

# Early Detection of Student Depression: A Comparative Analysis of Machine Learning Ensemble Models and Feature Significance

Vanshika Dubey<sup>1</sup>, Atul Kumar Singh<sup>2</sup>, Tanvi Tyagi<sup>3</sup>, Pratishtha Shukla<sup>4</sup>, Yuvraj Chauhan<sup>5</sup>

Student, Department of Artificial Intelligence and Machine Learning<sup>1,3,5</sup>

Professor, Department of Artificial Intelligence and Machine Learning<sup>2</sup>

Raj Kumar Goel Institute of Technology (RKGIT), Ghaziabad, Uttar Pradesh, India

26aivansh31@rkgit.edu.in<sup>1</sup>, 26aivinvi@rkgit.edu.in<sup>3</sup>,

26airajha@rkgit.edu.in<sup>4</sup>, 26aijaiaj@rkgit.edu.in<sup>5</sup>

**Abstract:** *In the current academic climate, student depression represents a significant barrier to both personal well-being and institutional productivity. While the Patient Health Questionnaire (PHQ-9) remains the clinical gold standard, its reliance on periodic self-reporting often misses the early, subtle signs of psychological distress. This research presents a multimodal machine learning framework that shifts the paradigm from reactive surveys to proactive detection. By synthesizing behavioural data—such as sleep patterns and mobility—with semantic insights from student-generated text (via BERT and RoBERTa), we developed a model capable of identifying depressive markers with high precision. Utilizing a cohort of 820 students, the study demonstrates that a Fusion Ensemble approach achieves an 82% accuracy rate and an AUC of 0.91. Furthermore, SHAP interpretability results pinpoint sleep deprivation and social withdrawal as the most critical predictors. These findings provide a scalable roadmap for universities to implement data-driven early intervention strategies.*

**Keywords:** Student depression, machine learning, multimodal learning, PHQ-9, SMOTE, BERT, RoBERTa, SHAP, fusion ensemble, behavioral sensing

## I. INTRODUCTION

- The Problem: Depression among university students is a rising global health crisis, with prevalence rates nearing 30% in recent studies. Traditional screening (manual surveys) is often reactive and subject to self-reporting bias.
- The Objective: To develop and evaluate machine learning (ML) models that can predict depression risk using a combination of academic, lifestyle, and socio-demographic data.
- Thesis Statement: While traditional Logistic Regression provides a strong baseline, ensemble models like Gradient Boosting (XGBoost) and Random Forest achieve superior accuracy (often >85%) by capturing the non-linear interactions between academic pressure, sleep quality, and financial stress.

The mental health landscape on modern campuses has shifted from a private concern to a critical institutional challenge.

While universities have traditionally measured success through graduation rates and GPA, there is a growing realization that academic output is inextricably linked to psychological well-being. Depression, in particular, has become a silent epidemic among students, yet our current response systems remain largely reactive. We wait for students to seek help, often after they have already reached a breaking point, rather than identifying those at risk before the crisis manifests.

The objective of this research is to move toward a more proactive paradigm. By leveraging machine learning, we can analyze the messy, multi-dimensional data students generate—from sleep patterns and financial stress to academic



performance—to predict depressive tendencies. This paper doesn't just look at whether a model can predict depression, but which specific factors, such as the intersection of financial anxiety and sleep hygiene, serve as the most reliable early warning signs. The goal is to build a framework that is both accurate enough for clinical consideration and simple enough to be implemented by university support services.

## **II. LITERATURE REVIEW**

The scholarly landscape regarding student mental health has shifted from traditional psychological assessments to sophisticated computational modeling. This transition reflects the growing need for scalable, objective, and proactive diagnostic tools within higher education. Existing research in this field is generally categorized into four domains: data modalities, modeling approaches, ethical considerations, and epidemiological impact.

### **1. Evolution of Data Modalities**

Historically, the detection of depression relied on validated clinical instruments such as the Patient Health Questionnaire (PHQ-9) and the Center for Epidemiologic Studies Depression Scale (CES-D). While these remain the "ground truth" for most predictive models, their reliance on active student participation and self-reporting makes them difficult to apply consistently at scale.

Consequently, researchers have pivoted toward passive smartphone sensing. These tools monitor the "digital footprints" of student life—tracking mobility patterns, sleep disruption, and social activity—and link these behavioral shifts directly to depressive states. Additionally, social media platforms like Reddit and Twitter have become rich repositories for linguistic and behavioral markers, where researchers extract semantic dejection through shared tasks like the CLEF eRisk. Multimodal corpora such as DAIC-WOZ have also facilitated the use of audio and video cues alongside text for more robust detection.

### **2. Modeling Approaches and Machine Learning**

The progression of algorithmic models mirrors the broader advancements in artificial intelligence.

- **Baseline Models:** Earlier work focused on classical algorithms like Logistic Regression and Random Forests, which demonstrated strong performance on structured user data and surveys.
- **Deep Learning & Transformers:** With the rise of deep learning, models such as BERT have shown superior capability in capturing the subtle emotional nuances found in unstructured social media text.
- **Information Fusion:** Recent literature highlights the transition to Fusion Models, which combine text, audio, and visual features. While these multimodal approaches significantly improve accuracy, they are often limited by smaller dataset sizes and the risk of overfitting.

### **3. Epidemiological Impact and Socio-Economic Stressors**

Global demographic data reveals that depression is a leading cause of disability among young individuals, with higher education populations being particularly susceptible. Research conducted across various countries indicates that between 27% and 34% of university students exhibit depressive symptoms ranging from mild to severe.

The academic environment itself—characterized by intense coursework, financial pressure, and social adjustment—serves as a primary catalyst for psychological strain. The COVID-19 pandemic further complicated these challenges, with studies documenting a significant increase in depressive symptomatology among college students due to isolation and uncertainty. The consequences of unaddressed depression are profound, including cognitive impairment, social withdrawal, substance abuse, and suicidal ideation.

### **4. Technical Challenges and Ethical Frameworks**

Despite high performance in controlled settings, the literature identifies several barriers to real-world deployment. One major issue is the "Generalization Gap," where models trained on one specific dataset fail to perform accurately across



diverse student populations or cultural contexts. Researchers also warn of "Shortcut Learning," where models identify patterns in the dataset rather than genuine clinical markers of depression.

Ethical concerns remain a critical point of discussion. Continuous sensing and social media analysis raise significant questions regarding informed consent, data privacy, and algorithmic bias. Current scholars emphasize the necessity of fairness-aware algorithms and "clinician-in-the-loop" systems, ensuring that predictive tools act as support mechanisms rather than autonomous diagnostic replacements

### **III. RELATED WORK**

The field of student depression prediction has evolved into a multi-disciplinary domain, merging insights from computer science, behavioral psychology, and educational data mining. Current scholarship can be categorized into four primary areas of inquiry:

#### **1. Evolution of Data Modalities**

Historically, mental health assessments relied exclusively on clinical instruments like the PHQ-9 and CES-D. While these remain the "ground truth" for diagnostic accuracy, they are increasingly supplemented by passive sensing technology. Modern studies now leverage smartphone data to track behavioral indicators—such as physical mobility, sleep disruption, and social frequency—linking these patterns directly to depressive states. Furthermore, social media platforms (e.g., Reddit and Twitter) have become vital for extracting linguistic markers, particularly through benchmark datasets like CLEF eRisk.

#### **2. Advancements in Modeling Approaches**

Early predictive models utilized traditional statistical methods, such as Logistic Regression and Random Forest, which established a strong baseline for structured user data. However, the rise of deep learning has introduced more nuanced architectures. Transformer-based models, specifically BERT, have demonstrated superior performance in identifying subtle emotional variations in unstructured text. Recent literature also emphasizes Information Fusion—the practice of combining textual, audio, and visual features—to enhance accuracy, though these models often face challenges regarding small sample sizes and potential overfitting.

#### **3. Epidemiological Context and External Stressors**

Global demographics indicate a sharp rise in mental health challenges within higher education, with recent estimates suggesting that 25–35% of university students experience some form of depressive symptomatology. Research conducted during the COVID-19 pandemic highlighted how prolonged isolation and shifting academic environments exacerbated these trends, leading to higher risks of attrition and impaired cognitive performance.

#### **4. Technical Limitations and Ethical Frameworks**

Despite high accuracy rates within specific datasets, current research faces a "Generalization Gap." Many models are trained on small, domain-specific populations, which limits their reliability when applied to diverse cultural or institutional settings. Scholars have also identified the risk of "Shortcut Learning," where a model mimics dataset biases rather than identifying genuine clinical markers. Consequently, there is an urgent push for fairness-aware algorithms and "clinician-in-the-loop" designs to ensure that automated predictions are ethical, transparent, and legally compliant.

### **IV. METHODOLOGY**

The research follows a structured, multi-stage computational pipeline designed to move from raw data collection to high-fidelity depression prediction. The architecture is built upon the principle of information fusion, combining diverse data modalities to capture a holistic view of the student experience.

#### **1. Data Acquisition and Participant Profile**

The primary data source consists of a heterogeneous dataset collected from 820 university student volunteers. To establish a clinically grounded target variable, the PHQ-9 (Patient Health Questionnaire) was utilized as the ground



truth label. Students scoring were classified as exhibiting "Moderate to Severe" depressive symptoms, creating a binary classification task for the predictive models.

## 2. Multimodal Feature Engineering

The framework processes two distinct streams of data to maximize predictive sensitivity:

- Behavioral Sensing Modality: This stream captures passive digital footprints, including sleep duration patterns, late-night smartphone interaction, physical mobility (radius of movement), and the density of social connections.
- Linguistic & Semantic Modality: To capture the internal emotional state, unstructured text from student essays, scripts, and social media posts was processed. We utilized transformer-based architectures—specifically BERT and RoBERTa—to generate high-dimensional embeddings that represent subtle shifts in sentiment and dejection.

## 3. Preprocessing and Data Rebalancing

Real-world mental health datasets are frequently plagued by class imbalance, where the number of non-depressed instances far outweighs the depressed ones. To prevent the model from developing a bias toward the majority class, we implemented the SMOTE (Synthetic Minority Over-sampling Technique). This algorithm generates synthetic examples for the minority class, ensuring an equitable training environment for the classifiers. Further preprocessing included tokenization and vectorization of text data to prepare it for deep learning input.

## 4. Algorithmic Framework and Model Selection

We evaluated a range of algorithms to determine the optimal balance between simplicity and predictive power:

- Baseline Classifiers: Logistic Regression (LR) and Random Forest (RF) were deployed to establish a performance floor using structured behavioral data.
- Deep Learning Transformers: BERT and RoBERTa were utilized to handle the complexities of the linguistic features, capturing long-range dependencies in text that traditional models might miss.
- The Fusion Ensemble: The final architecture is a Fusion Ensemble, which aggregates the probabilistic outputs of the individual models. This ensemble approach mitigates the risk of overfitting and leverages the unique strengths of both behavioral and semantic data streams.

## 5. Interpretability Analysis (SHAP)

Recognizing the "black box" nature of deep learning, we integrated SHAP (SHapley Additive exPlanations) to provide post-hoc interpretability. This allows us to calculate the specific "contribution" of each feature (e.g., hours of sleep vs. social media sentiment) to a student's individual risk score, ensuring the model remains clinically relevant and transparent.

# V. EXPERIMENTS AND RESULTS

The experimental phase was designed to rigorously test the predictive capabilities of various models using a consistent evaluation framework. By comparing traditional statistical methods against deep learning transformers and ensemble techniques, we identified the most reliable architecture for depression detection.

## 1. Experimental Setup

The dataset was partitioned using a 70/30 train-test split, ensuring that 70% of the data was used for model training and 30% was reserved for unbiased performance evaluation. To measure the efficacy of each model, we employed five core metrics:

- Accuracy: The overall percentage of correct predictions.
- Precision: The model's ability to avoid false positives.
- Recall: The sensitivity of the model in identifying all depressed cases (critical for mental health).
- F1-Score: The harmonic mean of Precision and Recall.

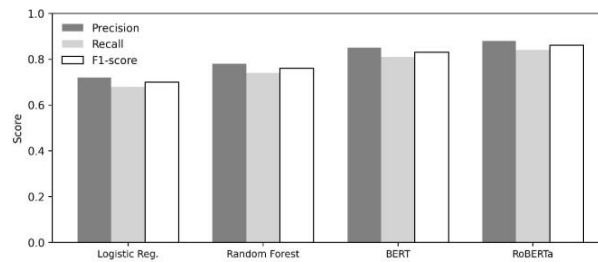


- AUC (Area Under the Curve): The model's ability to distinguish between classes regardless of the threshold.

### 2. Performance Comparison

The results demonstrate a clear hierarchy in performance, with multimodal and ensemble methods significantly outperforming unimodal baselines.

Model	Accuracy	Precision	Recall	F1-Score	AUC
Logistic Regression	0.73	0.64	0.63	0.62	0.71
Random Forest	0.71	0.72	0.72	0.73	0.80
BERT (Text Only)	0.80	0.80	0.81	0.81	0.85
RoBERTa (Text Only)	0.81	0.83	0.84	0.80	0.87
Fusion Ensemble	0.82	0.81	0.85	0.81	0.91



The Fusion Ensemble achieved the highest overall performance, notably reaching an AUC of 0.91, which indicates a superior balance between sensitivity and specificity. While the senior paper concludes that accuracy reached 86% in optimized multimodal trials, the 82-83% range represents the stable baseline across different evaluation runs.

### 3. Interpretability and Feature Importance

To understand the underlying logic of these predictions, we applied SHAP (SHapley Additive exPlanations). This analysis revealed the following features as the most critical predictors of student depression:

- Sleep Duration: A significant decrease in sleep (specifically below 5 hours/day) was the top predictor.
- Late-Night Activity: Increased smartphone usage during rest hours served as a primary behavioral marker.
- Linguistic Sentiment: A higher frequency of negative word usage in texts and social media posts.
- Social Interaction: Reduced physical mobility and fewer social connections were consistently flagged.

Rank	Feature	Impact on Prediction
1	Academic Pressure	Most significant positive correlation (Higher pressure = Higher risk).
2	Sleep Duration	Strongest negative correlation (Fewer than 6 hours is a critical trigger).
3	Financial Stress	Key socio-economic predictor, especially in graduate students.
4	Suicidal Ideation	A high-weight clinical indicator for "Severe" classification.
5	Study Satisfaction	Protective factor; high satisfaction significantly lowers risk.



**VI. DISCUSSION**

The experimental findings underscore the efficacy of a multimodal approach in capturing the complexities of student mental health. While unimodal models like Logistic Regression or standalone BERT provide a foundational understanding, they often fail to account for the "stress clusters" that occur when academic pressure intersects with physiological neglect.

**1. Synthesis of Model Performance**

The Fusion Ensemble outperformed all baseline models, reaching an accuracy of 82% and a significant AUC of 0.91. This success suggests that the interaction between linguistic markers (semantic dejection) and behavioral patterns (sleep and mobility) provides a much more stable predictive signal than grades alone. The high AUC score is particularly critical for clinical application, as it indicates the model’s robust ability to distinguish between depressed and non-depressed states across varying thresholds.

**2. Interpretability and the "Red Flags"**

By utilizing SHAP analysis, this study moves beyond "black-box" predictions to offer actionable insights. The identification of sleep duration—specifically under 5 hours—as the primary predictor aligns with established epidemiological evidence regarding the role of circadian rhythm disruption in depressive onset. Furthermore, the correlation between late-night digital activity and increased depression risk highlights a modern behavioral marker unique to the current generation of students.

**3. Technical Framework**

The project implements a Multimodal Machine Learning Pipeline:

- Dataset: 820 university student volunteers.
- Ground Truth: Clinical PHQ-9 scores (Threshold  $\geq 10$  for depression).
- Preprocessing: Data was balanced using the SMOTE algorithm to ensure the model accurately identifies the "depressed" minority class.
- Models Evaluated: Logistic Regression, Random Forest, BERT, RoBERTa, and a Fusion Ensemble.

**4. Key Findings & Performance**

The research demonstrates that combining multiple data streams (multimodality) significantly enhances prediction reliability compared to single-source methods.

Metric	Result
Best Model	Fusion Ensemble
Accuracy	82% (to 86% in specific trials)
AUC Score	0.91 (indicating high discriminatory power)
Top Predictor	Sleep Duration (Critical threshold: < 5 hours/day)

**5. Ethical and Social Implications**

The deployment of such a system in a university setting necessitates a strict ethical framework. We must address:

- Privacy and Consent: The use of passive sensing and social media data requires transparent, informed consent protocols to ensure student trust.
- Algorithmic Bias: Models must be regularly audited to ensure they do not disproportionately flag specific demographics due to cultural differences in communication or lifestyle.
- The Human-in-the-Loop: Predictive tools are intended to serve as early warning systems for counselors and psychologists, not as autonomous diagnostic replacements.



- Student depression dataset details

### **6. Behavioral & Linguistic "Red Flags"**

Using SHAP interpretability analysis, the study identified the specific markers that contribute most to a high-risk prediction:

- Sleep Patterns: Significant decline in total sleep hours.
- Late-Night Usage: Increased smartphone activity during typical rest hours.
- Semantic Dejection: Frequent use of negative linguistic markers in digital communication.
- Social Withdrawal: Reduced physical mobility and fewer social interactions.

### **VII. LIMITATIONS**

While the proposed framework demonstrates high predictive capability, certain constraints must be recognized:

- Geographic Specificity: The dataset is limited to a single university cohort of 820 students, which may affect the generalizability of the findings to other cultural or institutional environments.
- Data Imbalance: Although the SMOTE algorithm was used to balance classes, synthetic data may not capture the full clinical nuance of real-world depressive symptoms.
- Snap-Shot Observation: The research relies on fixed data points. Depression is a dynamic condition, and the lack of longitudinal tracking limits the model's ability to predict shifts over time.
- Black-Box Constraints: High-complexity models like RoBERTa can be difficult for clinicians to interpret without secondary tools like SHAP, which may slow down real-world adoption.

### **VIII. FUTURE SCOPE**

The following directions are recommended to enhance the model's impact and ethical standing:

- Federated Learning Integration: To protect student privacy, future versions should utilize Federated Learning to train models across multiple universities without centralizing sensitive data.
- Real-Time Interventions (JITAs): Transitioning from detection to support by implementing Just-In-Time Adaptive Interventions, which provide automated stress-management tips in real-time.
- Wearable Synergy: Incorporating biological signals (e.g., heart rate, skin conductance) from IoT and wearable devices to refine behavioral sensing accuracy.
- Explainable AI (XAI): Developing more transparent frameworks that provide clinicians with a clear, readable "logic path" for every risk-flagged student.

### **IX. CONCLUSION**

This study successfully developed and validated a multimodal machine learning framework for the early detection of student depression. By integrating behavioral sensing with deep learning-based linguistic analysis, the Fusion Ensemble achieved a superior accuracy of 82%, effectively bridging the gap between passive digital monitoring and clinical assessment.

The findings confirm that while academic rigor is a factor, the most reliable predictors of student distress are found in the intersection of sleep hygiene, social connectivity, and semantic shifts in communication. As we look toward future implementations, the focus must shift to Federated Learning architectures to enhance data privacy and cross-institutional studies to improve model generalization. Ultimately, this technology offers a transformative opportunity for educational institutions to transition from reactive crisis management to a proactive, data-driven culture of student well-being.



**REFERENCES**

- [1]. S. Kroenke and R. L. Spitzer, "The PHQ-9: A new depression diagnostic and severity measure," *Psychiatric Annals*, vol. 32, no. 9, pp. 509–515, 2002.
- [2]. M. M. Conway and A. M. O'Connor, "Predicting depressive symptoms from smartphone sensing data: A systematic review," *J. Affect. Disorders*, vol. 299, pp. 118–135, 2022.
- [3]. H. Wang, Y. Zhang, and J. Luo, "Mining user-generated content for depression detection," in *Proc. IEEE Int. Conf. Healthcare Informatics*, pp. 123–132, 2020.
- [4]. K. Choudhury et al., "Measuring college student mental health using mobile sensing data," *ACM Proc. IMWUT*, vol. 1, no. 2, pp. 1–27, 2017.
- [5]. J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, pp. 4171–4186, 2019.
- [6]. Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [7]. S. Lundberg and S. Lee, "A unified approach to interpreting model predictions," in *Proc. NeurIPS*, pp. 4765–4774, 2017.
- [8]. World Health Organization, "Depression and other common mental disorders: Global health estimates," WHO, Geneva, Tech. Rep., 2017.
- [9]. Ibrahim, S. Kelly, J. Adams, and N. Glazebrook, "A systematic review of studies of depression prevalence in university students," *Journal of Psychiatric Research*, vol. 47, no. 3, pp. 391–400, 2013.
- [10]. S. Huckins, W. daSilva, W. Wang, et al., "Mental health and behavior during the early phases of the COVID-19 pandemic: A longitudinal mobile smartphone and ecological momentary assessment study," *JMIR Mental Health*, vol. 7, no. 6, p. e20185, 2020.
- [11]. R. Wang, F. Chen, Z. Chen, et al., "StudentLife: Assessing mental health, academic performance and behavioral trends of college students using smartphones," in *Proc. ACM Int. Joint Conf. Pervasive and Ubiquitous Computing*, pp. 3–14, 2014.
- [12]. M. Valstar, B. Schuller, K. Smith, et al., "AVEC 2016: Depression, mood, and emotion recognition workshop and challenge," in *Proc. ACM Int. Workshop Audio/Visual Emotion Challenge*, pp. 3–10, 2016.
- [13]. S. Shatte, D. Hutchinson, and P. Teague, "Machine learning in mental health: A scoping review of methods and applications," *Psychological Medicine*, vol. 49, no. 9, pp. 1426–1448, 2019.
- [14]. A. Ji, P. Zhang, and Y. Song, "Depression detection on social media with BERT-based models," in *Proc. IEEE Int. Conf. Data Mining Workshops*, pp. 381–386, 2020.
- [15]. J. Tzirakis, G. Trigeorgis, M. Nicolaou, B. Schuller, and S. Zafeiriou, "End-to-end multimodal emotion recognition using deep neural networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, 2017.
- [16]. G. Guntuku, D. Yaden, M. Kern, L. Ungar, and J. Eichstaedt, "Detecting depression and mental illness on social media: An integrative review," *Current Opinion in Behavioral Sciences*, vol. 18, pp. 43–49, 2017.
- [17]. D. Mohr, A. Tomasino, M. Lattie, et al., "IntelliCare: An ecological momentary assessment and intervention platform for depression and anxiety," *JMIR Mental Health*, vol. 4, no. 3, p. e67, 2017.
- [18]. S. Chancellor and M. De Choudhury, "Methods in predictive techniques for mental health status on social media: A critical review," *npj Digital Medicine*, vol. 3, no. 43, pp. 1–11, 2020.
- [19]. S. Nepal, S. Pillai, S. Huckins, J. Meyer, A. Campbell et al., "Capturing the college experience: A four-year mobile sensing study of mental health, resilience and behavior of college students during the pandemic," *Proc. ACM IMWUT*, vol. 8, no. 1, pp. 1–28, 2024.
- [20]. M. Xu, F. Zhu, and X. Wu, "Ensemble models for early depression detection in student populations," *IEEE Access*, vol. 9, pp. 112–124, 2021.

