CrossMark

# Suicide and self-harm prediction based on social media data using machine learning algorithms

Abdulrazak Yahya Saleh [a,1,*], Fadzlyn Nasrini Binti Mostapa [a,2]

[a] Faculty of Cognitive Sciences and Human Development, Universiti Malaysia Sarawak, Malaysia
[1] abdulrazakalhababi@gmail.com; [2] 69662@siswa.unimas.my
* Corresponding Author

ABSTRACT

Online social networking (SN) data is a context and time rich data stream that has showed potential for predicting suicidal ideation and behaviour. Despite the obvious benefits of this digital media, predictive modelling of acute suicidal ideation (SI) remains underdeveloped at now. In combined with robust machine learning algorithms, social networking data may provide a potential path ahead. Researchers applied a machine learning models to a previously published Instagram dataset of youths. Using predictors that reflect language use and activity inside this social networking, researchers compared the performance of the out-of-sample, cross-validated model to that of earlier efforts and used a model explanation to further investigate relative predictor relevance and subject-level phenomenology. The application of ensemble learning approaches to SN data for the prediction of acute SI may reduce the complications and modelling issues associated with acute SI at these time scales. Future research is required on bigger, more diversified populations to refine digital biomarkers and assess their external validity with more rigor.

## 1. Introduction

The term "social media" refers to an online platform that enables social interaction, such as Facebook, Twitter, YouTube, and Instagram [1]. [2] assert that social media is fundamentally defined by three critical concepts which are cognition, communication, and cooperation. These concepts endow social media with a variety of various forms of sociality, including information, facts, and knowledge, activities, relationships, communities, and partnerships [2]. The numerous forms of sociality that social media supports allowing the platform to serve a variety of critical roles for its users, notably communication, relationship creation and maintenance, and thus a platform for providing knowledge [3]. These features are highly valued by social media users, particularly young adults. Youth are classified as ardent social media users in a variety of research conducted globally [3].

Malaysian youth are likewise reported to be avid social media users. While there is no debate about the use of social media [4], it is also necessary to exercise caution regarding the medium's hazards. Spam hoaxes, cyberbullying, online harassment, and sexting are just a few of the dangers [1]. Youth who use social media risk privacy assaults and depression [1]. These dangers can be avoided if youth possess an

acceptable level of social media competency. Presently, social media platforms such as Facebook and Instagram have become the primary sources of information for assisting people of modern society in adjusting to their new lifestyle. Indeed, the global population of social media users is predicted to reach 3.02 billion by 2021, and one of the important uses of social media may be to encourage healthy lifestyles and to enhance people's management of their health status [5]. In Malaysia specifically, a 2018 poll performed by the Malaysian Communications and Multimedia Commission (MCMC) discovered that many Malaysian internet users, particularly young users, shared content online via social media (61.8 percent). According to the survey, 97.3 percent of Malaysians use Facebook, making it the country's most popular social networking site, followed by Instagram, which has a user base of 57.0 percent [6]. And though the overall number of Facebook users in Malaysia is greater than the number of Instagram users, a January 2019 social media users' trend indicated that Instagram was more popular among Malaysian young users aged 18 to 24 years old, with 31.9 percent compared to 24.4 percent Facebook young users [7].

With machine learning and deep learning algorithms, this study will collect, and preprocess collected data sets from the social media. There will be a classifier that is a few of machine learning algorithms. The more variables a classifier can learn from, the greater, and the random forest evaluates the relative value of each characteristic in predicting the probability of any particular post displaying suicidal and self-harm purpose.

## 2. Method

### 2.1. Data Preparation

The data preparation phase is vital to the development of this study. This phase assembles datasets on suicide and self-harm on social media to meet specific needs. This phase is divided into two sections: data collection and data processing. The Fig. 1 below provides context for this period.
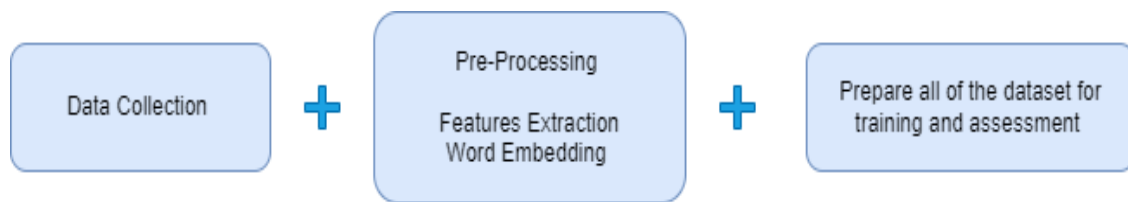


**Fig. 1.** Data Preparation

### 2.2. Datasets

To identify suicidal ideation, this study will develop classification models using an Instagram social media dataset in which users can express themselves through text, posts, links, or voting mechanism postings. They converse via comment threads related to each post [8]. [9] created the dataset for this investigation, which comprises of a list of suicide-indicative and non-suicidal posts. To protect users' privacy, personal information is substituted for a unique ID. Due to the users' proclivity for engaging in a variety of sub-Instagram's, each group is comprised of an equal number of messages generated from a variety of themes. This dataset includes only numeric values. As shown in Fig. 2, The table contains 59 columns and 28 rows, with each row representing a distinct sort of information.

**Fig. 2.** The summary of the dataset

## 2.3. Data Pre-Processing

Pre-processing raw posts before learning word embedding is called pre-processing, and it involves removing redundant features from an input text. Filters applied to Instagram posts transform raw data into a format that can be understood by learning models. NLTK [10] is used in our study to pre-process the dataset before it is used for training. We begin by concatenating the titles and content of the posts. The original dataset is cleansed of redundant sentences. To tokenize the Instagram posts, we first filter and convert the raw data into a manageable format. A single whitespace is then used to replace all URL addresses as well as contractions and redundant white spaces. All newline symbols and punctuation marks have been omitted because they could cause erratic results if left in place. Posts can be saved in lowercase and as separate text files by doing this. Our lemmatization process is designed to ensure that word endings will not be thrown away, resulting in meaningless word pieces like stemming. It is better to turn them into dictionary-related word lemmas. Word embedding can begin now that the data has been cleaned.

This study will utilize NumPy, pandas, matplotlib, and sklearn as its basic libraries. This study will utilize the Pandas library to import and analyze the dataset. The NumPy library facilitates the study's use of arrays. Before training, the data must be transformed to a NumPy array. The Matplotlib software will facilitate data visualization for this project. This study will utilize the Pandas library to import and analyze these datasets. The NumPy library facilitates the study's use of arrays. Before training, the data must be transformed to a NumPy array. The Matplotlib library will aid in the visualization of data.

The CSV-formatted dataset will be imported into this study using the Pandas package. CSV is an abbreviation for Comma-Separated Values. Fig. 3. is a view of the statistical info in the dataset.

```
# viewing statistical info about dataset
dataset.describe()
```

| | UserID | Demographics | Gender | Age | School-type | Acute_Suicidal Thoughts | Lifetime_Suicidal_Thoughts | Lifetime_Suicide_Attempt | Lifetime_NSSI | LIWC_Int |
|---|---|---|---|---|---|---|---|---|---|---|
| count | 59.000000 | 0.0 | 50.00000 | 50.000000 | 50.000000 | 55.000000 | 55.000000 | 55.000000 | 55.0 | |
| mean | 30.000000 | NaN | 0.86000 | 16.740000 | 1.320000 | 0.454545 | 0.945455 | 0.509091 | 1.0 | |
| std | 17.175564 | NaN | 0.35051 | 1.225744 | 0.767716 | 0.502519 | 0.229184 | 0.504525 | 0.0 | |
| min | 1.000000 | NaN | 1.00000 | 16.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 1.0 | |
| 25% | 15.500000 | NaN | 1.00000 | 16.000000 | 1.000000 | 0.000000 | 1.000000 | 0.000000 | 1.0 | |
| 50% | 30.000000 | NaN | 1.00000 | 16.000000 | 1.000000 | 0.000000 | 1.000000 | 1.000000 | 1.0 | |
| 75% | 44.500000 | NaN | 1.00000 | 17.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.0 | |
| max | 59.000000 | NaN | 1.00000 | 22.000000 | 4.000000 | 1.000000 | 1.000000 | 1.000000 | 1.0 | |

8 rows × 28 columns

**Fig. 3.** Viewing statistical info of the dataset

The dataset may occasionally include duplicate values. These duplicate values are unnecessary, so they must be discarded. There could be missing values in the dataset. A dataset with missing values cannot be used to train the model by the Python. Therefore, it must be determined if the dataset contains missing values. It is possible to replace missing values in numerical data with the mean, mode, or median of the column containing the missing value. As a result, it can preserve certain data required for the model. The mean is generally favored. It may also replace missing values with the value immediately preceding or after it in the same column. Fig. 4. Shows the handling of missing values with mean.

```
#using mean
dataset['Mean_likes_Instagram'].fillna(int(dataset['Mean_likes_Instagram'].mean()), inplace=True)

# checking the number of missing data
dataset.isnull().sum()

UserID                              0
Demographics                        0
Gender                              0
Age                                 0
School-type                         0
Acute_Suicidal Thoughts             0
Lifetime_Suicidal_Thoughts          0
Lifetime_Suicide_Attempt            0
Lifetime_NSSI                       0
LIWC_Interviews                     0
Wordcount_Interview                 0
Pronoun_Interviews                  0
Affect_Interviews                   0
Negativeemotion_Interviews          0
Cognitivemechanism_Interviews       0
FRE_German_Interviews               0
LIWC_Captions                       0
Wordcount_Captions                  0
Pronoun_Captions                    0
Affect Captions                     0
```

**Fig. 4.** Handling missing values with mean

The splitting of the dataset into training and test sets is an additional essential step in data preprocessing. A portion of the dataset will be used to train the model in this study. The remaining portion of the dataset will be utilized to test the model's performance on previously unseen data. Feature scaling places all data within the same range and scale. We do not want one characteristic of the dataset to overshadow another. In this study, feature scaling will be conducted using standardization. Standardization places all the values between -3 and 3. The StandardScaler class from the preprocessing module of the sklearn package will be utilized in this research.

### 2.4. Proposed Models

To identify the existence of suicide and self-harm in Instagram, this work will evaluate the strengths of some machine learning algorithms like: SVM, RF, KNN, DT, LR, GB and Naïve Bayes architectures and use these models for the classification of the chosen text data. The suggested models extract the features of the input text sentences and improve the outcomes of the classification accuracy.

### 2.5. Word Embedding Layer

Word embedding is a collection of language modelling and feature learning techniques used in natural language processing. It is an input layer of the selected SVM, RF, KNN, DT, LR, GB and Naïve Bayes models that converts the words to a vector representation in real-valued space. When words from the vocabulary are embedded, they tend to map onto a particular vector space of real numbers in a low-dimensional space [11]. The models are fundamentally based on unsupervised training of distributed representations for supervised problem solving [12]. In this part, we will use Word2vec [11], a shallow model composed of two neural layers trained to rebuild a word context or current words from their surrounding window of word vectors defined by embedding layer index numbers that are transformed

into the d-dimensional embedding vector Xt Rd via pre-training Word2Vec [11]. As noted in equation, d is the dimension of the word vector with an input text. At this stage, the tth word in the text is denoted by Xt Rd, where d is the word embedding vector and T is the text's length.

### 2.6. Baseline

This study is undertaken by a performance comparison of the suggested learning model against the baseline models to provide a fair comparative analysis to other competing models. Text characteristics (TF-IDF, Bag of Words, Statistical Features) are retrieved and fed into two classic machine learning algorithms (Support Vector Machine and Random Forest) using Word2vec embedding techniques. We use Scikit-learn to implement machine learning techniques [13]. The Support Vector Machine (SVM) is a supervised learning model that examines data and detects patterns for categorization [14]. It's frequently used in text classification [15], and it's shown to perform well in mental health tasks [16]. By establishing a hyperplane in a high-dimensional space, the SVM technique will be used to solve problems that are linearly and non-linearly separable in a lower space. The SVM approach, which has been shown to function well with succinct and categorical data, will be used to assess the effectiveness of word embeddings. Random Forest (RF) is an ensemble approach in which numerous weak classifiers are combined into a single strong classifier [17]. For binary class classification issues, RF is commonly employed [17]. The frequency (TF-IDF) approach is frequently utilized in the fields of information retrieval and text mining. It determines the frequency with which a word appears in a text; it picks significant words and eliminates low-importance terms for further text analysis [18]. Bag of Words (BOW) is an algorithm that lists the words in a text together with their word counts. The count of each word is utilized to construct a feature vector that will be used to summarize the material further [19]. The number of tokens, words, phrases, and their length are retrieved from the postings using statistical characteristics [20]. The CNN baseline network structure used for text categorization is comparable to the CNN model developed.

### 2.7. Evaluation

This study will employ assessment measures such the accuracy of estimates (Acc.) Equation [21] and the F-score (F1) Equation (15), which consists of precision (P) and recall (R) to evaluate the baseline using our proposed deep learning classification approach. It is based on a confusion matrix that incorporates data from each test sample prediction result. Precision estimates the number of positively identified samples; recall approximates the proportion of correctly identified positive samples; F1 Equation [22] score is a harmonic average of precision and recall; precision estimates the number of positively identified samples; recall approximates the proportion of correctly identified positive samples. The greater the F1 score, the closer the two values are. The number of true positive predictions (TP), true negative predictions (TN), false-positive predictions (FP), and false-negative predictions (FN) are among the assessment measures [23].

## 3. Results and Discussion

According to the data presented in Fig. 5, out of the seven classifiers that were tested, SVM achieved the highest accuracy rate at 97.47 percent, whereas Decision Tree had the lowest accuracy rate of 85.01 percent. It is worth noting that the range between the highest and lowest accuracy rates is relatively significant, which suggests that the choice of classifier could have a significant impact on the performance of the overall system.
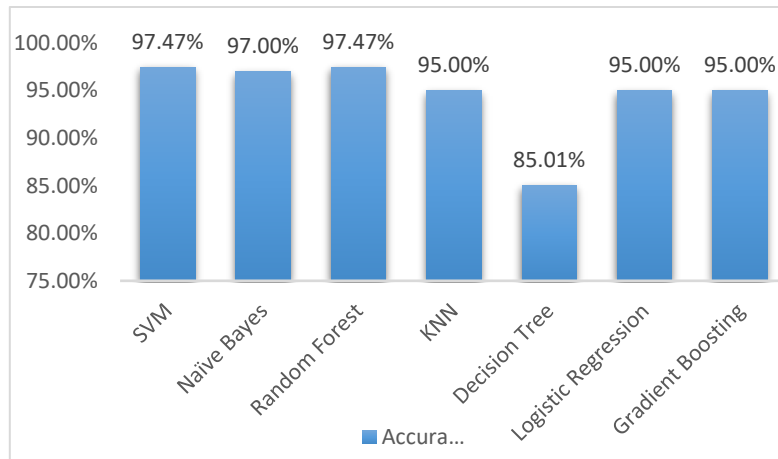
**Fig. 5.** Values of Accuracy

Based on the results presented in Fig. 6, it has been demonstrated that Random Forest (RF) outperformed all other classifiers in terms of precision, with a value of 97.54 percent. On the other hand, Decision Tree (DT) achieved the lowest precision value of 85.17 percent among all the tested classifiers. Precision is a crucial performance metric in machine learning, as it measures the percentage of correctly predicted positive instances out of all the instances that were predicted as positive. Therefore, a higher precision value indicates that the classifier was able to correctly identify positive instances more accurately, and with fewer false positives. The significant difference between the precision values of RF and DT suggests that the former may be a more suitable classifier for the given problem than the latter. However, it is important to note that other performance metrics, such as recall and F1 score, should also be considered when selecting a classifier, as they provide a more comprehensive evaluation of the model's overall performance.
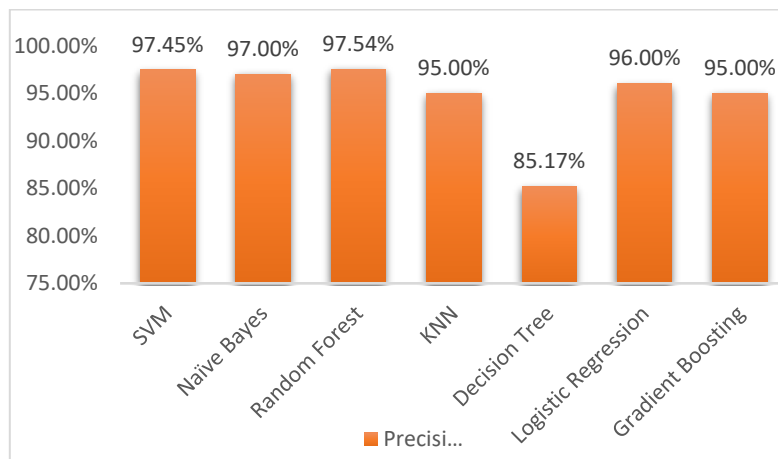


**Fig. 6.** Values of Precision

Fig.7 displays the results of the evaluation of several classifiers, and it has been observed that Naive Bayes (NB) achieved the highest value of recall at 98.00 percent, while the lowest value of recall was obtained by Decision Tree (DT) at 84.85 percent. Recall is a performance metric that measures the ability of a classifier to correctly identify all positive instances in a dataset, also known as sensitivity. A high recall value indicates that the classifier is more capable of correctly identifying positive instances, and hence, is less likely to miss any relevant information. The significant difference between the recall values obtained by NB and DT indicates that NB is a more suitable classifier for the given problem than

DT, in terms of identifying positive instances. However, it is important to consider other performance metrics, such as precision and F1 score, as they provide a more comprehensive evaluation of the classifier's overall performance. Therefore, based on the findings presented in Fig. 7, it can be concluded that NB is the most appropriate classifier in terms of recall, while DT is the least effective classifier for the given problem.
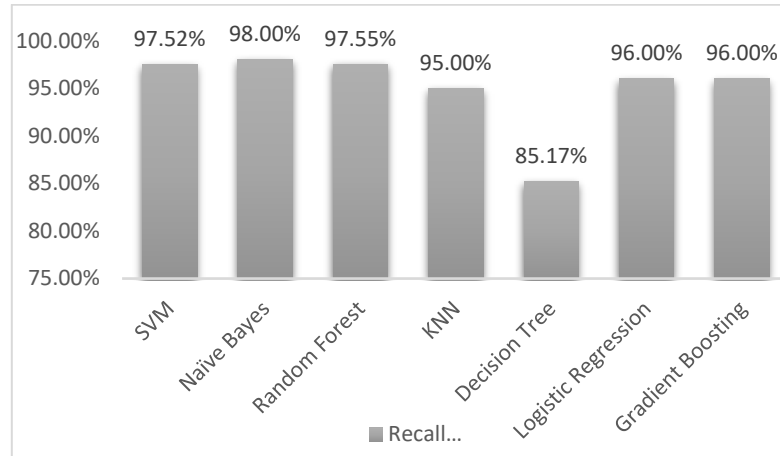


**Fig. 7.** Values of Recall

Fig. 8 presents the F1 Score values achieved by several classifiers, and it has been observed that Random Forest (RF) achieved the highest F1 Score at 97.53 percent, while the lowest F1 Score was obtained by Decision Tree (DT) at 84.92 percent. F1 Score is a performance metric that provides a balanced evaluation of the precision and recall of a classifier. It considers both false positives and false negatives, and a higher F1 Score indicates that the classifier has achieved a better balance between precision and recall. The significant difference between the F1 Score values obtained by RF and DT suggests that RF is a more suitable classifier for the given problem, as it achieves a better balance between precision and recall. However, it is essential to consider other performance metrics, such as accuracy, precision, and recall, in addition to the F1 Score, to obtain a more comprehensive evaluation of the classifier's overall performance. Therefore, based on the results presented in Figure 8, it can be concluded that RF is the most effective classifier in terms of the F1 Score, while DT is the least effective classifier for the given problem.
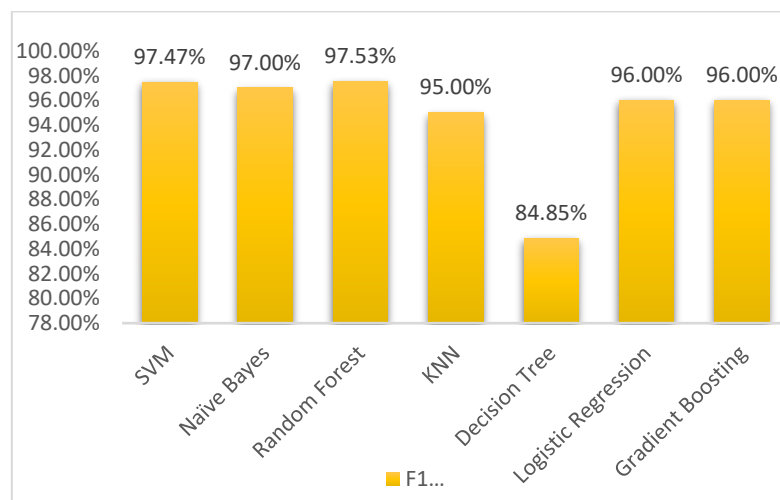


**Fig. 8.** Values of F1 Score

## 4. Conclusion

This research aimed to evaluate the effectiveness of machine learning models in predicting suicide and self-harm data on social media networks. The study employed seven different algorithms to assess the performance of the proposed approach in terms of classification accuracy, precision, recall, and F1 score. Furthermore, the models utilized independently generated features to enhance the prediction of suicidal and self-harming behaviors.

The results of the experiments demonstrated that the proposed approach met the requirements for the suicide and self-harm dataset, with satisfactory performance in terms of accuracy, precision, recall, and F1 score. However, the study recognizes that the investigation into the prediction model based on machine learning is still in its early stages, with considerable scope for improvement. Moving forward, there is potential for the proposed model to be refined and optimized to enable more precise forecasting of suicidal and self-harming behaviors. This includes exploring the various ways in which these behaviors can manifest and incorporating more diverse datasets. Additionally, further research can be conducted to investigate optimization methodologies that could enhance the performance of the classification system.

In conclusion, this study provides a valuable contribution to the field of machine learning-based prediction models for suicidal and self-harming behaviors on social media networks. The results suggest that the proposed approach shows promise and can be further improved to enhance the accuracy and effectiveness of suicide and self-harm prediction systems.

## Acknowledgment

## References

[1]     G. S. O'Keeffe and K. Clarke-Pearson, "The Impact of Social Media on Children, Adolescents, and Families," *Pediatrics*, vol. 127, no. 4, pp. 800–804, Apr. 2011, doi: https://doi.org/10.1542/peds.2011-0054.

[2]     D. Trottier and C. Fuchs, *Social Media, Politics and the State*. Routledge, p. 264, 2014, doi: 10.4324/9781315764832.

[3]     N. Ismail, "Young People ' s Use of New Media through Communities of Practice," *Malaysian J. Commun.*, no. October 2014, 2014, [Online]. Available at: https://www.researchgate.net/publication/280721106_Young_people's_use_of_new_media_through_communities_of_practice.

[4]     A. Kaizerman-Dinerman, N. Josman, and D. Roe, "The use of Cognitive Strategies among People with Schizophrenia: A Randomized Comparative Study," *Open J. Occup. Ther.*, vol. 7, no. 3, pp. 1–12, Jul. 2019, doi: 10.15453/2168-6408.1621.

[5]     M. Soheylizad and B. Moeini, "Social Media: An Opportunity for Developing Countries to Change Healthy Behaviors," *Heal. Educ. Heal. Promot.*, vol. 7, no. 2, pp. 57–58, Apr. 2019, doi: 10.29252/HEHP.7.2.57.

[6]  Malaysian Communications and Multimedia Commission, "Internet users survey 2018: Statistical brief number twenty-three," *Malaysian Commun. Multimed. Comm.*, pp. 1–39, 2018, [Online]. Available: https://www.mcmc.gov.my/skmmgovmy/media/General/pdf/Internet-Users-Survey-2018.pdf.

[7]  "Instagram Policy Changes on Self-Harm Related Content," 2019. [Online]. Available at: https://about.instagram.com/blog/announcements/supporting-and-protecting-vulnerable-people-on-instagram.

[8]  K. Kumar *et al.*, "Precursor feeding studies and molecular characterization of geraniol synthase establish the limiting role of geraniol in monoterpene indole alkaloid biosynthesis in Catharanthus roseus leaves," *Plant Sci.*, vol. 239, pp. 56–66, Oct. 2015, doi: 10.1016/j.plantsci.2015.07.007.

[9]  S. Ji, C. P. Yu, S. Fung, S. Pan, and G. Long, "Supervised Learning for Suicidal Ideation Detection in Online User Content," *Complexity*, vol. 2018, pp. 1–10, Sep. 2018, doi: 10.1155/2018/6157249.

[10] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit.* p. 463, 2009. [Online]. Available at: https://books.google.co.id/books?hl=en&lr=&id=KGIbfiiP1i4C&oi=fnd&pg=PR5&dq=Natural+Language+Processing+with+Python:+Analyzing+Text+with+the+Natural+Language+Toolkit&ots=Y4HlE1JDI2&sig=RzwFQTMKWvwBLjls5vS8aOa5wcU&redir_esc=y#v=onepage&q=Natural Language P.

[11] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient Estimation of Word Representations in Vector Space," in *1st International Conference on Learning Representations, ICLR 2013 - Workshop Track Proceedings*, Jan. 2013, pp. 1–12. [Online]. Available at: https://arxiv.org/abs/1301.3781.

[12] D. Sisk, "Simulation: Learning By Doing Revisited," *Gift. Child Q.*, vol. 19, no. 2, pp. 175–180, Jun. 1975, doi: 10.1177/001698627501900225.

[13] F. Pedregosa *et al.*, "Scikit-learn: Machine Learning in Python," *Environ. Health Perspect.*, vol. 127, no. 9, pp. 2825–2830, Sep. 2019, doi: 10.1289/EHP4713.

[14] C. D. Manning, P. Raghavan, and H. Schütze, "Introduction to Information Retrieval," *Cambridge University Press.*, p. 544, 2008. https://www-nlp.stanford.edu/IR-book/.

[15] T. Joachims, "Text categorization with Support Vector Machines: Learning with many relevant features," *Springer, Berlin, Heidelberg*, 1998, pp. 137–142, doi: 10.1007/BFb0026683.

[16] M. De Choudhury, "Role of social media in tackling challenges in mental health," in *Proceedings of the 2nd international workshop on Socially-aware multimedia*, Oct. 2013, pp. 49–52, doi: 10.1145/2509916.2509921.

[17] Y. Freund and R. E. Schapire, "A brief introduction to boosting," *J. Japanese Soc. Artif. Intell.*, vol. 2, no. 5, pp. 1401–1406, 1999, [Online]. Available at: https://www.yorku.ca/gisweb/eats4400/boost.pdf.

[18] M. Ikonomakis, S. Kotsiantis, and V. Tampakas, "Text Feedback Classification using Machine Learning Techniques," *Proc. 2nd Int. Conf. Edge Comput. Appl. ICECAA 2023*, vol. 4, no. 8, pp. 933–939, 2023, doi: 10.1109/ICECAA58104.2023.10212398.

[19] A. Fiori, "Innovative Document Summarization Techniques," *IGI Global*, pp. 1-341, 2014, doi: 10.4018/978-1-4666-5019-0.

[20] E. Okhapkina, V. Okhapkin, and O. Kazarin, "Adaptation of Information Retrieval Methods for Identifying of Destructive Informational Influence in Social Networks," in *2017 31st International Conference on Advanced Information Networking and Applications Workshops (WAINA)*, Mar. 2017, pp. 87–92, doi: 10.1109/WAINA.2017.116.

[21] T. Mikolov, M. Karafiˊa, L. Burget, J. "Honza" ˇCernockˊ, and S. Khudanpur, "Recurrent neural network based language model," *Proc. Annu. Meet. Assoc. Comput. Linguist.*, no. September, pp. 8093–8104, 2010, doi: 10.18653/v1/2020.acl-main.720.

[22] R. C. Hsiung, "A Suicide in an Online Mental Health Support Group: Reactions of the Group Members, Administrative Responses, and Recommendations," vol. 10, no. 4, pp. 495–500, Aug. 2007, doi: 10.1089/CPB.2007.9999.

[23]   T. Basu and C. A. Murthy, "Effective Text Classification by a Supervised Feature Selection Approach," in *2012 IEEE 12th International Conference on Data Mining Workshops*, Dec. 2012, pp. 918–925, doi: 10.1109/ICDMW.2012.45.