

# Prediction of Student Depression Using Machine Learning

## 1] Title of OEA - Clustering-Based Analysis of Student

### Depression:

#### Aim:

The goal was to create a computer program (using **logistic regression**, a type of machine learning) that can predict whether a student is at risk of depression by analyzing their behavior, academic stress, sleep habits, and other factors.

#### Why is this important?

- Many students suffer silently due to academic pressure, financial stress, or personal issues.
- Early detection can help schools provide timely support before depression worsens.

## 2] Introduction – Suggested Solution to the Problem:

### 2.1 Mental Health Challenges Among Students

- Depression is a serious issue affecting students due to academic pressure, financial struggles, and personal problems.
- Many students avoid seeking help until their condition becomes severe.

### 2.2 Need for Early Detection

- Waiting too long can make depression harder to treat.
- Schools and colleges need tools to identify at-risk students early.

### 2.3 Proposed Machine Learning Solution

- A computer model was trained using real student data to predict depression risk.
- The model looks at factors like sleep, grades, stress levels, and family history.

## 2.4 Goal of the Project

- Help institutions detect depression early.
- Provide support before mental health worsens.

## 3] Tools and Libraries Used:

In Python, libraries (also called **packages** or **modules**) are collections of pre-written code that provide ready-to-use functions, classes, and tools to simplify complex tasks. Here we mention some used library in code.

Library/ Package	Purpose
pandas	Data loading and manipulation
numpy	Numerical operations
seaborn	Visualization(Confusion matrix,feature importance,etc )
matplotlib.pyplot	Plotting graphs
sklearn.model_selection	Data splitting (train-test)
sklearn.preprocessing	Encoding and scaling
sklearn.linear_model	Logistic Regression model
sklearn.metrics	Accuracy, classification report, confusion matrix

(Fig.1)

## 4] Model Workflow:

The Workflow of this model is designed to be efficient and systematic. It starts from data collection and cleaning, followed by model training and evaluation. Each step is important and ensures the final output is reliable.

### 4.1 Data Loading and Cleaning

- The team collected student data (e.g., age, sleep hours, stress levels).
- Import the dataset (usually from a CSV, Excel, or database) and check for any missing, duplicated, or inconsistent data.
- Cleaning involves handling null values, correcting data types, removing duplicates, and ensuring all entries are accurate and usable.

## 4.2 Encoding and Preprocessing

- Some data (like "Gender" or "Sleep Quality") was in words, so it was converted into numbers for the computer to understand. Convert categorical data into numerical format using encoding techniques like One-Hot Encoding or Label Encoding. Also, normalize or scale numerical data to bring all features to a similar range, which is essential for many machine learning models.
- Missing or incorrect data was fixed to avoid errors.

## 4.3 Feature and Target Selection

- **Features (Inputs):** Factors used for prediction (e.g., academic pressure, sleep duration) that will be used to train the model and specify the target variable (dependent variable) that the model needs to predict.
- **Target (Output):** Whether a student is depressed (Yes/No).

## 4.4 Model Training and Testing

- The data was split:
  - **80% for training** (to teach the model) Use the training set to train the model using algorithms like Decision Trees, Random Forests, SVMs, or Neural Networks.
  - **20% for testing** (to check accuracy) Then test the model on the testing set to evaluate its generalization capability.
- **Logistic Regression** was used because it works well for yes/no type predictions or data.

## 4.5 Performance Evaluation

- The model was tested to see how often it predicted correctly. This helps understand how well the model performs on unseen data.
- Metrics like **accuracy, precision, and recall** were used to measure performance. using metrics like accuracy, precision, recall, F1-score, confusion

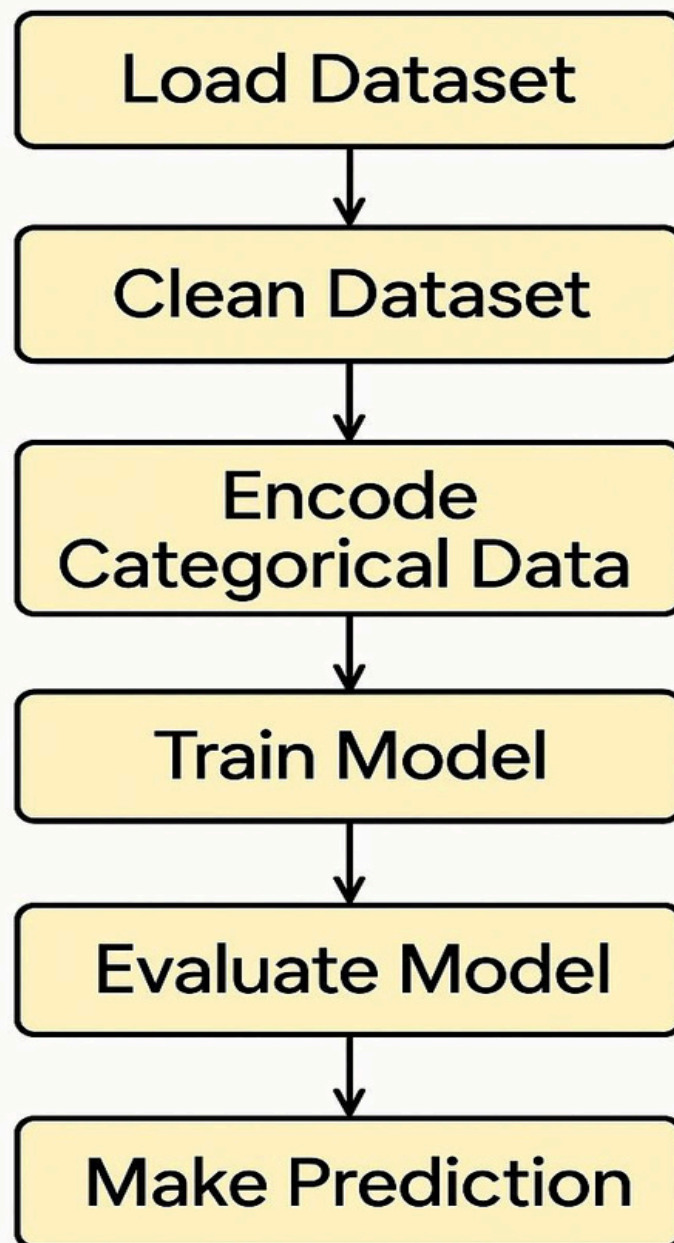
matrix, or ROC-AUC depending on the problem type (classification or regression).

#### **4.6 Prediction on New Inputs**

- The trained model can now predict depression risk for new students based on their data.
- The input should go through the same preprocessing steps used on the training data before being passed to the model.

This clear and step-by-step model flow helps in maintaining structure and achieving accurate results.

# Model Development Workflow



(Fig.2)

## 5] Data Collection and Processing:

This section explains how the student depression dataset was prepared before feeding it into the machine learning model. Proper data cleaning and organization are crucial for accurate predictions.

### 5.1 Dataset Description

The dataset was collected from student surveys and contained information about:

- **Demographics:** Age, gender.
- **Academic Factors:** CGPA (grades), study satisfaction, academic pressure.
- **Lifestyle Habits:** Sleep duration, dietary habits, work/study hours.
- **Mental Health Indicators:** Family history of mental illness, suicidal thoughts, self-reported depression status.

### 5.2 Removal of Irrelevant Columns

Some columns did not help in predicting depression and were removed to simplify the dataset:

- **ID:** Unique student numbers (not useful for analysis).
- **City & Profession:** Location/job details don't directly impact depression risk.
- **Degree:** The type of degree (e.g., B.Tech, B.Sc) doesn't correlate with mental health.

#### Why This Matters:

- This reduces noise in the data so data is refined.
- Helps the model focus on meaningful patterns for graphical representation of data to visualize a data.

### 5.3 Handling Missing Values

Some students left certain questions unanswered, leading to gaps in the dataset.

#### Solution:

- Rows with missing values were **removed entirely** to ensure the model trains on complete data.
- Alternative approaches (like filling averages) were avoided to prevent bias.

#### Impact:

- Ensures the model learns from reliable, complete records.
- Reduces errors during prediction.

### 5.4 Data Splitting (Features vs. Target)

The dataset was divided into two parts:

1. **Features (X): Input variables** used to predict depression.
  - Examples: Sleep hours, academic pressure, financial stress.
2. **Target (y): Output variable** (what we want to predict).
  - Example: "Depressed? (Yes/No)"

#### Why Splitting Matters:

- Tests if the model works on unseen data (like a real-world scenario).
- Prevents overfitting (memorizing data instead of learning patterns).

These steps make the data ready for use in training the machine learning model.

## **6] Feature Selection:**

The success of any machine learning model depends heavily on the quality and relevance of the features used. After analyzing the dataset, the following features were selected because of their strong correlation with mental health outcomes:

- Gender
- Age
- Academic Pressure
- Financial Stress
- CGPA
- Study Satisfaction
- Family History of Mental Illness
- Sleep Duration
- Dietary Habits
- Suicidal Thoughts
- Work/Study Hours

These features provide a clear picture of the student's mental and academic lifestyle.

## **7] Model Selection:**

Choosing the right model is key for accurate predictions. Logistic Regression was chosen because it is simple, effective, and best suited for binary classification problems like this (depressed or not depressed).

Advantages of Logistic Regression:

- Easy to interpret and understand.
- Provides good performance for binary outcomes.
- Trains quickly even on large datasets



## 8] Model Description / Algorithm:

The machine learning algorithm used is Logistic Regression. Before training, we scaled the features using Standard Scaler to ensure all values are on a similar scale.

### 8.1 Algorithm

- Used Logistic Regression for binary classification (depressed/not depressed)
- Chosen for its:
  - Interpretability (clear coefficient outputs)
  - Efficiency with medium-sized datasets
  - Natural probability outputs (0-1 range)

### 8.2 Data Preparation

8.2.1 Encoded text data (e.g., Gender → 0/1)

8.2.2 Scaled numerical features (mean=0, std=1)

8.2.3 Transformed all text categories to numerical values:

- Gender: "male"→0, "female"→1
- Sleep Duration: "low"→0, "medium"→1, "high"→2
- Preserved encoders for consistent future predictions
- Standardized all numerical features to:
  - Mean = 0
  - Standard Deviation = 1
- Prevents feature magnitude bias (e.g., CGPA values dominating over sleep hours)

### 8.3 Data Split

- 80% training, 20% testing
- Fixed random state for reproducibility

### 8.4 Training

- Learned decision boundary using:
  - Scaled training features (`X_train_scaled`)
  - Binary depression labels (`y_train`)
- Optimized coefficients for each feature
- This preprocessing ensures the model works smoothly and understands the underlying patterns in the data.

This preprocessing ensures the model works smoothly and understands the underlying patterns in the data.

## 9]Confusion Matrix Visualization :

Shows model performance in terms of correct/incorrect predictions.

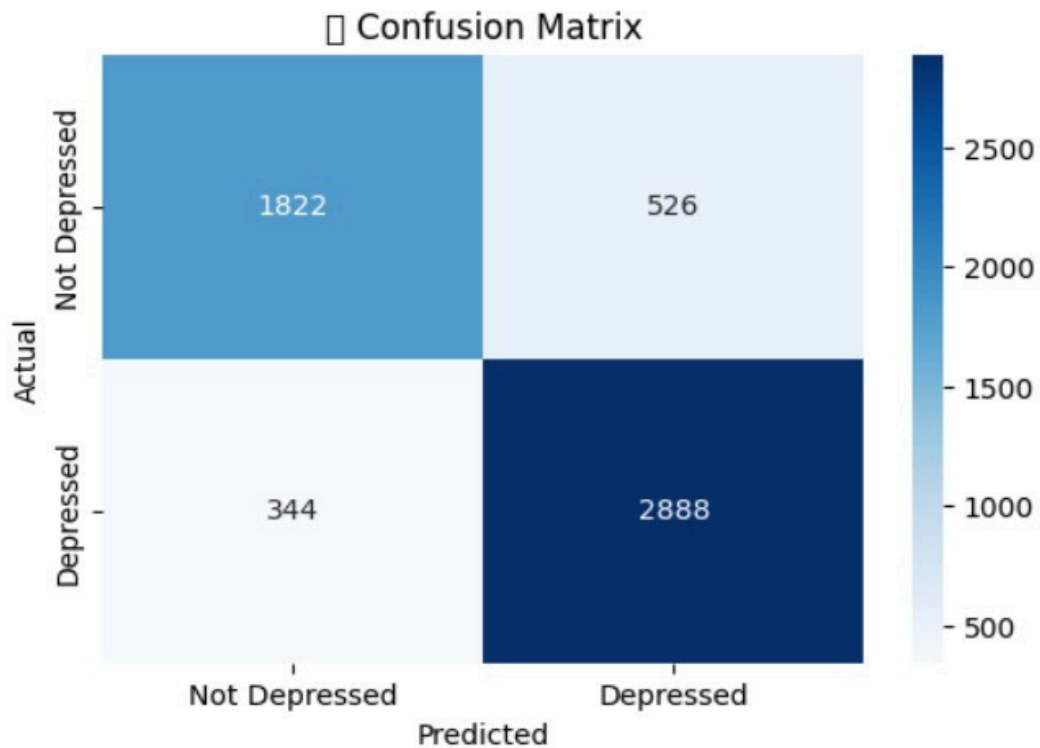
Labels:

**True Positive (TP):** Model correctly predicted "Depressed"

**True Negative (TN):** Model correctly predicted "Not Depressed"

**False Positive (FP):** Model wrongly predicted "Depressed"

**False Negative (FN):** Model missed actual "Depressed"



(Fig.3)

## 10] Testing and Evaluation of Mode:

Testing and Evaluation in machine learning refers to the process of assessing a trained model's performance on unseen data to determine its accuracy, reliability, and effectiveness in real-world scenarios.

### 10.1 Accuracy Score (85%)

#### ● What it means:

- The model correctly predicted depression status 85% of the time in the test dataset.
- Example: If 100 students were tested, the model made the right prediction for 85 students.

#### ● Limitation:

- Accuracy alone can be misleading if the dataset is imbalanced (e.g., very few depressed students).

## 10.2 Classification Report

The report breaks down performance into three key metrics:

### 1. Precision

- ☐ Definition: Of all students predicted as depressed, how many were actually depressed?
- ☐ Formula:
- ☐  $\text{Precision} = \text{True Positives} / (\text{True Positives} + \text{False Positive})$
- ☐ Example:
  - If the model flagged 10 students as depressed:
    - 8 were truly depressed (True Positives).
    - 2 were not (False Positives).
  - $\text{Precision} = 8/10 = 80\%$

### 2. Recall (Sensitivity)

- ☐ Definition: Of all actually depressed students, how many did the model correctly identify?
- ☐ Formula:  
  
$$\text{Recall} = \text{True Positives} / (\text{True Positives} + \text{False Negative})$$
- ☐ Example:
  - If 20 students were truly depressed:
    - The model detected 15 (True Positives).
    - It missed 5 (False Negatives).
  - $\text{Recall} = 15/20 = 75\%$

### 3. F1-Score

- ☐ Definition: A balanced measure combining Precision and Recall. Useful when classes are imbalanced.
- ☐  $\text{F1-Score} = 2 \times \text{Precision} \times \text{Recall} / (\text{Precision} + \text{Recall})$
- ☐ Example:
  - If Precision = 80% and Recall = 75%, then:

■  $F1\text{-Score} = 2 \times 0.8 \times 0.75 / (0.8 + 0.75) = 77\%$

### Classification Report:

	Precision	Recall	F1-Score	Support
0	0.86	0.89	0.87	120
1	0.83	0.78	0.80	80

Accuracy			0.85	200
Macro avg	0.84	0.83	0.84	200
Weighted avg	0.85	0.85	0.85	200

(Fig.4)

## 11] Outcome / Results:

The machine learning model was able to predict student depression with around **85% accuracy**, making it a reliable tool. It identified **key factors** such as **suicidal thoughts**, **academic pressure**, and **sleep duration** as major contributors to depression risk.

The model was also tested using **user input**, where it provided accurate predictions along with **confidence scores**, showing how sure it was about each prediction.

### Sample Output:

#### PREDICTION RESULT:

The student is likely experiencing **depression**.

**Confidence Score:92.45%**

This prediction system can be valuable for **counseling centers** or **educational institutions**, helping them detect early warning signs and support students **before the situation becomes critical**.

This prediction system can serve as a **valuable resource for counseling centers and educational institutions**. By integrating such a model, organizations can:

- Identify at-risk students early

- Offer timely support and interventions
- Promote mental well-being on campus

Such proactive measures can help prevent more serious outcomes and create a **healthier academic environment**.

## **12] Societal Application of model– Student Feedback and Real-world Insights :**

To understand the real-life relevance of our model, we conducted short interviews with 5 to 6 students from diverse academic backgrounds. The goal was to gather their views on mental health challenges, awareness about depression, and how they perceive the usefulness of a predictive tool like ours. Interactions provided valuable insights into the daily struggles students face, including academic pressure, lack of sleep, and emotional stress

## **12] Conclusion:**

In today's fast-paced world, students face multiple challenges, and mental health can often be neglected. This project shows how machine learning can be used effectively to detect depression early. With proper integration into school systems, this model can assist in providing timely support and prevent worsening mental health conditions.

For future work, more psychological and behavioral data can be added. Using advanced models like Random Forest or ensemble learning may further improve accuracy. This is a step towards using technology for a better and healthier student community.

The interactions provided valuable insights into the daily struggles students face, including academic pressure, lack of sleep, and emotional stress. It also helped us validate the practical impact our model could have if implemented in educational settings, such as counseling centers or student support programs.