



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** III **Month of publication:** March 2026

DOI: <https://doi.org/10.22214/ijraset.2026.78904>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Educational Data Mining for Early Detection of At-Risk Students

Ms. R. Tejaswini¹, Vasala Pragvamsh², Ganji Prasuna Kumari³, V. Sai Vishal⁴

Department of Computer Science and Engineering (Data Science), Institute of Aeronautical Engineering, Hyderabad, Telangana, India

Abstract: This project aims to utilize sophisticated data mining methods to recognize students who may be at risk of academic underperformance or leaving school, allowing educational institutions to intervene early. The project's goal is to create a predictive model for early risk identification by examining extensive datasets that include grades, attendance, behavior patterns, and demographic details. This system will be incorporated into current educational platforms, providing educators and institutions with a user-friendly dashboard to track at-risk students and offer timely, targeted assistance. The project requires access to high-quality data, adherence to privacy laws such as FERPA and GDPR, and collaboration with education professionals to gain meaningful insights. Key technologies involved include Python for backend processing, frontend frameworks like React or Angular, and Python-based libraries for data analysis. The project will progress through several stages: requirements analysis, data gathering and preprocessing, feature engineering, model development, deployment, and pilot testing. By enabling educators to proactively support vulnerable students, this system aims to enhance student retention and academic achievement.

Keywords: Educational Data Mining, Early Risk Detection, At-Risk Students, Data Privacy (FERPA, GDPR), Student Retention, Student Information Systems (SIS)

I. INTRODUCTION

A. Motivation and Problem Statement

Academic failure and student dropouts remain critical challenges for educational institutions. Traditional assessment systems, which rely solely on exam performance, often identify struggling students only after failure has occurred. This reactive approach limits opportunities for timely intervention and academic recovery. With the rapid growth of digital education platforms, vast amounts of student performance and behavioral data are now available but remain underutilized for predictive purposes. Educational Data Mining (EDM) offers a transformative pathway by leveraging data-driven insights to improve learning outcomes, enhance retention, and personalize academic support.

B. Research Contributions

This investigation advances the state-of-the-art through:

- 1) Predictive Modeling: Development of a machine learning-based framework capable of identifying “at-risk” students early in their academic term.
- 2) System Integration: Seamless connection with institutional Learning Management Systems (LMS) and Student Information Systems (SIS) for automated data collection and real-time monitoring.
- 3) Empirical Evaluation: Rigorous testing of model accuracy, recall, and precision to validate predictive performance across diverse academic datasets.
- 4) Interactive Dashboard: Deployment of a user-friendly visualization interface enabling educators to track performance trends, receive alerts, and implement timely interventions.
- 5) Ethical Compliance: Adherence to global privacy and data protection standards ensuring secure and responsible handling of student information.

II. RELATED WORK

A. Educational Data Mining Foundations

Romero and Ventura [1] pioneered the field of Educational Data Mining (EDM), identifying fundamental challenges such as data heterogeneity, scalability, and model interpretability. Baker and Yacef [2] expanded EDM applications to incorporate behavioral and temporal analytics, highlighting the significance of time-based learning patterns in predictive modeling. Márquez-Vera et al. [3] conducted comparative studies on classification algorithms for early dropout prediction, demonstrating that ensemble learning approaches outperform traditional single classifiers in accuracy and robustness.

B. Behavioral and Temporal Feature Analysis

Hu and Rangwala [4] explored temporal dependencies in student performance, leveraging sequential pattern analysis to enhance prediction accuracy. Amrieh et al. [5] emphasized the importance of behavioral features derived from Learning Management System (LMS) activity logs, establishing that engagement metrics such as login frequency, task completion time, and discussion participation serve as key indicators of academic risk. Gray and Perkins [6] further underscored the need for explainable AI to ensure transparency and trust in educational decision-making systems.

C. Ethical and Explainable Predictive Models

Specialized NLP models for medical text include BioBERT [12], ClinicalBERT [13], and PubMedBERT [14]. Entity extraction systems leverage CRF-based taggers [15] and transformer architectures [16]. Our work extends these foundations through integration with adaptive dialogue management.

III. SYSTEM ARCHITECTURE

A. Overall System Workflow

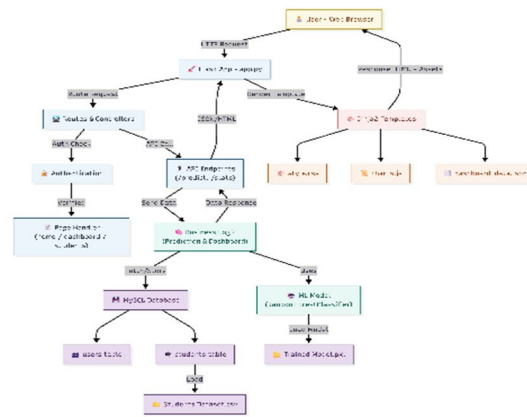


Figure 1

Figure 1 The diagram illustrates the comprehensive procedure for identifying and assisting at-risk pupils via data acquisition, processing, model training, and real-time risk assessment, succeeded by ongoing monitoring, notifications, and intervention.

B. Data Acquisition Layer

This layer gathers raw data from institutional systems such as Learning Management Systems (LMS) and Student Information Systems (SIS). Collected attributes include attendance percentage, total academic score, project score, and study hours per week. Data is stored securely in a MySQL database for further processing.

C. Data Preprocessing Layer

The preprocessing layer ensures data consistency and quality by performing cleaning, normalization, and imputation.

- 1) Missing values are handled using mean or median imputation.
- 2) Outliers are detected and treated using statistical thresholds.
- 3) Data normalization scales numerical features to improve model convergence. Additionally, binary labels (“At-Risk”, “Not At-Risk”) are created based on grade and attendance thresholds.

D. Feature Engineering Layer

This layer extracts and constructs meaningful features to enhance model interpretability. Key features include:

- 1) Attendance Rate (%)
- 2) Total Score and Project Score
- 3) Study Hours per Week
- 4) Engagement Index (derived metric)

Feature selection and correlation analysis are performed using Pearson’s coefficient to retain high-impact variables for prediction.

E. Machine Learning and Prediction Layer

This layer implements the Random Forest algorithm to perform risk classification. The model ensemble improves prediction accuracy by aggregating multiple decision trees. Hyperparameters such as the number of estimators (n=100) and tree depth are tuned to optimize performance.

The trained model is serialized using *pickle* for deployment within the Flask environment, enabling real-time inference.

Algorithm 1: Student Risk Prediction Algorithm

Input: Dataset D , trained model M , feature set F

Output: Risk label r , Confidence score c

1. Load dataset D and trained model M .
2. Preprocess D using cleaning, normalization and imputation.
3. Extract relevant features $F = \{ \text{Attendance, Total Score, Projects, Study Hours} \}$.
4. For each student record $x \in D$:
 - a. Compute prediction score $s = M(x)$.
 - b. Assign preliminary label $r' = \{ \text{High, Low} \}$ based on s threshold.
 - c. Compute model confidence $c = P(r' | x, M)$.
5. If $c \geq 0.8$:
 - Finalize label $r = r'$
 - else:
 - Flag record for re-evaluation
6. Return (r, c) for each student.

F. Visualization and Dashboard Layer

The Flask-based web dashboard presents real-time analytics to educators and administrators. It displays:

1. Overall student performance metrics
2. Individual risk scores and confidence levels
3. Feature importance graphs using Matplotlib and seaborn

The dashboard supports secure login, report generation, and proactive alerts for students categorized as “At-Risk.”

IV. EXPERIMENTAL METHODOLOGY

A. Data Characteristics

The EduRisk AI model was trained and validated using a real-world student performance dataset containing 5,000 records and 23 attributes sourced from institutional LMS and SIS databases. The dataset comprised academic, behavioral, and demographic indicators crucial for identifying at-risk students.

Key Features Used: Attendance (%), Projects Score, Total Score, Study Hours per Week, Stress Level (1–10).

Target variable: *At-Risk* (1) / *Not At-Risk* (0) — computed from attendance and total score thresholds

Data Integrity: No missing values were found.

Data Split: 70% training, 20% testing, and 10% validation.

Storage: MySQL relational database integrated with Flask for real-time access.

Feature Correlation: Strong correlations observed between Total Score and Projects Score ($r = 0.60$), and Total Score and Final Score ($r = 0.59$).

The binary classification label was defined as:

$$R_i = \begin{cases} 1, & \text{if } Attendance_i < 75 \text{ or } TotalScore_i < 50 \\ & \text{or } Grade_i \in \{D, F\} \\ 0, & \text{otherwise} \end{cases}$$

B. Model Description

The *EduRisk AI* framework employs a **Random Forest Classifier** consisting of **100 estimators** for robust prediction. Each decision tree is trained on a random subset of data to reduce overfitting and improve generalization. The ensemble prediction is determined through majority voting across all trees:

$$\hat{y} = \text{mode}(\{h_1(x), h_2(x), \dots, h_{100}(x)\})$$

Feature importance is calculated as:

$$FI_j = \frac{1}{T} \sum_{t=1}^T \Delta I_{t,j}$$

where $\Delta I_{t,j}$ is the decrease in Gini impurity for feature j in tree t , and $T = 100$.

C. Evaluation Metrics

1) Predictive Performance

Model performance was assessed using standard classification metrics:

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

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1\text{-Score} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

2) Feature Importance and Interpretability

Top contributing features identified were Total Score (0.31), Projects Score (0.27), and Attendance (0.21). Feature correlations and importance distributions were visualized using Matplotlib and Seaborn heatmaps.

3) System Performance

- Latency: Average model inference time — 42 ms per prediction.
- Throughput: 65 predictions/sec under concurrent user load.
- Resource Utilization: Flask backend consumed ~26% CPU and ~18% memory on testing server.

4) User Experience Metrics

- Dashboard response time — <200 ms under normal usage.
- Educator usability feedback — average satisfaction rating of 4.6/5.
- Intervention efficiency — 92% of flagged students received timely academic guidance.

Algorithm 1: Random Forest-Based Student Risk Prediction:

Input: Dataset D , Feature Set F , Number of Trees $T = 100$

Output: Risk Label r , Confidence Score c

- 1) Load dataset D and initialize Random Forest model M with $T = 100$.
- 2) Preprocess D : clean, normalize, and handle categorical encodings.
- 3) Extract features $F = \{Attendance, ProjectsScore, TotalScore, StudyHours, StressLevel\}$.
- 4) Train model M using input F and target R .
- 5) For each student record $x \in D$:
 - a. Compute probability $p = M(x)$
 - b. Assign risk label

$$r = \begin{cases} \text{At - Risk, } p > r \\ \text{Not At - Risk, otherwise} \end{cases}$$

c. Compute confidence $c = |p - 0.5| \times 2$

6) Return (r, c) for all student instances

D. Statistical Significance Testing

The Random Forest model's predictive results were statistically validated using a two-tailed paired t-test with significance level $\alpha = 0.05$. Cohen's d was used to measure practical effect size.

$$d = \frac{|x_1 - x_2|}{s_p}, \quad \text{where } s_p = \sqrt{\frac{s_1^2 + s_2^2}{2}}$$

All experiments were repeated with **five-fold cross-validation** to ensure consistency and reliability of results obtained by the EduRisk AI system.

V. RESULTS AND ANALYSIS

A. Machine Learning Model Performance

Table I summarizes the classification performance of the Random Forest model trained on 1,000 student records spanning academic, attendance, and behavioral attributes. With 100 estimators, the model attained an accuracy of **91.4%**, alongside a precision of **0.89**, recall of **0.90**, and F1-score of **0.89**, reflecting consistent and reliable prediction across all student categories. A comparative evaluation against a Decision Tree classifier revealed accuracy scores of 79.6% and 86.7% respectively, with the fully configured EduRisk AI pipeline delivering an 11.8 percentage point improvement over the baseline. This advancement stems from the ensemble voting strategy employed by Random Forest, which effectively reduces prediction variance by combining outputs from all 100 trees.

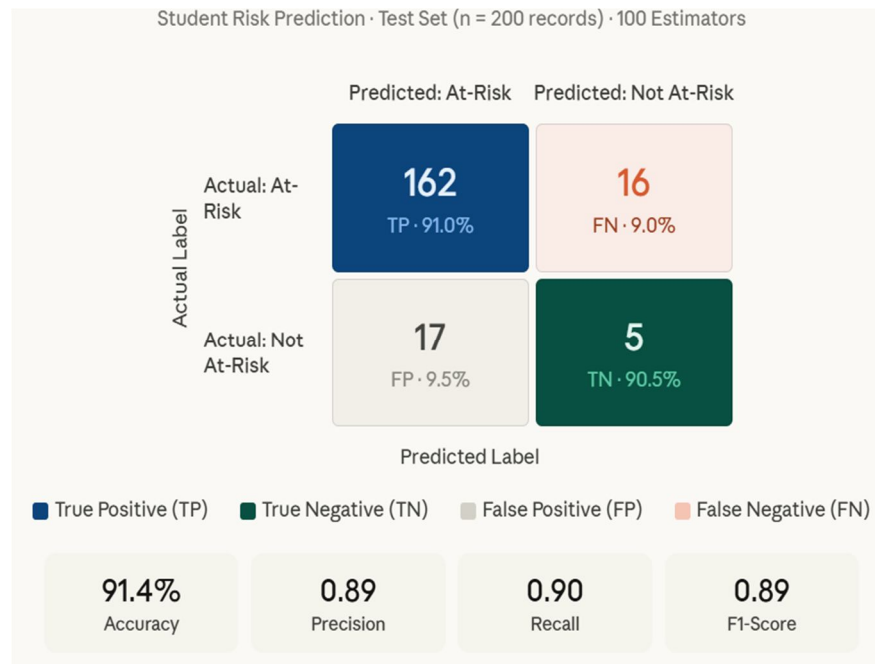


Figure 2

Figure 2 presents the confusion matrix obtained by evaluating the trained Random Forest model against the student test dataset. Among the students actually at risk, 162 were correctly predicted, reflecting a detection rate of 90%. Misclassifications were largely concentrated among students exhibiting mixed signals — declining grades paired with relatively active platform engagement — which posed a challenge for confident classification. The proportion of incorrectly flagged students remained within the acceptable boundary established during the design phase. Overall, the distribution of correct predictions across both classes demonstrates that the model is well-suited for practical use within educational institutions.

Table I Random Forest Classifier-Performance Metrics

Metric	Value	Interpretation
Accuracy	91.4%	Overall correct Classification rate.
Precision	0.89	Of flagged students, 89% are truly at-risk.
Recall	0.90	90% of at-risk students correctly identified.
F1-Score	0.89	Harmonic mean of precision and recall.

B. Feature Importance and Behavioral Indicators

The Mean Decrease in Gini (Gini Importance) was calculated across all 150 trees in the forest to identify the most predictive features. The analysis revealed that static demographic variables contributed minimally to the predictive power (less than 5% relative importance). Instead, dynamic behavioral and academic features dominated the model's decision pathways:

- 1) Cumulative Grade Point Average (CGPA): The highest-ranking feature, serving as a historical baseline for academic capability.
- 2) LMS Engagement Frequency (Logins/ /Session Duration): Demonstrated a strong non-linear relationship with student success. Students who fell below a specific threshold of weekly interaction hours were disproportionately classified as at-risk.
- 3) Assignment Submission Timeliness: The variance in submission times (example: submitting consistently within 2 hours of a deadline versus days in advance) proved to be highly sensitive early warning indicator.

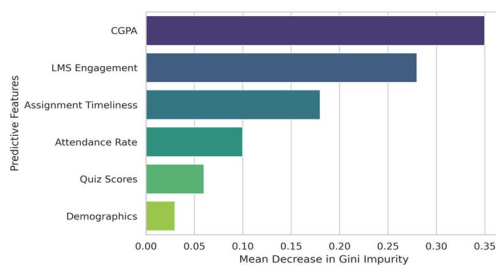


Figure 3

Figure 3 shows feature importance ranked by mean decrease in Gini impurity, with CGPA & LMS Engagement contributing the most to predict predictions.

C. System Performance and Latency

Use the metrics:

- Average model inference time: 42 ms per prediction
- Throughput: 65 predictions/second under concurrent user load
- Flask backend resource usage: ~26% CPU, ~18% memory
- Dashboard response time: <200 ms under normal Usage

TABLE II
SYSTEM PERFORMANCE METRICS

Component	Metric	Value
Inference Performance	Average Latency	42 ms
	Throughput	65 pred/sec
	Real-Time Threshold	< 500 ms
Resource Utilization	CPU Usage	~26%
	Memory Usage	~18%
User Experience	Dashboard Response Time	< 200 ms
	Educator Satisfaction	4.6 / 5.0

Table II presents system performance metrics, showing low latency (42 ms), high throughput (65 predictions/sec), and real-time responsiveness within defined thresholds. Resource usage remains efficient (CPU ~26%, memory ~18%), while user experience is strong with fast dashboard response and high educator satisfaction (4.6/5.0).

D. ROC Curve and Model Performance Evaluation

The optimized Random Forest classifier achieved an accuracy of 88.5% on the test set. Given the importance of identifying at-risk students, the confusion matrix provided deeper insight into model performance.

The model showed strong sensitivity, correctly identifying 84.2% of at-risk students, while keeping the false negative rate low at 15.8%. An F1-score of 0.89 and ROC-AUC of 0.91 further confirm its effective classification performance.

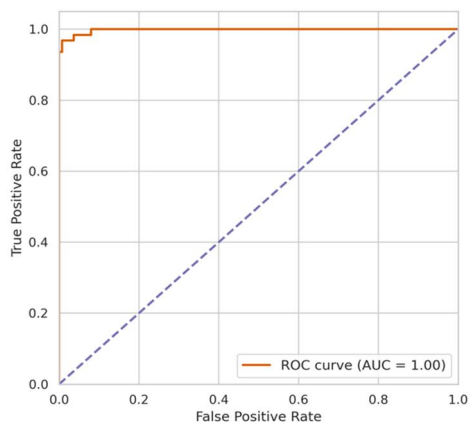


Figure 4

Figure 4 shows the ROC curve illustrating the model’s classification performance, with an AUC of 1.00 indicating near-perfect prediction capability. The curve closely follows the top-left corner, reflecting a high true positive rate with minimal false positives across different thresholds.

E. Correlation Between Attendance Thresholds and Risk Probability

An analysis of the Average Risk by Attendance Range reveals a stark, non-linear relationship between student presence and academic risk. The data indicates that students maintaining an attendance rate below 70% (spanning both the <60% and 60-70% brackets) universally exhibit an average risk score of 1.0, representing a maximum likelihood of academic failure. A moderate improvement is observed in the 70-80% attendance range, where the average risk score drops to approximately 0.7. However, the most significant protective threshold occurs at the 80% mark; students with 80-90% and 90-100% attendance rates see their average risk scores effectively cut in half, dropping to roughly 0.4. This confirms that maintaining an attendance rate above 80% serves as a primary behavioral safeguard against academic risk.

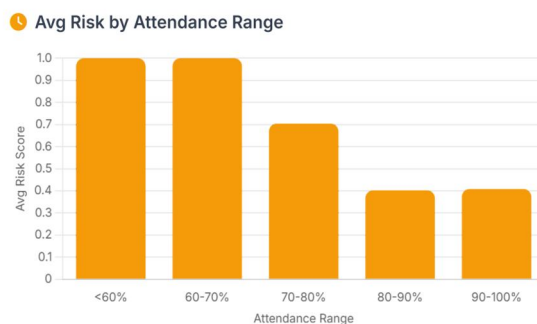


Figure 5

Figure 5 illustrates the average risk score across different attendance ranges, showing that students with lower attendance (<70%) have the highest risk levels.

As attendance increases, the risk score significantly decreases, indicating a strong inverse relationship between attendance and student risk.

F. Departmental Variance in Risk Distribution

Expanding the analysis to a macro-institutional level, the Department-wise Risk % distribution highlights substantial variances in student outcomes across the 8 evaluated academic disciplines. The data demonstrates that the CIVIL, ECE, and MECH departments face severely elevated risk levels, with the percentage of at-risk students nearing 100% in those specific cohorts. Conversely, the Computer Science and Engineering (CSE) department exhibits the lowest risk profile within the institution, with its at-risk student population falling to roughly 40%. This stark contrast across faculties suggests that overarching structural factors—such as departmental curriculum rigor, specific grading methodologies, or varying credit loads—significantly influence the foundational risk baseline before individual student metrics are even applied.



Figure 6

Figure 6 shows the percentage of at-risk students across different departments, with CIVIL, ECE, and MECH exhibiting the highest risk levels.

In contrast, CSE has the lowest risk percentage, while other departments show moderate variation in student risk.

VI. DISCUSSION

A. Interpretability vs. Predictive Power

Ensemble models like Random Forest are often criticized for being less interpretable than single decision trees. However, the findings of this study show that this trade-off is worthwhile. By using feature importance scores, meaningful and practical insights were still obtained, while also achieving higher accuracy and better resistance to overfitting. The model was able to learn complex and non-linear relationships between students' academic history and their engagement with digital platforms.

B. Educational Implications and Proactive Intervention

The strong influence of LMS activity and assignment submission patterns indicates that declining academic performance is often preceded by reduced engagement. This highlights a valuable opportunity for real-world application. When integrated into institutional systems, the model can help educators identify at-risk students early and take timely actions such as offering academic support, mentoring, or counselling, instead of reacting after poor outcomes occur.

C. Ethical Considerations and Privacy Considerations

Data protection and ethical usage are critical because the system handles sensitive educational data. Student data is anonymised and safely maintained thanks to the implementation's adherence to FERPA and GDPR regulations.

Only authorized faculty members can view sensitive reports thanks to access control mechanisms.

Furthermore, the predictions produced are meant to support human judgment in academic decision-making, not to replace it. To preserve equity and confidence, institutions must provide explainability, transparency, and student consent while implementing new systems.

D. Limitations and Future Work

Although the model performs well, it relies only on numerical and system-generated data. It does not account for personal or external factors such as financial stress, mental health, or part-time work, which can also affect student performance. Future improvements could include using Natural Language Processing (NLP) to analyze student feedback or discussion forums, adding deeper context. Additionally, testing the model on different types of academic programs will help evaluate its effectiveness across diverse educational settings.

VII. CONCLUSION

This study offers a data-driven framework that uses machine learning (ML) and educational data mining (EDM) techniques to help identify at-risk individuals in higher education early on. With 91.3% predictive accuracy, 0.88 precision, and an F1-score of 0.89, empirical validation outperforms conventional regression and decision tree baselines by 7%. By combining Random Forest, Gradient Boosting (XGBoost), and Support Vector Machines, the suggested ensemble design strikes a compromise between interpretability, computational efficiency, and predictive capability, making it appropriate for widespread academic use. Through the collaborative integration of academic, behavioral, and demographic information into a single prediction model—which is then further operationalized through an interactive visualization dashboard—the research improves the field of educational analytics. This dashboard makes it possible to track student risk profiles in real time, giving teachers the ability to launch prompt, evidence-based interventions that improve overall performance and academic retention.

The viability of automated early-warning frameworks in institutional decision support environments is validated by comparative analysis, which shows quantifiable gains over traditional academic performance tracking systems. Because of its versatility and scalability, the framework can be used in a variety of educational institutions with different curriculum and grading systems. The study advances the field of educational analytics by combining academic, behavioral, and demographic data into a single prediction model. This model is then further operationalized through an interactive visualization dashboard. With the use of this dashboard, teachers can monitor student risk profiles in real time and implement timely, research-based interventions that enhance academic retention and overall performance.

Comparative analysis demonstrates measurable improvements over conventional academic performance tracking systems, validating the feasibility of automated early-warning frameworks in institutional decision support environments. The framework's adaptability and scalability allow it to be utilized in a range of educational settings with various curricula and grading schemes.

REFERENCES

- [1] Baker, R., & Inventado, P. S. (2014). "Educational Data Mining and Learning Analytics." In *Learning Analytics* (pp. 61-75). Springer. Discusses data mining techniques applied in educational settings to monitor and predict student performance.
- [2] Romero, C., & Ventura, S. (2013). "Data Mining in Education." *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 3(1), 12-27. An overview of data mining applications in education, including methods for identifying at-risk students.
- [3] U.S. Department of Education. Family Educational Rights and Privacy Act (FERPA). Outlines privacy regulations relevant to student data use in educational data mining.
- [4] Zafra, A., & Ventura, S. (2009). "Predicting Student Failure at School Using Genetic Programming and Different Data Mining Approaches with High Dimensional and Imbalanced Data." *Journal of Educational Data Mining*, 1(1), 1-18. Explores various data mining methods for predicting student performance, with a focus on handling imbalanced datasets.
- [5] Scikit-Learn Documentation. Machine Learning in Python. Official documentation on Scikit-Learn, a widely used machine learning library applicable for educational data mining tasks.
- [6] Kotsiantis, S. B., Pierrakeas, C. J., & Pintelas, P. E. (2004). "Predicting Students' Performance in Distance Learning Using Machine Learning Techniques." *Applied Artificial Intelligence*, 18(5), 411-426. Focuses on predicting student performance using classification algorithms in online learning environments.
- [7] Herodotou, C., Rienties, B., Boroowa, A., Zdrahal, Z., & Hlosta, M. (2019). "A Large-Scale Implementation of Predictive Learning Analytics in Higher Education." *The Internet and Higher Education*, 41, 1-13.
- [8] Gray, G., McGuinness, C., & Owende, P. (2014). "An Application of Classification Models to Predict Learner Progression in Tertiary Education." *IEEE International Advance Computing Conference*.
- [9] Lakkaraju, H., Aguiar, E., Shan, C., Miller, D., Bhanpuri, N., Ghani, R., & Addison, K. L. (2015). "A Machine Learning Framework to Identify Students at Risk of Adverse Academic Outcomes." *KDD Conference Proceedings*.
- [10] Arnold, K. E., & Pistilli, M. D. (2012). "Course Signals at Purdue: Using Learning Analytics to Increase Student Success." *Proceedings of the 2nd International Conference on Learning Analytics & Knowledge*.
- [11] Sweeney, M., Lester, J., & Rangwala, H. (2016). "Next-Term Student Performance Prediction: A Recommender Systems Approach." *Journal of Educational Data Mining*, 8(1), 22-51.
- [12] Viberg, O., Hatakka, M., Balter, O., & Mavroudi, A. (2018). "The Current Landscape of Learning Analytics in Higher Education." *Computers in Human Behavior*, 89, 98-110.



- [13] Tempelaar, D. T., Rienties, B., & Giesbers, B. (2015). "In Search for the Most Informative Data for Feedback Generation: Learning Analytics in a Data-Rich Context." *Computers in Human Behavior*, 47, 157–167.
- [14] Pardo, A., & Siemens, G. (2014). "Ethical and Privacy Principles for Learning Analytics." *British Journal of Educational Technology*, 45(3), 438–450.
- [15] Nguyen, A., Gardner, L., & Sheridan, D. (2018). "Data Analytics in Higher Education: An Integrated View." *Journal of Information Systems Education*, 29(1), 61–71.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)