Project for the Degree of Bachelor of Science in CSE

# Mental Health Analysis and Prediction Using Machine Learning Approach

## M. Shafwan Jarif
## ID: B180305054

Department of Computer Science and Engineering

Jagannath University

Dhaka - 1100, Bangladesh

May, 2024

# Project Title

Mental Health Analysis and Prediction Using Machine Learning Approach

### Supervised by

## Dr. Mohammed Nasir Uddin

## Professor

## Jagannath Unibersity, Dhaka

**Submitted to the Department of Computer Science and Engineering of Jagannath University in partial fulfillment of the requirements for the degree of B.Sc. in CSE**

### Thesis/Project Evaluation Committee:

**Examiner 1** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

**Examiner 2** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

**Examiner 3** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

**Examiner 4** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

# Project Approval

## Mental Health Analysis and Prediction Using Machine Learning Approach

Student's Name:M. Shafwan Jarif

ID : B180305054

We the undersigned, recommend that the project completed by the student (s)

listed above, in partial fulfillment of B.Sc. in CSE degree requirements, be

accepted by the Department of Computer Science and Engineering, Jagannath

University for deposit

### Supervisor Approval

.............................

Dr. Mohammed Nasir Uddin

Professor, CSE, JnU

### Departmental Approval

.............................

Professor Dr. Uzzal Kumar Acharjee

Chairman

Department of CSE

**Jagannath University**

**Dhaka - 1100, Bangladesh**

*Dedicated to Our parents and Teachers*

# Abstract

Mental health, an integral part of overall wellbeing include disorders that interfere with the normal flow of a person's personal, social and professional life. Due to which predictive methods are needed for its efficient diagnosis at present. The World Health Organization (WHO) estimates that one in four people globally suffers from mental health conditions. Moreover, from 1999 to 2019, the global prevalence of mental health issues increased by 50 percent was found from a report of the Global Burden of Disease Study. According to the World Health Organization (WHO), mental health illnesses or issues are characterised by a mix of aberrant thoughts, feelings, behavior in day-to-day activities, and interpersonal connections.. In the initial stage of this research project, we analyzed exploratory data. These analyzes provide important insights into the distribution and interrelationships of variables and assist in the identification of potential predictors and risk factors associated with mental health outcomes as well as aiding in the identification of potential predictors. Determining the distribution of the data, analyzing the correlation matrix of the variables, and using visualization techniques like box and count plots are some of the important tasks involved in data analysis. 4 distinct ML models are applied to predict mental health outcomes: Multilayer Perceptron (MLP), XGBoost, Naive Bayes (NB) and Random Forest. Among all of these models, the MLP and Random Forest models return the highest accuracy (Both having 90 percent).

**Key words: Mental Health, Correlation Matrix, Box Plots, Count plots, Multilayer Perceptron (MLP), XGBoost, Naive Bayes (NB).**

# Acknowledgement

# Table of Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

## 1.1 Introduction

Mental health problems are common worldwide including changes in mood, personality, inability to cope with daily problems or stress, withdrawal from friends and activities, and so on. In 2010 mental health problems were the leading causes of years lived with disability (YLDs) worldwide with depressive and anxiety disorders among the most frequent disorders [1]. Polanczyk et al. estimated a worldwide prevalence of mental disorders in children and adolescents of 13.4 % [2]. In the last years, their prevalence has increased even further [3], [4]. Dealing with mental health issues can have an impact not just, on the individuals directly involved but also on their families and the wider community. Nearly 1 billion people worldwide live with a mental disorder [1]. Timely treatment can prevent exacerbating the symptoms that lead to such crises and subsequent hospitalization [5]. But most of the time, people have already reached an emergency stage by the time they seek professional assistance for mental health issues. Corrective action cannot be particularly important in these situations. Consequently, managing a patient's mental health concerns requires effective early detection of mental health issues. Because machine learning can detect mental health disorders with high accuracy and speed, it can be highly useful in this regard. Additionally, it has recently been noted that a sizable percentage of students have a variety of issues linked to mental health issues. The majority of students grumble about the high level of stress they experience in their university lives, including feelings of anxiety and depression, especially towards the end of the semester [6]. The level of stress increases as the learning process pro-

gresses due to the need to balance assessments, workload, and examinations [3]. Other factors may also affect students' mental health. Students may face a high risk of developing mental health problems due to family issues, uncertainties about their future careers, financial troubles and difficulties arising out of living away from home [4]. Balancing between life at university and other demands or needs can also lead the students to face the risk of developing mental health problems [1]. Students experiencing symptoms of mental health problems have claimed that they are not receiving any treatments and would not seek help to address their emotional troubles. These students do not place any importance on their predicament as their peers also experience similar symptoms, and thus they see this as something common in their university lives [7]. However, some of them are aware of the need for proper treatment, but they lack the courage to seek help and worry too much about other people's perceptions [6]. They fear that the stigma of being diagnosed with mental health problems may lead to discrimination or prejudice by society, and they worry about the negative impact of being labelled sick, overly emotional or crazy [8]. Therefore, a new, safe and reliable strategy to solve human health problems is urgently needed. The use of machine learning (ML) methods has shown potential in times for enhancing the prediction and analysis of health. ML models provide the ability to analyze amounts of data identify patterns and pinpoint indicators linked to different mental health outcomes. By applying machine learning techniques such, as Random Forest, XGBoost, Naive Bayes (NB) Multilayer Perceptron (MLP) and others we can develop models that accurately detect mental health issues. In this research project, we explored the application of machine learning in analyzing and predicting health. We start by conducting an examination of the data to uncover connections, between factors to understand the elements that impact mental health results. Subsequently we employ machine learning models to forecast mental health issues and evaluate their effectiveness in recognizing individuals who may need support.

## 1.2    Motivation

This thesis is motivated by a sincere desire to work together with advanced technologies and mental health treatment. Conventional approaches to mental health diagnosis and treatment frequently have limitations that result in incorrect diagnoses and inadequate care for those who require it. We tried to close these gaps by utilizing machine learning to create more intelligent and proactive approaches to comprehending and anticipating mental health issues. Our ultimate goal is to use data-driven insights to improve people's quality of life, provide better tools for intervention to healthcare providers, and ultimately help build a more efficient and compassionate mental health care system.

## 1.3    Problem Statement

The issue revolves around correctly predicting mental health problems. Subjective evaluations of the topic in our hand make it difficult to diagnoses mental health issues properly. Machine learning models face challanges predicting mental disorders for several factors such as: data quality, dynamic nature of mental health issues, diversity of symptoms. On the other hand, the traditional methods of treating mental health problems face barrier due to a good number of reasons. The lack of trained professionals, a patient's reluctance to disclose the symptoms and bias can lead to misdiagnoses. Traditional approaches often address mental health issues reactively after they have already become apparent, which may be less effective in preventing the progression of these issues.

## 1.4    Objectives

To address the above issues, this project aims to illustrate how Mental Health issues may be efficiently predicted, as well as analyzed properly using Machine Learning algorithms and techniques. To achieve this goal, we attempted to build machine

learning models that can correctly predict mental health issues and analyze them with the help of EDA. Our developed models have adapted the following strategies:

- Derived Feature Creation: Two derived features named "CGPA Midpoint" and "Study Year" were engineered from the "Cgpa" and "Your Current Year of Study" columns respectively . These two derived features were created with an aim to capture more precise representation of the data.

- Binary feature creation: A new binary feature named 'Total Mental Health Issues' is added which is basically a combination of "Depression", "Anxiety" and "Panic Attack" columns of the dataset. This process summarizes the values of these binary indicators for each respondent, providing an overall representation of the mental health problems they may be experiencing.

- To detect, mental disorders accurately and efficiently, the predictive models worked on the data based on the discussed feature selection technique.

## 1.5    Contributions of the Project

In order to understand the function of machine learning in the diagnosis of mental disorders and other mental health-related concerns, we read a number of research articles and scientific journals. The primary contributions of the research work are discussed next:

- In order to obtain a sense of the dataset, we engaged in exploratory data analysis (EDA). We closely examined variables such as age, education level, year in school, CGPA, and several mental health markers. An intriguing discovery revealed that the participants ranged in age from 18 to 24 years old, with an average age of approximately 20.53 years.

- We looked at the relationships between numerical from the dataset, and visualized these interactions using boxplots and histograms to help make sense of

these linkages and help understand the links between these variables.

- We used machine learning techniques like Random Forest, XGBoost, Multi-layer Perceptron (MLP), Naive Bayes to build predictive models. These models were trained on prepared data to predict the likelihood and severity of mental health issues among students. To assess their effectiveness, we examined classification reports, which gave us details on precision, recall, and F1-score for each class. Our findings showed that different algorithms performed differently in predicting mental health issues.

## 1.6    Organization of the Project

- **Chapter 1 Introduction:** This chapter include opening remarks of Mental Health problems and challenges that comes with dealing with mental health problems. It also includes the motivation behind this work, the problem statement, contribution and objective of this research work.

- **Chapter 2 Background Study:** This section presents the theoretic background and role of EDA techniques and Machine Learning algorithms..

- **Chapter 3 Related Works:** This chapter analyzes different works, limitations with existing methods in relation to the issue and the comparison of the final outcomes of different models.

- **Chapter 4 Proposed Methodology:**This section holds the most importance as it contains the explanation of the structure of the suggested model.

- **Chapter 5 Result Analysis:** The parameters of the result analysis and a comparison of the outcomes from various machine learning models are included in this chapter.

- **Chapter 6 Conclusion and Future Work:** This chapter concludes the dissertation indicating the limitations and future works.

# Chapter 2

# Background Study

## 2.1 Depression

Depression is characterized by constant sadness, loss of interest or excitement, feelings of guilt or low self-worth, disturbed sleep, loss of appetite, fatigue, and inability to concentrate [1]. As to the National Institute of Mental Health, clinical depression, sometimes known as depression, is a grave mood disorder characterized by intense symptoms that impact an individual's mood, thoughts, and ability to carry out everyday tasks [9]. Depression can cause pain to the person suffering from the ailment and the people around them. It can be a serious health concern as it may lead to suicide [10]. Here's a key observation: we should remember that the understanding of the definition of depression is crucial since understanding its dimensions and definition will help us understand it better. Depression, sometimes referred to as major depressive disorder, clinical depression, or severe depression, is a common but dangerous mood disease, according to the National Institute of Mental Health. It causes severe symptoms that affect how a person feels, thinks, and handles daily activities, such as sleeping, eating, or working. And it's from this point on that we need to distinguish between despair and irritation or the sense of frustration. The National Institute of Mental Health states that a patient's symptoms must persist for at least two weeks in order for a diagnosis of depression to be made. Furthermore, although practically everyone experiences melancholy, inattention, self-defeating thoughts, and other similar emotions occasionally, these feelings cannot be readily classified as depression. The National Institute of Mental Health has also offered assistance in this instance by providing a list of symptoms of depression; we will go into that in a

later section. Depression is a time-dependent subject. Patients might be categorized according to the length of time they have experienced depression. Furthermore, it is evident that the patient is more affected by it over time. There are several types of depression, and a few of them are persistent depressive disorder, postpartum depression and psychotic depression [11]. Persistent depressive disorder, also known as dysthymia, is a state of low mood that lasts for at least two years [9]. A person who is diagnosed with persistent depressive disorder may have major depressive episodes along with periods of less severe symptoms, but signs must last more than two years in order to be considered persistent depressive disorder [9].Psychotic depression is characterized by severe depression combined with psychosis, such as the sensation of disturbingly erroneous ideas or unsettling phenomena perceived by the sufferer that are not perceived by others. The symptoms of psychotic depression typically have a grim "theme," such as delusions of guilt, poverty, or illness [13]. The development of another mood disorder, seasonal affective disorder (SAD), generally happens in winter months when less natural sunlight is available [13]. Winter depression always appears and disappears at the same time each year. It is typically accompanied by social isolation, excessive sleep, and weight gain. When bipolar disordered, a person experiences powerful mood episodes that alternate between an extreme low that satisfies the primary depressive symptoms and an excessive high, commonly referred to as mania, when the individual is either euphoric or irritable. A less severe form of mania is known as hypomania [12].

## 2.2 Anxiety Disorder

Anxiety disorders are characterized by overwhelming worry and fear, especially when confronted with problems or decision-making [13], [14]. Symptoms of intense uneasiness, worry, and excessive dread impact the lives of those with anxiety disorders. Other symptoms, such as heart palpitations, breathing issues, excessive sweating, tremors, or nausea, may also manifest in uncomfortable settings [14]. Anyone can have anxiety problems since they are not limited to any particular ailment or age

group [15]. Four main forms of anxiety disorders have been identified by the National Institute of Mental Health: social anxiety disorder, panic disorder, generalized anxiety disorder, and illnesses connected to phobias. While everyone experiences occasional anxiety, those who suffer from anxiety disorders describe their feelings as being trapped in a never-ending tornado of dread and terror. It's not enough to simply feel a little strange sometimes; anxiety disorders can cause major disruptions to day-to-day functioning. They have a way of interfering with even the most basic duties and can seriously damage relationships with friends, family, and coworkers. According to WHO an estimated 4 percent of the global population currently experience an anxiety disorder. In 2019, 301 million people in the world had an anxiety disorder, making anxiety disorders the most common of all mental disorders. Despite the availability of very effective therapies for anxiety disorders, only roughly one in four individuals in need (27.6percent) receive any kind of care.

## 2.3 MLP

Multilayer Perceptron (MLP) is an Artificial Neural Network (ANN) architecture widely used in machine learning. It is an architecture that is composed of multiple layers of neurons. Here, a neuron is connected to each neuron in its adjacent layer, also this architecture has an activation function that introduces non-linearity. There are basically three main layers in MLP: The input layer, the output layer, and the hidden layer. The input layer is responsible for receiving data, and the output layer is responsible for functions like classification, prediction, etc. Hidden layers are mainly used for feature extraction. Also, hidden layers are used to train the model. Using the back propagation approach, MLP is trained. Here, weight value and loss function are two key ideas. The difference between the predicted and obtained outputs is calculated using the loss function. The model is trained and the loss function is optimized in the back propagation approach by changing the weight values from the output layer to the input layer.
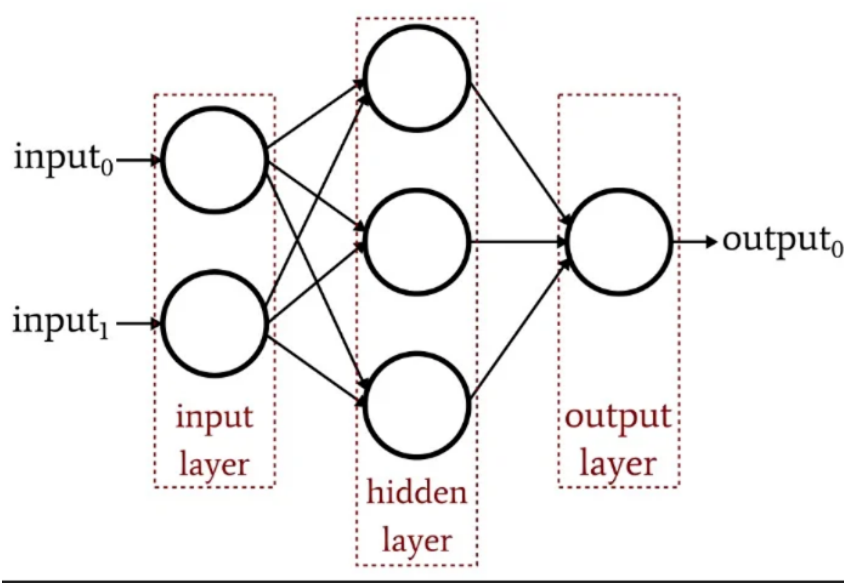
Figure 2.1: General Structure of MLP

## 2.4 XGBoost

Extreme Gradient Boosting, abbreviated as XGBoost, is a distributed gradient boosting library used in machine learning. It is called an ensemble machine learning technique due to its ability to combine the predictions of multiple models to give an accurate prediction. Another way to put it is that the XGBoost technique takes the predictions from several poor models and, by combining them, produces a nearly flawless prediction result, increasing the prediction's accuracy and robustness. Also this process performs well when dealing with large datasets. A notable aspect of this technique is that it has built-in parallel processing support, which takes a reasonable amount of time to train the model. The working method of XGBoost is based on Decision Tree. To enhance the outcomes, techniques such as gradient boosting are applied. The first step in the working procedure of XGBoost is to make a prediction with the objective of fitting the training dataset. Predicted values and observed values are used in the calculation of residuals. A decision tree is constructed based on the similarity score. The next step is to calculate the log of odds and probabilities for classification. Here, the strategy used in the decision tree is repeated many times

Figure 2.2: General Structure of XGBoost

until an efficient result is obtained and in this case each subsequent tree learns from the previous tree and changes the assigned weights.

## 2.5   Naive Bayes

The Naive Bayes method is a set of supervised machine learning algorithms based primarily on the application of Bayes' theorem. The "Naive" part in its name refers to an assumption, where each pair of features is assumed to be conditionally independent. Here, we'll assume that the dependent feature vector is X1 and the class variable is y.The relationship is defined by the "Bayes" theorem as:

$$P(y|X_1, \ldots, X_n) = \frac{P(y) \cdot P(X_1, \ldots, X_n|y)}{P(X_1, \ldots, X_n)} \tag{2.1}$$

Assumption of the Naive conditional independence:

$$P(X_i|y, X_1, \ldots, X_{i-1}, \ldots, X_n) = P(X_i|y) \tag{2.2}$$

After simplification, the classification rule can be stated as follow:

$$P(y|X_1, \ldots, X_n) \propto P(y) \prod_{i=1}^{n} P(X_i|y) \tag{2.3}$$

$$\hat{y} = \arg \max_{y} P(y) \prod_{i=1}^{n} P(X_i|y) \tag{2.4}$$

In practical scenarios, the Naive Bayes approach has demonstrated good performance, despite its reliance on oversimplified assumptions for categorization. Our research revealed that this technique's accuracy was effective. It is quicker than other techniques because of its straightforward implementation and comparatively low processing cost.

## 2.6   Random Forest

Random Forest is a widely used machine learning algorithm. Its strong classification and regression performance have contributed to its rise in popularity. It employs decision trees, just as our previously covered XGBoost algorithm, and condenses several outcomes into a single result. Here, the fundamental distinction is in how the decision tree is created and used. The Random Forest algorithm builds numerous decision trees continuously in parallel and combines the outcomes of each one, whereas the XGBoost algorithm builds a sequential decision tree and attempts to enhance it by leveraging the previous decision tree's findings at each iteration. Every decision tree is trained independently in the Random Forest method. For every decision tree in this instance, a random subset of features is chosen to train the tree.

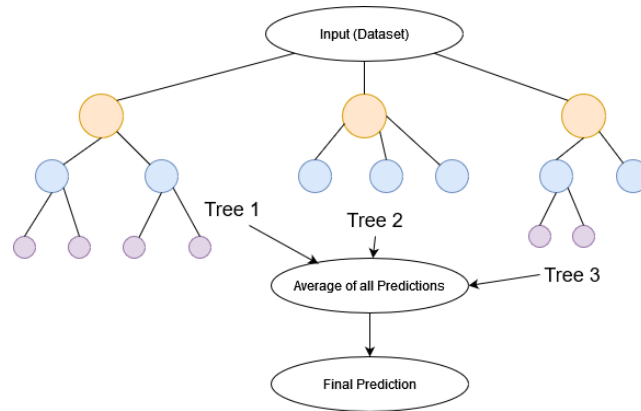Figure 2.3: General Structure of Random Forest

This algorithm uses two final steps in making predictions. First, an average of the predictions obtained from each tree is taken. Then, in a final step, the Random Forest algorithm uses these averages to produce a final prediction. Random Forest is known for its high accuracy, robustness, and ability to handle large datasets with high dimensionality.

# Chapter 3

# Related Works

Several studies have investigated the application of various machine learning models for predicting mental health outcomes using different sets of variables. Sofianita Mutalib [1] explored the predictive capabilities of Decision Trees (DT), Multi-Layer Perceptron (MLP), Naive Bayes (NB), Support Vector Machine (SVM), and Logistic Regression (LR) models utilizing predictors such as gender, age, program, part, cumulative grade point average (CGPA), and financial support. Their findings demonstrated varied accuracies across models, with Decision Trees achieving the highest accuracy at 84.44 percent.

Ryan C. McCabe et al. [4] investigated the predictive performance of Random Forest (RF), Support Vector Machine (SVM), XGBoost, Neural Networks (NN), and Logistic Regression (LR) models using symptom-related variables. Although the accuracies varied, Random Forest and SVM models exhibited comparable performances, both achieving accuracies above 70 percent.

Fadhluddin Sahlan et al. [16] focused on a reduced set of features including age, gender, physical health, and personal traits to predict mental health outcomes using Decision Trees (DT), K-Nearest Neighbors (KNN), and Support Vector Machine (SVM) models. However, they reported relatively low accuracies ranging from 44.00 percent to 64.00 percent.

Ashley A Sabourin et al. [17] evaluated Random Forest (RF), Naive Bayes (NB), Support Vector Machine (SVM), and K-Nearest Neighbors (KNN) models using the Perceived Stress Scale (PSS) questionnaire. Their study revealed promising performance with SVM achieving the highest accuracy of 85.71 percent.

Fenfen Ge et al. [18] examined Random Forest (RF), Decision Trees (DT), Multi-Layer Perceptron (MLP), and Support Vector Machine (SVM) models to predict mental health outcomes based on variables such as age, Mini-Mental State Examination (MMSE) score, neurological conditions, depression (GDS), and MoCA test. Their results indicated exceptionally high accuracies, particularly with RF and SVM models achieving accuracies of 100 percent and 99.5 percent, respectively.

Furthermore, A. A. Choudhury et al. [19] assessed K-Nearest Neighbors (KNN), Random Forest (RF), and Support Vector Machine (SVM) models using a questionnaire comprising 55 questions. Their findings highlighted the effectiveness of SVM with an accuracy of 80.20 percent.

V. Laijawala et al. [20] examined the performance of Decision Trees (DT), Random Forest (RF), and Naive Bayes (NB) models utilizing survey data. While Decision Trees achieved the highest accuracy at 82.02 percent, all models demonstrated competitive performances.

Chekroud et al. [21], Gradient Boosting (GB), Naive Bayes (NB), Logistic Regression (LR), and Multi-Layer Perceptron (MLP) models were employed to predict depression levels in a large cohort of patients. Despite Gradient Boosting demonstrating relatively low accuracy, other models exhibited competitive performances, with Naive Bayes achieving the highest accuracy at 79.6 percent.

K. M. Mitravinda et al. [22] compared the performance of various models including Logistic Regression (LR), Decision Trees (DT), Random Forest (RF), K-Nearest Neighbors (KNN), XGBoost, and Gradient Boosting (GB). Their results demonstrated the effectiveness of XGBoost and GB models, both achieving accuracies above 90 percent.

Tate et al. [23] investigated the predictive capabilities of Random Forest (RF), Support Vector Machine (SVM), Neural Networks (NN), Logistic Regression (LR), and XGBoost models using data from the child and adolescent Twin Study in Sweden. While the accuracies varied, all models displayed competitive performances, with Random Forest achieving the highest accuracy at 79.00 percent.

Table 3.1: Summary of Related Works

| Author | Variables | Accuracy | Limitations |
| --- | --- | --- | --- |
| Sofianita Mutalib [1] | Gender, Age, Program, CGPA, Financial Support | DT: 84.44, MLP: 80.00, NB: 74.81, SVM: 82.22, LR: 82.96 | Reduced Feature Set |
| Ryan C. McCabe, et al. [4] | Symptoms (Impulsivity, Inattention, Emotional, OD), Dysfunctions, Difficulties | RF: 73.90, SVM: 73.60, NN:70.50 ,LR: 70.00, XG-Boost: 69.20 | Low accuracy |
| Fadhluddin Sahlan, et al. [16] | Age, Gender, Physical health, Personal traits | DT: 64.00, KNN: 59.00, SVM: 44.00 | Low accuracy |
| Ashley Sabourin, et al. [17] | Perceived Stress Scale (PSS) questionnaires | DT: 99.4, SVM: 99.5, RF: 100, MLP: 98.70 | |
| Fenfen Ge, et al. [18] | Age, MMSE Score, Neurological condition, Depression, MoCA test | DT: 99.4, SVM: 99.5, RF: 100, MLP: 98.70 | Non-standardized measures |
| A. A. Choudhury, et al. [19] | Questionnaire consisting of 55 questions | RF: 75.00, SVM: 80.20, KNN: 60.00 | Use of limited dataset |
| V. Laijawala, et al. [20] | Survey questions | DT: 82.02, RF: 79.3, NB: 78.7 | Non specific mental illness targeted |
| Chekroud, et al. [21] | 1949 patients with level 1 of depression | GB: 64.6, NB: 79.6, LR: 72.4, MLP: 77.8 7 | GB method with low accuracy |
| Tate, et al. [23] | 7638 twins from the child and adolescent Twin Study in Sweeden | RF: 73.9, SVM: 73.6, NN: 70.05, LR: 70.00, XG-Boost: 79.00 | The performances of the models are very close to each other |

# Chapter 4

# Proposed Methodology

The section Proposed Methodology mostly addresses the project work's systematic approach. The architecture of the suggested methodology for analysis and prediction of mental health is shown in Figure 4.1; data collection, data preprocessing, exploratory data analysis, feature engineering, modeling, and model evolution are the key components of our suggested approach. Initially, the steps data collection and preprocessing were completed. To help with the next critical stage of feature engineering, exploratory data analysis was carried out in order to comprehend the distribution and interactions between the datasets. Next, a few machine learning models were developed with the intention of making prediction. Using the attributes accuracy, recall, F1 score, and support, we compare each model's performance in the model evaluation stage, which is the final step in our proposed approach, analyzing the performance of these constructed models. All things considered, the methodology emphasizes a thorough and moral approach to examining mental health issues.

## 4.1 Data Collection

The dataset used in this project work titled "Student Mental Health" was sourced from Kaggle. The dataset includes mental health-related factors for students. Gender, age, education, marital status, year of study at the moment, and CGPA are the variables. It is known from the dataset's source that a survey was used to gather the data. This dataset was chosen for analysis in our research project with the aim of creating prediction models.
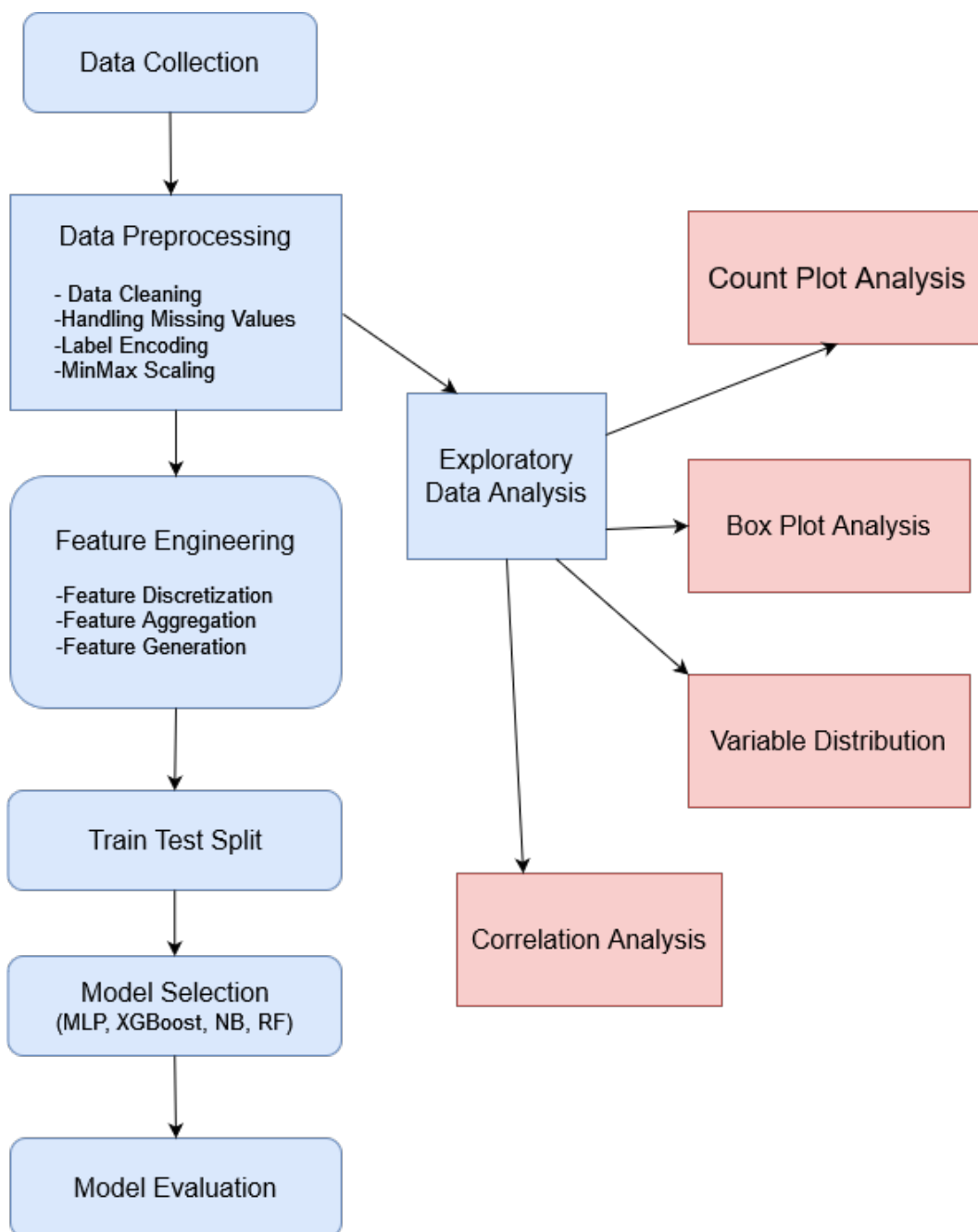
Figure 4.1: Proposed Methodology

## 4.2    Dataset Description

There are 130 occurrences in the dataset we used for this research work. As per the poll used to create this dataset, 66 individuals have been diagnosed with depression, whereas 35 persons are stated to be devoid of the disorder. Regarding anxiety, the findings indicate that 67 individuals are positive, meaning they experience anxiety, while 34 individuals do not experience anxiety. Two target variables and eight variables total in this dataset.

| Gender | Age | Course | Study Year | CGPA | Marital status | Depression | Anxiety |
|---|---|---|---|---|---|---|---|
| Female | 18 | Engineering | year 1 | 3.00 - 3.49 | No | Yes | No |
| Male | 21 | Islamic edu. | year 2 | 3.00 - 3.49 | No | No | No |
| Male | 19 | BIT | Year 1 | 3.00 - 3.49 | No | Yes | No |
| Female | 22 | Laws | year 3 | 3.00 - 3.49 | Yes | Yes | No |
| Male | 23 | Mathematics | year 4 | 3.00 - 3.49 | No | No | No |
| Male | 19 | Engineering | Year 2 | 3.50-4.00 | No | No | Yes |
| Female | 19 | HRM | Year 2 | 2.50 - 2.99 | No | No | No |
| Male | 18 | Psychology | Year 1 | 3.50-4.00 | No | No | Yes |
| Female | 20 | Engineering | year 2 | 3.00-3.49 | No | Yes | Yes |
| Female | 24 | Accounting | Year 4 | 3.00-3.49 | No | No | No |
| Female | 18 | BRM | Year 1 | 3.50-4.00 | No | No | Yes |
| Male | 19 | BIT | Year 2 | 2.50-2.99 | No | Yes | Yes |
| Male | 18 | Marine | year 1 | 3.50-4.00 | No | No | Yes |

Table 4.1: Used Dataset for Our Project Work

## 4.3    Data Preprocessing

Data preparation stage is crucial to ensuring that the dataset is appropriate for training machine learning models. This research project's preprocessing procedures

comprised several important procedures. The dataset was first loaded using the pandas library, which made data manipulation and analysis easy. Having clean and noiseless data is essential when working with ML models. This involves addressing missing values and maintaining consistent, error-free, and null-value-free data in the dataset. The model's training process is greatly impacted by each of these processes, which in turn affects the model's prediction accuracy. To increase the dataset's correctness as much as feasible was our aim.

Data Cleaning: Data cleaning commenced with the systematic identification and rectification of any inconsistencies, errors, or extraneous data points within the dataset. Duplicate entries were identified and removed, typographical errors were corrected, and outliers were addressed, thereby fortifying the dataset's reliability and coherence. This preliminary phase laid the foundation for robust and dependable analysis.

Handling Missing Values: Handling missing values constituted a critical facet of data preprocessing. Through meticulous examination, missing data points were discerned and subsequently addressed using prudent strategies embedded within the code. Numerical features were imputed with mean or median values, while categorical features were imputed with mode values, ensuring the dataset's completeness and coherence while mitigating potential biases.

Label Encoding: In order to carry out the label encoding in this research work, we employed the sklearn.preprocessing library. With the use of this library's fit transform method, numerical data can be transformed from categorical values. By assigning a unique integer value to every category that exists inside a column, this encoding process effectively transforms categorical data into a format that can be used for mathematical operations and model training. In order to employ the numerical representation of the category values from the original dataset, these

modified data were later reintegrated into the data frame.

MinMax Scaling: MinMax Scaling is a well received feature engineering method. It reduces a dataset's features to a predetermined range, usually between 0 and 1. The primary benefit of the MinMax Scaler is in its ability to maintain the original distribution's structure while bringing the values within the intended range. The MinMax scaller method also made advantage of the sklearn preprocessing library. This approach made sure that no variable dominated the others and that every variable contributed equitably to the analysis.This technique plays a crucial role in enhancing the performance and robustness of machine learning models by improving convergence speed and reducing the sensitivity to the scale of input features. Additionally, MinMax scaling aids in optimizing the performance of algorithms that rely on distance measures or gradient descent for optimization. Through the implementation of MinMax scaling, our research endeavors to improve the efficacy and reliability of predictive models in analyzing mental health issues, facilitating more accurate assessments and interventions.

## 4.4 Feature Engineering

Feature engineering involves gathering valuable features from a dataset, or collection of data. Raw data is transformed into relevant data using this strategy. The main goal of any model in machine learning is to make predictions with a good accuracy depends on efficient feature engineering .It increases machine learning's capacity for prediction and assists in revealing hidden patterns in the data. Feature engineering involve tasks of selecting important features, modifying the features of the original dataset and generating new features. We have attempted to use efficient feature engineering in our research work. The steps in feature engineering our research work have included are: Feature Discretization, Feature Aggregation and Feature Interaction.

Feature Discretization: Feature discretization involves converting continuous variables into categorical or ordinal variables, thereby enabling the models to capture non-linear relationships and underlying patterns more effectively. Feature discretization was systematically executed for the "Age" variable to categorize respondents into distinct age groups, thereby enriching the dataset with more informative features. Implemented through a systematic binning process, the continuous "Age" variable underwent transformation into discrete age groups, each encapsulating a range of ages with similar characteristics or patterns related to mental health outcomes. This binning process was meticulously orchestrated, ensuring the creation of meaningful and interpretable age groups tailored to the research work's context.

Feature Aggregation: The process of feature aggregation emerged as a fundamental technique aimed at consolidating multiple indicators related to mental health into composite features, thereby enhancing the predictive capacity and interpretability of the machine learning models. Feature aggregation involves the amalgamation of individual variables or attributes to create new, more informative features. In our

implementation part, feature aggregation was systematically executed to synthesize binary indicators of depression and anxiety into a unified metric representing the overall mental health status of each respondent. Implemented through a meticulous combination of binary variables, the feature aggregation process entailed summing or combining the presence or absence of depression, anxiety, and panic attacks to generate a composite feature denoting the cumulative burden of mental health issues.

Feature Generation: Feature generation is a critical aspect of feature engineering in machine learning. In order to enhance a machine learning model's performance, it entails generating new features or variables from already-existing data. With the help of these additional characteristics, the data should be able to hold more information, have more representational capacity, and help the model comprehend the underlying patterns. In this research work, new feature generation was conducted to augment the dataset with additional informative variables, thereby enriching the predictive capacity of the machine learning models. This process involved creating novel features derived from existing ones, aimed at capturing nuanced patterns and relationships within the data. Total Mental Health Issues and CGPA Midpoint were generated by aggregating and transforming relevant attributes. Total Mental Health Issues was derived from the sum of individual columns indicating the presence of depression, anxiety, and panic attacks, offering a holistic perspective on mental health status. Similarly, CGPA Midpoint was computed by mapping the range of CGPA values to their respective midpoints, providing a more nuanced representation of academic performance. By introducing these novel features, our research seeks to enhance the predictive accuracy and interpretability of the models, enabling a deeper understanding of the factors influencing mental health outcomes.

## 4.5   EDA (Exploratory Data Analysis)

Exploratory Data Analysis is the process of understanding the relationships, traits and patterns among data. It is also regarded as the initial stage of modeling. The distribution of the data inside a dataset can be found. Additionally, choosing prospective predictive factors is another usage for EDA. In this research project, we conducted exploratory data analysis using the Pandas package. The structure and dimensions of the dataset can be easily ascertained with the help of the Pandas library. Data visualization, which allows for the learning of data-related insights, is one of the main components of EDA. With the use of box plots, count plots, histograms, and other visualization tools, analytical figures are created once the data relationship is well comprehended and the data has been properly examined. In general, it can be said that through exploratory data analysis, a clear understanding of the data can be obtained.

In our project work, the distribution of the variables is first ascertained using exploratory data analysis, and then the data is visualized using a histogram, box plot, or count plot. Additionally, the correlation matrix illustrates how the variables are related to one another. Within the domain of mental health analysis, the impacts of several characteristics on mental well-being are illustrated, perhaps yielding practical insights.

## 4.6   Exploring the Distribution of Variables

The following figures in page 25 shows the distribution of the variables "Age" and "CGPA" attributes. "Age" and "CGPA" provides a deep insight into the decision making process based on the dataset. These attributes were fundamental from the analytical context as they were important predicting factors.

### 4.6.1 Distribution of Age

The distribution of age among the participants in the study reveals a pattern reflective of the wider demographic landscape.
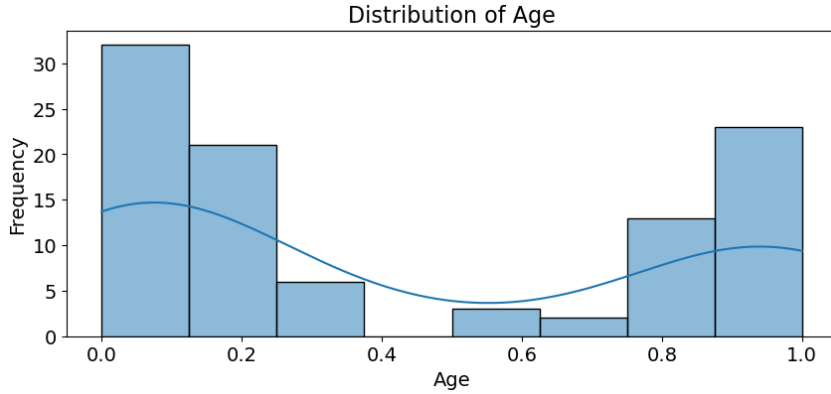


Figure 4.2: Distribution of Age

As depicted in the histogram, the age distribution appears to follow a fairly normal distribution, with the majority of participants falling within the range of 20 to 25 years. This age range aligns with the typical demographic profile of undergraduate and graduate students, who often constitute the primary cohort in academic studies.

### 4.6.2 Distribution of CGPA

The distribution of Cumulative Grade Point Average (CGPA) offers a glimpse into the academic performance landscape of the participants. As evidenced by the box plot, the distribution of CGPA showcases considerable variability, with a median value indicating the central tendency of the dataset. The interquartile range (IQR) highlights the spread of CGPA values within the middle 50 percent of the distribution, illustrating the diversity in academic achievement levels among the participants.
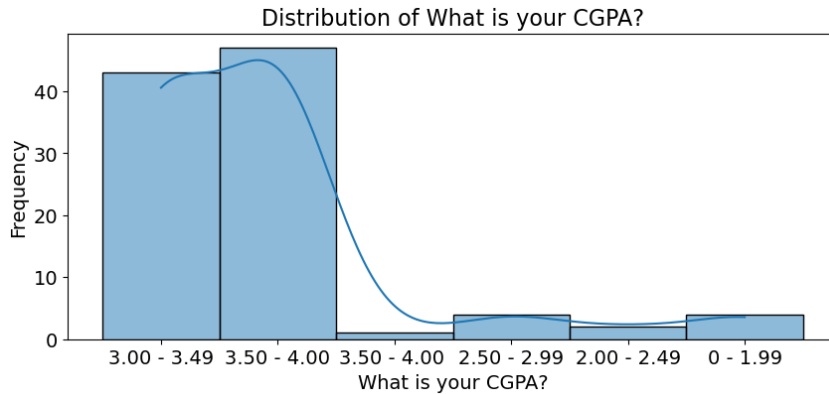
Figure 4.3: Distribution of CGPA

## 4.7 Correlation Matrix of Variables

A correlation matrix provides a comprehensive overview of the relationships between variables in a dataset. It quantifies the degree and direction of linear association between pairs of variables, ranging from -1 to 1. A value close to 1 indicates a strong positive correlation, while a value close to -1 indicates a strong negative correlation.
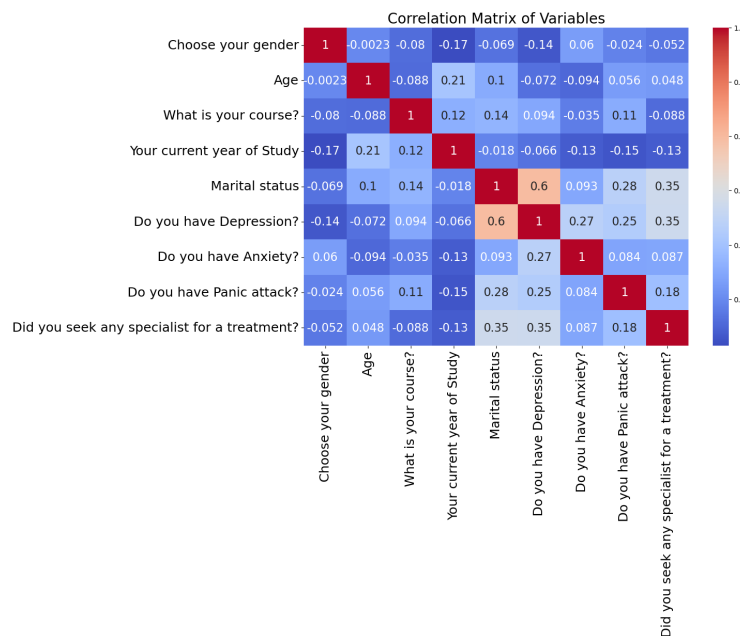


Figure 4.4: Correlation Matrix of Variables

By analyzing the correlation matrix, it is possible to gain insights into the strength and nature of the relationship between variables, which is essential for understanding patterns and dependencies within the dataset.

## 4.8 Boxplot Analysis

In this section, we utilize boxplots to explore the relationship between students' anxiety, depression levels, and CGPA. Boxplots are graphical representations of the distribution of data that provide insights into central tendency, variability, and potential outliers within a dataset.
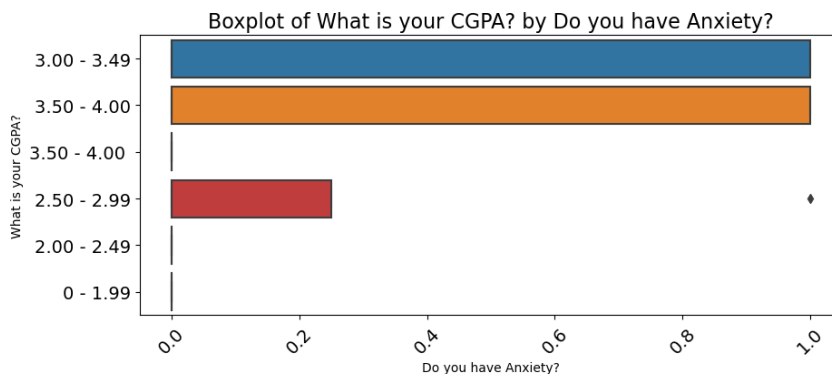


Figure 4.5: Anxiety Vs CGPA Boxplot

The boxplot comparing Anxiety scores with CGPA reveals valuable insights into the potential relationship between students' anxiety levels and their academic performance. The boxplot displays the distribution of CGPA for different levels of anxiety and represents the interquartile range (IQR), with the middle line denoting the median CGPA. Here, the median of "CGPA" decreases or shows variations across different anxiety levels, it suggests a potential correlation between anxiety and academic performance. Similarly, the boxplot comparing Depression scores with CGPA aims to uncover any discernible patterns between students' depression levels and their academic achievements. The boxplot illustrates the distribution of CGPA concerning different levels of depression. Similar to the Anxiety Vs CGPA boxplot, it
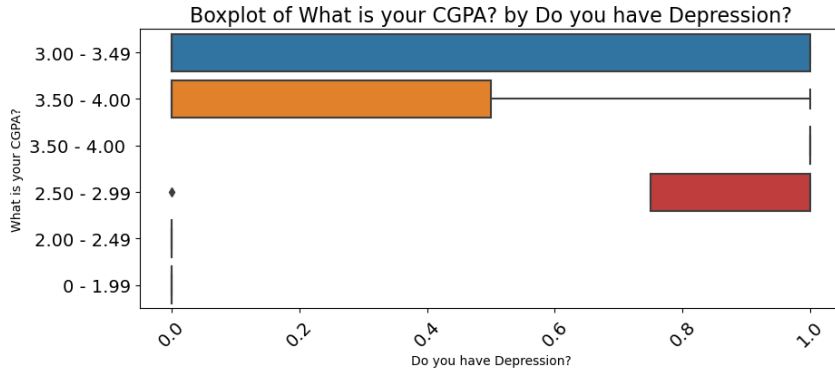
Figure 4.6: Depression Vs CGPA Boxplot

showcases the IQR, median CGPA, whiskers, and potential outliers. Here, variations in median of "CGPA" across depression levels may suggest a correlation between depression and academic performance.

## 4.9    Train Test Split

Train-test split is a fundamental practice in machine learning used to evaluate the performance of predictive models. It involves dividing the available dataset into two subsets: the training set and the testing set. The training set is utilized to train the model on historical data, allowing it to learn patterns and relationships between input features and target variables. On the other hand, the testing set serves as an independent dataset to assess the model's generalization capability by evaluating its performance on unseen data. The train-test split is conducted to simulate the model's performance in real-world scenarios where it encounters new instances beyond the training data. By evaluating the model on unseen data, it provides insights into its ability to make accurate predictions on future observations. This process helps to assess the model's robustness, reliability, and potential for deployment in practical applications. In our implementation, the train-test split was conducted using the train test split function from the sklearn model selection module. This function randomly divided the dataset into two subsets: the training

set and the testing set. The data was partitioned based on a specified ratio, typically 80 percent for training and 20 percent for testing. The random sampling ensured that both sets adequately represented the overall distribution of the data, enabling unbiased evaluation of the model's performance. The training set was used to train the machine learning models, while the testing set was reserved for assessing their predictive accuracy on unseen data. This separation facilitated robust evaluation and estimation of the models' generalization ability.

## 4.10 Model Selection

The selected models for our project work are: MLP, XGBoost, Naive Bayes and Random Forest. In our research project, MLP was considered as a candidate algorithm due to its inherent flexibility, adaptability, and capability to learn from intricate patterns in the data. Its ability to model complex relationships between input features and mental health outcomes was particularly relevant to our research objectives. Furthermore, MLP's scalability and effectiveness in handling large datasets made it a viable choice for analyzing mental health issues among university students. Overall, MLP emerged as a promising algorithm for our predictive modeling tasks, offering the potential to uncover meaningful insights and patterns in mental health data through its sophisticated learning capabilities. The second model, XGBoost offers several advantages, including its ability to handle large datasets efficiently, feature importance estimation, and regularization techniques to prevent overfitting. It utilizes a gradient boosting framework, which sequentially adds new models to correct errors made by existing ones, thereby improving overall predictive accuracy. Moreover, XGBoost incorporates advanced optimization techniques, such as parallelization and tree pruning, to enhance computational speed and model performance. In our research project, XGBoost was selected as a candidate algorithm due to its versatility, scalability, and proven success in various real-world applications. Its flexibility in handling both structured and unstructured data, along with its ability to capture complex relationships within the data, made it well-suited for our mental

health prediction task. Furthermore, its robustness against overfitting and its ability to provide interpretable insights into feature importance were crucial considerations in our model selection process. After this comes the Naive Bayes model. It calculates the probability of a given instance belonging to a particular class based on the joint probability of its features. One of the key advantages of Naive Bayes is its ability to handle high-dimensional data with a relatively small number of training samples. By leveraging probability theory, Naive Bayes calculates the likelihood of each class based on the features present in the data, allowing it to make predictions swiftly and accurately. Moreover, Naive Bayes is robust to irrelevant features and noise in the data, making it suitable for datasets with diverse and noisy attributes. Lastly, we have used the Random Forest model for our project work. Random Forest is highly interpretable, as it provides measures of feature importance based on the information gain or Gini impurity reduction achieved by each feature across all trees in the ensemble. This feature importance analysis enables researchers to identify the most influential variables driving the predictive performance of the model, facilitating insights into the underlying relationships within the data.

# Chapter 5

# Result Analysis

The performance of models for predicting mental health is the main topic of this chapter. Here, the models' accuracy, precision, recall, and F1 score are assessed. Furthermore displayed are the models' performance matrices and confusion matrices.

## 5.1 Confusion Matrix

Confusion matrices were constructed to visualize the performance of the models. A confusion matrix is a table that summarizes the performance of a classification algorithm by displaying the counts of true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions. One of the most crucial performance parameters that the confusion matrix assesses is accuracy. It is necessary to define accuracy while talking about confusion and classification matrices.

Accuracy: The percentage of correctly identified cases relative to the total number of examples is known as accuracy. It is computed by dividing the total number of predictions made by the sum of true positive (TP) and true negative (TN) forecasts.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \tag{5.1}$$

Maintaining a machine learning model's accuracy as near to 100 percent as feasible is the aim while working with one; the closer the accuracy is to 100 percent, the better the model is thought to have performed.

Precision: Precision defines how well a categorization model predicts favorable outcomes. The ratio of true positive (TP) predictions to all positive predictions (including false positive (FP) and true positive (TP) predictions) is its definition.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (5.2)$$

The precision of a model is its capacity to prevent false positives—cases that were mistakenly labeled as positive.

Recall: A classification model's recall—also referred to as sensitivity or true positive rate—is a statistic that assesses how well it can distinguish genuine positives from false positives in the dataset.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (5.3)$$

Recall quantifies the percentage of accurately detected positive examples (TP) and cases that were mistakenly labeled as negative (false negatives, FN) among all genuine positive instances.

F1 Score: F1 score is a comprehensive metric that evaluates how well a classification model balances recall and precision. It provides a single numerical value that sums up the model's performance by taking recall and precision into account at the same time. When there is an imbalance between the classes or when the effects of false positives and false negatives are different, the F1 score is especially useful.

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5.4)$$
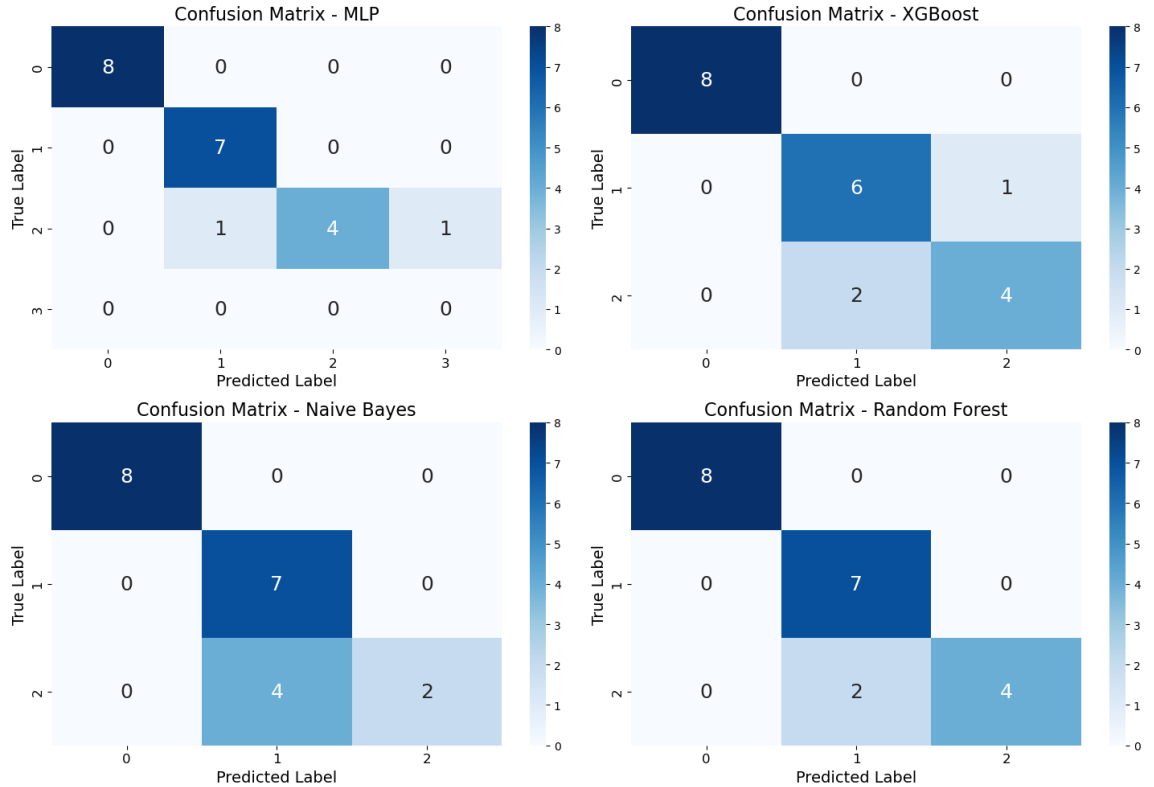
Figure 5.1: Confusion Matrices of Different Models

A higher F1 score, which ranges from 0 to 1, denotes better model performance. The model performs best at 1 and worst at 0. Perfect precision and recall, devoid of false positives or false negatives, are indicated by an F1 score of 1.

Support: The number of real instances of each class in the dataset is referred to as support. Regardless of how accurately or wrongly the model identified an instance, it shows the overall number of examples falling into that class. Understanding the distribution of classes in the dataset requires an understanding of support, which can shed light on the relative prevalence of certain classes.
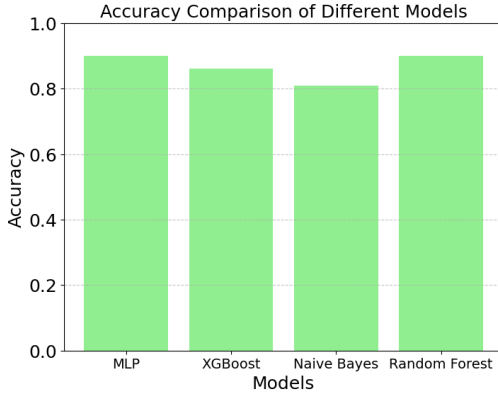
## 5.2 Comparison of Different Models

In our research project, 4 different models (MLP, XGBoost, Naive Bayes and Random Forest) were implemented to predict mental health problems. In terms of Accuracy, Precision, Recall and F1 Score the models were evaluated. The MLP
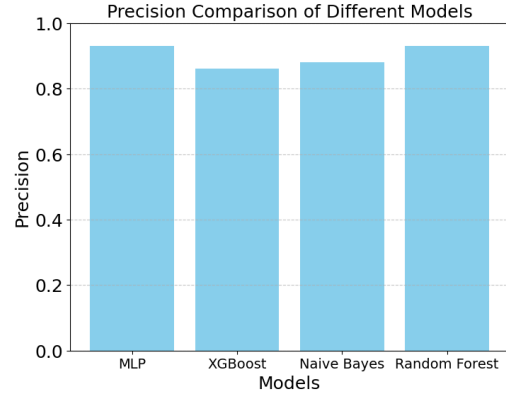
| Model | Accuracy (%) | Precision | Recall | F1 Score |
|---|---|---|---|---|
| MLP | 90 | 0.93 | 0.90 | 0.90 |
| XGBoost | 86 | 0.86 | 0.84 | 0.86 |
| Naive Bayes | 81 | 0.88 | 0.78 | 0.78 |
| Random Forest | 90 | 0.93 | 0.90 | 0.90 |

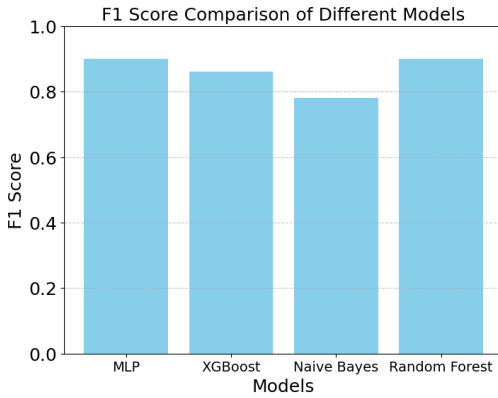Table 5.1: Performance Comparison of Different Models

model demonstrated strong performance in terms of precision, recall, and F1-score metrics. It also achieved a high level of accuracy, demonstrating accurateidentification and comprehensive instance capture.XGBoost outperformed MLP in terms of overall accuracy, albeit it did so somewhat less well. while attaining ideal recall and precision. For classes 1 and 2, it demonstrated somewhat reduced recall and precision, which led to somewhat lower F1-scores as compared to MLP. Even with its simplicity, Naive Bayes performed admirably, especially when it came to recall and precision for class 1. Its performance in class 2, however, was significantly worse, leading to a lower F1-score and overall accuracy. With great recall and precision in every class, Random Forest performed similarly to MLP, producing excellent F1-scores and total accuracy. Overall, the findings indicate that MLP and Random Forest models perform well across a range of measures and are well-suited for predicting mental health disorders based on the features that are provided. They also achieve high accuracy.
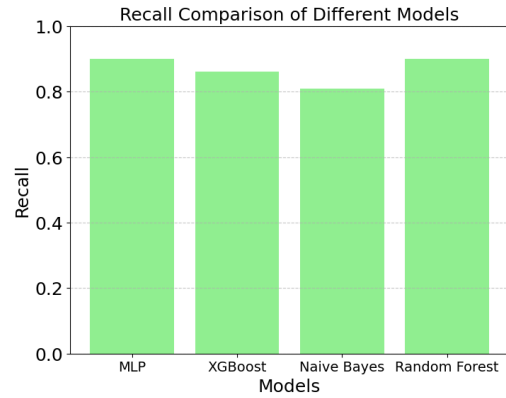
(a) Accuracy of Different Models


(b) Precision of Different Models


(c) F1 Score of Different Models


(d) Recall Score of Different Models

Figure 5.2: Overall Performance Comparison of Different Models

## 5.3 Comparative Analysis

Table (5.2) demonstrates the comparison between our work and related works in terms of accuracy where the MLP Model was used. Table (Table 5.3) demonstrates the accuracy comparison of our work to related works for the XGBoost Model. Table 5.4 and Table 5.5 demonstrates the comparison of our work to related works for the Naive Bayes and Random Forest Model respectively.

| Work Ref. | Related Work Accuracy | Our Accuracy |
|---|:---:|:---:|
| Sofianita Mutalib [1] | 80.00 | |
| Ashley A Sabourin, et al [17] | 98.70 | 90.00 |
| Chekroud, et al. [21] | 77.8 | |

Table 5.2: Accuracy Comparison of the MLP Model in Related Works

| Work Ref. | Related Work Accuracy | Our Accuracy |
|---|:---:|:---:|
| Ryan C. McCabe, et al. [4] | 69.20 | |
| Tate, et al. [23] | 79.00 | 86.00 |

Table 5.3: Accuracy Comparison of the XGBoost in Related Works

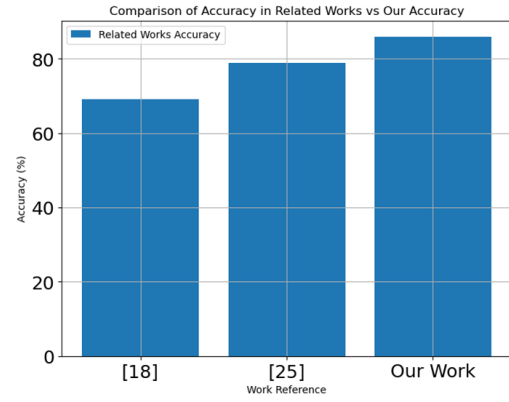| Work Ref. | Related Work Accuracy | Our Accuracy |
|---|:---:|:---:|
| Sofianita Mutalib [1] | 74.81 | |
| Ashley A Sabourin, et al. [17] | 71.42 | 81.00 |
| V. Laijawala, et al. [20] | 78.7 | |
| Chekroud, et al. [21] | 79.6 | |

Table 5.4: Accuracy Comparison of the Naive Bayes Model in Related Works

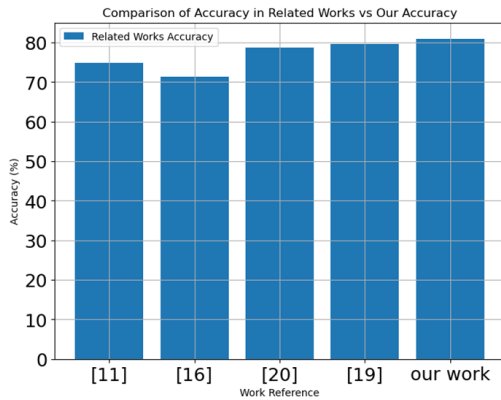| Work Ref. | Related Work Accuracy | Our Accuracy |
|---|:---:|:---:|
| Ashley A Sabourin, et al. [17]] | 83.33 | |
| Ryan C. McCabe, et al. [4] | 73.90 | 90.00 |
| K. M. Mitravinda1, et al. [22] | 91.41 | |
| V. Laijawala, et al. [20] | 79.3 | |

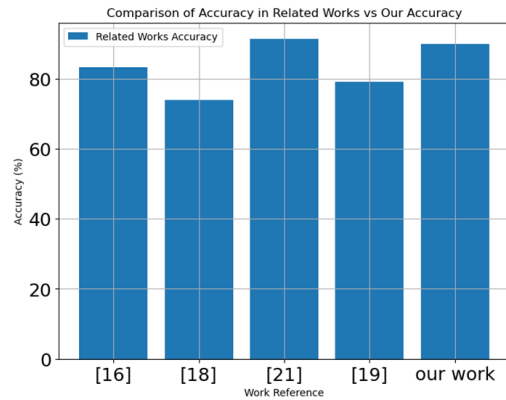Table 5.5: Accuracy Comparison of the Random Forest Model in Related Works

(a) Accuracy Comparison of the MLP model (Related Works Vs Our Work)

(b) Accuracy Comparison of the XGBoost model (Related Works Vs Our Work)

(c) Accuracy Comparison of the Naive Bayes model (Related Works Vs Our Work)

(d) Accuracy Comparison of the Random Forest model (Related Works Vs Our Work)

Figure 5.3: Comparative Analysis of Different Models

# Bibliography

[1] S. Mutalib, "Mental health prediction models using machine learning in higher education institution," in *Turkish Journal of Computer and Mathematics.* TURCOMAT, 2021, pp. 1782–1792.

[2] Polanczyk and Salum, "Annual research review: A meta-analysis of the worldwide prevalence of mental disorders in children and adolescents," in *J. Child Pyschol. Psychiatry*, 2015, pp. 56(3): 345–365.

[3] e. a. M McLafferty, CR Lapsley, "Mental health, behavioural problems and treatment seeking among students commencing university in northern ireland," in *J. Child Pyschol. Psychiatry.* PLOS ONE, 2017, pp. 1–14.

[4] R. Vidourek and M. Burbage, "Positive mental health and mental health stigma: A qualitative study assessing student attitudes," in *Mental Health and Prevention*, 2019, pp. 13:1–6.

[5] D. Olfson and M. S.C, "Trends in mental health care among children and adolescents," in *New England J. Med*, 2018, pp. 372 (21):2029–2038.

[6] R. Ahuja and A. Banga, "Mental stress detection in university students using machine learning algorithms," in *Procedia Computer Science*, 2019, p. 349–353.

[7] T. H. e. a. S Shannon, G Breslin, "Predicting student-athlete and non-athletes' intentions to self-manage mental health: Testing an integrated behaviour change model," in *", Mental Health and Prevention*, 2019, pp. 13:92–99.

[8] R. Vidourek and M. Burbage, "Positive mental health and mental health stigma: A qualitative study assessing student attitudes," in *Mental Health and Prevention*, 2019, pp. 13:1–6.

[9] N. I. of Mental Health (NIMH) (n.d.)., in *"Anxiety Disorder"*, viewable at: https://www.nimh.nih.gov/health/topics/anxiety-disorders/index.shtml.

[10] W. H. O. (WHO), "Depression and other common mental disorders - global health estimation," in *"Obstetrics and Gynecology"*, 2017, pp. 48, 1, 56–60.

[11] T. JM., "Time period and birth cohort differences in depressive symptoms in the us," in *Social Indic. Res.*, 2015, pp. 121(2): 437–454.

[12] D. L. Whiteford H.A, "Global burden of disease attributable to mental and substance use disorder," in *Global Burdern of Disease Study*, 2013, pp. 1575–1586.

[13] C. B. G Andrews and P. Boyce, "Royal australian and new zealand college of psychiatrists clinical practice guidelines for the treatment of panic disorder, social anxiety disorder and generalised anxiety disorder," in *Australian and New Zealand Journal of Psychiatry*, 2018, pp. 1109–1172.

[14] R. Vidourek and M. Burbage, "Positive mental health and mental health stigma: A qualitative study assessing student attitudes," in *Mental Health and Prevention*, 2019, pp. 13:1–6.

[15] E. Kvarnstrom, "The dangers of mental health misdiagnosis," in *Bridge to Recovery*, 2017, viewable at: https://www.bridgestorecovery.com/blog/the-dangers-of-mental-health-misdiagnosis-why-accuracy-matters/.

[16] F. H. N. Fadhludluddin Sahlan and M. H. A. Zamzuri, "Prediction of mental health among university students," in *International Journal on Perceptive Cognitive Computing*, 2021, pp. 7(1): 85–91.

[17] N. M. JC Prater and A. Sabourin, "Assessment of mental health in doctor of pharmacy students. currents in pharmacy teaching and learning," in *Assessment of mental health in doctor of pharmacy students. Currents in Pharmacy Teaching and Learning*, 2019.

[18] W. Z. Yuan, J Zhang and M. Yuan, "Identifying predictors of probable post-traumatic stress disorder in children and adolescents with earthquake exposure: A longitudinal study using a machine learning approach," in *Journal of Affective Disorders*, 2020, pp. 483–493.

[19] K. A. Choudhury and Chakrabarty, "Predicting depression in bangladeshi undergraduates using machine learning," in *SN Computer Science*, 2022, p. Vol.:(0123456789).

[20] H. J. Laijawala, A. Aachaliya and V. Pinjarkar, "Classification algorithms based mental health prediction using data mining," 2020, pp. 1174–1178.

[21] R. J. Z. . M. Chekroud and Z. Shehzad, "Cross-trial prediction of treatment outcome in depression: a machine learning approach," in *Ae Lancet Psychiatry*, 2017, pp. 243–250.

[22] D. N. Mitravinda and G. Srinivasa, "Mental health in tech: Analysis of workplace risk factors and impact of covid-19," in *SN Computer Science*, 2022, p. Vol.:(0123456789).

[23] H. L. Tate, McCabe and K. Halkola, "Predicting mental health problems in adolescence using machine learning techniques," in *PLoS One*, 2020, pp. vol. 15, no.4, p. e0 230 389.