Chen, Chetan Prasad, Chen Wang,
Rachit Rastogi, Sumit Mukhija

School of Computer Science and Statistics, Trinity College Dublin
{chenc1, cprasad, wangc5, rrastogi, mukhijas}@tcd.ie

April 14, 2020

Abstract

A large number of war veterans are diagnosed with post traumatic stress disorder every day. This warrants a concern to evaluate the mental health of the war veterans. Existing works show that mental health changes caused by wars can be reflected in lexical and non-lexical features of the social media texts. In order to detect and compare those changes we collected data from 208 soldiers' tweets (n=6,61,342) with 280 civilians' tweets (n=6,00,173). We examined them with lexical and non-lexical aspects. Elements and attributes of tweets are processed and included as non-lexical features. Sentiment and emotion analysis is made on the corpora. We had a close look at the results with discussion and identified the differences between civilians and soldiers in both positive and negative directions.

Keywords: Twitter, tweet, soldier, veteran, lexical, SentiWordNet, EmoLex, lexicon, sentiment, emotion

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 users. In this research, we intend to analyze and compare the tweets of 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 of both veterans/soldiers and civilians. The accounts selected for the research all meet the criteria of having a minimum of 50 tweets and do not belong to any organisation. The tweets collected are processed using SentiWordNet to recognize the polarity of the tweets, we also count the number of words, the number of cuss words used. The data extracted from the processing of veteran/soldier tweets are compared with those of the civilians.

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 take a deeper look into the results and present 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 Soldiers

To make the 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 can easily use the collected dataset and machine learning techniques for sentiment analysis.

Linguistic contents that users post 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 civilians. In this paper, veterans and civilians 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. However, this paper will concentrate on analyzing the impact of the war on veterans' mental state by comparing the tweets posted by soldiers with the tweets generated by civilians. In addition, comprehensive sentiment analysis of veterans will be summarized.

2.2 Sentiment and Emotion Analysis on Social Media

Sentiment analysis has been applied in various fields such as the research about consumer behaviour was in the field of product marketing and the analysis of voter bias. The advances in Natural Language Processing and linguistic research have led to the development of different methods of sentiment analysis.

Nowadays, people tend to use social media such as Twitter to post tweets and express their opinions and emotions. Most of the tweets generated from Twitter accounts are public by default and easy to obtain. Also, the tweets are short (limited to 140 characters) and often appear with spelling mistakes and slangs. A tweet often comes with other features like spreading tweets (retweet) from other accounts. The factors mentioned above make Twitter a viable source to carry out sentiment analysis.

Generally speaking, the common approaches for sentiment analysis consist of two parts, including the machine learning techniques based and the lexicon-based. Analyzing users' social activities and calculating linguistic features of user-generated texts are the core of the machine learning algorithm. Compared with machine learning-based method, the lexicon based method is more direct and straightforward. In sentiment analysis, the typical task is finding the polarities of the given texts. The text content is either positive, negative or neutral. Lexicon based method could recognize and analyze the words about sentiment and other emoticons and hashtags which are associated with the sentiment.

Therefore, the sentiment lexicons are adopted for matching the words from tweets, thus analyzing and determining the polarities of the corpus.

Azizan et al. (2019) performed sentiment analysis on Twitter data about movie review tweets using R and lexicon-based method. They found that the lexicon-based method is more effective than the machine learning based method under the same calculation cost. Ray and Chakrabarti (2017) used a dictionary-based method and analyzed the results at aspect level and document level to predict the public's sentiment using tweets about product review.

SentiWordNet and SenticNet, as open lexicons resources, have been developed in recent years. Sentiwordnet is a lexical resource which scores a text on three premises object, positivity, negativity and objectivity. It is an open-source software which is free to use and helps in extracting the sentiment of the text. Due to the high accuracy, the SentiWordNet 3.0 (Baccianella et al., 2010) will be used as lexicons in this paper.

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

3 Experiments and Results

We believe the tweets originate from soldiers and veterans can have differences with civilians on both lexical and non-lexical aspects. The experiments are divided into two parts. We first summarized the features on the large aspects to compare the difference between corpora from soldiers/veterans and civilians. Then we calculated sentiment and emotion metrics based on lexicons.

3.1 Data Collection

We use TWINT (Poldi, 2017) as our data collection tool. All the data can be accessed publicly so there are no ethical considerations.

Numerous strategies considered in an attempt to procure data that belonged to the war veterans. Transcripts of podcasts and YouTube videos involving accounts of wars from the veterans, books that were written by the ex-servicemen, the public dataset that had diaries, and letters of first world war soldiers were a few sources. However, as an inference, all of these sources were highly specific to the negative aspects and impacts of war and eventually would add bias to the data.

On account of being a platform that is widely used by a large number of service-men and the civilians, Twitter was selected as the platform to extract data. Several verified Twitter pages linked to the US Army were manually analysed and a few profiles of the veterans were used as an initial set. But, we finally used the data from a verified page with the name IAVA. This is said to be the most significant association speaking to the new age of vets specifically from The United States. The profiles with a minimum of 50 tweets were only considered eligible. The final set had Twitter profiles with as few as 57 tweets and as many as 65,000 tweets. The succeeding veterans were picked by scouring through the followers of the veterans in the original set based on some keywords like the army, us-army, military in their bio.

A collection from 208 veterans profile was performed which was used for our final set of experiments. It helped us to have a dataset set of 6,61,342 tweets that were used in our analysis.

Similarly, for collection the data of civilians, we opted to chooses pages such as Netflix, USA official page, The US open and more. Here also the same guidelines were followed as in the case of finding the profiles of vets, but the only difference was that, here we were selection profiles based on keywords not similar to “army”, “us-army”, “military” in their bio. A final list of 280 civilians usernames was selected with variations in age, sex, and the number of tweets. This ultimately added 6,00,173 tweets by the general civilians and was used in our processing.

Once we had all the data that is from both the veterans and the civilians, we merged all the data to two CSV file so that is easier to process the data that we have extracted. This merged data for both were further used in our pre-processing.

3.2 Data Process and Analysis

3.2.1 Lexical & Non-Lexical Feature Summarization

Once the data is gathered and saved in the CSV format, we then start pre-processing the gathered data after which we do our analysis where we come up with a set of results to prove our hypothesis. Sentiment analysis can be broadly categorized into two kinds, based on the type of output the analysis generates. Under our processing, we are trying to label text to be under “curses” - or bad words. Take into account the tweets of all the selected veterans from our data and then run an analysis to gather information from there. Not only this, but information such as the total counts of retweets, replies, likes, emojis, URLs, mentions and total words are calculated. Special elements in tweets can be referred to Table 1. This is one of the analysis that we are done from the collected data.

This identical analysis is then run on tweets by ordinary people or civilians data as well. Finally, we will categorise the difference in the count of “curse” words to check the two sets of results which will ultimately help us identify an individual being in the state of depression.

The initial step we took for pre-processing our data were:

- Only selecting tweets that were in English.
- Applying the rule to only select textual data.

It was then followed by removal of: (1) Retweets numbers; (2) Tokenizing.

Furthermore, the emojis were given a textual form so that we can get information from that part as well. This is because in this everchanging world the use of emojis has increased. These converted emojis were counted. Subsequently, profanity was predicted in the tweets.

The step to remove punctuation followed by Tokenizing which is the way toward separating a goliath string into a rundown of words. NLTK (Bird et al., 2009), a python library is used for this process. Stopwords, where pull-out, as they do not change any meaning of a sentence hence, can be ignored. Tuples were generated with each word and part of speech. Finally, the counts were extracted from the tweets of both the veterans and the civilians. And was then was compared with each other.

An abstract of the process is illustrated in Figure 1.

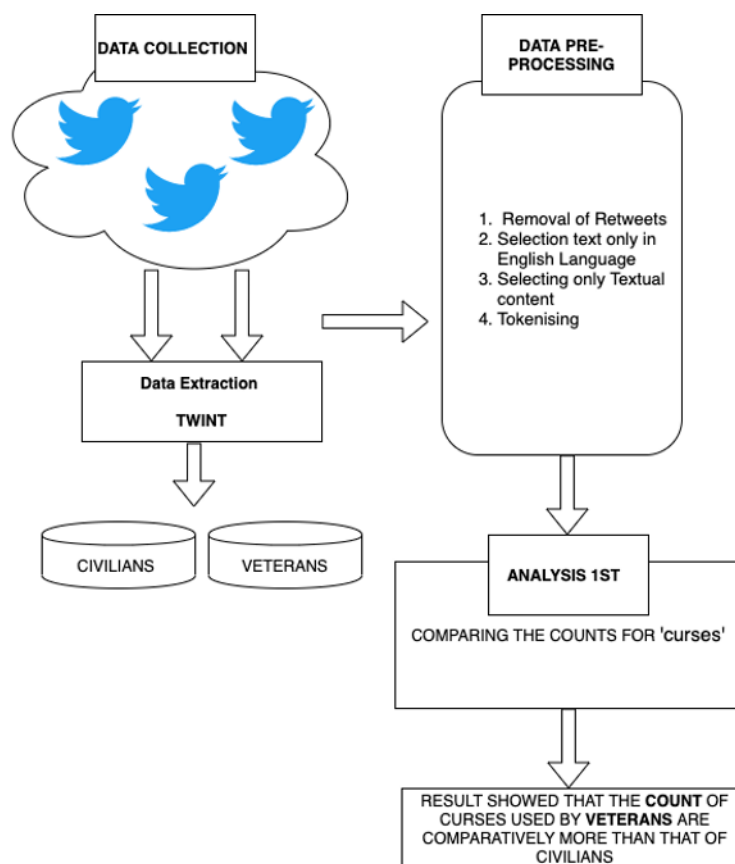


Figure 1: Process of lexical & non-lexical feature summarization

Table 1: Elements to handle when preprocessing tweets

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	Stopwords	it I as

Here is the data that was collected from the tweets:

Table 2: Summarized features from soldier and civilian corpora

	Soldiers n=208	Civilians n=280
tweets	661342	600173
retweets	869272	24102124
replies	333204	5919866
likes	2735134	123551386
emojis	159583	69284
urls	250026	248854
mentions	226533	325691
words	51528493	44355619
curses	28249	8983

3.2.2 Sentiment & Emotion Analysis

Tweets are filtered and only tweets with texts originate from users themselves remain, which means the likes and retweets are filtered.

The corpora are then pre-processed to remove elements mentioned in Table 1.

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).

We use lexicons to score the words in our corpora. SentiWordNet is used for sentiment polarity analysis and NRC Word-Emotion Association Lexicon (EmoLex) (Mohammad and Turney, 2013) based on the model of Plutchik’s wheel of emotions (Plutchik, 2003) (see Figure 2, with additional Positiveness and Negativeness) is for emotion analysis.

Once the POS tags of words are generated. We search the synonyms of words in SentiWordNet to determine the scoring for positiveness, negativeness and objectiveness by calculating means among synonyms. Meanwhile, EmoLex is used to perform emotion analysis on 10 emotions. Scores of one tweet are generated calculating the means of the scores of all the words after preprocessing.

The result data obtained by applying SentiWordNet is shown in Table 3. The result produced by using EmoLex is shown in Table 4.

We also counted adjectives with top 100 frequencies in soldiers and civilians corpora, for we think that adjectives have more subjective meanings than verbs, nouns, etc. We discovered some words with more “political” meanings appear to be different in the lists of two corpora. The list of adjectives is shown in Table 5.

Table 3: Results of sentiment analysis using SentiWordNet

		Valid Cnt.	Valid Len.	Positive.	Negative	Objective.
Soldiers n=208	Mean	3179.54*	16.450	257.58×10^{-4}	196.10×10^{-4}	3371.7×10^{-4}
	Std.	5041.70	6.6427	78.586×10^{-4}	64.386×10^{-4}	750.49×10^{-4}
Civilians n=280	Mean	2143.66*	14.293	262.65×10^{-4}	177.39×10^{-4}	3530.5×10^{-4}
	Std.	5286.12	5.2067	87.432×10^{-4}	64.786×10^{-4}	720.09×10^{-4}

4 Discussion

The entire idea for our research was to check the mental condition of veterans to that of to a normal civilian and eventually compare the sentiment traits. Our research involved 488 people that included the civilians as well as the war veterans, shows that there is definitely an effect of war that can be reflected by the word selection and other traits they use while tweeting. One of the analysis shows that the number



Figure 2: Plutchik's wheel of emotions

of "curses" used in by the war veterans is much greater than that used by ordinary civilians. The word count for curse word by war veterans accounts to be 28,249 which is much greater than that curse word count used by civilians which is at 8983. On an average, Veterans employed profanity 3 times more than a civilian did.

From Table 3 we cannot get much inference on the sentiment part, instead we find that soldiers are more likely to post lengthy tweets (see the numbers with *).

We can see from Table 4 that corpus of soldiers' tweets has more "negative" emotions like Disgust, Fear, Anger and Sadness. The corpus of soldiers' tweets is judged as negative on the whole. While civilians' corpus tends to be more positive, with better metrics on Surprise, Anticipation, and Joy. One interesting emotion is Trust, from Plutchik's wheel of emotions 2 we can infer that Trust is a kind of emotion related to submission, acceptance and admiration, which is related to soldiers' loyalty obeying the commands. While Surprise is related to disapproval and distraction, which can somewhat indicate the quality of disorder among internet users.

With the list of adjectives (Table 5) we can see that soldiers are more involved in political topics. It might be that soldiers/veterans are more involved in political events fighting for their rights, also they usually have closer relations to governments and military.

5 Conclusion

The results that we received after extracting and processing the data for the soldiers and the civilians conclude that there exists an evident logical separation between the way a war veteran and a commoner

Table 4: Results of emotion analysis using EmoLex

Soldiers: n=208					
	Trust+	Anger+	Surprise	Joy	Positive.
Mean $\times 10^{-4}$	422.84	167.17	149.99	312.51	637.50
Std. $\times 10^{-4}$	154.94	83.143	62.113	151.64	213.55
	Disgust+	Fear+	Anticipat.	Sadness+	Negative.+
Mean $\times 10^{-4}$	122.09	193.43	295.57	149.34	339.61
Std. $\times 10^{-4}$	74.795	89.380	112.40	66.053	148.79
Civilians: n=280					
	Trust	Anger	Surprise+	Joy+	Positive.+
Mean $\times 10^{-4}$	399.72	132.03	163.72	349.44	650.63
Std. $\times 10^{-4}$	170.45	80.189	108.13	224.58	269.03
	Disgust	Fear	Anticipat.+	Sadness	Negative.
Mean $\times 10^{-4}$	98.934	163.78	330.09	131.82	283.00
Std. $\times 10^{-4}$	82.884	108.23	152.86	87.180	160.08

Table 5: List of “political” adjectives with rankings and frequencies

Word	Soldiers	Civilians	Word	Soldiers	Civilians
military	4.46 (17th)	-	dead	1.20 (77th)	-
american	3.85 (24th)	1.20 (79th)	human	1.17 (80th)	1.01 (91st)
political	2.12 (40th)	-	local	1.16 (81st)	1.26 (72nd)
medical	1.75 (47th)	-	democratic	1.15 (83rd)	-
public	1.56 (51st)	1.32 (68th)	illegal	1.14* (85th)	-
social	1.44 (60th)	1.78 (51st)	foreign	1.14* (85th)	-
sick	1.37 (64th)	-	poor	1.10 (90th)	-
personal	1.31 (71st)	1.21 (78th)	republican	1.07 (93rd)	-

tweets to the veterans when compared to normal people. They do show characteristics of using a higher amount of curse words along with longer messages just to name a few. With the high amount of data that we had along with the preprocessing and methodologies used we are certain that we were able to prove our hypothesis and come to a conclusion at the end of our research. Nevertheless, we feel that there is still a lot of work that exists in these areas and we hope that our findings enlighten other researchers to deep dive furthermore.

6 Future Works

Though we did try to answer our research question by applying two different analysis, there is always a scope of improvement. One area where we like to work further is to extract data depending upon the timestamp. We wanted to get extract data depending on times when the veterans returned from war and then compare it with their own tweets before the war. For the sentiment and emotion analysis, negators (e.g. not, isn't, won't) need to be considered for that they contribute strong and opposite effects towards polarity and word meaning. A cross comparison should be made as the work of Mohammad and Bravo-Marquez (2017) but on our corpora, to find better methods with lexicons to calculate polarities. These works can be done in future, and definitely worth exploring.

References

- Azizan, A., N. N. S. A. Jamal, M. N. Abdullah, M. Mohamad, and N. Khairudin (2019). Lexicon-based sentiment analysis for movie review tweets. In *2019 1st International Conference on Artificial Intelligence and Data Sciences (AiDAS)*, pp. 132–136. IEEE.
- Baccianella, S., A. Esuli, and F. Sebastiani (2010). Sentiwordnet 3.0: an enhanced lexical resource for sentiment analysis and opinion mining. In *Lrec*, Volume 10, pp. 2200–2204.
- Bird, S., E. Klein, and E. Loper (2009). Natural language processing with python.
- 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.
- Mohammad, S. M. and F. Bravo-Marquez (2017). Emotion intensities in tweets.
- 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.
- Plutchik, R. (2003). *Emotions and life: Perspectives from psychology, biology, and evolution*. American Psychological Association.
- Poldi, F. (2017). TWINT - twitter intelligence tool. <https://github.com/twintproject/twint>. Accessed: 2020-03-04.
- Ray, P. and A. Chakrabarti (2017). Twitter sentiment analysis for product review using lexicon method. In *2017 International Conference on Data Management, Analytics and Innovation (ICDMAI)*, pp. 211–216.
- 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.