# Text Analytics for Suicidal Ideation using Natural Language Processing Techniques

Yi Xiang Chin
*School of Computer Science*
*Universiti Sains Malaysia*
Penang, Malaysia
chinyixiang@student.usm.mY

Hui Ting Ling
*School of Computer Science*
*Universiti Sains Malaysia*
Penang, Malaysia
linghuiting@student.usm.my

*Abstract*— **Suicide rate is increasing at an alarming pace in worldwide. In online social networks, users can express their thoughts via texts. In this paper, we leverage Natural Language Processing (NLP) techniques and Machine Learning algorithms in our study. We explore the potential topic modeling for suicidal posts expressed on social networks, as well as correctly identify the emotions in order to make inferences from the posts. We use publicly available word-emotion NRC lexicons to assist in our study. First, we performed text cleaning, second, we build suicide psychological lexicons for emotions detection and third, we employ topic modeling to identify the suicidal ideation. To this end, our result shows that majority of suicide attempter are suffered in sadness and fear when writing the suicidal note. And, three topics are identified from the study, where topic 1 has significant differences in terms of emotion compositions.**

*Keywords*— *Suicide, Natural Language Processing (NLP), topic modeling, emotions detection, social networks, emotion lexicons*

## I. INTRODUCTION

Suicide is a significant mental health problem worldwide, its accounted 1.4% of deaths globally and this makes the 18th leading cause of death in 2016. [1] According to the World Health Organization [2] reported, estimated there are 800,000 died from suicide every year, and 79% are occurs in low-middle-income countries, however, the figure is usually assumed underestimated due to misclassification of the deaths. [3][4]

Suicide ideation is viewed as ones' has a suicide attempt tendency with a plans [5][6]. People who become suicidal and seek suicide as the only solution when they feel overwhelmed by life challenges. The risk factor believed to increase the impulsive suicide behavior such as history of substances abuse, access to firearms, difficult life events, isolation from others, history of mental illness, history of physical or sexual abuse, chronic illness, and past suicide attempts. [7]

In recent year, with the rise of the use of social media platform, it is become a "venting window" to provide a space to users to utter their suicidal thoughts while remain anonymous where suicide is society stigmatized openly discuss topics. The suicidal notes can be noticed by accompanied with a clear marker like killing themselves, their life has no purpose, feeling like a burden, feeling stuck, not wanting to exist, etc. [7] Suicidal notes always convey feeling, emotions, and behavior. [8] Some words have semantic core emotion of the feeling, for example, dejected and wistful associated with some amount of sadness. On the other hand, some words may not denote the effect but it associate some degree of the emotions, for example, failure and death are usually accompanied by sadness. [8]

Emotion detection is a subfield of sentiment analysis (SA), where SA core intent is to analyze the polarity either positive, negative, or neutral. [9] Emotion detection is sought to extract finer grain of emotions. Importance of identifying actual emotions rather than sentiment polarities, [10] For example, "I am crying tonight (sadness)" and "I am furious, I shall have my revenge (anger)" are classified under "negative polarity". However, the two messages convey different feelings.

Robert Plutchik was a psychologist studying emotions, [11] he created Plutchick's wheel of emotions classified primary emotions into eight elements that can be clearly distinguished, such as joy, anger, sadness, fear, disgust, surprise, anticipation and trust. [11]

National Research Council Canada (NRC) [12] has a collection of English words with sentiment and emotion. NRC Words-Emotion Association Lexicon associates with eight emotions (anger, fear, anticipation, trust, surprise, sadness, joy and disgust) and two sentiments (negative and positive). The annotations were manually done by crowdsourcing, [13] and entirely created by the expert of the National Research Council of Canada.

Natural Language Processing (NLP) is defined as the naturally occurring text driven by human language to address a range of objectives (information retrieval or artificial intelligence). The analysis can be from a set of rules or combined with machine learning algorithms. The NLP output should represent as close as possible to human output. [5]

Suicide could be prevent by early detection, treatment on acute emotional distress. [2] Using NLP approach is a new venture in suicidal prevention research. [5] In this paper, we present two novel NLP approaches in the project, we use the NRC list of lexicon to help us identify the emotions when people have suicidal thoughts through the notes, and detect the specific words used by the respective problem they had faced.

The structure of our paper as follows: Section II describes the problem statement. Section III discusses the related work. Section IV explains the methodology. Section V analyzes the analysis, and Section VI and VII for discussion and conclusion.

## II. PROBLEM STATEMENT

Talking about suicide can be a scary topic as suicide is a taboo in many societies to discuss openly. The prevention of suicide has not been adequately addressed due to lack of awareness of suicide is a major public health problem. [2]

Malaysia has seen a rise in suicide during Covid-19 pandemic. Befriender Malaysia received 4142 calls between March to May, where over one third of the cases were suicidal in nature. The government hotline received about 11791 calls between March to August. The figure could be underreported

as the act of attempted suicide is criminalized in Malaysia. [14]

The objective in this project is to use NLP techniques to pin the problems and perhaps measure up the preventions.

- What is the emotion when people have suicidal thoughts?

- What is the unbearable problem they had when people seek to end their life?

- What are the keywords to the respective topic in the suicidal notes?

- What is the greater emotion to the respective problem they had?

## III. RELATED WORK

There are several literatures on text mining methods used to identify suicidal behavior for research. Huang et al.[4] paper study on identify suicidal ideation by topic modeling, the proposed method was collect suicidal data and tagging into 12 warning signs of suicide. Then use word embedding technique to construct psychological lexicon and then employ a topic modeling approach to look for suicidal ideation. They also use the lexicon sentiment analysis based on Chinese sentiment dictionary (Hownet). Fernandes et al. [5] using NLP techniques to identify suicide ideation and suicide attempts in their clinical database to predict suicidal behavior, they use rule-based approach to manually classify the data as training set, a hybrid machine learning rule-based to identify the suicidal risk. M. Tadasse et al. [6] study is to detect suicide ideation in reddit social media forum to detect whether it is suicidal or non-suicidal post. They proposed using deep learning LSTM-CNN algorithm and machine learning based classification approaches, the dataset labeled suicide-indicative or suicide-non-indicative, next by feature extraction using word embedding then followed by deep learning classifier. Metzger et al. [15] employ seven standard machine learning techniques (Support Vector Machines, predictive association rules, decision trees, logistic regression, Naïve Bayes, random forest and neural networks) to classify suicidal ideation from suicide attempt from a list of demographic and clinical variables.
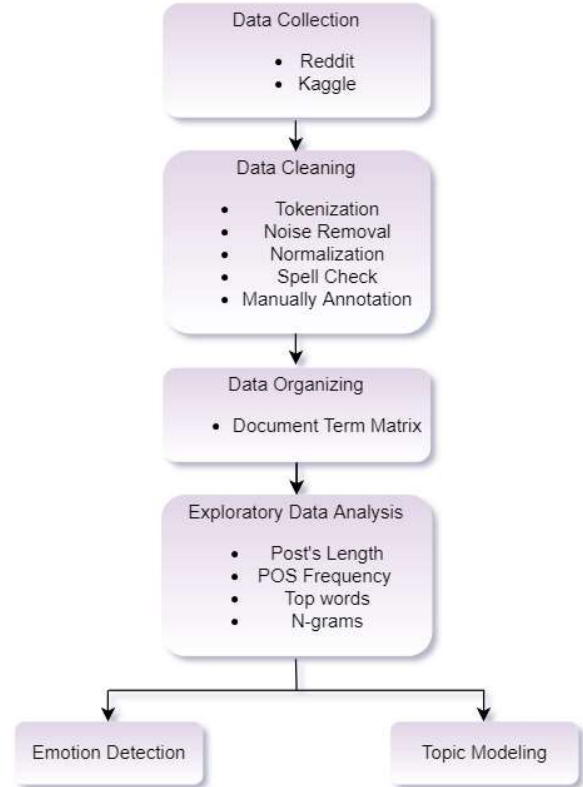
Yang et al [16] was study on emotion recognition in suicide notes, a Conditional Random Field (CRF) model for identifying emotion cues at the token level, and three different machine learning models, Naïve Bayes, Maximum Entropy and Support Vector Machine (SVM) for emotion sentence level. Razek and Frasson [17] proposed dominant meaning technique to recognize emotions, the system has two stages, training stage and classification stage, they use ISEAR dataset to form dominant meaning by hierarchy trees. Then use it for incoming examples for emotion classification. Kusen et al. [10] study was use three open available word-emotion lexicon resources such as NRC, DepecheMood and EmoSentiNet to identify emotion in social media posts, they use ISEAR dataset and survey response as ground golden set, then extract emotion from texts and compare the performances of different emotion-lexicon. NRC and DepecheMood gave comparable results. Acheampong et al. [18] use BERT- based approach which the output is put into Bi-LSTM classifier for predicting emotion. They use the ISEAR dataset as a golden sample. Ghosh et al. [19] build an emotion corpus named CEASE, the study are manually annotate sentences by 15 emotions class, then evaluate the

result by the ensemble architecture with three deep learning models, Convolutional Neural Network (CNN), Gated Recurrent Unit (GRU) and Long Short Term Memory (LSTM).

## IV. METHODOLOGY

The purpose of this study is to detect the emotion and topic modeling for the suicide ideation. We use various NLP and text mining techniques in our study. Figure 1 shows a general overview of our framework. It depicts methodology flow of the project.

FIGURE 1. METHODOLOGY FRAMEWORK



### A. Data Collection

To collect suicidal posts, we retrieve data from two sources. We used the Beautiful Soup package to web-scraping the posts from a subreddit social media platform: Suicidal_Thoughts where users can express their thoughts via texts. And also we downloaded ready suicide notes dataset from Kaggle. After combining the data, a total 1490 posts in original size. We perform the first stage of cleaning, removing the empty posts and removing the duplicate posts.

### B. Data Cleaning

For second stage text data cleaning, we use NLTK library to handle the text preprocessing tasks.
- Tokenization

For many NLP tasks, we need to access each word in a string. We first use nltk's .lower() to convert the texts into lower cases. To access each word, we tokenize the text into

individual tokens, we use nltk's word.tokenize() function to tokenize the string of text returned into a list of words.

- Cleaning (Noise Removal)

Then we further treatment on the unwanted information such as punctuation, non-alphabetical tokens or special characters and empty token posts. We use the .sub() method in Python's regular expression (re) library to clean the noise. And also, we remove one-token posts which do not make sense for text analysis.

- Normalization

Next, we further process through normalization by removing stop-words and lemmatization. We use nltk's stopwords.word("English") function to remove the default stop-words. In lemmatization, it reduces the inflected words properly ensuring the root word belongs to the language. For example, "run", "running" and "ran" belong to the lemma word of "run". We use the WordNet Lemmatizer package to lookup lemmas words.

However, after lemmatization has been done, we still find there is not clean enough such as spelling error or contraction of words like "idk", "tbh". Therefore, we further extract the words and indicate them as non-dictionary words. To further process the non-dictionary words, we first use SnowballStemmer("English") package to stem the words then apply Speller() to perform spell check and autocorrect. With that, we managed to autocorrect 80% of the non-dictionary words. We then manually annotate and correct for the rest of the 20%. We inspect the high-density non-dictionary word, then 1. Manually correction and update, 2. Update the certain word directly to the dictionary, 3. Assign the word as stop-word.

After the text is cleaned, we proceed to organize the texts into Document Term Matrix (DTM). We use scikit-learn's CountVectorizer for this task, DTM every row represents a different document (post) and every column represents a different word.

After that we save the after clean data and document term matrix as a pickle file.

*C. Exploratory Data Analysis*

In this portion, we are interested in deep-diving to see how many words appear in a post, what is the top word or specific phrase, POS tag usage, which terms are likely to co-occur.

- Post's Length

A note is to express the thoughts via text in the whole of the suicide ideation, hence the length of notes is critical to find the context of the sentence and classified into correct topic and emotion detection.

We load the cleaned data pickle file and count the token in each post. To see the distribution of length we plot the values in histogram.

- Part-of-Speech (POS) tagging Frequency

A note is formed by several POS to become a notes. Here, we are interested to see which POS tagging is likely to be written in the suicidal notes. For POS tagging, there are eight main categories that are: Nouns, Pronouns, Verbs, Adjectives, Adverbs, Conjunctions, Interjections, and Propositions. We use nltk's pos_tag for the tagging. We first group the Verbs into one category (VB, VBD, VBG, VBN, VBP, VBZ)

- Each-Post-Top-words

To understand each individual state-of-thoughts and what it is same in common. We extract 30 top words from each post to see what token occurs frequently in each post and how likely the same token will have the same in common in different posts.

- N-grams

N-grams in this study is not about to check the probability of words co-occurrence and information retrieval. The purpose of N-grams are likely to be the same as top-words, Fwe segmented it into uni-gram, bi-gram and tri-gram to see what are the words they are likely to express and this is an important message to detect suicidal ideation.

- Wordcloud

Wordcloud is an image composed of words used in a particular text, in which the size of each word indicate its frequency and importance. [20] We generate wordcloud to help on identify the important keyword for POS tagging, uni-gram and bi-gram.

*D. Emotion Detection*

Emotions play vital roles in the existence or complete make-up of an individual. They provide information on the current state and well-being. Emotion detection is important to understand the state of emotions of the people when they write the suicidal notes. We utilize the NRC lexicon in our study, however we do not include the positive emotions because suicidal notes itself is a negative-state in nature, for the positive words like "happy", "good", "and grateful" are considered passive-aggressive behavior. Therefore, we only pick negative emotions such as "Anger", "Sadness", "Fear", "Disgust" and "Surprise" for this study. We match annotated lexicons with each suicidal note, then count and group the same emotion group then summarize the dominant emotions from each post. The dominant emotions are summarized using two ways: by count or by score. For summarization by count, each token of the post is determined it's emotion using the highest score, then the dominant emotion of the post is determined by the highest count of the token's emotion. For summarization by score, scores of 5 emotions of each token in a post is added up, then dominant emotion is determined by the highest score.

*E. Topic Modeling*

Topic modeling is an unsupervised machine learning technique that is capable of scanning a set of documents, detecting word and phrase patterns within them, and automatically clustering word groups and similar expressions that best characterize a set of documents. [21] We based on POS tagging approach and tried a few attempts to find the best topic model. We use the Gensim python library for topic modeling. We first transpose the document-term-matrix into term-document matrix, then put it into Gensim format to create Gensim sparse matrix and Gensim corpus. Gensim also requires a dictionary of all terms and their respective location in the term-document matrix. Several parameters can be specified when training a topic model, including number of topics, alpha, beta and number of passes. As topic modelling is unsupervised, it is difficult to determine which model to choose. For each corpus (Full Text, Nouns, Nouns + Adj, Nouns + Verbs), a base model is trained with 5 topics, 50

passes, and default alpha and beta. The base model is then evaluated using the topic coherence measurement [22]. Then, by using 2 passes, the best combination of number of topics, alpha and beta is tested. Then, the best model for each corpus is selected, then trained using 100 passes and evaluated again. The best model is chosen from the 4 models. Each post is then labeled with the topic using the topic model.

## V. ANALYSIS

In this section, we present the analysis results and outputs from what we discussed in methodology.

### A. Data Exploratory Analysis

● Post's Length

We have a total 1283 posts after data preprocess has been done. Figure 2 shows the length of each post can range from 2 tokens to over 800 tokens a post, where 94% of the posts are less than 200 tokens.

FIGURE 2. POST'S LENGTH

● Part-of-Speech (POS) tagging Frequency

Figure 3 is frequency of POS tag used in suicidal notes, the top five of POS tags are Nouns, followed by Verbs, Adjectives, Adverbs, and Interjections. And Figure 4 is a wordcloud of the top five POS.

FIGURE 3. FREQUENCY OF POS TAG USE IN SUICIDAL NOTES

FIGURE 4. POS TAGGING

As Figure 4 showed, people are likely to express themselves in Nouns. That being said, noun terms are a keyword or a hint of the problem they faced such as "friend", "family", "life", "thing", "school", etc. And Verbs are the feeling and thoughts they expressed such as swear word "fuck", "tired", "kill", "trying", etc. Adjectives are semantic role to the nouns, to the surprise, "happy" word are occur quite often, but believed that it is passive aggressive behavior, majority are still surrounded by negative words such as "suicide", " bad", "hard", "sad", and etc.

● Wordcloud
Each-Post-Top-words

Figure 5 depicts the 30 top words in each post, "Feel", "know", "want", "life" are the common word occur in each post. Whereas "flesh", "flew", "flight", "fledged" are unique words in some post where it does not see in uni-gram wordcloud.

FIGURE 5. UNI-GRAM WORDCLOUD

Uni-gram

Figure 6 shows the uni-gram wordcloud. The wordcloud is worth a thousand words. Uni-gram and 30-top-words are quite similar. The "feel" is the highest freqency in this study, this means that people express their feelings a lot when writing their thoughts. Where "life", "want", "thing",

"people" are the keywords or problems they are facing in their life challenge.

Bi-gram

Figure 7 shows the bi-gram wordcloud. Suicidal behaviour could be impulsive due to some event change. They used to describe the event in "last day", "last night", "obvious reason" why they have suicidal thoughts.

## B. Emotion Detection

FIGURE 8. POST EMOTIONS BY COUNT



Figure 8 shows the percentage of post emotions by count. From the pie chart, we can observed that suicide attempter are emotionally in "Sadness" and "Fear" when expressed in a suicidal posts which takes up 74% of the posts. Anger emotion is 17%, and self-disgust is 3%, surprisingly, the emotion "Surprise" is 6% more than the emotion "Disgust". Figure 9 shows the percentage of post emotions by score. Can be observed that the rank of emotions in the posts are still the same, but the percentage is different. Sadness and Fear become more dominant compared to by count.

This observation shows that emotions are complex and subtle. Every word can carry a different proportion of emotions which contributes to the whole semantic of the text.

FIGURE 9: POST EMOTIONS BY SCORE



## C. Topic Modeling

FIGURE 10: TOPIC MODELS AND COHERENCE SCORES

| Corpus | Tuning | Number of Topics | Alpha | Beta | Number of Passes | Coherence Score |
|---|---|---|---|---|---|---|
| Full | No | 5 | - | - | 50 | 0.7560 |
| Full | Yes | 3 | symmetric | 0.01 | 100 | 0.7910 |
| Nouns | No | 5 | - | - | 50 | 0.7375 |
| Nouns | Yes | 3 | 0.31 | 0.01 | 100 | 0.7567 |
| Nouns + Adjective | No | 5 | - | - | 50 | 0.7685 |
| Nouns + Adjective | Yes | 3 | 0.01 | 0.61 | 100 | 0.7791 |
| Nouns + Verbs | No | 5 | - | - | 50 | 0.7479 |
| Nouns + Verbs | Yes | 3 | 0.61 | 0.01 | 100 | 0.7622 |

Figure 10 shows a table of comparison among different topic models. Different corpus is generated by filtering tokens with specific POS tags. Tuning is done on number of topics, alpha, and beta values to determine the best combination of that corpus. The highest coherence score is achieved by full text corpus, 3 topics, symmetric alpha, and beta of 0.01 with coherence score of 0.7910.
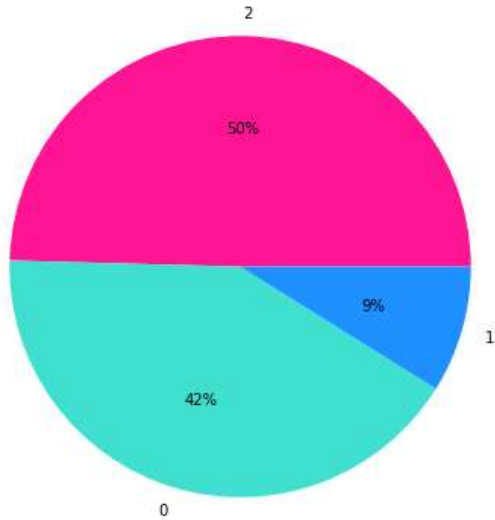
From the table, can be seen that grid search tuning does improve the coherence score of the topic model, but mostly less than 0.05 improvement.

The topics generated from the modelling are as follows:
[(0,
 '0.026*"magic" + 0.019*"rob" + 0.013*"base" + 0.012*"seekers" + 0.011*"fading" + 0.011*"blacklist" + 0.009*"aide" + 0.008*"theyre" + 0.008*"cite" + 0.008*"rot"'),

(1,
 '0.015*"magic" + 0.013*"gifted" + 0.012*"rob" + 0.011*"theyre" + 0.011*"monitor" + 0.010*"base" + 0.009*"fading" + 0.008*"incarnate" + 0.007*"fiend" + 0.007*"insist"'),
 (2,
 '0.024*"blacklist" + 0.021*"base" + 0.019*"magic" + 0.018*"rob" + 0.017*"seekers" + 0.016*"aide" + 0.012*"barking" + 0.011*"fiend" + 0.011*"everlasting" + 0.009*"chase"')]

There is no obvious difference in the topics to relate to the post intention or reason of suicide. A lot of the factors are overlapping in the topic. This might be due to the highly similar nature of suicidal posts.

FIGURE 11: SUICIDAL POSTS BY TOPICS



From Figure 11, topic 0 and topic 2 takes up 92% of the posts. This can also be related to the highly similar nature of topic 0 and topic 1 in the word factors. It is possible that there's only 2 topics available in the posts, but 2 topics are not considered in the grid search tuning.

### D.   Relating Topics to Emotions

At the last part of the project, we tried to relate results of topic modeling to the emotion detection.

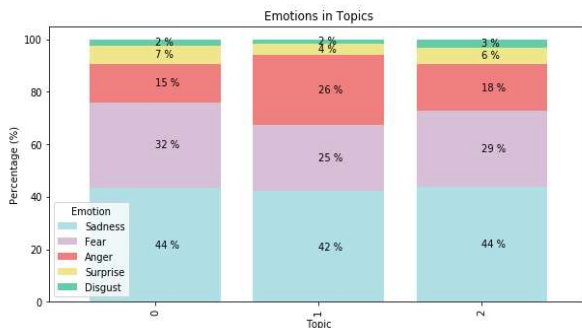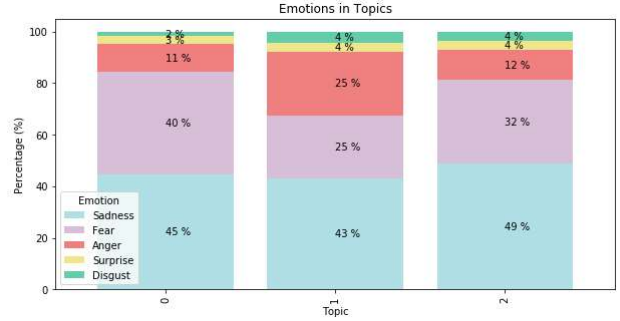FIGURE 12: EMOTIONS BY COUNT IN TOPICS



FIGURE 13: EMOTIONS BY SCORE IN TOPICS



From Figure 12 and Figure 13, we can observe a slight difference in proportion of emotions in each topics using by count or by score methods. Labelling emotion by score shows higher sensitivity of topic difference compared to labelling by count. However, we can observe that topic 1 is only occupied 9% in the suicidal posts, but the emotion "Anger" is higher than topic 0 and topic 2 about 10%, therefore, we cannot ignore this topic as a topic.

## VI.   DISCUSSION

### A.   Challenges

One of the biggest challenges we faced in this project is in the text cleaning process. In this project, we used Speller autocorrect and identifying non-dict word and manual correction. However, from the final results of the topic modelling, we can still observe some words that escaped the check, for example "theyre". The suggestions to overcome this is to use a more specific dictionary to identify words that are out of context.

Besides, as an unsupervised learning method, topic modelling is difficult to be evaluated and tuned. As what happened in this project, the results of the tuning is not understandable from a subjective view.

### B.   Limitations

The results of this study are subjected to the limited data source that we managed to gather. By combining web scraped data with Kaggle data sources, we only managed to gather a total of 1490 posts. This sample size is very small and might not represent the population.

The sampling of the suicidal posts is also subjected to the medium used. As all the data collected are digitally available online, this limits study target to populations that use social media as the medium to express their final thoughts.

### C.   Discussion

Through this study, we have found that text analytics techniques are still very subjective and difficult to fully automate. One example is in the labelling of emotion. Different ways of annotating emotions will result in different outcomes.

Same thing also happens to the labelling of POS tags. As POS is a grammatical annotation, it is mostly done before the removal of stopwords and before tokenization to preserve the sentence structure. However, in this study, we have done the POS after text cleaning. This is because we are not very concerned about the function of the words in the sentence.

POS tagging in this case works more as a classification method and becomes a characteristic of the token itself.

*D.     Future work*

Similar to Emotion analysis, word embedding can be used in the future for similar topics. With a well defined dictionary, word embedding can help us analyze the texts with a numerical approach.

Further analysis can also be done if the age of the post owner can be obtained. This can help identify the problems faced by different age groups.

## VII. CONCLUSION

This paper present the comprehensive methods and analysis on the suicidal ideation. We observed that when people having suicidal thoughts, they are majority in "Sadness" and "Fear" of emotion-in-state compare to other negative emotions "Anger", "Disgust" or "Surprise". From the exploratory analysis, the high frequency of words express that are "feel", "life", "people", "friend", "family" could be the key abruptly change and unbearable problems they could not stand anymore. Three topics can be distinguish with the unique keywords from topic modelling. The dominant emotion for topics respectively are still "Sadness", and "Fear" occupied the highest, however emotion "Anger" in topic 1 is greater than topic 0 and topic 2 about 10%. In the end of this paper, we also discussed the challenges, limitations, discussion and future work.

## REFERENCES

[1]     "Mental Health and Substance Use," *World Health Organization*.    [Online].    Available: https://www.who.int/teams/mental-health-and-substance-use/suicide-data.

[2]     "Suicide," *World Health Organization*. [Online]. Available:    https://www.who.int/news-room/fact-sheets/detail/suicide.

[3]     C. Katz, J. Bolton, and J. Sareen, "The prevalence rates of suicide are likely underestimated worldwide: why it matters," *Soc. Psychiatry Psychiatr. Epidemiol.*, vol. 51, no. 1, pp. 125–127, 2016.

[4]     X. Huang, X. Li, L. Zhang, T. Liu, D. Chiu, and T. Zhu, "Topic model for identifying suicidal ideation in Chinese microblog," *29th Pacific Asia Conf. Lang. Inf. Comput. PACLIC 2015*, pp. 553–562, 2015.

[5]     A. C. Fernandes, R. Dutta, S. Velupillai, J. Sanyal, R. Stewart, and D. Chandran, "Identifying Suicide Ideation and Suicidal Attempts in a Psychiatric Clinical Research Database using Natural Language Processing," *Sci. Rep.*, vol. 8, no. 1, pp. 1–10, 2018.

[6]     M. M. Tadesse, H. Lin, B. Xu, and L. Yang, "Detection of suicide ideation in social media forums using deep learning," *Algorithms*, vol. 13, no. 1. 2020.

[7]     "Suicide Warning Signs," *PSYCOM*. [Online]. Available: https://www.psycom.net/suicide-warning-signs.

[8]     S. M. Mohammad, "Word affect intensities," *arXiv*, 2017.

[9]     F. A. Acheampong, C. Wenyu, and H. Nunoo-Mensah, "Text-based emotion detection: Advances, challenges, and opportunities," *Eng. Reports*, vol. 2, no. 7, pp. 1–24, 2020.

[10]     E. Kušen, G. Cascavilla, K. Figl, M. Conti, and M. Strembeck, "Identifying emotions in social media: Comparison of word-emotion lexicons," *Proc. - 2017 5th Int. Conf. Futur. Internet Things Cloud Work. W-FiCloud 2017*, vol. 2017-Janua, pp. 132–137, 2017.

[11]     Wikipedia, "Robert Plutchik." [Online]. Available: https://en.wikipedia.org/wiki/Robert_Plutchik.

[12]     "Sentiment and emotion lexicons," *Government of Canada*.     [Online].     Available: https://nrc.canada.ca/en/research-development/products-services/technical-advisory-services/sentiment-emotion-lexicons.

[13]     "NRC Word-Emotion Association Lexicon," *NRC*. [Online].     Available: https://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm.

[14]     "Suicide cases on the rise amid pandemic," *The Malaysia    Reserve*.    [Online].    Available: https://themalaysianreserve.com/2020/11/11/suicide-cases-on-the-rise-amid-pandemic/.

[15]     M. H. Metzger, E. Poulet, and Q. Gicquel, "Use of emergency department electronic medical records for automated epidemiological surveillance of suicide attempts : a French pilot study," no. June, pp. 1–10, 2016.

[16]     H. Yang, A. Willis, A. De Roeck, and B. Nuseibeh, "A Hybrid Model for Automatic Emotion Recognition in Suicide Notes," *Biomed. Inform. Insights*, vol. 5s1, p. BII.S8948, 2012.

[17]     M. A. Razek and C. Frasson, "Text-Based Intelligent Learning Emotion System," *J. Intell. Learn. Syst. Appl.*, vol. 09, no. 01, pp. 17–20, 2017.

[18]     F. A. Acheampong and H. Nunoo-mensah, "Recognizing Emotions from Texts using a BERT-based Approach," no. November, 2020.

[19]     S. Ghosh, A. Ekbal, and P. Bhattacharyya, "CEASE , a Corpus of Emotion Annotated Suicide notes in English," no. May, pp. 1618–1626, 2020.

[20]     "How Do You Create a Word Cloud," *ProWritingAid*.     [Online].     Available: https://prowritingaid.com/art/425/What-the-Heck-is-a-Word-Cloud-and-Why-Would-I-Use-One.aspx#:~:text=A word cloud can help,want to use the most.&text=A word cloud can also,%2C themes%2C and plot points.

[21]     "Introduction to Topic Modeling," *MonkeyLearn*. [Online].     Available: https://monkeylearn.com/blog/introduction-to-topic-modeling/#:~:text=Topic    modeling    is    an unsupervised,characterize a set of documents.

[22]     S. Kapadia. "Evaluate Topic Models: Latent Dirichlet Allocation (LDA)." Towards Data Science Inc. https://towardsdatascience.com/evaluate-topic-model-in-python-latent-dirichlet-allocation-lda-7d57484bb5d0 (accessed 2020-12-20, 2020).