

# Detection of Suicidal Twitter Posts

Fatima Chiroma<sup>1</sup>, Mihaela Cocea<sup>1</sup>, and Han Liu<sup>2</sup>

<sup>1</sup> University of Portsmouth, Portsmouth, UK,  
fatima.chiroma@port.ac.uk, mihaela.cocea@port.ac.uk

<sup>2</sup> Cardiff University, Cardiff, UK  
LiuH48@cardiff.ac.uk

**Abstract.** As web data evolves, new technological challenges arise and one of the contributing factors to these challenges is the online social networks. Although they have some benefits, their negative impact on vulnerable users such as the spread of suicidal ideation is concerning. As such, it is vital to fine tune the approaches and techniques in order to understand the users and their context for early intervention. Therefore, in this study, we measured the impact of data manipulation and feature extraction, specifically using N-grams, on suicide-related social network text (tweets). We propose a diversified ensemble approach (multi-classifier fusion) to improve the detection of suicide-related text classification. Four machine classifiers were used for the fusion: Support Vector Machine, Random Forest, Naïve Bayes and Decision Tree. The results of our proposed approach have shown that the multi-classifier fusion has improved the detection of suicide-related text and, also, that Support Vector Machine has shown some promising results when dealing with multi-class datasets.

**Keywords:** Ensemble Learning, Suicide-related Tweets, Text Classification

## 1 Introduction

The ability of machine learning to automatically detect and uncover patterns in data for the purpose of prediction as well as to enhance decision making has, in the last decade, led to the rapid increase in its application across diverse areas including Healthcare [15], Finance [31] and Law Enforcement [11], to name a few. Law enforcement and other intelligence organizations have used machine learning approaches and techniques to explore large databases efficiently [11]. Unfortunately, the emergence of the online social networks (often referred to as social media) has created a new source for immense and uncontrollable data generation, in addition to transforming the way crime and victimisation is understood, experienced and committed [19].

This type of crimes, such as hate crime which is a criminal offence that is prejudicially targeted towards someone based on a personal characteristic [6, 14], make up around two percent of crimes based on the notifiable reports recorded by the

UK Home Office [14]. However, the percentage is believed to be higher as this type of crime is largely unreported to the police [8]. Additionally, hate crimes often follow hate speech and over the last decade hate speech online has significantly increased [4, 24]. Online hate speech such as misogyny and the spread of suicidal ideation may impact vulnerable social network users as they are at potential risk of harming themselves due to the information they receive [9, 12].

Furthermore, several studies have shown the correlation between social media and suicidal behaviour [17, 27, 30]. Hence, there is a need to understand social network users and the contents they post for the enhancement of existing approaches and techniques, for possible interventions, as well as keeping up with the evolving web. Therefore, in this study, we use Twitter posts that were extracted using suicide-related search terms to propose an approach based on text classification, a machine learning task, to investigate whether employing ensemble learning would lead to an improved detection of suicidal risk from social media text. According to [7], Twitter is a logical source for suicide-related communications as users are more likely to deindividualize and express themselves emotionally while other studies [9, 10] have shown that people are more likely to seek for support through social networks, such as Twitter, than professional help due to anonymity and concerns of social stigmatization.

Detection of suicidal risk from social media text using automatic techniques, such as text classification, has only recently started to be explored, with only few studies [9, 12, 13, 20] reported. In this paper, we further investigate the performance of classifiers, while exploring ensemble learning. In particular, we are investigating the influence of ensemble learning, i.e. the use of several classifiers, on classification performance in comparison with using individual classifiers.

The rest of the paper is organized as follows: Section 2 describes the background and related work on social network suicide-related communications and machine learning, focusing specifically on text classification; Section 3 provides details of the proposed approach including the experimental process; in Section 4 the results obtained are presented and discussed; and Section 5 draws the conclusions which include a summary of contributions of the paper and future directions.

## 2 Background and Related Work

In this section, background related to text classification is covered, as well as related work on suicidal ideation on social media.

### 2.1 Text Classification

Classification is one of the most prominent machine learning tasks where the category of an unseen instance is judged [12] and it typically involves employing an algorithm to build a model to identify an instance’s category. Furthermore,

classification relating to text is referred to as text classification and it can also be defined as the assigning of a pre-defined class to a textual instance in a dataset [3]. Although the concept of classification, regardless of type, seems quite straightforward, it is complex and cannot guarantee an accurate classification for unseen instances, especially when dealing with real world data which contains many irrelevant, noisy and redundant features [18].

Furthermore, several studies have identified some problems relating to the misclassification of data which includes class imbalance [2, 21, 23]. Class imbalance is the insufficient representation of the target or minority class in a dataset [23], and most machine learning algorithms do not consider the underlying distribution of a dataset, generally leading to good performance on the detection of majority classes (as the algorithm has more instances to learn from) and poor performance on the minority classes (as the algorithm may not have been exposed to sufficient information to learn reliable patterns).

## 2.2 Suicide and Online Social Networks

Text classification relating to suicide-related communications is still in its infancy stage; as such, the research in this area is limited. Some studies have been carried out using text classification to try and detect social network users that are at risk of suicide. An example of such a study is [1], where they used machine learning to identify risk factors relating to suicide from Twitter conversations and they found a strong correlation between geographical suicide rates and Twitter data, with an accuracy of approximately 63%.

Additionally, in their study, [9] used machine classifiers to classify suicide-related Twitter communications. Their baseline experiment achieved an F-measure of 0.702 for all their (seven) classes, however, they further improved the results to 0.728 by applying an ensemble learning approach. Another study was carried out by [12], where they conducted a baseline experiment to measure the performance of popular machine classifiers in distinguishing suicide-related communications from Twitter. They acquired these dataset from [9] and used Decision Tree, Naive Bayes, Random Forest and Support Vector Machine for the classification. Their result showed an F-measure of up to 0.778 was achieved by the Decision Tree classifier.

## 3 Experimental Approach

We propose an approach based on ensemble learning, to investigate whether it would lead to an improved detection of suicide-related communications. Fig. 1 provides an overview of the experimental approach which consists of four stages: Data Preparation, Feature Preparation, Individual Classification and Ensemble

Classification. Furthermore, this experiment was carried using Knime<sup>3</sup>, an open-source data analytics tool.

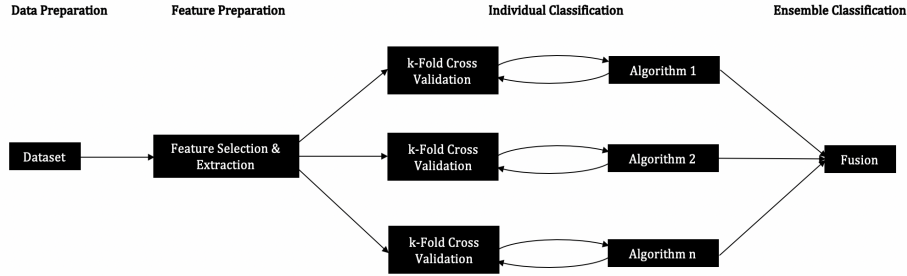


Fig. 1. The Experimental Approach

### 3.1 Data Preparation

This is the initial stage of the approach where a dataset is cleaned and partly transformed for modelling. It comprises of the data collection, data manipulation and pre-processing.

**Data Collection:** Twitter is a good source for suicide-related communications as stated in Section 1. Therefore, in our studies, we used 2,000 suicide-related communication tweets from [9], which were collected from Twitter using the Twitter Streaming Application Programming Interface (API). They used search keywords from reported news and lexicon of terms such as *don't want to exist* and *Kill myself* which were derived from known suicide websites for the collection. These were further annotated by four human annotators from CrowdFlower<sup>4</sup>, a crowd sourcing online service, into seven suicide categories as shown in Table 1. These categories were developed by [9] with expert researchers in the area of suicide to best capture people's general representation when communicating on suicide topics. Additionally, following established methods [9, 27], we also discarded tweets that have less than 75% annotator agreement leaving a total of 1064 tweets. The tweets are organised into several datasets for experimentation, as outlined below.

**Data Manipulation:** The data can be organised in different datasets reflecting different labeling schemes derived from the initial seven labels, which can be manipulated and categorized based on the level of similarity or dissimilarity. For example, class 1 is suicide and class 2 is the flippant reference to suicide,

<sup>3</sup> <https://www.knime.com>

<sup>4</sup> <http://www.crowdflower.com>

**Table 1.** Instances Per Class (Adapted from [9])

Class	Type	Description	Instances
1	Suicide	Possible suicidal intent	159
2	Flippant	Un-serious reference to suicide	133
3	Campaign	Suicide petitions	158
4	Support	support or information	178
5	Memorial	Condolences or memorial	142
6	Reports	Suicide reports excluding bombing	165
7	Other	None of the above	129
<b>Total: 7</b>	-	-	<b>1064</b>

**Table 2.** Data Manipulation Distribution (Adapted from [12])

Class	Dataset	Class name	Raw	Processed	Total
1	Binary-class	Suicide	159	156	289
2		Flippant	133	133	
1	Three-class	Suicide	159	156	1060
2		Flippant	133	133	
3, 4, 5, 6, 7		Non-suicide	772	771	
1	Seven-class	Suicide	159	156	1060
2		Flippant	133	133	
3		Campaign	158	158	
4		Support	178	178	
5		Memorial	142	142	
6		Reports	165	165	
7		Other	129	128	

however classes 3 to 7 are about suicide in other contexts, i.e. not in the context of a person considering the possibility of committing suicide. Table 2 describes the data manipulation distributions – although all the datasets are from the same original data, the classes, size and complexity have changed for the resulting datasets.

**Pre-processing:** When dealing with real world textual datasets, especially social network user generated datasets, pre-processing is vital. These datasets contain noise and redundant information; therefore, the use of pre-processing techniques lead to the removal of irrelevant features and reduces the vector space, thereby improving classification performance [5, 16]. For this study, standard and established pre-processing methods [9, 16, 25] were applied, which include the removal of stop words, words containing numbers, punctuations, URLs and non-ASCII characters, Part-of-speech tagging and reducing terms to their stem for redundancy reduction (these reduced the number of instances from 1064 to 1060).

### 3.2 Feature Preparation

Subsequent to pre-processing, the Bag-of-words representation, which ignores the syntactic and semantic information and treats the text as a collection of words [25], was used to extract relevant unigrams (i.e. individual words) based on their term frequencies, while the  $n$ -gram approach (i.e. each feature is a term including  $n$  words) was used to extract terms between 2 and 5 words, also based on term frequencies. However, we found that the use of 4-grams and 5-grams de-

**Table 3.** NGrams Per Dataset

Dataset	1 Gram	1 - 3 Gram
Binary-class	599	2,075
Three-class	2,223	10,990
Seven-class	2,223	10,990

grades the classification performance due to almost non-existent representation of the target classes (i.e. suicide as well as flippant). Consequently, we report only the experiments using 1 to 3-grams and the number of features for each of the three datasets are displayed in Table 3. Additionally, dimensionality reduction is typically applied to textual data, as it is known that high dimensional text data hinders classifiers' performance [22, 29]. However, in this study, no dimensionality reduction technique has been applied as the number of features are small and dimensionality reduction is usually applied to datasets with hundreds of thousands of features [28].

### 3.3 Individual Classification

This phase consists of the training and evaluation of the individual classifiers, where we use two sets of features: 1 Gram and 1-3 Gram. The machine learning algorithms that are used to train both sets of features were chosen based on their performance from previous studies [12]. These classifiers are Decision Tree (DT), Naïve Bayes(NB), Random Forest (RF) and Support Vector Machine (SVM). This was done for each of the datasets, i.e. the binary-class, three-class and seven-class datasets. Additionally, this (individual classification) phase will be used as the baseline for comparison with the ensemble classification (Multi-classifier fusion) results – please refer to Fig. 1 for these phases.

**Training Setup:** For training, stratified sampling was used, which preserves the original class distribution for the training and test data. 10-fold cross-validation was used given the relatively small size of the data, especially for the binary-class dataset.

**Evaluation:** Evaluation as a process allows us to determine the extent to which an objective has been attained [26]. For text classification purposes, the standard classification metrics, such as Precision, Recall and F-measure are used to measure the performance. Accuracy is not typically used for measuring the performance of text classification, especially for data with class imbalance, as the higher results of the majority class (which is often not the one of interest) gives the false impression of a good performance. Consequently, the F-measure is preferred to accuracy as it can reflect how the overall performance is affected due to the low performance for some classes.

### 3.4 Ensemble Classification

The ensemble classification stage uses the ensemble learning approach, which involves combining (fusing) the outputs of each classifier through techniques such as majority voting or algebraic formulas (e.g. weighted sum). In our approach, we use majority voting but excluded Random Forest as part of the ensemble, as it is already an ensemble method and will provide another point of comparison i.e. between the two ensembles.

## 4 Results and Discussion

This study investigated whether the use of ensemble learning will improve the detection of suicide-related communications on social media, and more specifically, from Twitter. The experimental investigation in this paper builds on previous work by [9], as well as our own previous studies [12] by using the same dataset, pre-processing techniques and machine classifiers for this study. Additionally, the results from this study are presented in two categories: (a) the individual classification results which comprises of the results for the individual classifiers and (b) the ensemble classification results, the results for the individual classifiers by applying the ensemble learning approach.

### 4.1 Individual Classification

We report the results for each of the datasets, i.e. binary-class, three-class and seven-class; we present the overall results (see Fig. 2), as well as per class, as we are interested in the performance of the *Suicide* and *Flippant* classes in particular, as these are the ones reflecting suicide risk. The results from the binary-class dataset indicate an F-measure between 0.411 to 0.776 was achieved for the 1 Gram whereas a similar but lower result of 0.380 to 0.771 (see Fig. 2) was achieved for the 1-3 Gram. The suicide class has a higher F-measure of up to 0.80 than the flippant class (see Table 4) in both 1 Gram and 1-3 Gram, which was achieved by SVM. The lowest performing classifiers for the binary-class dataset is NB with an F-measure of 0.38 and 0.41, respectively; Furthermore, both NB and RF seem to have performed well on the suicide class, but performed (very) poorly on the flippant class.

In addition, the performance of the classifiers for the three-class dataset (see Fig. 2 and Table 5) and seven-class dataset (see Fig. 2, Table 6 and 7) has varying performance ranging from 0.00 to 0.90 for the three-class and 0.04 to 0.74 for the seven-class. NB had the worst performance for all the dataset categories i.e. binary-class, three-class and seven-class. Furthermore, DT and SVM are the two highest performing classifiers in this phase, however their performance varies depending on the dataset. For instance, DT had the highest F-measure for the binary-class, whereas SVM had the best performance for the three-class and seven-class dataset; which may imply that SVM performs better with larger and/or multi-class datasets.

**Table 4.** Individual Classification Results: Binary-class

Classifier	Measure	1 Gram		1-3 Gram	
		<i>Suicide</i>	<i>Flippant</i>	<i>Suicide</i>	<i>Flippant</i>
DT	Recall	0.731	0.820	0.757	0.789
	Precision	0.826	0.722	0.821	0.719
	F-measure	0.776	0.768	0.788	0.753
NB	Recall	0.942	0.075	0.988	0.023
	Precision	0.544	0.526	0.562	0.600
	F-measure	0.690	0.132	0.717	0.043
RF	Recall	0.917	0.444	0.982	0.203
	Precision	0.659	0.819	0.610	0.900
	F-measure	0.767	0.576	0.753	0.331
SVM	Recall	0.821	0.729	0.822	0.707
	Precision	0.780	0.776	0.781	0.758
	F-measure	0.800	0.752	0.801	0.732

**Table 5.** Individual Classification Results: Three-class

Classifier	Measure	1 Gram			1-3 Gram		
		<i>Suicide</i>	<i>Flippant</i>	<i>Non-suicide</i>	<i>Suicide</i>	<i>Flippant</i>	<i>Non-suicide</i>
DT	Recall	0.603	0.316	0.883	0.604	0.293	0.901
	Precision	0.537	0.512	0.848	0.551	0.506	0.865
	F-measure	0.568	0.391	0.865	0.576	0.371	0.882
NB	Recall	0.506	0.053	0.970	0.485	0.023	0.983
	Precision	0.840	0.280	0.794	0.891	0.250	0.816
	F-measure	0.632	0.089	0.874	0.628	0.041	0.892
RF	Recall	0.615	0.015	0.914	0.538	0.000	0.948
	Precision	0.568	0.182	0.801	0.615	0.000	0.819
	F-measure	0.591	0.028	0.854	0.574	0.000	0.879
SVM	Recall	0.641	0.226	0.921	0.675	0.278	0.929
	Precision	0.637	0.536	0.838	0.640	0.544	0.878
	F-measure	0.639	0.317	0.877	0.657	0.368	0.903



**Table 6.** Individual Classification (1 Gram) Results: Seven-class

Datasets	DT			NB			RF			SVM		
	P	R	F	P	R	F	P	R	F	P	R	F
Suicide	0.67	0.54	0.60	0.50	0.87	0.63	0.88	0.43	0.57	0.74	0.58	0.65
Campaign	0.63	0.59	0.61	0.24	0.91	0.38	0.67	0.68	0.67	0.64	0.67	0.66
Flippant	0.42	0.35	0.38	0.05	0.54	0.10	0.17	0.48	0.25	0.33	0.41	0.37
Support	0.55	0.68	0.60	0.98	0.21	0.35	0.77	0.47	0.58	0.70	0.71	0.71
Memorial	0.36	0.30	0.33	0.08	0.41	0.13	0.23	0.39	0.29	0.47	0.36	0.41
Reports	0.35	0.45	0.39	0.10	0.55	0.17	0.26	0.60	0.37	0.42	0.54	0.47
Other	0.28	0.41	0.33	0.16	0.69	0.25	0.44	0.66	0.53	0.49	0.56	0.52

**Table 7.** Individual Classification (1 - 3 Gram) Results: Seven-class

Datasets	DT			NB			RF			SVM		
	R	P	F	R	P	F	R	P	F	R	P	F
Suicide	0.71	0.46	0.56	0.48	0.91	0.63	0.76	0.50	0.61	0.79	0.54	0.64
Campaign	0.73	0.68	0.70	0.41	0.99	0.58	0.57	0.79	0.66	0.72	0.77	0.74
Flippant	0.35	0.39	0.37	0.02	0.60	0.04	0.12	0.52	0.20	0.29	0.41	0.34
Support	0.66	0.84	0.74	0.99	0.27	0.42	0.95	0.38	0.54	0.74	0.75	0.74
Memorial	0.46	0.32	0.38	0.14	0.88	0.24	0.12	0.83	0.22	0.52	0.37	0.44
Reports	0.31	0.43	0.36	0.06	0.83	0.11	0.21	0.68	0.32	0.33	0.55	0.41
Other	0.36	0.55	0.44	0.26	0.95	0.41	0.35	0.86	0.49	0.50	0.55	0.53

## 4.2 Ensemble Classification

The results show that each combination performed differently on each dataset. Interestingly, combining DT, NB and SVM gave a higher performance than combining only DT and SVM even though NB has the lowest performance amongst all the classifiers. Additionally, there is an improved performance when the ensemble learning is applied except in two cases (1 Gram for seven-class and 1-3 Gram for three-class) where SVM has the higher performance. Also, the multi-classifier fusion has outperformed RF for all the datasets.

Additionally, in previous work [12], the feature extraction methods used were the bag-of-words and document frequency to generate only unigrams however in this study we further explored the use of bigrams and trigrams. Furthermore, applying the ensemble approach gave an improved performance compared with the individual classifiers for all the datasets, except in two cases mentioned earlier. The combination of DT, NB and SVM gave the best performance for all the datasets while removing NB from the ensemble combination marginally deteriorates the performance of all the datasets except the binary-class dataset.

## 5 Conclusion and Future Direction

In this paper, we investigated whether employing ensemble learning would lead to an improved detection of suicidal risk from social media text. Although there

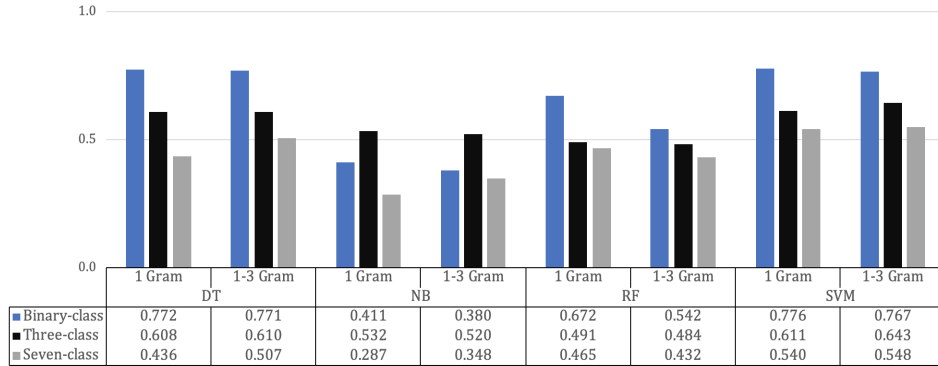


Fig. 2. Training results (F-measure) per dataset

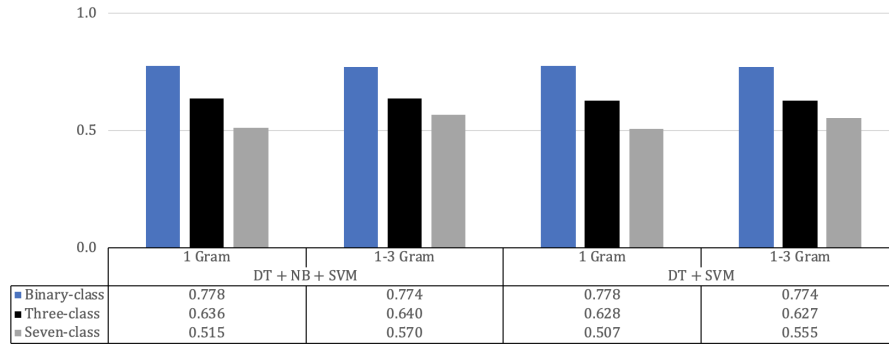


Fig. 3. Ensemble Classification: F-measure

is evidence of classification improved, the improvement when compared to the worst performing classifier is significant but it is not significant when compared to the best performing classifiers. Consequently, the use of ensemble learning with majority voting improves the prediction of suicide-related text, but only marginally. In future work, we will investigate alternative ways of fusing the classifiers and their influence on performance.

## 6 Acknowledgement

This is an independent research that is supported by the Petroleum Development Technology Fund (PTDF) and the Department of Health Policy Research Programme (Understanding the Role of Social Media in the Aftermath of Youth Suicides, Project Number 023/0165). The views expressed in this publication are those of the authors and not necessarily those of PTDF or Department of Health.

## Bibliography

- [1] Abboute A, Boudjeriou Y, Entringer G, Azé J, Bringay S, Poncelet P (2014) Mining twitter for suicide prevention. In: International Conference on Applications of Natural Language to Data Bases/Information Systems, Springer, pp 250–253
- [2] Ali A, Shamsuddin SM, Ralescu AL (2015) Classification with class imbalance problem: a review. *Int J Advance Soft Compu Appl* 7(3):176–204
- [3] Allahyari M, Pouriyeh S, Assefi M, Safaei S, Trippe ED, Gutierrez JB, Kochut K (2017) A brief survey of text mining: Classification, clustering and extraction techniques. *Proceedings of KDD Bigdas, Halifax, Canada, August 2017* p 13
- [4] Banks J (2010) Regulating hate speech online. *International Review of Law, Computers Technology* 24(3):233–239
- [5] Barbosa L, Feng J (2010) Robust sentiment detection on twitter from biased and noisy data. In: *Proceedings of the 23rd international conference on computational linguistics: posters*, pp 36–44
- [6] Burnap P, Williams ML (2014) Hate speech, machine classification and statistical modelling of information flows on twitter: Interpretation and communication for policy decision making. In *Proceedings of IPP 2014* pp 1–18
- [7] Burnap P, Williams ML (2015) Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and decision making. *Policy & Internet* 7(2):223–242
- [8] Burnap P, Williams ML (2016) Us and them: identifying cyber hate on twitter across multiple protected characteristics. *EPJ Data Science* 5(1):11
- [9] Burnap P, Colombo W, Scourfield J (2015) Machine classification and analysis of suicide-related communication on twitter. In: *Proceedings of the 26th ACM conference on hypertext & social media, ACM*, pp 75–84
- [10] Cavazos-Rehg PA, Krauss MJ, Sowles S, Connolly S, Rosas C, Bhargava M, Bierut LJ (2016) A content analysis of depression-related tweets. *Computers in Human Behavior* 54:351–357
- [11] Chen H, Chung W, Xu JJ, Wang G, Qin Y, Chau M (2004) Crime data mining: a general framework and some examples. *computer* 37(4):50–56
- [12] Chiroma F, Liu H, Cocea M (2018) Text classification for suicide related tweets. In: *2018 International Conference on Machine Learning and Cybernetics (ICMLC), IEEE*, vol 2, pp 587–592
- [13] Colombo GB, Burnap P, Hodorog A, Scourfield J (2016) Analysing the connectivity and communication of suicidal users on twitter. *Computer communications* 73:291–300
- [14] Corcoran H, Smith K (2016) Hate crime, England and Wales, 2015/16. URL [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/559319/hate-crime-1516-hosb1116.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/559319/hate-crime-1516-hosb1116.pdf)

- [15] Dipnall JF, Pasco JA, Berk M, Williams LJ, Dodd S, Jacka FN, Meyer D (2016) Fusing data mining, machine learning and traditional statistics to detect biomarkers associated with depression. *PloS one* 11(2):1–23
- [16] Haddi E, Liu X, Shi Y (2013) The role of text pre-processing in sentiment analysis. *Procedia Computer Science* 17:26–32
- [17] Jashinsky J, Burton SH, Hanson CL, West J, Giraud-Carrier C, Barnes MD, Argyle T (2014) Tracking suicide risk factors through Twitter in the US. *Crisis* 35(1):51–59
- [18] Li J, Cheng K, Wang S, Morstatter F, Trevino RP, Tang J, Liu H (2017) Feature selection: A data perspective. *ACM Computing Surveys (CSUR)* 50(6):94
- [19] McGovern A, Milivojevic S (2016) Social media and crime: the good, the bad and the ugly. URL <https://theconversation.com/social-media-and-crime-the-good-the-bad-and-the-ugly-66397>
- [20] O’Dea B, Wan S, Batterham PJ, Calear AL, Paris C, Christensen H (2015) Detecting suicidality on twitter. *Internet Interventions* 2(2):183–188
- [21] Picek S, Heuser A, Jović A, Bhasin S, Regazzoni F (2018) The curse of class imbalance and conflicting metrics with machine learning for side-channel evaluations. *IACR Transactions on Cryptographic Hardware and Embedded Systems* 2019(1):209–237
- [22] Rehman A, Javed K, Babri HA, Saeed M (2015) Relative discrimination criterion—a novel feature ranking method for text data. *Expert Systems with Applications* 42(7):3670–3681
- [23] Sagi O, Rokach L (2018) Ensemble learning: A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 8(4):e1249
- [24] Schmidt P (2018) Human rights online. URL [http://www.inach.net/wp-content/uploads/2018/05/INACH\\_HumanRightsOnline.pdf](http://www.inach.net/wp-content/uploads/2018/05/INACH_HumanRightsOnline.pdf)
- [25] Schütze H, Manning CD, Raghavan P (2008) Introduction to information retrieval, vol 39. Cambridge University Press
- [26] Steele SM (1970) Program evaluation—a broader definition. *Journal of Extension* 8(2):5–17
- [27] Sueki H (2015) The association of suicide-related Twitter use with suicidal behaviour: A cross-sectional study of young internet users in Japan. *Journal of Affective Disorders* 170(September 2014):155–160
- [28] Tang B, He H, Baggenstoss PM, Kay S (2016) A bayesian classification approach using class-specific features for text categorization. *IEEE Transactions on Knowledge and Data Engineering* 28(6):1602–1606
- [29] Uğuz H (2011) A two-stage feature selection method for text categorization by using information gain, principal component analysis and genetic algorithm. *Knowledge-Based Systems* 24(7):1024–1032
- [30] Won HH, Myung W, Song GY, Lee WH, Kim JW, Carroll BJ, Kim DK (2013) Predicting National Suicide Numbers with Social Media Data. *PLoS ONE* 8(4), DOI 10.1371/journal.pone.0061809
- [31] Yao J, Zhang J, Wang L (2018) A financial statement fraud detection model based on hybrid data mining methods. In: 2018 International Conference on Artificial Intelligence and Big Data (ICAIBD), IEEE, pp 57–61