

Transparent insights: explainable AI with machine learning classifiers for early stage of depression classification

S. M. Rakibul Islam¹, Shaykh Yunus¹, Rashiduzzaman Shakil², Fatema Tuz Johora^{1,4}, Aditya Rajbongshi¹, Sujon Chandra Sutradhar³

¹Department of Educational Technology and Engineering, University of Frontier Technology, Bangladesh, Gazipur, Bangladesh

²Department of Computer Science and Engineering, Daffodil International University, Daffodil Smart City (DSC), Dhaka, Bangladesh

³Department of General Education, University of Frontier Technology, Bangladesh, Gazipur, Bangladesh

⁴Department of Software Engineering, Daffodil International University, Daffodil Smart City (DSC), Dhaka, Bangladesh

Article Info

Article history:

Received Nov 2, 2025

Revised Mar 4, 2026

Accepted Mar 29, 2026

Keywords:

Depression classification

Explainable artificial intelligence

Machine learning

Shapley additive explanations

Synthetic minority over sampling technique

ABSTRACT

Depression is a widespread mental health condition characterized by enduring feelings of persistent sadness, loss of interest, and impaired daily functioning. Untreated depression can result in significant implications, such as academic failure, social isolation, and even suicide. This study presents a machine learning (ML)-based framework for classifying depression severity among university students using the Zahir depression scale dataset, comprising 478 responses categorized into mild, moderate, severe, and profound depression. In order to address the issue of class imbalance, we utilized the synthetic minority over sampling technique (SMOTE) on the dataset. In addition, seven different ML algorithms are employed to classify the severity of depression, and each algorithm's efficiency is determined by four performance evaluation metrics. Among the applied ML classifiers, extra tree classifier outperformed with an average accuracy of 97.85% and 95.75% precision, 95.76% recall, and 95.75% F1-score. To enhance interpretability, the shapley additive explanations (SHAP) method was integrated to identify influential features, providing transparency and insight into the model's decision process. The proposed framework demonstrates that combining explainable artificial intelligence (XAI) with traditional ML can support healthcare professionals in early depression screening and data-driven mental health interventions.

This is an open access article under the [CC BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.



Corresponding Author:

Aditya Rajbongshi

Department of Educational Technology and Engineering, University of Frontier Technology, Bangladesh

Gazipur 1750, Bangladesh

Email: aditya0001@uftb.ac.bd

1. INTRODUCTION

Depression is a pervasive mental disorder marked by persistent sadness and diminished interest in daily activities. At present, depression affects people of all ages, as it is considered one of the primary reasons for illness and disability. Psychological factors, such as negative thought patterns, learned helplessness, and a history of trauma can exacerbate vulnerability [1]. Thus, the impact of depression on mental well-being was transcended. It can manifest physically, leading to sleep disturbances, changes in appetite, and negative effects on the immune system [2]. Moreover, the chance of developing long-term health issues such as heart disease, diabetes, and obesity increases owing to depression [3]. According to the World Health Organization (WHO), depression affects over 264 million people globally and disrupts daily

life and wellbeing [4]. More than 300 million people are estimated to suffer from depression worldwide, and the rate is higher in low and middle-income countries such as Bangladesh [5].

University students represent a unique group of individuals who typically undergo critical transition due to their physical development. A survey from 2021 to 2022 revealed an alarming rise in depression, particularly among students facing immense academic pressure, with 44% of students across 133 US college campuses reporting depressive symptoms [6]. As a result, more mental health issues of higher levels are being reported now. Among pre-university students in Bangladesh, depression and anxiety rates of 44% and 27%, respectively, have been reported [7], but these rates have escalated to 52% and 58%, respectively [8]. Currently, this is a worldwide issue that transcends geographical boundaries. According to research, there was concern about the rise in student depression across various Asian countries in 2020, with India reporting a prevalence of 31.9% among medical students [9]. In addition, the WHO forecasts that over 54 million individuals in China suffer from depression, whereas approximately 41 million people suffer from anxiety disorders [10]. In Bangladesh, a research investigation found a disturbing rate of suicidal ideation and attempts among young adults, highlighting the need for improved mental health support services [11]. An estimated 25% of the world's population suffer from mental health diseases, whereas approximately 7 million Bangladeshis suffer from anxiety and depression. Suicide is the most concerning consequence of depression. The WHO anticipated that 703,000 individuals would die annually from suicide worldwide in 2020 [12]. Numerous factors, namely academic or non-academic pressures such as socioeconomic, environmental, cultural, and psychological attributes, may cause stress among youths [13].

The traditional diagnosis of depression relies on clinical treatment, and self-administered measures are subjective and time-consuming. With recent developments in both medical research and technology, machine learning (ML) has gradually supplanted more conventional methods for diagnosing and treating depression. ML can process vast amounts of data, potentially leading to earlier detection, more accurate diagnoses, and improved treatment outcomes [14]. Studies using ML algorithms can analyze various data sources, including clinical records, social media activity, and speech patterns, to identify individuals at risk for depression [15]. Conventional ML models are particularly complicated, as are ensemble approaches that frequently act as “black boxes.” Because of this lack of transparency, it is difficult to understand the prediction generation process [5], [16]. Consequently, explainable artificial intelligence (XAI) techniques offer insight into how models obtain their findings. For example, shapley additive explanations (SHAP) values can emphasize which clinical factors contribute the most to a specific diagnosis, making the decision-making process more transparent.

This study bridges the gap between predictive accuracy and model interpretability by integrating SHAP-based explainability into depression severity classification using structured psychological survey data. It also attempts to provide insight into the multiple aspects of depression to bring attention to this significant public health concern. By conducting a thorough examination of the factors that contribute to depression, the effects it has, and the possible remedies.

2. LITERATURE REVIEW

In order to address the depression, many researchers studied and proposed their methodology to find out the level for ensuring proper guidance for a patient. Li [17] utilized four distinct ML models, highlighting the urgent worldwide issue of depression and suicide and underlining the crucial requirement for proactive mental health interventions. In this study they used datasets provided by “Suicide Watch” and “Depression” on the Reddit platform and employed a deep neural network (DNN) model to classify the level of depression that achieved 95% accuracy. Biradar and Totad [18] highlight the growing utilization of Twitter and other social networking sites (SNS) as platforms for expressing emotions and daily activities, providing important data for psychological research. Researchers have used a back propagation neural network (BPNN) to classify tweets as either indicative of depression or not.

Almars *et al.* [19] collected a dataset from Twitter, which consisted of around 6000 tweets, for analysis of the people's depression level. Similarly, Zulfiker *et al.* [20] also worked in the similar domain. They used synthetic minority oversampling technique (SMOTE) to eliminate the class imbalance problem in their dataset. After that, the researchers employed three ML algorithms: adaptive boosting (AdaBoost), support vector machine (SVM), and logistic regression (LR) to classify depression patients, and AdaBoost outperformed the others with an accuracy of 92.56%.

Haque *et al.* [21] introduced the Boruta feature selection algorithm and techniques, along with random forest (RF), to identify depression in children and adolescents aged 4-17 years old. Among applied models, RF demonstrated superior performance, with an accuracy of 95%. Similarly, Priya *et al.* [22] focused on accurately classifying depression by employing five ML models, where Naïve Bayes (NB) outperformed

with 85.5% accuracy. Sau and Bhakta [23] used an RF classifier for the anxiety level of a generic patient, where the model achieved the highest prediction accuracy of 89%.

Islam *et al.* [24] focused on the classification of depression, mainly focusing on four factors: emotional process, linguistic, and temporal. A feature-selection-based method presented by Alghowinem *et al.* [25] focuses on the importance of feature selection and the performance and limitations of existing models. This methodology was conducted based on 38 different features from 902 patients.

In order to predict the hypertension for US people, Lee and Kim [26] suggest an ML-based framework. Among applied methods, SVM had the highest accuracy (0.771) for classifying depression. By employing psychological data, Gupta *et al.* [27] used it for depression detection. They utilized the SMOTE technique to address the issue of class imbalance. Nemesure *et al.* [28] classify generalized anxiety disorder (GAD) and major depressive disorder (MDD) using ML. To conduct their research, they used 4,184 patients' data, whereas the biomedical and demographic features were 59. Magboo and Magboo [29] accessed a publicly available dataset from GitHub which consists of 604 sample. To anticipate depression, eight different ML models were used, where the LR algorithm gained 91% accuracy.

3. METHOD

This section comprises three distinct subsections: data acquisition, data preprocessing, and the applied model with performance evaluation metrics. In the subsection of data acquisition, we describe the collection of raw data, while preprocessing focuses on cleaning and preparing the data for analysis. All the modeling approaches and evaluation criteria used to assess performance are explained in a subsection. Figure 1 illustrates the overall working process of the proposed methodology, detailing each step from data acquisition to the final model outcome.

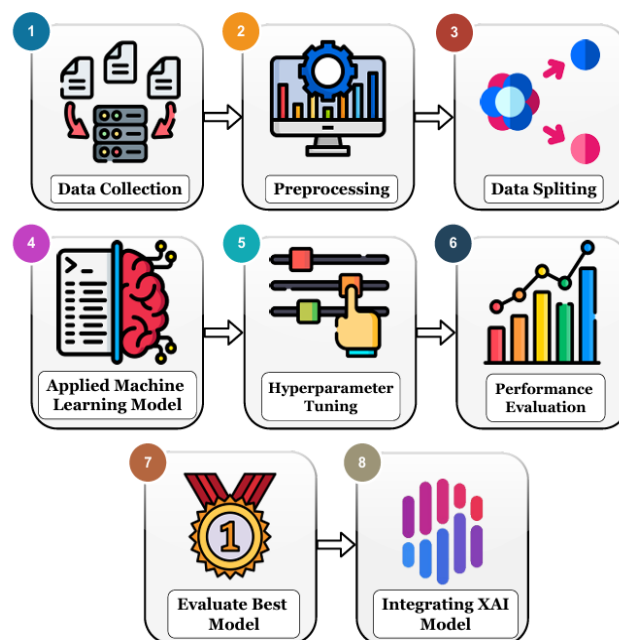


Figure 1. Workflow of the proposed depression classification system

3.1. Data acquisition

To conduct this research, we utilized the “Zahir Depression Scale” questionnaire, which was collected from the National Institute of Mental Health (NIMH), Dhaka, Bangladesh [30]. Based on these questionnaires, data was collected through a Google form among Bangladeshi students from different regions but primarily focused on Bangladeshi university students. It consisted of approximately 478 individuals with 27 distinct attributes. 289 participants were found to have mild depression, and 73, 50, and 66 were found to have moderate, severe, and profound depression, respectively. Table 1 consists of attributes and description of the data. After that, according to the “Zahir Depression Scale”, Table 2 provide an overview of several classes of depression depending on percentiles and corresponding scores on depression scale.

Table 1. Overview of the dataset attributes

| SI No. | Attribute name | SI No. | Attribute name |
|--------|----------------------------|--------|--------------------------|
| 1 | Name | 15 | Incompetence and failure |
| 2 | Age | 16 | Guilt |
| 3 | Gender | 17 | Low self-confidence |
| 4 | Upset | 18 | Emptiness |
| 5 | Joyless | 19 | Attention deficit |
| 6 | To cry | 20 | Making decisions |
| 7 | Turmoil | 21 | Slowing work |
| 8 | Lack of interest | 22 | Social interaction |
| 9 | Worthlessness | 23 | Weakness and fatigue |
| 10 | Hopelessness | 24 | Appetite changes |
| 11 | Death wish | 25 | Weight changes |
| 12 | Suicide | 26 | Sleep changes |
| 13 | Sense of loss | 27 | Sex interest |
| 14 | Pain of unmet expectations | | |

Table 2. Severity of depression based on the corresponding score of depression scale

| Percentiles | Corresponding scores on the depression scale | Severity of depression |
|-------------|--|------------------------|
| 25 | 0-34 | Mild depression |
| 50 | 35-43 | Moderate depression |
| 75 | 44-50 | Severe depression |
| 100 | 51 and above | Profound depression |

3.2. Data preprocessing

Data preprocessing is the crucial part of research. The dataset initially contained unprocessed records that required preprocessing to ensure data consistency and reliability. During this phase, various preprocessing techniques were implemented to adequately prepare the data. Initially, we checked for missing values in our dataset. However, there were no missing values. Second, priority-based encoding mapping was utilized to assess the priority mapping of each column. Third, we classified the age into four groups: 18–21, 22–24, 25–28, and 29–34 years. Further, we used label encoding for the target column named depression level, where ‘Mild’: 0, ‘Moderate’: 1, ‘Severe’: 2, and ‘Profound’: 3. Data normalization is the process of rescaling attributes to have a mean of 0 and a variance of 1 [31]. We applied a standard scaler to normalize all features to a uniform scale while preserving the original range of values.

Balancing the dataset enhances model training by preventing bias towards a certain class. Hence, to mitigate the risk of improper model training, we employed the borderline-SMOTE to equalize the distribution of data in the training dataset [32]. In this strategy, oversampling produces misleading data points from a minority class [33], [34]. Figure 2(a) illustrates the imbalanced dataset, and Figure 2(b) depicts the balanced dataset before and after employing SMOTE. We divided the dataset into two distinct categories: training (80%) and testing (20%) sets. Out of 478 students, 382 samples were used for training and 96 samples for testing.

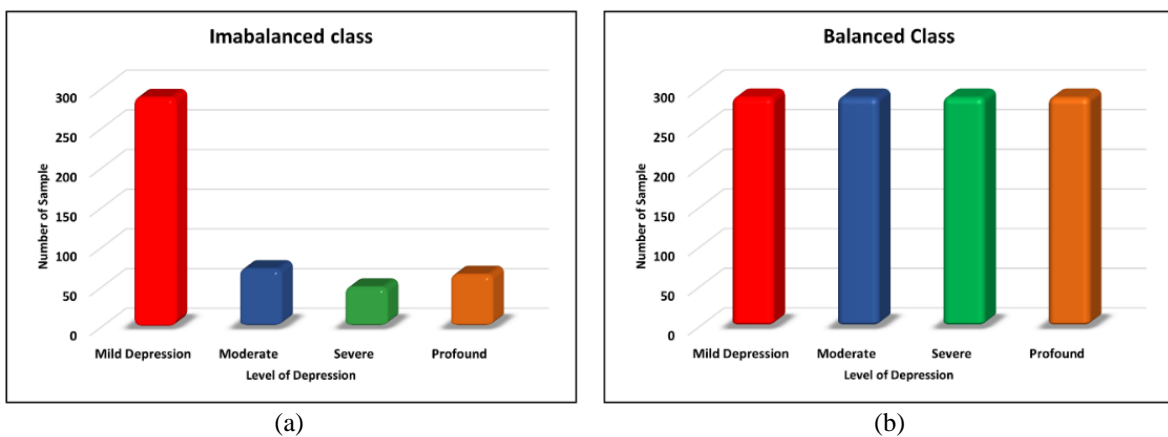


Figure 2. Visualization of depression severity class distribution showing: (a) imbalanced classes and (b) balanced classes after applying SMOTE

3.3. Applied ML models and performance evaluation metrics

In this study, we have implemented seven ML classifiers, including LR [33], SVM [34], extra tree [35], AdaBoost [36], gradient boosting (GB) [37], RF [38], and K-nearest neighbor (KNN) [39], and those are trained with the training dataset. Performance evaluation metrics are used in a study to evaluate the effectiveness of applied classifiers. At this stage, four performance metrics, namely accuracy, precision, recall, and F1-score, have been used to determine the efficiency of the applied ML algorithm.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F1\ Score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (4)$$

4. RESULTS

To conduct this research, all experiments were conducted in Jupyter Notebook using Python 3.9.18. ML models were implemented with scikit-learn 1.2.2, and result visualization was performed using matplotlib 3.7.3. In this study, we have developed an automation system that can precisely classify depression level utilizing the ML algorithms to help medical professionals diagnose depression more quickly and accurately, improving the depression prediction outcomes.

For evaluating the model performance, a confusion matrix was generated for each ML model. Figures 3 and 4 presented the class-wise confusion matrix evaluation values for imbalanced and balanced data. For each confusion matrix, 0: mild, 1: moderate, 2: severe, and 3: profound depression. In Figure 3(a), the LR classifier in the “Mild” class has true positive (TP), true negative (TN), false positive (FP), and false negative (FN) values of 59, 34, 2, and 1, respectively. The SVM in Figure 3(b) successfully predicts the highest number of TP (58) value for the mild class, while in Figures 3(c) and (d) for extra tree and AdaBoost, the TP values are 59 and 57 for mild class. For the Figure 3(e) GB classifier, the profound class TP and FN values are gradually 8 and 3. In Figure 3(f) RF, and Figure 3(g) KNN predict the similar number of TP values, 60, for the mild class, though their TN values are different, 26 and 29.

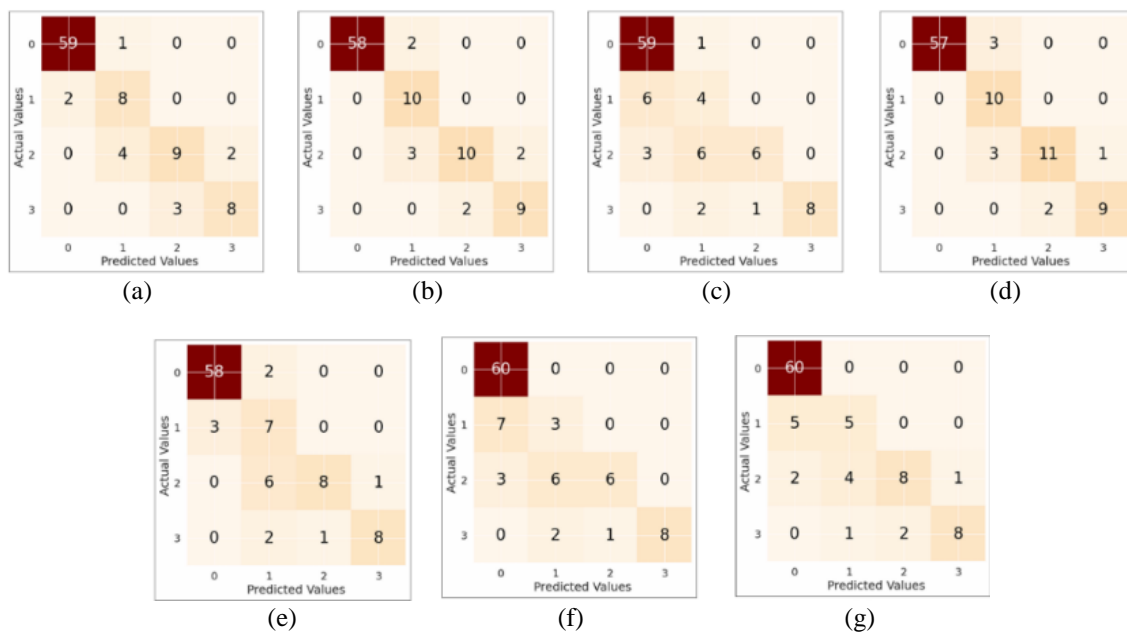


Figure 3. Confusion matrices for ML models on imbalanced data: (a) LR, (b) SVM, (c) extra tree, (d) AdaBoost, (e) GB, (f) RF, and (g) KNN

Figure 4 depicts confusion matrix values after solving the data imbalance issue using SMOTE. In Figure 4(a), the LR for the moderate class, TP, and TN values are 53 and 164. The SVM algorithm in Figure 4(b) for the mild class has the TP, TN, FP, and FN values of 49, 176, 1, and 6 respectively. The severe class in Figure 4(c) for extra tree classifier has 60 (TP) and 1 (FP). The similar number of TP (52) was evaluated for severe and profound classes in the AdaBoost classifier in Figure 4(d). For the GB classifier in Figure 4(e), the profound class achieved 53 (TP) and 2 (FN), while in Figure 4(f), the RF classifier, the profound class recorded 53 TP and 0 FN. In Figure 4(g) the KNN classifier in yielded TP counts of 47, 60, 62, and 52 for the mild, moderate, severe, and profound classes, respectively.

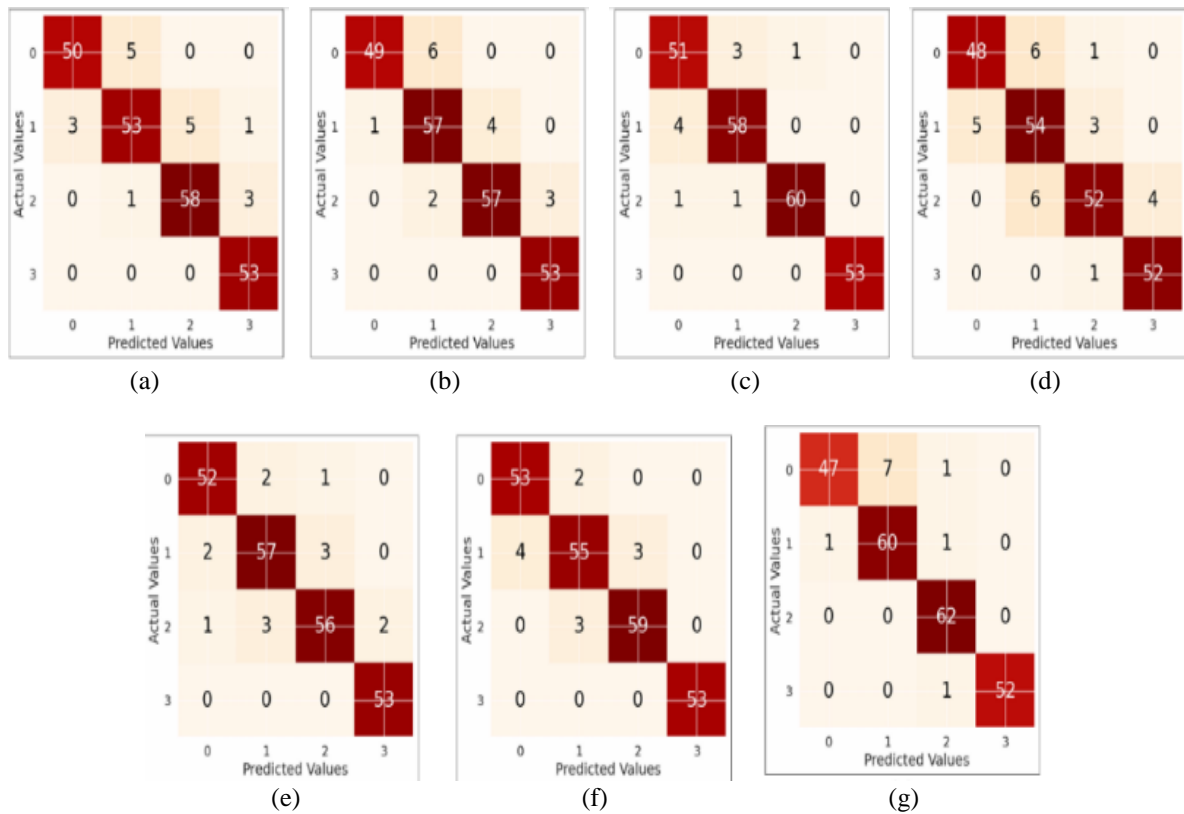


Figure 4. Confusion matrices for ML models on balanced data: (a) LR, (b) SVM, (c) extra tree, (d) AdaBoost, (e) GB, (f) RF, and (g) KNN

4.1. Depression class-wise performance of the applied ML algorithms

After estimating the confusion matrix value, we evaluated the performance result of 4 classes for depression classification that is presented in Tables 3 and 4, respectively. Among the classifiers in Table 3, the SVM outperformed, the accuracy was 97.92%, 94.79%, 92.71% and 95.83% of mild, moderate, severe and profound class respectively. In the extra tree, the profound class achieves the highest accuracy of 96.88% and 100% precision, whereas the moderate class achieves the lowest precision of 30.77%. In addition, for RF, the accuracy of the moderate class was 84.38% and the recall of the mild class was 100%.

Additionally, after applying the SMOTE technique, the performance evaluation metrics for each classifier was improved. In Table 4 the accuracy in KNN algorithm, 99.57%, 98.71%, 96.12% and 96.12% for Profound, Severe, Moderate and Mild depression class. In the RF the Moderate class accuracy is 94.83%, and the Recall of Mild class become 96.36%. The accuracy of severe class for LR, extra tree, GB, and KNN are 96.12%, 98.71%, 95.69% and 98.71%, whereas the recall values become 93.55%, 96.77%, 90.32% and 100% respectively. The AdaBoost gains highest precision 92.86% and lowest recall is 83.87% among their four classes and the accuracy, precision, recall and F1-score of profound class are 97.84%, 92.86%, 98.11%, and 95.41% respectively.

Table 3. Depression class-wise performance of the applied ML algorithms on imbalanced data

| Algorithm | Class | Accuracy (%) | Precision (%) | Recall (%) | F1-score (%) |
|------------|----------|--------------|---------------|------------|--------------|
| LR | Mild | 96.88 | 96.72 | 98.33 | 97.52 |
| | Moderate | 92.71 | 61.54 | 80 | 69.57 |
| | Severe | 90.62 | 75 | 60 | 66.67 |
| | Profound | 94.79 | 80 | 72.73 | 76.19 |
| SVM | Mild | 97.92 | 100 | 96.67 | 98.31 |
| | Moderate | 94.79 | 66.67 | 100 | 80 |
| | Severe | 92.71 | 83.33 | 66.67 | 74.07 |
| | Profound | 95.83 | 81.82 | 81.82 | 81.82 |
| Extra tree | Mild | 89.58 | 86.76 | 98.33 | 92.19 |
| | Moderate | 84.38 | 30.77 | 40 | 34.78 |
| | Severe | 89.58 | 85.71 | 40 | 54.55 |
| | Profound | 96.88 | 100 | 72.73 | 84.21 |
| AdaBoost | Mild | 96.88 | 100 | 95 | 97.74 |
| | Moderate | 93.75 | 62.5 | 100 | 76.92 |
| | Severe | 93.75 | 84.62 | 73.33 | 78.57 |
| | Profound | 96.88 | 90 | 81.82 | 85.71 |
| GB | Mild | 94.79 | 95.08 | 96.67 | 95.87 |
| | Moderate | 86.46 | 41.18 | 70 | 51.85 |
| | Severe | 91.67 | 88.89 | 53.33 | 66.67 |
| | Profound | 95.83 | 88.89 | 72.73 | 80 |
| RF | Mild | 89.58 | 85.71 | 100 | 92.31 |
| | Moderate | 84.38 | 27.27 | 30 | 28.57 |
| | Severe | 89.58 | 85.71 | 40 | 54.55 |
| | Profound | 96.88 | 100 | 72.73 | 84.21 |
| KNN | Mild | 92.71 | 89.55 | 100 | 94.49 |
| | Moderate | 89.58 | 50 | 50 | 50 |
| | Severe | 90.62 | 80 | 53.33 | 64 |
| | Profound | 95.83 | 88.89 | 72.73 | 80 |

Table 4. Depression class-wise performance of the applied ML algorithms on balanced data

| Algorithm | Class | Accuracy (%) | Precision (%) | Recall (%) | F1-score (%) |
|------------|----------|--------------|---------------|------------|--------------|
| LR | Mild | 96.55 | 94.34 | 90.91 | 92.59 |
| | Moderate | 93.53 | 89.83 | 85.48 | 87.6 |
| | Severe | 96.12 | 92.06 | 93.55 | 92.8 |
| | Profound | 98.28 | 92.98 | 100 | 96.36 |
| SVM | Mild | 96.98 | 98 | 89.09 | 93.33 |
| | Moderate | 94.4 | 87.69 | 91.94 | 89.76 |
| | Severe | 96.12 | 93.44 | 91.94 | 92.68 |
| | Profound | 98.71 | 94.64 | 100 | 97.25 |
| Extra tree | Mild | 96.12 | 91.07 | 92.73 | 91.89 |
| | Moderate | 96.55 | 93.55 | 93.55 | 93.55 |
| | Severe | 98.71 | 98.36 | 96.77 | 97.56 |
| | Profound | 100 | 100 | 100 | 100 |
| AdaBoost | Mild | 94.83 | 90.57 | 87.27 | 88.89 |
| | Moderate | 91.38 | 81.82 | 87.1 | 84.38 |
| | Severe | 93.53 | 91.23 | 83.87 | 87.39 |
| | Profound | 97.84 | 92.86 | 98.11 | 95.41 |
| GB | Mild | 97.41 | 94.55 | 94.55 | 94.55 |
| | Moderate | 95.69 | 91.94 | 91.94 | 91.94 |
| | Severe | 95.69 | 93.33 | 90.32 | 91.8 |
| | Profound | 99.14 | 96.36 | 100 | 98.15 |
| RF | Mild | 97.41 | 92.98 | 96.36 | 94.64 |
| | Moderate | 94.83 | 91.67 | 88.71 | 90.16 |
| | Severe | 97.41 | 95.16 | 95.16 | 95.16 |
| | Profound | 100 | 100 | 100 | 100 |
| KNN | Mild | 96.12 | 97.92 | 85.45 | 91.26 |
| | Moderate | 96.12 | 89.55 | 96.77 | 93.02 |
| | Severe | 98.71 | 95.38 | 100 | 97.64 |
| | Profound | 99.57 | 100 | 98.11 | 99.05 |

4.2. Overall performance evaluation results of each algorithm

Figure 5 depicts the overall performance of each algorithm to classify the four distinct types of depression. Figure 5(a) represents the overall performance result of imbalanced data, where the LR, SVM, extra tree, AdaBoost, GB, RF and KNN achieve 93.75%, 95.31%, 90.11%, 95.31%, 92.19%, 90.11% and 92.19% average accuracy respectively. AdaBoost and SVM outperformed compared to another model.

On the other hand, Figure 5(b) represents the overall result for balanced data. The extra trees algorithm achieved the highest average accuracy rate of 97.85%, with a precision of 95.75%, recall of

95.76%, and an F1-score of 95.75%. KNN and RF algorithm closely followed with average accuracy rates of 97.63% and 97.41%, respectively.

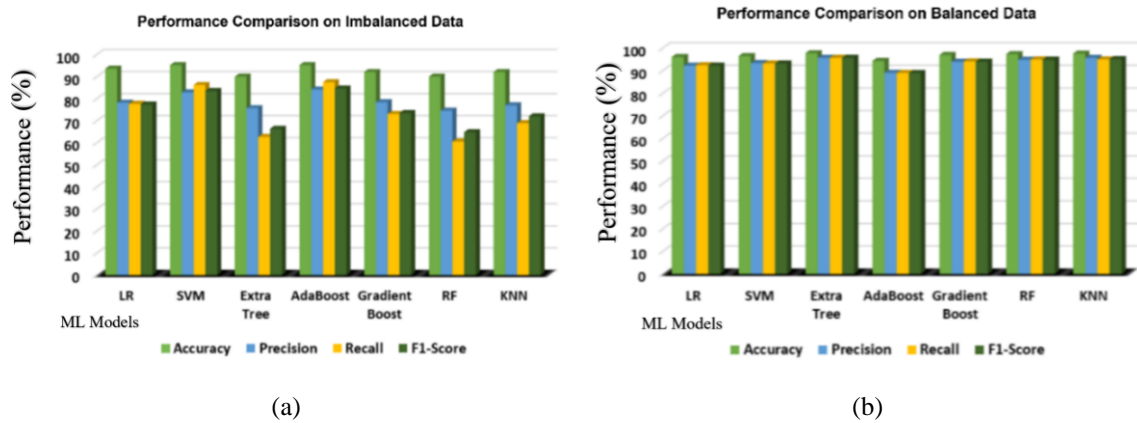


Figure 5. Performance comparison of ML models on (a) imbalanced and (b) balanced data

As traditional ML models work like a black box, to make it transparent, we employed SHAP to make this more precise for depression classification. In this study, we evaluate feature importance using SHAP with the extra trees classifier. Figure 6 depicts the synopsis of SHAP, which visualizes the effect of each feature on the model’s predictions. The X-axis represents the SHAP value, where higher positive values indicate a stronger tendency towards predicting early-stage depression, while lower negative values indicate a lower likelihood. The features are ranked by importance, with the most influential features appearing at the top. The plot reveals that the “turmoil” feature exhibits the highest positive SHAP values, suggesting that higher levels of turmoil significantly contribute to the prediction of early-stage depression. On the other hand, the “sex interest” feature predominantly shows negative SHAP values, indicating that lower values of sex interest are associated with a decreased probability of early-stage depression.

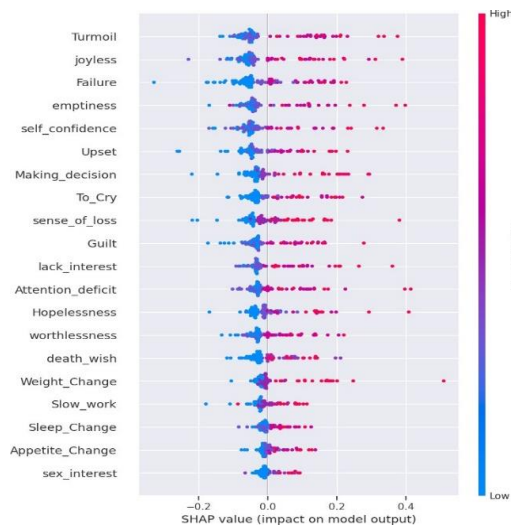


Figure 6. Summary plot of SHAP utilizing extra tree classifier

Numerous researchers have previously engaged in depression classification utilizing artificial intelligence; thus, we provide a comprehensive study alongside their findings in Table 5. Many of them examine depression by monitoring their social media habits and determining whether they are depressed or not. Among these studies, most of the researchers concentrated on identifying the binary classification of depression, whereas we classified four categories based on 27 distinct criteria. In addition, XAI is integrated to make this classification transparent.

Table 5. Comparative analysis of existing studies and the proposed methodology

| Authors | Data frequency | Total attributes | No of class | XAI | Highest accuracy |
|-------------------------------|-----------------------------------|------------------|----------------|-----|--------------------|
| Li [17] | 202374 tweets data | - | Binary | × | LR: 84.15% |
| Biradar Totad [18] | 61,400 tweets of different types | 8 | Binary | × | BPNN: 81.46% |
| Almars [19] | 6000 tweets | 4 | Binary | × | Bi-LSTM: 083% |
| Zulfiker <i>et al.</i> [20] | 604 participants | 6 | Binary | × | Adaboost: 92.56% |
| Haque <i>et al.</i> [21] | 6,310 children | 12 | Binary | × | RF: 95% |
| Priya <i>et al.</i> [22] | 348 participants | 7 | Multiclass | × | Naïve Bayes: 85.5% |
| Sau and Bhakta [23] | 510 geriatric patients | 7 | Binary | × | RF: 89% |
| Islam <i>et al.</i> [24] | 7145 comments, | 4 | Binary | × | - |
| Alghowinem <i>et al.</i> [25] | 902 behavioral cues | 8 | Binary | ✓ | - |
| Lee and Kim [26] | 8,628 adults with hypertension | 6 | Binary | × | SVM: 77.1% |
| Gupta <i>et al.</i> [27] | Sentiment_140: 1.6 million tweets | 4 | Binary | × | LSTM: 83% |
| Nemesure <i>et al.</i> [28] | 4,184 undergraduate students | 3 | - | ✓ | - |
| Magboo and Magboo [29] | 604 instances | 5 | Binary | ✓ | LR: 91% |
| Proposed methodology | 478 participants | 27 | Multiclass (4) | ✓ | Extra tree: 97.85% |

5. CONCLUSION

This study developed an explainable ML framework for classifying the severity of depression among university students in Bangladesh. We determine the prevalence and intensity of four unique levels of depression among individuals in the age range from 18 to 34. Our research findings indicate that a substantial number of students polled experienced moderate and severe levels of severity of depression. To carry out this study, we utilize seven ML classifiers to classify the level of depression. After utilizing the SMOTE technique, the performance of each algorithm surged. The extra trees classifier achieves the highest average accuracy of 97.85%, outperforming six other algorithms, while the integration of the SHAP technique provides transparent reasoning behind model predictions. The analysis reveals that turmoil was the most influential feature positively correlated with depression, whereas sex interest showed the strongest negative contribution. These insights demonstrate the potential of XAI in supporting early depression detection and fostering data-driven mental health interventions. Future work will focus on expanding the dataset to include diverse demographic groups and integrating multimodal features such as behavioral and textual data to further improve predictive performance and interpretability.





REFERENCES

- [1] R. Sarwar and A. Asmat, "Mental health expert's perspective on risk and protective factors of suicide ideation in Patients with OCD and depression," *BMC Psychiatry*, vol. 25, no. 1, 2025, doi: 10.1186/s12888-024-06404-9.
- [2] American Psychological Association, "Depression." Accessed: Apr. 16, 2025. [Online]. Available: <https://www.apa.org/topics/depression>
- [3] A. M. Delamater, A. Guzman, and K. Aparicio, "Mental health issues in children and adolescents with chronic illness," *International Journal of Human Rights in Healthcare*, vol. 10, no. 3, pp. 163–173, 2017, doi: 10.1108/ijhrh-05-2017-0020.
- [4] World Health Organization (WHO), "Depressive disorder (depression)." Accessed: Apr. 11, 2025. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/depression>
- [5] World Health Organization (WHO), "Depression and Other Common Mental Disorders: Global Health Estimates." Accessed: Apr. 13, 2025. [Online]. Available: <https://www.who.int/publications/i/item/depression-global-health-estimates>
- [6] K. Bowe, "College students and depression: A guide for parents," Speaking of Health, Mayo Clinic Health Systems. Accessed: Apr. 11, 2025. [Online]. Available: <https://www.mayoclinichealthsystem.org/hometown-health/speaking-of-health/college-students-and-depression>
- [7] M. A. H. Bhuiyan, M. D. Griffiths, and M. A. Mamun, "Depression literacy among Bangladeshi pre-university students: Differences based on gender, educational attainment, depression, and anxiety," *Asian Journal of Psychiatry*, vol. 50, p. 101944, 2020, doi: 10.1016/j.ajp.2020.101944.
- [8] S. Hossain, A. Anjum, M. E. Uddin, M. A. Rahman, and M. F. Hossain, "Impacts of socio-cultural environment and lifestyle factors on the psychological health of university students in Bangladesh: A longitudinal study," *Journal of Affective Disorders*, vol. 256, pp. 393–403, 2019, doi: 10.1016/j.jad.2019.06.001.
- [9] S. Sarkar, R. Gupta, and V. Menon, "A systematic review of depression, anxiety, and stress among medical students in India," *Journal of Mental Health and Human Behaviour*, vol. 22, no. 2, pp. 88–96, 2017, doi: 10.4103/jmhbb.jmhbb_20_17.
- [10] World Health Organization (WHO), "Mental health." Accessed: Apr. 13, 2025. [Online]. Available: <https://www.who.int/china/health-topics/mental-health>
- [11] A. Irish and N. S. Murshid, "Suicide ideation, plan, and attempt among youth in Bangladesh: Incidence and risk factors," *Children and Youth Services Review*, vol. 116, p. 105215, 2020, doi: 10.1016/j.childyouth.2020.105215.
- [12] World Health Organization (WHO), "Suicide," 2025. Accessed: Apr. 13, 2025. [Online]. Available: <https://www.who.int/health-topics/suicide>
- [13] A. R. Arusha and R. K. Biswas, "Prevalence of stress, anxiety and depression due to examination in Bangladeshi youths: A pilot study," *Children and Youth Services Review*, vol. 116, p. 105254, 2020, doi: 10.1016/j.childyouth.2020.105254.
- [14] M. M. Tadesse, H. Lin, B. Xu, and L. Yang, "Detection of Depression-Related Posts in Reddit Social Media Forum," *IEEE Access*, vol. 7, pp. 44883–44893, 2019, doi: 10.1109/access.2019.2909180.
- [15] F. Cacheda, D. Fernandez, F. J. Novoa, and V. Carneiro, "Early Detection of Depression: Social Network Analysis and Random




- Forest Techniques,” *Journal of Medical Internet Research*, vol. 21, no. 6, p. e12554, 2019, doi: 10.2196/12554.
- [16] J. Wawira Gichoya, L. G. McCoy, L. A. Celi, and M. Ghassemi, “Equity in essence: a call for operationalising fairness in machine learning for healthcare,” *BMJ Health & Care Informatics*, vol. 28, no. 1, p. e100289, 2021, doi: 10.1136/bmjhci-2020-100289.
- [17] Y. Li, “Depression and Suicide Risk Prediction Based on Machine Learning Models,” *Journal of Education, Humanities and Social Sciences*, vol. 15, pp. 302–307, 2023, doi: 10.54097/ehss.v15i.9312.
- [18] A. Biradar and S. G. Totad, “Detecting Depression in Social Media Posts Using Machine Learning,” 2019, *Springer Singapore*. doi: 10.1007/978-981-13-9187-3_64.
- [19] A. M. Almars, “Attention-Based Bi-LSTM Model for Arabic Depression Classification,” *Computers, Materials & Continua*, vol. 71, no. 2, pp. 3091–3106, 2022, doi: 10.32604/cmc.2022.022609.
- [20] M. S. Zulfiker, N. Kabir, A. A. Biswas, T. Nazneen, and M. S. Uddin, “An in-depth analysis of machine learning approaches to predict depression,” *Current Research in Behavioral Sciences*, vol. 2, p. 100044, 2021, doi: 10.1016/j.crbeha.2021.100044.
- [21] U. M. Haque, E. Kabir, and R. Khanam, “Detection of child depression using machine learning methods,” *PLOS ONE*, vol. 16, no. 12, p. e0261131, 2021, doi: 10.1371/journal.pone.0261131.
- [22] A. Priya, S. Garg, and N. P. Tigga, “Predicting Anxiety, Depression and Stress in Modern Life using Machine Learning Algorithms,” *Procedia Computer Science*, vol. 167, pp. 1258–1267, 2020, doi: 10.1016/j.procs.2020.03.442.
- [23] A. Sau and I. Bhakta, “Predicting anxiety and depression in elderly patients using machine learning technology,” *Healthcare Technology Letters*, vol. 4, no. 6, pp. 238–243, 2017, doi: 10.1049/htl.2016.0096.
- [24] M. R. Islam, M. A. Kabir, A. Ahmed, A. R. M. Kamal, H. Wang, and A. Ulhaq, “Depression detection from social network data using machine learning techniques,” *Health Information Science and Systems*, vol. 6, no. 1, 2018, doi: 10.1007/s13755-018-0046-0.
- [25] S. Alghowinem, T. Gedeon, R. Goecke, J. F. Cohn, and G. Parker, “Interpretation of Depression Detection Models via Feature Selection Methods,” *IEEE Transactions on Affective Computing*, vol. 14, no. 1, pp. 133–152, 2023, doi: 10.1109/taffc.2020.3035535.
- [26] C. Lee and H. Kim, “Machine learning-based predictive modeling of depression in hypertensive populations,” *PLOS ONE*, vol. 17, no. 7, p. e0272330, 2022, doi: 10.1371/journal.pone.0272330.
- [27] S. Gupta, L. Goel, A. Singh, A. Prasad, and M. A. Ullah, “Psychological Analysis for Depression Detection from Social Networking Sites,” *Computational Intelligence and Neuroscience*, vol. 2022, pp. 1–14, 2022, doi: 10.1155/2022/4395358.
- [28] M. D. Nemesure, M. V. Heinz, R. Huang, and N. C. Jacobson, “Predictive modeling of depression and anxiety using electronic health records and a novel machine learning approach with artificial intelligence,” *Scientific Reports*, vol. 11, no. 1, 2021, doi: 10.1038/s41598-021-81368-4.
- [29] V. P. C. Magboo and M. S. A. Magboo, “Important Features Associated with Depression Prediction and Explainable AI,” 2022, *Springer International Publishing*. doi: 10.1007/978-3-031-14832-3_2.
- [30] M. Z. Uddin, “Phenomenology of depression and revision of depression scale in Bangladesh,” Ph.D. dissertation, University of Dhaka, 2022.
- [31] R. Bharti, A. Khamparia, M. Shabaz, G. Dhiman, S. Pande, and P. Singh, “Prediction of Heart Disease Using a Combination of Machine Learning and Deep Learning,” *Computational Intelligence and Neuroscience*, vol. 2021, no. 1, 2021, doi: 10.1155/2021/8387680.
- [32] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “SMOTE: Synthetic Minority Over-sampling Technique,” *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002, doi: 10.1613/jair.953.
- [33] H. Jiang *et al.*, “Detecting Depression Using an Ensemble Logistic Regression Model Based on Multiple Speech Features,” *Computational and Mathematical Methods in Medicine*, vol. 2018, pp. 1–9, 2018, doi: 10.1155/2018/6508319.
- [34] C.-T. Wu, D. G. Dillon, H.-C. Hsu, S. Huang, E. Barrick, and Y.-H. Liu, “Depression Detection Using Relative EEG Power Induced by Emotionally Positive Images and a Conformal Kernel Support Vector Machine,” *Applied Sciences*, vol. 8, no. 8, p. 1244, 2018, doi: 10.3390/app8081244.
- [35] E. O. Ogunseye, C. A. Adenusi, A. C. Nwanakwagwu, S. A. Ajagbe, and S. O. Akinola, “Predictive Analysis of Mental Health Conditions Using AdaBoost Algorithm,” *ParadigmPlus*, vol. 3, no. 2, pp. 11–26, 2022, doi: 10.55969/paradigmplus.v3n2a2.
- [36] V. Arun, P. V., M. Krishna, A. B.V., P. S.K., and S. V., “A Boosted Machine Learning Approach For Detection of Depression,” 2018 *IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 41–47, Nov. 2018. doi: 10.1109/ssci.2018.8628945.
- [37] H. AlSagari and M. Ykhlef, “Quantifying Feature Importance for Detecting Depression using Random Forest,” *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 5, 2020, doi: 10.14569/ijacsa.2020.0110577.
- [38] M. R. Islam, A. R. M. Kamal, N. Sultana, R. Islam, M. A. Moni, and A. ulhaq, “Detecting Depression Using K-Nearest Neighbors (KNN) Classification Technique,” 2018 *International Conference on Computer, Communication, Chemical, Material and Electronic Engineering (IC4ME2)*, pp. 1–4, Feb. 2018. doi: 10.1109/ic4me2.2018.8465641.
- [39] M. Squires *et al.*, “Deep learning and machine learning in psychiatry: a survey of current progress in depression detection, diagnosis and treatment,” *Brain Informatics*, vol. 10, no. 1, 2023, doi: 10.1186/s40708-023-00188-6.

BIOGRAPHIES OF AUTHORS






S. M. Rakibul Islam     completed his B.Sc. in Educational Technology and Engineering from University of Frontier Technology, Bangladesh. His research passions lie in the realms of machine learning, deep learning, explainable AI, blended learning and educational technology with a strong focus on curriculum design and development. He is actively engaged in research collaborations with Bangladeshi professors and has worked on projects such as Alumnex (a digital transformation for alumni engagement), and future-focused curriculum development in Bangladesh. He can be contacted at email: smrakibulislam34@gmail.com.






Shaykh Yunus    is currently pursuing a B.Sc. in Educational Technology and Engineering from University of Frontier Technology, Bangladesh. He has a keen interest in data science and machine learning, continuously exploring innovative applications in these fields. His passion lies in leveraging technology to enhance education and develop intelligent systems. He actively engages in research and projects related to data-driven decision-making, artificial intelligence, and modern educational methodologies. With a strong foundation in engineering and technology, he aspires to contribute to advancements in AI-driven learning systems and data-centric solutions. He can be contacted at email: shaykhyunus2000@gmail.com.






Rashiduzzaman Shakil    received his bachelor of science (B.Sc.) degree in the Department of Computer Science and Engineering at Daffodil International University, Bangladesh, in 2023. He predominantly works on machine learning, deep learning, and image processing and biomarker analysis in the application of the health domain. He has been working as a collaborator on a research project with researchers in Bangladesh and Malaysia. His numerous research papers were published in the prestigious journals (Scopus) and conferences (Scopus). He can be contacted at email: rashiduzzaman.diucse@gmail.com.






Fatema Tuz Johora    completed her B.Sc. in Educational Technology and Engineering from the University of Frontier Technology, Bangladesh, achieving the distinction of graduating first in her department. She is currently serving as a lecturer in the Department of Computer Science and Engineering at Daffodil International University. She has a deep passion for exploring cutting-edge technologies and modern education systems. Her current research interests include machine learning, explainable AI, computer vision, e-learning, curriculum development, and educational technology. She can be contacted at email: meem8494@gmail.com.



Aditya Rajbongshi    is currently working at University of Frontier Technology, Bangladesh as Assistant Professor in the Department of Educational Technology and Engineering. His areas of interest in research are computer vision, image processing, machine learning, and deep learning. He has had over 30 research articles accepted for publication in journals and conferences worldwide. He reviewed research articles for several international journals and conferences. He can be contacted at email: aditya0001@uftb.ac.bd.



Sujon Chandra Sutradhar    was born in 1996 in Dohar, Dhaka. He graduated from the University of Dhaka with a B.Sc. and an M.Sc. (Thesis) in Applied Mathematics. He is currently employed with University of Frontier Technology, Bangladesh as a mathematics lecturer. Prior to this he was a lecturer of mathematics in the department of EEE at South-East University. Besides that, He worked as a faculty (Adjunct) in BRAC university and National Institute of Textile Engineering and Research for a semester. He is linked to one funded UGC research project and several research articles and has one journal article published in the area of option pricing. His areas of interest include data analysis, machine learning, time series analysis, stochastic modeling in finance, and related fields. He can be contacted at email: sujon0001@uftb.ac.bd.