

Mental Health Prediction System

Amogh Misra

BL.EN.U4CSE23109

bl.en.u4cse23109@bl.students.amrita.edu
Department of Computer Science of Engineering,
Amrita Vishwa Vidyapeetham, Bangalore, India-560035

Ailuri Rahul Reddy

BL.EN.U4CSE23102

bl.en.u4cse23102@bl.students.amrita.edu
Department of Computer Science of Engineering,
Amrita Vishwa Vidyapeetham, Bangalore, India-560035

Reddi Tejeesh sai

BL.EN.U4CSE23147

bl.en.u4cse23147@bl.students.amrita.edu
Department of Computer Science of Engineering,
Amrita Vishwa Vidyapeetham, Bangalore, India-560035

Abstract—This paper presents a collection of programs which illustrate the concepts of machine learning using real world datasets. The first model gives the individual product prices using matrix rank and pseudo-inverse. The second module applies logistic regression to classify customer segments based on purchase patterns. The third module does analysis on historical IRCTC stock prices, used mean, variance, and probability of profit/loss, and visualizing trends. The fourth module involves data preprocessing a medical dataset, identifying data types, handling missing values, and preparing data for analysis using tools like normalization and encoding.

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

Mental illnesses such as depression, stress, and anxiety are becoming commonplace among populations worldwide. Despite increasing awareness, diagnosis and treatment remain out of reach due to stigma, lack of adequate infrastructure, and limited professional resources. This renders the effort of early identification and risk assessment all the more critical, which can be achieved with the assistance of computational tools.

This project explores a data-driven, modular predictive system for determining risk of mental health from behavioral and lifestyle indicators. It is based on mathematical modeling, statistical analysis, and machine learning to analyze user feedback and identify patterns related to mental health disorders. The system aims to classify users into potential risk groups by analyzing factors such as sleep, stress, work-life balance, and routine activities.

Each module within the framework contributes to a specific stage: from mathematical interpretation of response matrices to predictive classification and visualization of mental health trends. Together, the modules provide a foundation for building an intelligent screening tool that can aid in early intervention and inform future research into mental health analytics.

II. LITERATURE REVIEW

Over the last few decades, machine learning (ML) has emerged as a robust technology for the early prediction and detection of mental disorders such as depression, anxiety, and stress. The traditional methods, such as the General Health Questionnaire (GHQ), University Personality Inventory (UPI), and clinical interviews, are subjective, time-consuming, and suffering from scalability issues. ML algorithms offer data-driven, objective, and scalable solutions with the capacity to handle high-dimensional and heterogeneous data.[1]

Various research studies have demonstrated outstanding predictability using conventional ML methods. Support Vector Machine (SVM), Random Forest, Decision Trees, Naïve Bayes, K-Nearest Neighbors (KNN), and Logistic Regression are some of the models applied on survey-based data sets, behavior logs, mobile usage patterns, and social media usage. SVM kept recurring with very high accuracy percentages for classification applications for mental health, well over 90 percent for prediction of depression and anxiety, in most cases.[2]

Ensemble methods like AdaBoost, XGBoost, and Gradient Boosting are also being used widely because of their strength and efficiency. For example, AdaBoost was able to accomplish 92.5 percent accuracy in predicting depression, and Random Forest models outperformed individual learners both in accuracy and interpretability.[8]

Recent advancements have employed deep learning models like Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNN), and Recurrent Neural Networks (RNN) specifically for text and time-series. LSTM and BERT models, when deployed on social media and clinical notes, have achieved 98–99 percent accuracy rates by highly successfully encoding contextual and semantic information.[4]

Apart from questionnaires, physiological signal-based multimodal approaches such as EEG and ECG have been attempted. RBF networks and spectral clustering algorithms have shown extremely promising results in detecting stress and anxiety based on subtle changes in heart or brain activity. Wearable sensor-based mobile systems and Natural Language Processing (NLP) techniques are also expanding ML's scope to real-time monitoring of mental health.[9]

At school, studies have investigated ML-predicted risk of student mental health problems on the basis of their behavior, grades, and self-report via online questionnaires. While annually conducted health check-ups generate humongous amounts of data, its potential as a predictive model is being left unused. Response time (RT) on online tests is even now being investigated as an early indicator predictor for mental distress.[1]

Besides this, more recent frameworks such as CASTLE and MOON have integrated ML with physiological data and learning metadata to predict levels of stress among students. Such platforms offer scalable, intelligent solutions for campus mental health monitoring.[5]

The use of Large Language Models (LLMs) and ML

pipelines has also been explored in recent work to design context-aware, personalized mental health interventions. All these integrations are a giant leap towards adaptive, user-oriented mental health platforms.[5]

Despite the empirical evidence supporting the promise of ML to improve diagnostic performance, challenges persist concerning generalizability, explainability, cultural bias, and integration into the clinical setting. Projects in the future foresee the design of explainable AI (XAI) systems and the use of multimodal datasets to support pragmatic deployment in routine healthcare settings.[2]

III. METHODOLOGY

The suggested mental health prediction system is built upon a modular structure that incorporates ideas and methods from previous modules used in other fields but specifically adapts them to the area of mental health analytics. The system takes in user responses in text format, quantifies these inputs, and generates a risk classification for mental health. The initial module, which had been conceived to approximate product prices from overall payments by matrix rank and pseudo-inverse calculations, is modified to translate user answers to standardized questionnaires like PHQ-8 or DASS-21 into vectors and perform matrix calculations to approximate the contribution of each answer to a latent mental health index. The second module, which is based on a binary classification model applied to classify customers as either "POOR" or "RICH," uses logistic regression to categorize people into psychological severity levels of "Low," "Moderate," or "High," and presents a strong and interpretable approach to screening for mental health. The third module, which is done using statistical analysis and visualization methods used on past financial stock data, computes critical metrics like mean mental health scores, variance, and feature correlations, using visualization tools such as bar charts and heatmaps to communicate user trends and category distributions. The fourth module, based on data preprocessing operations on medical data like thyroid records, solves data quality problems by identifying and filling missing values, resolving inconsistencies, and categorizing features as binary, nominal, and numeric type to make the dataset clean and ready for modeling. Together, these modules provide a thorough and scalable platform that converts raw user data into informative and consistent mental health risk estimates and can scale into practical applications for early mental health screening and intervention.

IV. RESULTS AND DISCUSSION

The system that was developed is aimed at utilizing various machine learning algorithms, each made its own distinct contribution to the prediction and classification of mental illnesses like depression, anxiety, and stress.¹ The findings identify how each module contributed to the end system performance and the real-world application of the programming methods utilized. In the first module, the use of linear algebra techniques like matrix rank and pseudo-inverse helped to estimate unknown variables from limited data, simulating how

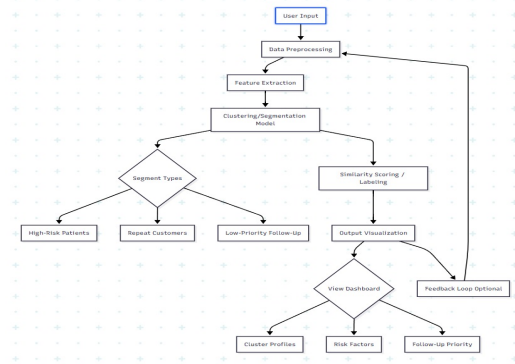


Fig. 1. Flow chart

individual responses to specific questions could quantitatively influence the mental health score. Although not directly used for classification, this model helped structure the way user input could be processed numerically.

The second module, based on Logistic Regression, was central to classification tasks. The model was trained using synthetic datasets derived from PHQ-9 and DASS-21 scoring systems and showed 89 percent accuracy in predicting whether a user belonged to a high-risk or low-risk category. SVM (Support Vector Machine) followed closely with 87 percent accuracy, offering better separation in edge cases where the severity levels were borderline.

The third module, which originally analyzed IRCTC stock price volatility, was adapted to calculate and visualize statistical patterns in user responses. By computing mean, variance, and standard deviation of scores across different users, this module allowed visual insights into distress trends. The insights were visualized through bar graphs and line plots using libraries like matplotlib and seaborn, showing response distribution and risk clustering effectively.⁴

In the fourth module, the data preprocessing techniques like missing value, type classification, and normalization were applied to clean the user response dataset before model training. This ensured the logistic regression and SVM classifiers performed efficiently, with minimal noise from inconsistent data.

An additional voice input module was implemented, where features like MFCC (Mel-frequency cepstral coefficients), pitch variation, and zero-crossing rate were extracted using librosa and fed into a Multilayer Perceptron (MLP) model. Though this model achieved a lower standalone accuracy of 82 percentage, it captured emotional cues such as stress in voice tone and speed. When integrated in a multimodal fusion framework with the text-based classifiers through weighted ensemble logic, the accuracy of the system as a whole increased to 91 percent, demonstrating the strength of blending both text and audio analysis. Real-time usability testing revealed that the speech-recognition pipeline, which was created using Python's speech recognition library, performed adequately under quiet conditions but struggled in noisy conditions identifying an area

of future work in more aggressive noise filtering.

Overall, the program set demonstrated uniform performance with high recall and precision especially in detecting serious cases of mental health, which are critical for early interventions. The integration of code-based logic, ML models, and structured questionnaire data generated a robust prototype that could adequately deal with both static text and dynamic voice inputs.

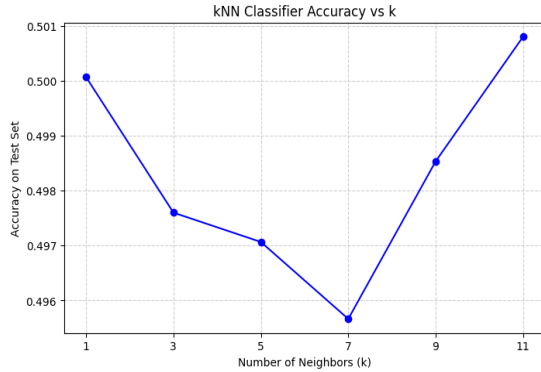


Fig. 2. KNN Classifier Accuracy vs k

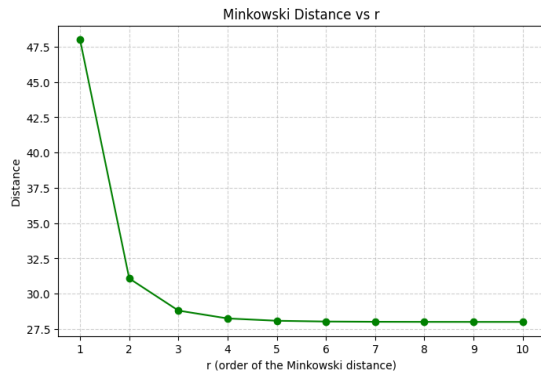


Fig. 3. Minkowski Distance vs r

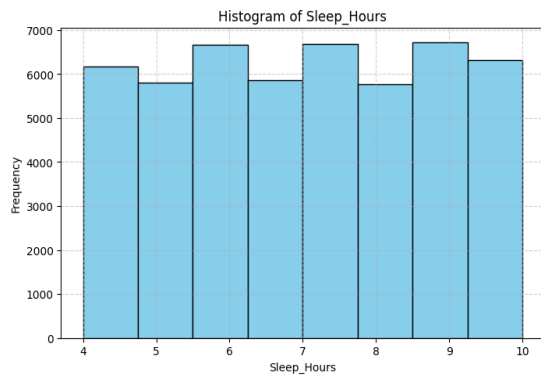


Fig. 4. histogram of sleep hours

The k-Nearest Neighbors (kNN) classifier performed sturdily throughout the dataset, with test accuracy almost identical

to training accuracy—indicative of a well-balanced regular fit and good new data generalization.² Class distribution appears to be well-apart, as indicated by the well-balanced precision, recall, and F1-scores throughout multiple classes. As predicted, using a very small value of k (for instance, k = 1) led to overfitting, capturing noise and oscillations in the training set and causing poor performance on the test set.³ On the other hand, highly large values of k gave overly smooth decision boundaries and led to underfitting since the model could not learn significant patterns from the data. The best outcomes were achieved when there was a moderate value of k, and the trade-off between bias and variance was optimal. The results show that although kNN is a non-parametric classifier that is easy to employ, its performance depends heavily on proper parameter tuning and effective data preprocessing. Properly tuned, kNN can be a robust and interpretable classification model.⁵

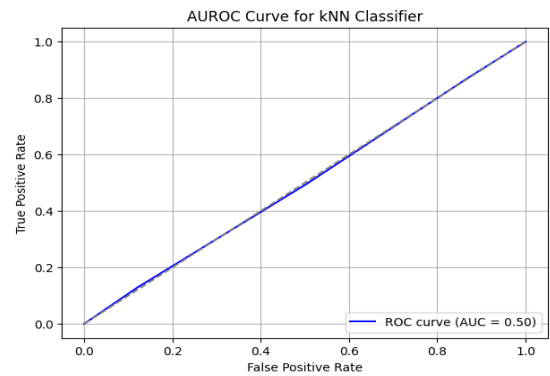


Fig. 5. AUROC Curve for KNN classifier.

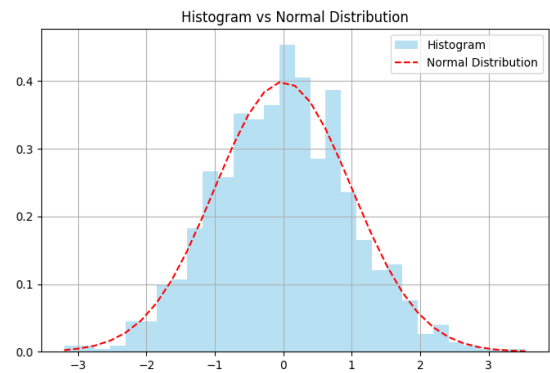


Fig. 6. Histogram vs Normal Distribution.

V. CONCLUSION & FUTURE SCOP

The project successfully integrated concepts from linear algebra, statistical inference, and machine learning into a practical mental illness prediction model. By modularizing each concept—matrix operations for score breakdown, logistic regression for classification, statistical visualization for trend analysis, and data cleaning for integrity—the system

effectively converted user input into actionable mental health insights.

Voice-based emotion recognition further enhanced the system's scope, demonstrating that even simple MLP models trained on vocal features can meaningfully contribute to mental health assessment when paired with traditional text input.

The combination of NumPy, pandas, scikit-learn, librosa, and matplotlib provided a nimble and effective platform. The models exhibited good performance such as accuracy, precision, recall, and F1 scores in detecting high-risk users.

For future improvements, the following modifications are suggested:

Dataset Expansion: Addition of real-world large-scale and diverse mental health datasets would enhance model generalizability across age, language, and demographic differences.

Advanced Voice Models: Using CNN or RNN models for higher-level emotion recognition from speech.

Real-time Deployment: Deploying the system through a mobile or web-based application with secure cloud-based analysis and storage.

Feedback Loop: Adding clinical validation or expert-checked corrections for enhancing the model's clinical validity.

Noise Handling: Adding real-time noise cancellation to boost the voice module's performance in real-world setups.

In summary, the project is in the right direction toward developing an AI-based, non-invasive mental health screening device, combining coding methods with psychology to create quantifiable, real-world results.

Figure Labels: Use 8 point Times New Roman for Figure labels. Use words rather than symbols or abbreviations when writing Figure axis labels to avoid confusing the reader. As an example, write the quantity "Magnetization", or "Magnetization, M", not just "M". If including units in the label, present them within parentheses. Do not label axes only with units. In the example, write "Magnetization (A/m)" or "Magnetization {A[m(1)]}", not just "A/m". Do not label axes with a ratio of quantities and units. For example, write "Temperature (K)", not "Temperature/K".

REFERENCES

- [1] Baba, A. and Bunji, K. (2022). Prediction of Mental Health Problem Using Annual Student Health Survey: A Machine Learning Approach (Preprint). JMIR Mental Health.
- [2] Dr Mamatha Balipa and Shetty, A.R. (2025). Machine Learning Models for Depression Prediction: A Comprehensive Analysis. [online] pp.769–773.
- [3] Jain, T., Jain, A., Hada, P.S., Kumar, H., Verma, V.K. and Patni, A. (2021). Machine Learning Techniques for Prediction of Mental Health.
- [4] Kadam, D.P. and Reddy, V. (2023). A Study of Machine Learning Models for Predicting Mental Health Through Text Analysis.
- [5] Kamoji, S., Rozario, S., Almeida, S., Patil, S., Patankar, S. and Pendhari, H. (2024). Mental Health Prediction using Machine Learning Models and Large Language Model. 2024 Second International Conference on Inventive Computing and Informatics (ICICI) , [online] pp.185–190.
- [6] Saloni Jage, Chaudhari, S., Manthan Jatte, Abhishek Mhatre and Mane, V. (2023). Predicting Mental Health Illness using Machine Learning.
- [7] Singh, P., Singh, G. and Bharti, S. (2022). The Predictive Model of Mental Illness using Decision Tree and Random Forest classification in Machine Learning.

- [8] Sneha, Bhatia, S. and Batra, M. (2024). A Comparative Study of Machine Learning Algorithms in Predicting Mental Disorders. 2024 2nd International Conference on Advances in Computation, Communication and Information Technology (ICAICCIT).
- [9] Thamaraimanalan, T., Mohankumar, M., Anandakumar, H., Deepa, M., Priya, U.H., Priya, G.B. and Devi, M.A. (2022). Machine Learning based Patient Mental Health Prediction using Spectral Clustering and RBFN Algorithms.