Singapore Management University

## Institutional Knowledge at Singapore Management University

# Data mining approach to the detection of suicide in social media: A case study of Singapore

Jane H. K. SEAH
*Singapore Management University*, jane.seah.2016@smu.edu.sg

Kyong Jin SHIM
*Singapore Management University*, kjshim@smu.edu.sg

# Data Mining Approach to the Detection of Suicide in Social Media: A Case Study of Singapore

Jane H. K. Seah
School of Information Systems
Singapore Management University
Singapore
jane.seah.2016@sis.smu.edu.sg

Kyong Jin Shim
School of Information Systems
Singapore Management University
Singapore
kjshim@smu.edu.sg

*Abstract*—**In this research, we focus on the social phenomenon of suicide. Specifically, we perform social sensing on digital traces obtained from Reddit. We analyze the posts and comments in that are related to depression and suicide. We perform natural language processing to better understand different aspects of human life that relate to suicide.**

*Keywords—suicide, suicide detection, social media, data mining*

## I. INTRODUCTION

Suicide is a global phenomenon in all regions of the world. The World Health Organization reports [1] that some 800,000 lives are lost annually due to suicide. A recent article [5] reports that the global suicide rate in developed countries rose since 1999 and by the year 2010, suicide became the major reason behind unnatural death. Singapore is not an exception. As population ageing trend continues, more cases of elderly suicide are being reported in Singapore [2]. Also, the number of suicide-related deaths is reported to be far exceeding that of transport-related accidents [3]. The same report shows that especially for the younger age group of 10 years old to 29 years old, suicide is known to be the leading cause of their death. Just last year, a suicide case involving a young student [6] was reported in Singapore. According to the article [6], extreme stress caused by exams is suspected to have caused the suicide.

With over 900 active social media platforms globally as of today, platforms such as Facebook, Twitter and Reddit have seen a drastic increase in active conversations about various topics and issues globally. Such digital traces provide valuable insights into human emotional states and feelings. While offline methods such as hotlines and in-person counselling exist to capture human emotional states and feelings, social media platforms provide rich digital traces that can be analyzed to detect suicidal intentions using data mining techniques.

## II. DATASET

In this study, we investigate depression and suicide-related conversations in Singapore. We performed data crawling from Reddit using their API. After obtaining an API key, we wrote Python scripts to search and crawl posts and comments from Subreddit "/r/Singapore." We use PRAW [7], a popular Python package and a Reddit API wrapper. Using the pre-defined and customized list of key terms related to depression and suicide, the Python scripts then used these terms to query Reddit (Figure 1). The crawled data are stored as CSV files (Figure 2).



```python
dir_name = word_file[:-4] + '-reddit-singapore'
create_project_dir(dir_name)
subreddit = reddit.subreddit('Singapore')
for word in words:
    path = os.path.join(dir_name, word+'.csv')
    if not os.path.isfile(path):
        subreddit_query = subreddit.search(word)

        if subreddit_query:
            print(word)

            topics_dict = {
                "title":[],
                "score":[],
                "id":[],
                "url":[],
                "comms_num": [],
                "created": [],
                "body": [],
                "author_name": [],
                #    "reports_num":[]
            }

            for submission in subreddit_query:
                topics_dict["title"].append(submission.title)
                topics_dict["score"].append(submission.score)
                topics_dict["id"].append(submission.id)
                topics_dict["url"].append(submission.url)
                topics_dict["comms_num"].append(submission.num_comments)
                topics_dict["created"].append(get_date(submission.created))
                topics_dict["body"].append(submission.selftext[:-3])
                topics_dict["author_name"].append(submission.author.name)
                #    topics_dict["reports_num"].append(submission.num_reports)

            # Load into dataframe
            topics_data = pd.DataFrame(topics_dict)

            # Write to csv
            topics_data.to_csv(path, index=False)
```

Fig. 1.   Reddit Crawling Script Written in Python



Fig. 2.   Depression and Suicide-Related Posts and Comments Crawled from Reddit

Figure 3 shows a depression-related post in Reddit that our script crawled. 385 posts and over 21,000 comments mentioning depression and suicide-related terms were retrieved from Reddit. The data span the Year 2010 through the Year 2018. We used the suicide-related terms from a previous study [4], which includes terms indicative of depression and suicide (e.g. "feel useless", "hang myself", known suicide methods, prescription drug names, etc.).
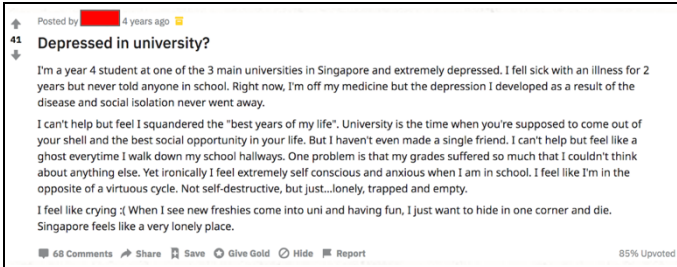


Fig. 3.   A Depression-Related Reddit Post (the user's name is masked).

Prior to programmatically performing data cleaning and transformation, our human researchers manually inspected the crawled data to ensure relevancy of the data. Some 5% of the data contained Chinese characters (as Chinese is widely spoken in Singapore). We engaged human translators to translate Chinese phrases and words into English. Initially, we used freely available online translation services, however, the quality of the translation was deemed poor. A majority of the translated content turned out to be irrelevant to our topic (e.g. comments related to recent political elections), and hence, we removed it from our analysis.

Our researchers also spent a substantial amount of time expanding our corpus with acronyms or terms previously unknown. For instance, "pes" or "physical employment status" was mentioned in some of the posts. This term is specific to Singapore's National Service. Physical Employment Status (PES) is a required physical fitness exam for National Service registrants, enlisted soldiers and reservists. As it appeared along with depression and suicide indicating terms, we took note of this and added it to our corpus.

We further enhanced this corpus with new terms such as new drug names associated with treatment of depression, other suicide methods, and mostly local terms such as "imh" (Institute of Mental Health in Singapore) and terms related to Singapore's National Service. We included these new terms in our customized corpus based on term frequency in the data that our study obtained.

Further data cleaning and transformation involved removal of certain stop words, Unicode characters, and hyperlinks. Upon manual data cleaning and inspection, next we performed stemming, lemmatization, spell checking prior to performing analyses.

## III. FINDINGS & ANALYSIS

We performed topic modeling using Latent Dirichlet Allocation (LDA) [8] which automatically discovers topics. We used Python Gensim package [9] for topic modelling.



Fig. 4.   Some of the Topics Discovered by LDA

Figure 4 and Figure 5 show the topics the LDA algorithm discovered from our Reddit dataset.



Fig. 5.   More Topics Discovered by LDA

Conversations about treatment, clinics and costs are seen in a cluster (Figure 6). Terms such as "hospital", "counsellor", "doctor", "medication", "imh", and "therapy" appear in this topic.
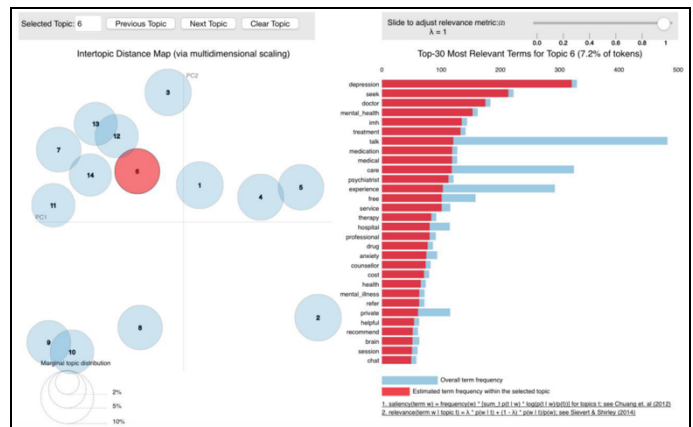


Fig. 6.   Results of Topic Modeling. Cluster showing treatment-related terms
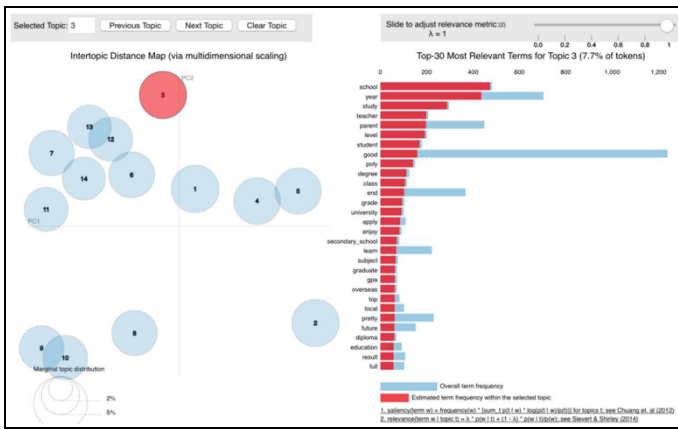
Fig. 7. Results of Topic Modeling. Cluster Showing Terms Related to Schools and Studies.
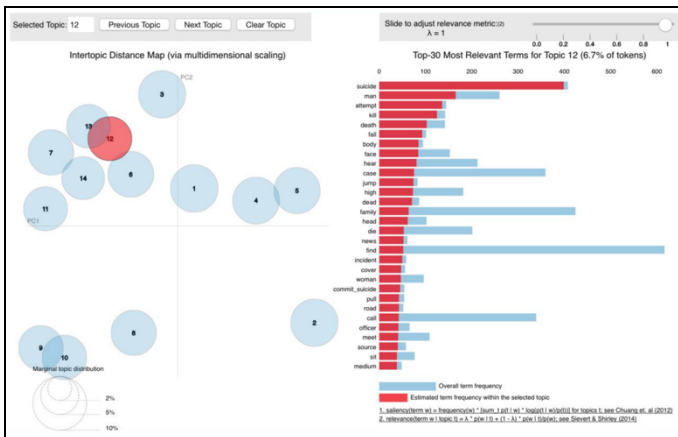


Fig. 8. Results of Topic Modeling. Cluster Showing Terms Related to Suicide Cases and Reports.

Figure 7 shows a topic related to schools and studies. And in Figure 8, we can see another topic with terms related to reports of suicide cases. Below, we contextually explain the topics we discovered using LDA algorithm (Figure 9).

| Topic | Keywords | Description |
|---|---|---|
| 1 | life, live, good, end, make, people, give, die, euthanasia, decide | Decision-making on life and death, euthanasia |
| 2 | feel, friend, thing, good, talk, find, understand, bad, relationship, mind | Friendships, support, relationships |
| 3 | school, year, study, teacher, parent, level, student, good, poly, degree | School life, stress, expectations |
| 4 | people, make, love, kid, wrong, stop, point, matter, choice, guess | Societal opinions |
| 5 | issue, problem, delete, give, question, op, point, fact, condition, social | Discussions of problems faced, comments on Reddit posts ("op" refers to the original poster) |
| 6 | work, job, good, find, hard, time, learn, thing, part, start | Career, stress, difficulties |
| 7 | people, make, change, country, place, world, kind, bad, great, singaporean | Society-at-large, attitudes |
| 8 | call, case, shit, fuck, police, guy, happen, report, girl, pretty | Cases which were reported to the police |
| 9 | post, hope, advice, comment, read, url, share, give, experience, hear | Encouraging comments, sharing experiences and advice |
| 10 | law, gay, society, child, accept, state, sex, man, religion, opinion | LGBT acceptance, religious beliefs, 377A |
| 11 | suicide, man, attempt, kill, death, fall, body, face, hear, case | Suicide attempts, explicit descriptions |
| 12 | family, parent, money, leave, move, dad, house, home, situation, mother | Familial problems |
| 13 | time, day, back, start, month, year, thing, long, remember, bad | Recollections of past events |
| 14 | depression, seek, doctor, mental_health, imh, treatment, talk, medical, medication, care | Treatment, seeking help ("imh" refers to the Institute of Mental Health) |

Fig. 9. Summary of 14 Topics Related to Depression and Suicide.

Our analysis reveals interesting insights into different aspects or themes concerning conversations about depression and suicide. Given a very large amount of textual content, topic modelling was able to quickly summarize and group social conversations about depression and suicide into 14 topics where each topic is uniquely identifiable by a set of keywords.

The discovered topics reveal different contexts in which suicidal thoughts or discussions occur. They may not necessarily be direct causes of suicide acts. However, the insights into the contexts in which conversations about depression and suicide occur can help all of us better understand where to look in order to find troubled individuals in need of desperate help.

## IV. CONCLUSION AND FUTURE DIRECTIONS

Our study demonstrates that a data mining approach allows for efficient and automated ways for detecting suicide in social media. In Singapore, there are two 24-hour suicide hotlines available: 1) Samaritans of Singapore and 2) SAF Hotline (for military personnel). A number of other non-emergency helplines are available such as Youth Line and Tinkle Friend. While these services remain critical in ensuring safety of people, we believe that tapping into digital traces such as public forums can help authorities proactively reach out to those in need of help.

Future directions of this proof-of-concept study include: 1) expanding the current corpus with local languages beyond the Chinese language, 2) working with local organizations to combine digital traces with offline data (e.g. suicide hotlines, counseling), and 3) expanding our analysis to other popular public forums beyond Reddit.

## REFERENCES

[1] Suicide data. (2018 November 11). Retrieved from http://www.who.int/mental_health/prevention/suicide/suicideprevent/en/

[2] Rashith, R. (2018 November 11). "Number of suicides committed by the elderly hits record high as Singapore population ages". Retrieved from https://www.straitstimes.com/singapore/more-than-1-in-3-suicides-committed-by-elderly-as-singapore-population-ages

[3] The Situation in Singapore. (2018 November 11). Retrieved from https://www.sos.org.sg/learn-about-suicide/quick-facts

[4] J. Jashinsky, S. H. Burton, C. L. Hanson, J. West, C. Giraud-Carrier, M. D. Barnes, and T. Argyle. Tracking suicide risk factors through twitter in the us. Crisis, 35(1):51–59, 2014.

[5] Taylor, D. Suicide as a Product of Modernization. Retrieved from https://medium.com/@dvlan/suicide-as-a-product-of-modernization-776e9823cfaf

[6] Fun, P. Lonely 12-Year-Old Girl Sadly Commits Suicide Because She Barely Passed Exams. Retrieved from https://www.worldofbuzz.com/lonely-12-year-old-girl-sadly-commits-suicide-barely-passed-exams/

[7] PRAW: The Python Reddit API Wrapper. https://praw.readthedocs.io/

[8] Blei, D. M., Ng, A. Y. and Jordan, M. I. (2003) Latent Dirichlet Allocation. The Journal of Machine Learning Research, 3, 993-1022.

[9] Gensim: https://radimrehurek.com/gensim/