

# Flamingos\_python@LT-EDI: An Ensemble Model to Detect Severity of Depression

Abirami P S, Amritha S, Pavithra M, C.Jerin Mahiba

Meenakshi Sundararajan Engineering College, Chennai

{abiabhi2712, pavithrameganathan15, amrithasenthil2001, jerinmahibha}@gmail.com

## Abstract

The prevalence of depression is increasing globally, and there is a need for effective screening and detection tools. Social media platforms offer a rich source of data for mental health research. The paper aims to detect the signs of depression of a person from their social media postings wherein people share their feelings and emotions. The task is to create a system that, given social media posts in English, should classify the level of depression as ‘not depressed’, ‘moderately depressed’ or ‘severely depressed’. The paper presents the solution for the Shared Task on Detecting Signs of Depression from Social Media Text at LT-EDI@RANLP 2023. The proposed system aims to develop a machine learning model using machine learning algorithms like SVM, Random forest and Naive Bayes to detect signs of depression from social media text. The model is trained on a dataset of social media posts to detect the level of depression of the individuals as ‘not depressed’, ‘moderately depressed’ or ‘severely depressed’. The dataset is pre-processed to remove duplicates and irrelevant features, and then, feature engineering techniques is used to extract meaningful features from the text data. The model is trained on these features to classify the text into the three categories. The performance of the model is evaluated using metrics such as accuracy, precision, recall, and F1-score. The ensemble model is used to combine these algorithms which gives accuracy of 47.7% and the F1 score is 0.262. The results of the proposed approach could potentially aid in the early detection and prevention of depression for individuals who may be at risk.

## 1 Introduction

Depression is a mood disorder that causes a persistent feeling of sadness and loss of interest. It is also called major depressive disorder or clinical depression, it affects how one feel, think and behave and can lead to a variety of emotional

and physical problems. Detecting depression is important since it has to be observed and treated at an early stage to avoid severe consequences. It aims to detect the signs of depression of a person from their social media postings wherein people share their feelings and emotions. Social media platforms have emerged as a valuable source of data for mental health research. They provide an opportunity to study the language and behaviour patterns of individuals, which may reveal early signs of depression. Given social media postings in English, the system should classify the signs of depression into three labels namely “not depressed”, “moderately depressed”, and “severely depressed”. This project aims to develop a machine learning model using Support Vector Machine (SVM), Random forest and Naive Bayes algorithm to detect signs of depression from social media text.

## 2 Related Works

A depression recognition method for college students [3] proposed by Yan Ding et al., had used a deep integrated support vector machine (DISVM) algorithm to classify the input data, and finally realize the recognition of depression. Wolohan et al., (2018) created a dataset based on Reddit posts in which users were assigned to one group: depressed or control. A transformers approach to detect depression in social media [6] by Malviya et al., had analyzed the posts using the linguistic inquiry and word count tool (LIWC). Findings of the Shared Task on Detecting Signs of Depression from Social Media [7] had used a variety of technologies from traditional machine learning algorithms to deep learning models. ScubeMSEC@LT-EDI-ACL2022: Detection of Depression using Transformer Models [8] by S, Sivamanikandan et al., used different transformer models like DistilBERT, RoBERTa and ALBERT to detect depression which had achieved a Macro F1 score of 0.337, 0.457 and

Text Data	Label
Like no matter how much sleep I get always fatigued.	not depression
All I wanna do right now is crawl out of myself. Does that make sense?	moderate
A few anxiety attacks and I'm ready to go	severe

Table 1: Example Instances

0.387 respectively. Early Detection of Depression from Social Media Data Using Machine Learning Algorithms [4] by G. Geetha et al., had used machine learning algorithms like Support Vector Machine (svm), Logistic Regression, Random Forest, Bayes Theorem for the early detection of depression. A machine learning based depression analysis and suicidal ideation detection system using questionnaires and Twitter [5] by Jain et al., had proposed a system for predicting the suicidal acts based on the level of depression using XGBoost classifier. Depression detection by analyzing social media posts of user [2] by Nafiz Al Asad, et al., had presented a structural model that identified users' depression level from their social media posts. This system had used SVM classifier and Naïve Bayes classifier. A machine learning approach to detect depression and anxiety using supervised learning [1] by A. Ahmed et al., had aimed to apply natural language processing on Twitter feeds for conducting emotion analysis focusing on depression. Detection of depression related posts in Reddit social media forum [12] by Tadesse et al. had implemented class prediction using support vector machine and Naive-Bayes classifier. Sentiment analysis from depression related user generated contents in social media texts by Ananna Saha et al. [9] had found the usage and effectiveness of the five different types of AI algorithms: Convolutional Neural Network, Support Vector Machine, Linear Discriminant Analysis, K Nearest Neighbor Classifier and Linear Regression on two datasets of anxiety and depression. Detection of major depressive disorder using signal processing and machine learning approaches [10] had used classification algorithms such as Logistic Regression, Support Vector Machine, and Naive-Bayes classifier for the process of classification. To check the accuracy and precision, ten-fold cross validation had been performed. [11] All the related works helps to know the research works carried out to detect depression from social media texts using different machine learning algorithms.

### 3 Dataset

The dataset used by the proposed approach consists of social media text posts that have been annotated with labels indicating the presence or absence of signs of depression. Specifically, the labels are 'not depressed', 'moderately depressed' or 'severely depressed'. The dataset has been collected from a variety of social media platforms and contains text data in English. Example texts with labels from the dataset are presented in Table 1. The dataset is divided into three parts: train data, development data, and test data. Train data and Development data consists of three columns namely PID, Text\_Data and Label. Each label column in the train data and development data consists of three different values namely 'not depressed', 'moderately depressed' or 'severely depressed' according to the text data which consists of the social media comments of an individual. There are about 7202 entries in the training dataset and 3246 entries in the development dataset. After removing the duplicates there are about 7202 entries in the training dataset and 3246 entries in the development dataset.

### 4 Solution

#### 4.1 Ensemble model to combine Random Forest Classifier, Naïve Bayes Classifier, and SVM

Ensemble learning helps to improve machine learning results by combining several models. This approach allows the production of better predictive performance compared to a single model. The main challenge is not to obtain highly accurate base models, but rather to obtain base models which make different kinds of errors. For example, if ensembles are used for classification, high accuracies can be accomplished if different base models misclassify different training examples, even if the base classifier accuracy is low.

The training set is fitted to the SVM classifier. To create the SVM classifier, we import SVC class from Sklearn.svm library.

```
classifier=SVC(kernel='linear', random_state=0)
classifier.fit(x_train, y_train)
```

## 4.2 Methods for Independently Constructing Ensembles

The different methods for constructing ensemble models are:

- Majority Vote
- Bagging and Random Forest
- Randomness Injection
- Feature-Selection Ensembles
- Error-Correcting Output Coding

**Majority Vote:** A voting ensemble involves summing the predictions made by classification models or averaging the predictions made by regression models. How voting ensembles work, when to use voting ensembles, and the limitations of the approach. How to implement a hard voting ensemble and soft voting ensemble for classification predictive modeling.

**Bagging and Random Forest:** Bagging is an ensemble algorithm that fits multiple models on different subsets of a training dataset, then combines the predictions from all models.

Random forest is an extension of bagging that also randomly selects subsets of features used in each data sample. Both bagging and random forests have proven effective on a wide range of different predictive modeling problems.

**Randomness Injection:** Random values in machine learning are derived by random number generators. To create random values, the generator is first initialized with a seed, a number that represents the starting point for the random number generation. The generator then creates random values from that starting point with a specific algorithm.

**Feature-Selection Ensembles:** The idea behind ensemble feature selection is to combine multiple different feature selection methods, taking into account their strengths, and create an optimal best subset. In general, it makes a better feature space and reduces the risk of choosing an unstable subset.

**Error-Correcting Output Coding:** The Error-Correcting Output Coding method is a technique that allows a multi-class classification problem to be reframed as multiple binary classification problems, allowing the use of native binary classification models to be used directly.

## 4.3 Types of Ensemble models

**Bagging:** Bagging (Bootstrap Aggregation) is used to reduce the variance of a decision tree. Suppose a set  $D$  of  $d$  tuples, at each iteration  $i$ , a training set  $D_i$  of  $d$  tuples is sampled with replacement from  $D$ . Then a classifier model  $M_i$  is learned for each training set  $D < i$ . Each classifier  $M_i$  returns its class prediction. The bagged classifier  $M^*$  counts the votes and assigns the class with the most votes to  $X$  (unknown sample).

**Stacking:** There are many ways to ensemble models in machine learning, such as Bagging, Boosting, and stacking. Stacking is one of the most popular ensemble machine learning techniques used to predict multiple nodes to build a new model and improve model performance. Stacking enables us to train multiple models to solve similar problems, and based on their combined output, it builds a new model with improved performance.

**Boosting:** Boosting is an ensemble method that enables each member to learn from the preceding member's mistakes and make better predictions for the future. Unlike the bagging method, in boosting, all base learners (weak) are arranged in a sequential format so that they can learn from the mistakes of their preceding learner.

Hence, in this way, all weak learners get turned into strong learners and make a better predictive model with significantly improved performance.

## 4.4 Creating the confusion matrix

To create the confusion matrix, we need to import the `confusion_matrix` function of the `sklearn` library. After importing the function, we will call it using a new variable `cm`. The function takes two parameters, mainly `y_true` (the actual values) and `y_pred` (the targeted value return by the classifier).

```
cm= confusion_matrix(y_test, y_pred)
```

## 5 Results and Discussions

### 5.1 Metrics

The metrics used during the experiments are accuracy, macro-averaged precision, macro-averaged recall and macro-averaged F1-score across all the classes. The macro-averaged F1-score was the main measure when evaluating solutions.

Model	Accuracy	Precision	Recall	F1- score
Random Forest Classifier	0.599	0.62	0.60	0.56
Naive Bayes Classifier	0.523	0.58	0.52	0.40
SVM Classifier	0.636	0.63	0.61	0.63
Ensemble Model	0.902	0.91	0.89	0.91

Table 2: Results of model on the development dataset

## 5.2 Performance

Table 2 shows the result of each model on the training dataset. Among the machine learning language models as shown in the Table 2 Ensemble model was the best in terms of accuracy (0.90) and F1-score (0.89). Sklearn SVC is the implementation of SVC provided by the popular machine learning library Scikit-learn. There are three datasets namely training dataset, development dataset and test dataset. The training dataset is pre-processed and is used to train the SVM model. Flask server is used to display the depression detection website. In the website the user will be prompted to enter a comment. The SVM classifier is used to detect depression and classify them as ‘not depressed’, ‘moderately depressed’ or ‘severely depressed’. The result is displayed on the next webpage. Confusion matrix is plotted to determine the performance of the classification models for a given set of test data. Figure 1 to 4 shows the confusion matrix for the three machine learning models used: Random Forest Classifier, Naive Bayes Classifier and SVM Classifier.

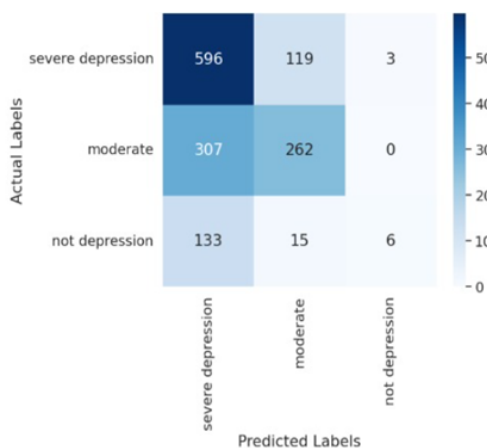


Figure 1: Random Forest

## 6 Conclusion

The proposed system aim to detect severity of depression from social media text using machine

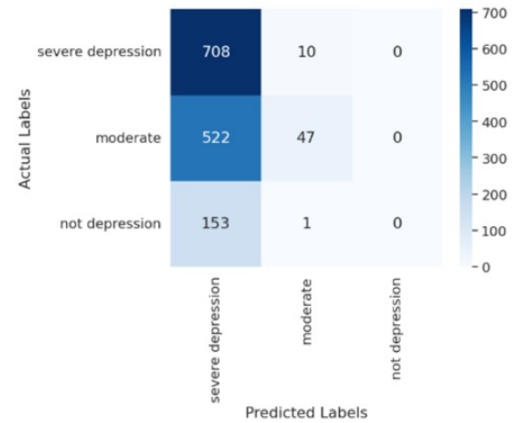


Figure 2: Naive Bayes

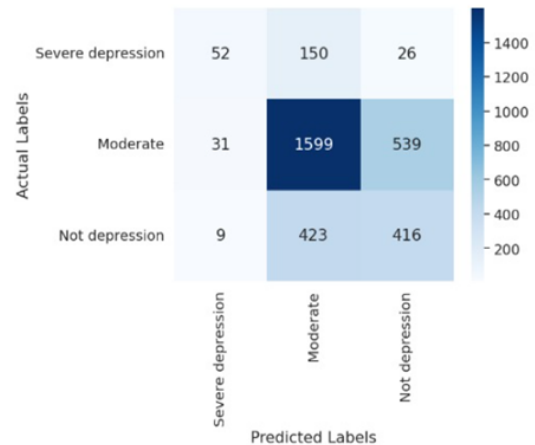


Figure 3: SVM

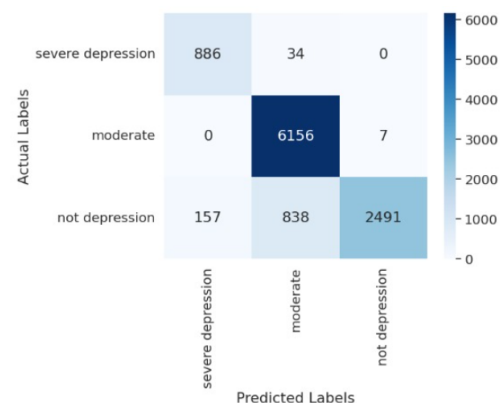


Figure 4: Ensemble model

learning techniques. The given data is pre-processed by cleaning, tokenizing, and removing stop words. Then, feature engineering techniques like TF-IDF are applied to represent the text data. Experimentation was conducted using three machine learning algorithms, namely SVM, Random Forest and Naive Bayes, and evaluated their performance using various metrics like accuracy, precision, recall, and F1-score. By using ensemble model the accuracy obtained is 47.7% and the F1 score is 0.262. Overall, the proposed system demonstrate the potential of machine learning techniques to detect severity of depression from social media text, which can aid in early detection and intervention of depression.

## References

- [1] Anamika Ahmed, Raihan Sultana, Md Tahmidur Rahman Ullas, Mariyam Begom, Md. Muzahidul Islam Rahi, and Md. Ashraful Alam. A machine learning approach to detect depression and anxiety using supervised learning. In *2020 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE)*, pages 1–6, 2020.
- [2] Nafiz Al Asad, Md. Appel Mahmud Pranto, Sadia Afreen, and Md. Maynul Islam. Depression detection by analyzing social media posts of user. In *2019 IEEE International Conference on Signal Processing, Information, Communication Systems (SPICSCON)*, pages 13–17, 2019.
- [3] Yan Ding, Xuemei Chen, Qiming Fu, and Shan Zhong. A depression recognition method for college students using deep integrated support vector algorithm. *IEEE Access*, 8:75616–75629, 2020.
- [4] G. Geetha, G. Saranya, K. Chakrapani, J. Godwin Ponsam, M. Safa, and S. Karpagaselvi. Early detection of depression from social media data using machine learning algorithms. In *2020 International Conference on Power, Energy, Control and Transmission Systems (ICPECTS)*, pages 1–6, 2020.
- [5] Swati Jain, Suraj Prakash Narayan, Rupesh Kumar Dewang, Utkarsh Bhartiya, Nalini Meena, and Varun Kumar. A machine learning based depression analysis and suicidal ideation detection system using questionnaires and twitter. In *2019 IEEE Students Conference on Engineering and Systems (SCES)*, pages 1–6, 2019.
- [6] Keshu Malviya, Bholanath Roy, and SK Saritha. A transformers approach to detect depression in social media. In *2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS)*, pages 718–723, 2021.
- [7] Kayalvizhi S, Thenmozhi Durairaj, Bharathi Raja Chakravarthi, and Jerin Mahibha C. Findings of the shared task on detecting signs of depression from social media. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 331–338, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [8] Sivamanikandan S, Santhosh V, Sanjaykumar N, Jerin Mahibha C, and Thenmozhi Durairaj. scubeMSEC@LT-EDI-ACL2022: Detection of depression using transformer models. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 212–217, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [9] Ananna Saha, Ahmed Al Marouf, and Rafayet Hosain. Sentiment analysis from depression-related user-generated contents from social media. In *2021 8th International Conference on Computer and Communication Engineering (ICCCE)*, pages 259–264, 2021.
- [10] Shahriar Saleque, Gul-A-Zannat Spriha, Rasheeq Ishraq Kamal, Rafia Tabassum Khan, Amitabha Chakrabarty, and Mohammad Zavid Parvez. Detection of major depressive disorder using signal processing and machine learning approaches. In *2020 15th IEEE Conference on Industrial Electronics and Applications (ICIEA)*, pages 1032–1037, 2020.
- [11] Kayalvizhi Sampath, Thenmozhi Durairaj, Bharathi Raja Chakravarthi, Jerin Mahibha C, Kogilavani Shanmugavadivel, and Pratik Anil booktitle = Rahood. Overview of the second shared task on detecting signs of depression from social media text.
- [12] Michael M. Tadesse, Hongfei Lin, Bo Xu, and Liang Yang. Detection of depression-related posts in reddit social media forum. *IEEE Access*, 7:44883–44893, 2019.