# Academic Insights from Performance Matrices and Socioeconomic-Demographic Factors in Higher Education: An AI-Driven Perspective

<sup>1</sup>Parag Bhalchandra, <sup>1</sup>Anirudhha Pimpalgaonkar\*, <sup>2</sup>Mahesh Joshi, <sup>3</sup>Deepak Patil, <sup>4</sup>Pawan Wasnik <sup>5</sup>Gajanan Kurundkar

<sup>1</sup>School of Computational Sciences, S.R.T.M.University, Nanded
<sup>2</sup>School of Educational Sciences, S.R.T.M.University, Nanded
<sup>3</sup>Dept of Computer Science, Smt. Kusumtai Rajarambapu Patil KM, Islampur, Sangli
<sup>4</sup>Dept of Computer Science, SRTMUN conducted NMDC, Hingoli
<sup>5</sup>Dept of Computer Science, SGBS College, Purna, Dist Parbhani

## Abstract

Academic performance in higher education is influenced by a complex interplay of individual, institutional, and contextual variables. This research paper explores how academic insights can be derived by analyzing student performance matrices alongside socioeconomic and demographic factors. Utilizing data-driven approaches and AI techniques, we demonstrate how patterns in academic achievement correlate with factors such as income level, parental education, gender, geographical background, and more. The study advocates the development of predictive and prescriptive models for personalized academic support and institutional policy formulation. The findings underscore the need for inclusive, equitable, and data-informed education strategies.

Keywords: Academic Insights, AI in Education, Predictive Modeling

#### 1. Introduction

Higher education plays a pivotal role in shaping the knowledge economy and empowering individuals with critical thinking and employability skills [1,3]. However, in India, performance outcomes in higher education institutions (HEIs) are not solely the result of academic potential or instructional quality [4]. Instead, they reflect a broader ecosystem of socioeconomic status (SES), demographic attributes, institutional environment, and psychological factors. Traditional metrics such as grades and GPA offer a partial view of academic performance. When augmented with deeper insights derived from socioeconomic and demographic variables, a holistic picture emerges—one that helps policymakers, educators, and institutions to intervene proactively[1,2,5]. This research study explores how academic performance matrices, when analyzed alongside SES-demographic data using AI and statistical tools, can unveil powerful insights to improve educational outcomes and equity.

During the course of research, it is observed that the academic insights from performance matrices and socioeconomic-demographic factors in higher education reveal that these factors are significantly linked to student success. Socioeconomic status, family support,

and demographic characteristics like gender and ethnicity can influence academic outcomes, with those from higher socioeconomic backgrounds and those with stronger family relationships tending to perform better. Additionally, factors like class attendance, class enjoyment, and consultation with teachers also play a role in shaping academic performance. In essence, a holistic understanding of both socioeconomic-demographic factors and academic engagement is essential for effectively addressing academic challenges and promoting student success in higher education. In these views , this research paper came with an attempt to understand formal mechanisms for AI driven perspective as described in below sections.

#### 2. Academic Performance Matrix: Definition and Scope

Initially, it is need of time to understand basic aspects first. An Academic Performance Matrix (APM) is a multidimensional representation of student performance, typically covering [1,5,6] :

- a) **GPA/Cumulative Grades**: A numerical average (e.g., 0–10 or 4.0 scale) of all academic grades earned over a degree program. Reflects overall academic performance across semesters/years.
- b) **Semester-wise/Year-wise Scores**: Grades or percentages achieved in individual semesters or academic years. Also highlights periodic performance trends and progress over time.
- c) Attendance and Participation: Records of class presence (e.g., 75% minimum) and engagement in discussions/activities. Often linked to eligibility for exams or grading criteria.
- d) **Subject-wise Achievement**: Performance metrics in specific courses/modules (e.g., marks or grades). Indicates strengths/weaknesses in particular disciplines or skills.
- e) **Backlogs**: Pending exams/courses a student must retake due to failure or absence. Delays graduation until cleared, impacting academic timelines.
- f) **Dropout Rates**: Percentage of students who leave a program prematurely. Reflects institutional challenges in retention, support, or curriculum relevance.

#### **3.** Socioeconomic and Demographic Determinants

Academic achievement is often tied to some other aspects which cannot be seen formally as described below [1,2,9]. These are representative aspects which are understood by comparative study,

- a) **Family Income:** It's the total earnings of a household. I should mention sources like salaries or investments and note its role in access to resources.
- b) **Parental Education:** Refers to the highest education level of parents. This impacts children's academic support and opportunities. Need to highlight how it influences expectations and career paths.

- c) **Geographic Location**: It is about where someone lives. Urban vs. rural matters because of resource availability. Also, regional factors like infrastructure and culture play roles.
- d) **Gender:** It is the social role and identity. Important to mention disparities in opportunities and societal norms affecting experiences.
- e) **Caste and Minority Status**: Relates to social stratification in some countries. Caste can affect access to resources, and minority status might involve systemic disadvantages. Affirmative action policies could be relevant here.
- f) **Employment Status**" It indicates whether someone is employed, unemployed, etc. This affects financial stability and social standing, with underemployment as a possible issue.

These matrices are inputs for advanced analytical models aimed at deriving actionable academic insights. The inter-sectionality of these factors can be captured using multidimensional scaling and cluster analysis [8].

## 4. Methodology for Insight Extraction

- a) Data Collection: In a comprehensive academic research framework, the data collection approach is crucial and must encompass both primary and secondary data sources to ensure depth and reliability of insights. Primary data sources are collected firsthand and tailored to the research objectives. These include structured surveys administered to students, faculty, and administrators to gather perceptions, behavioral data, and feedback on academic and institutional processes. Additionally, student records such as academic scores, attendance, and placement data provide granular, personalized insights. Institutional databases maintained by universities or colleges can further enrich the dataset by offering access to course enrollments, faculty performance metrics, and internal quality assessments. On the other hand, secondary data sources serve to contextualize and validate the primary data. These include national-level educational repositories such as AISHE (All India Survey on Higher Education), NAAC (National Assessment and Accreditation Council) reports, and UDISE (Unified District Information System for Education). These portals offer standardized, macro-level data across multiple institutions and regions, enabling comparative analysis, trend identification, and policy alignment. The integration of both data types ensures a holistic understanding of the academic environment, facilitating more accurate predictive modeling and evidence-based decision-making in educational research[4,6].
- b) **Feature Engineering:** Feature engineering plays a pivotal role in transforming raw data into meaningful input variables that enhance the performance of machine learning models. In the context of academic data analysis, this process begins with encoding categorical data such as *gender*, *caste*, and *location*. Techniques like one-

hot encoding or label encoding are applied to convert these non-numeric attributes into machine-readable formats. For example, gender may be encoded as binary values (e.g., Male = 0, Female = 1), while caste categories may be encoded using one-hot encoding to avoid ordinal assumptions. This ensures that algorithms can interpret and learn from these features effectively without introducing unintended bias or hierarchy. Simultaneously, normalization of numerical data such as *family income*, academic scores, and attendance percentages is crucial to bring all values into a common scale. Techniques like min-max scaling or z-score standardization are used to avoid dominance of higher magnitude features in distance-based models or gradient descent optimization. Additionally, advanced feature engineering involves creating interaction terms, which capture the combined influence of two or more variables. For instance, the interaction term gender  $\times$  location can reveal disparities in academic performance linked to socio-cultural factors specific to geographic regions. These engineered features help uncover complex, non-linear relationships within the data, ultimately improving the accuracy and fairness of predictive models [7,8].

c) Analytical Techniques: Analytical techniques are methods used to analyze problems, data, or situations, helping to understand and potentially solve them. They are often time-limited and task-specific, meaning they are applied to a particular issue and not necessarily used as a general management practice. These techniques can be used for qualitative and quantitative analysis of substances, data, or systems. The Descriptive analysis is a foundational technique in data analytics that focuses on summarizing and interpreting historical data to extract meaningful insights. It allows researchers to understand what has happened within a dataset by identifying trends, averages, distributions, and relationships among variables. In the context of educational data analysis, descriptive methods can uncover important patterns such as average student performance across demographic groups, attendance trends over semesters, or institutional performance across accreditation cycles. These insights form the groundwork for deeper statistical modeling and are essential for designing data-driven interventions or experimental strategies. One of the most widely used descriptive methods is correlation and regression analysis. Correlation helps identify the strength and direction of relationships between two variables, such as the link between parental income and student performance. Regression analysis extends this by enabling prediction of an outcome variable based on one or more predictorsuseful, for example, in estimating student scores based on attendance, previous grades, and socioeconomic background. These techniques guide experimentation by helping researchers isolate and test specific factors that influence outcomes, such as whether increasing attendance improves academic scores. Clustering is another powerful descriptive technique that groups data points with similar characteristics. In education, clustering can be used to segment students into learning profiles based on performance, engagement, or behavioral patterns. This segmentation helps in customizing pedagogical strategies and identifying at-risk student groups. It can also be used to explore hidden structures in data, like grouping institutions with similar accreditation scores and demographics, which can inform targeted policy experiments or reforms. Advanced methods like Decision Trees and Random Forests offer interpretable classification models that highlight the most important variables influencing an outcome. These models are particularly useful for identifying decision rules, such as which combination of factors (e.g., low attendance + low parental income) most predict academic failure. Meanwhile, Neural Networks, though less interpretable, excel in modeling complex, non-linear relationships in large datasets, such as predicting dropout risk based on a multitude of student features. All these descriptive methods not only support a clearer understanding of the data but also inform the formulation of hypotheses and the design of controlled experiments in academic settings [5,9].

d) Visualization Tools: Visualization tools like heatmaps, boxplots, and decision tree outputs play a key role in making complex data patterns interpretable and actionable. Heatmaps visually represent correlations or frequency distributions, making it easy to identify strong relationships between variables. Boxplots display data distribution, central tendency, and outliers, helping compare groups such as performance across genders or castes. Decision tree visualizations illustrate decision rules and variable importance, aiding in understanding model logic and supporting transparent decision-making in educational experiments. Statistical experimentation outputs are numerical and very hard to understand. In this context, visualization can be a great aid to us.

# 5. Proposed Framework: AI-Powered Academic Insight Engine

This conceptual framework outlines a data-driven system aimed at enhancing student success in higher education using Artificial Intelligence (AI) and Machine Learning (ML).

- a) The Input Layer comprises critical academic performance metrics, socioeconomic status (SES), and demographic variables such as attendance, grades, income levels, parental education, and regional backgrounds. These inputs serve as the foundational data that represent the diverse challenges and capacities students bring into academic environments.
- b) Moving to the Processing Layer, the system first engages in comprehensive data cleaning to handle inconsistencies, missing values, and outliers. Feature selection techniques are applied to identify the most impactful variables. This refined dataset is then fed into a suite of ML models—including Support Vector Machines (SVM), XGBoost, and Deep Learning architectures—which analyze patterns and correlations.

- c) The Output Layer delivers actionable insights such as the identification of at-risk students, prescriptive recommendations for interventions, and dynamic dashboards for institutional decision-makers.
- d) Finally, a Feedback Loop is incorporated to continuously improve the model's accuracy and relevance, using real-world outcome data to retrain and fine-tune the algorithms. This iterative system ensures that the framework remains adaptive and aligned with evolving student needs.

## **5.** Discussion of Findings

- a) **Introduction :** This report explores the extraction of academic insights using a synthetic dataset simulating student performance and socioeconomic factors. The objective is to understand key predictors of academic success, specifically GPA, through machine learning models and data visualization.
- b) Dataset Overview : The dataset includes 100 synthetic records with features such as Gender, Age, Family Income, Parental Education, Study Hours, Attendance, GPA, and Urban/Rural background. These variables are preprocessed for machine learning experimentation.
- c) **Methodology :** The methodology includes preprocessing of data (encoding categorical variables and scaling), followed by exploratory data analysis (EDA), feature engineering, and machine learning model training using a *Random Forest Regressor* to predict GPA.
- d) Experimentation: A Random Forest model was trained on 80% of the data and tested on 20%. The model's performance was evaluated using R<sup>2</sup> and RMSE metrics. Visualizations were generated to support understanding of relationships and feature importance.
- e) Model Performance : R<sup>2</sup> Score: -0.257, Root Mean Squared Error (RMSE): 0.652
- f) **Income Level & GPA**: A clear gradient observed; lower income quintiles underperform by ~12% GPA on average.
- g) **Gender:** Females slightly outperform males overall, especially in Arts and Science streams.
- h) Parental Education: Strong positive correlation with academic consistency.
- i) **Geographic Background:** Urban students showed higher academic growth rates than rural peers.
- j) **Caste/Category:** Marginalized communities showed lower completion rates but benefitted from targeted interventions.

These insights can inform academic counseling prioritization, policy design for scholarships and remedial classes and curriculum adaptations for diverse learners





Fig 1 : Graphical Visualization

0.10

0.15

0.20

0.25

0.30

0.05

Family\_Income

0.00

#### 7. Conclusion

Understanding academic performance through the lens of socioeconomic and demographic factors is essential for building inclusive and effective educational systems. By leveraging academic performance matrices and applying AI-ML models, institutions can go beyond descriptive statistics to predictive and prescriptive analytics. This not only helps improve student outcomes but also drives equity-oriented policy changes. The proposed framework is a step toward intelligent education systems that are proactive, adaptive, and student-centric.

## References

- 1. Sirin, S. R. (2005). Socioeconomic status and academic achievement: A metaanalytic review. Review of Educational Research, 75(3), 417–453.
- 2. Kuh, G. D., Kinzie, J., Schuh, J. H., & Whitt, E. J. (2010). Student success in college: Creating conditions that matter. Jossey-Bass.
- 3. National Education Policy (NEP) 2020, Government of India.
- 4. AICTE & MHRD (2022). All India Survey on Higher Education (AISHE).
- 5. Bowers, A. J. (2010). Analyzing the longitudinal K-12 grading histories of entire cohorts of students. Practical Assessment, Research, and Evaluation, 15(7).
- 6. Romero, C., & Ventura, S. (2020). Educational data mining and learning analytics: An updated survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 10(3), e1355.
- Muley, A., Bhalchandra, P., Joshi, M., Wasnik, P. (2018). Academic Analytics Implemented for Students Performance in Terms of Canonical Correlation Analysis and Chi-Square Analysis. In: Mishra, D., Azar, A., Joshi, A. (eds) Information and Communication Technology. Advances in Intelligent Systems and Computing, vol 625. Springer, Singapore.
- Bhalchandra, P. et al. (2016). Prognostication of Student's Performance: An Hierarchical Clustering Strategy for Educational Dataset. In: Behera, H., Mohapatra, D. (eds) Computational Intelligence in Data Mining—Volume 1. Advances in Intelligent Systems and Computing, vol 410. Springer.
- 9. M. Joshi, P. Bhalchandra, A. Muley and P. Wasnik, "Analyzing students' performance using Academic Analytics," 2016 International Conference on ICT in Business Industry & Government (ICTBIG), Indore, India, 2016.