

Utilizing Artificial Intelligence for the Identification of Students with Depression and Anxiety through Social Media Analysis

Taneesh Sebastian – 10th grade, Plymouth Christian Academy, Michigan.

Mentored by: Abdulla Kerimov, PhD in Geophysics from Stanford, Applied Scientist @ BP

Abstract

In response to the escalating concerns about the mental well-being of students, this research represents a significant stride in utilizing artificial intelligence (AI) to discern individuals experiencing depression and anxiety through their social media comments. The study employs a comprehensive suite of regression methods, including Logistic Regression, Decision Trees, Support Vector Classifier (SVC), Random Forest, and the Ridge model. The dataset, meticulously curated with 6982 comments sourced from Kaggle, undergoes processing, incorporating a conscientious split into training and testing sets for robust evaluation. The AI models exhibit remarkable performance metrics, with precision reaching an impressive 99%, F1 scores achieving 99% for normal and 95% for depressive comments, and accuracy and weighted averages standing at 99%. The study serves its immediate purpose of identifying students on social media and hints at broader applications across diverse contexts. Beyond specific demographics, the research underscores the potential of AI in addressing mental health challenges, offering insights that extend far beyond the confines of student populations, thereby contributing to the broader discourse of AI's impact on mental health.

Introduction

The adoption of multiple machine learning techniques in Artificial Intelligence research enhanced the ability to identify and address mental health concerns among students through analysis of their online interactions, using social media. The application of machine learning analysis allows for a systematic examination of social media content, providing an opportunity for early detection and intervention. By utilizing a combination of machine learning algorithms, one can analyze large volumes of data from social media platforms to identify patterns and indicators of depression.

The focus of this study was identifying depressive and normal social media comments using the machine learning model. The primary objective was to accurately decipher between these two categories to gain insights into the emotional tone of online interactions. To achieve this, the dataset consisting of social media comments was utilized for machine learning training and testing. This research has significant implications for understanding mental

health trends and detecting potential signs of depression through online platforms.

Dataset

“Students anxiety and depression dataset”, from Kaggle, was used for this project. The dataset comprises 6982 comments from students who have access to social media, labeled as depressive or normal (Table 1). Data underwent preprocessing using the Bag of Words (BoW) method, along with a Count Vectorizer. None of the features were excluded from the dataset as they were all relevant to this project, except for the unnecessary punctuation and stop-words that were excluded while processing the data.

I'm confused, I'm not feeling good lately. Every time I want to sleep, I always feel restless	1
Have you ever felt nervous but didn't know why?	1
I'm already in a bad mood and then my heart seems to be beating really fast... I'm really nervous. Is there something wrong???	1
Feeling happy today	0
This morning I smell really good	0
Dudeeee! What's up?	0

Table 1. Sample of dataset; The label “1.0” is for depressive comments, while “0.0” is for normal comments.

Because of imbalanced distribution of the depressive and normal comments in the dataset, stratified train-test split has been implemented in this study. It splits the dataset into train and test sets in a way that preserves the same proportions of examples in each class as observed in the original dataset. 33% of the data was split as testing data, while the remaining 67% was classified as training data.

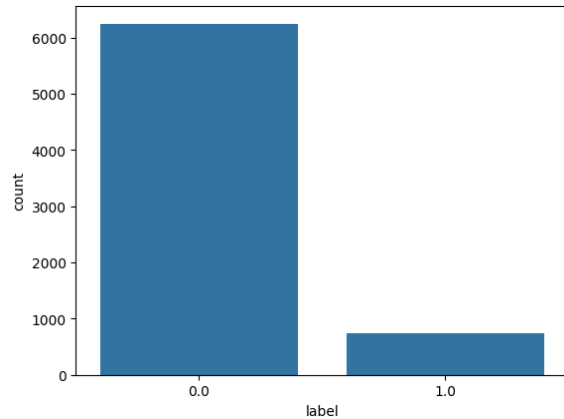


Figure 1. Imbalanced count distribution of labels. The label “1.0” is for depressive comments, while “0.0” is for normal comments.

Machine Learning Models

A combination of regression algorithms was used to analyze large volumes of data from social media platforms to identify patterns and indicators of depression. These techniques enable a more comprehensive understanding of the complex relationship between language usage in online communication and mental health. The application of multiple machine learning techniques provides a robust framework for training models that can accurately classify and detect signs of depression within social media comments made by students. This experimental approach of using multiple methods contributes to the development of innovative solutions that can support early intervention and mental health support services for the demographic. The intricacies of regression methods, including Ridge, Logistic, Decision Tree, SVC, and Random Forest, are meticulously used. The models I used each have their own limitations, which is why I used multiple. Hyperparameter tuning is explored as a means to enhance model performance. I used hyperparameter tuning on the Random Forest model because it was the best to use it on. In refining a random forest model through hyperparameter tuning, we adjust parameters such as `min_samples_split`, which determines the minimum number of samples required to split a node, `n_estimators`, which controls the number of trees in the forest, and `max_features`, governing the maximum number of features considered for splitting. Tuning involves exploring ranges and selecting from feature subset sizes

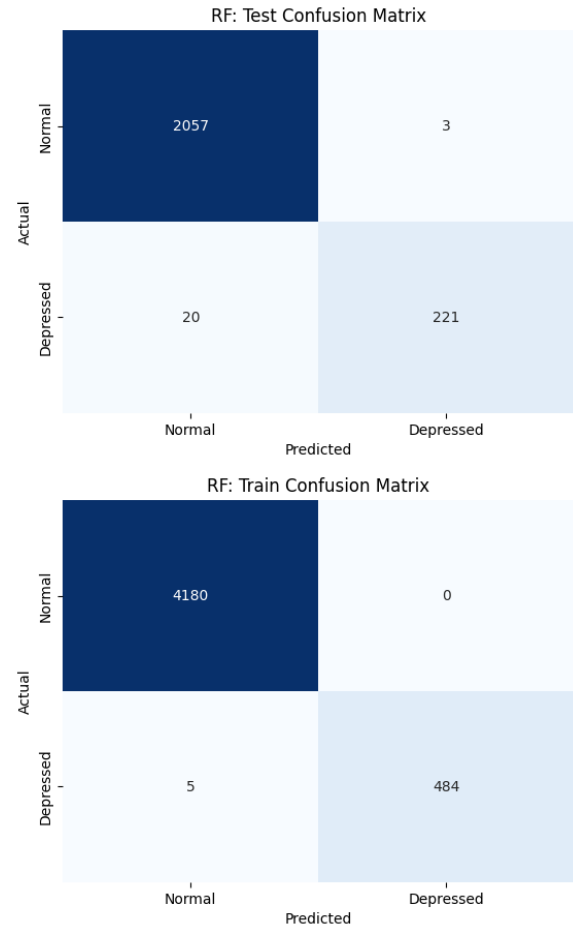


Figure 2. Confusion Matrices of train and test random forest models.

Results and Discussion

Through extensive research and analysis, the implementation of the model showcased remarkable results in accurately identifying the sensitivity of comments made by students on social media platforms. When classifying comments as either depressive or normal, the AI model demonstrates a precision rate of 99%. This level of accuracy highlights the potential for AI to effectively detect and understand the emotional state expressed within social media interactions. While these results are impressive, it is important to acknowledge that there may be certain sources of error in the identification process. One possible limitation is the AI's inability to consistently recognize comments as either depressive or normal due to various factors such as linguistic nuances, context, or unique expressions used by individuals. Some possible modifications could be processing different datasets with more content that could be used to expand the knowledge of the AI.

Conclusions

In conclusion, this research represents a significant advancement in leveraging artificial intelligence (AI) to discern signs of depression and anxiety through the analysis of social media comments. The comprehensive suite of regression methods, dataset curation, and thorough evaluation demonstrate the effectiveness of AI in identifying depressive and normal comments. The study not only serves its immediate purpose in the context of student mental well-being but also hints at broader applications across diverse populations. Beyond specific demographics, the research underscores the potential of AI to address mental health challenges, contributing valuable insights to the broader discourse on AI's impact in this domain. While the impressive results showcase the capability of AI in detecting emotional states through online interactions, it is essential to acknowledge potential limitations and consider further refinements to enhance the model's accuracy and robustness in recognizing nuances within language usage on social media platforms. It is important to explore other data preprocessing techniques, such as word embeddings, and their effect on the accuracy of the model. Also, oversampling techniques that generate synthetic samples from the minority class can be implemented to handle the imbalanced dataset.

	Precision	Recall	F1
0.0	99%	100%	99%
1.0	99%	92%	95%

Acknowledgments

I would like to thank Abdulla Kerimov, PhD in Geophysics from Stanford, Applied Scientist @ BP and Carlson Hoo, Master's in Data Science, Postgraduate Student at UKM for their continuous support and assistance.

References

- <https://carlson-hoo.medium.com/what-type-of-social-media-users-are-more-likely-to-have-social-anxiety-disorder-65194323f8e1>
- <https://www.kaggle.com/datasets/sahasourav17/students-anxiety-and-depression-dataset/data>