

Student Dropout Prediction: A Data-Driven Approach to Identifying At-Risk High School Students

ABSTRACT

This research explores the patterns and predictive factors of high school student dropout using a comprehensive dataset of US high school students. The study employs machine learning techniques, specifically Random Forest and XGBoost algorithms, to identify key indicators that contribute to dropout risk. Through data preprocessing, feature engineering, and model evaluation, the research revealed that academic performance metrics (GPA, attendance rates, and test scores), parental education levels, and social worker visits are significant predictors of dropout risk. The developed model achieved an accuracy of over 90%, demonstrating strong predictive capability. This paper provides insights that could help educational institutions develop targeted intervention strategies to support at-risk students and improve retention rates. The findings emphasize the importance of early detection and holistic support systems that address both academic and social factors contributing to student dropout.

TABLE OF CONTENTS

ABSTRACT	i
TABLE OF CONTENTS	ii
CHAPTER I: INTRODUCTION	1
1.1. Research background	1
1.2. Problem statements.....	1
1.3. Objectives of research.....	2
1.4. Significant of research	2
CHAPTER II: LITERATURE REVIEW	3
CHAPTER IV: METHODOLOGY	4
4.1. Dataset Collection	4
4.2. Data Preprocessing	4
4.3. Prediction Models	6
4.4. Hyperparameter Turning	7
4.5. Evaluation Metrics.....	7
CHAPTER V: RESULTS AND DISCUSSION	8
5.1. Exploratory Data Analysis	8
5.2. Results of Prediction Model	9
CHAPTER VI: CONCLUSION AND RECOMMENDATION.....	12
6.1. Conclusion	12
6.2. Recommendation.....	12
REFERENCES.....	14

CHAPTER I: INTRODUCTION

1.1. Research background

High school dropout remains a critical issue in the American education system, with far-reaching consequences for individuals and society. According to the National Center for Education Statistics, while the overall dropout rate has declined over the past decade, approximately 5.1% of high school students still fail to complete their education. This translates to hundreds of thousands of students annually who do not receive their high school diploma, significantly limiting their future educational and career opportunities.

The impacts of dropping out extend beyond the individual student. Research consistently demonstrates that high school dropouts face higher rates of unemployment, lower lifetime earnings, increased likelihood of incarceration, and poorer health outcomes. From a societal perspective, dropouts represent a substantial economic cost through reduced tax revenue, increased social services expenditure, and higher crime rates.

Traditional approaches to addressing dropouts have often been reactive rather than preventive, focusing on intervention after warning signs become apparent. However, with the advent of data science and machine learning techniques, there is tremendous potential to develop predictive models that can identify at-risk students before they drop out, enabling more timely and targeted interventions.

1.2. Problem statements

Despite the recognized importance of preventing student dropout, many educational institutions lack effective systems to identify students at risk before they leave school. Current approaches often rely on obvious warning signs such as prolonged absences or failing grades, which may appear too late for effective intervention. Additionally, there is limited understanding of the complex interplay of factors that contribute to dropout decisions, including academic performance, socioeconomic background, parental involvement, and school environment.

1. This research addresses several key questions:
2. What factors are most predictive of high school student dropout?
3. How accurately can machine learning models predict student dropout risk?
4. What interventions might be most effective based on the identified risk factors?
5. How can schools leverage data-driven approaches to implement early warning systems?

1.3. Objectives of research

The primary objectives of this research are:

1. To develop and evaluate machine learning models that can accurately predict high school student dropout risk based on available student data.
2. To identify the most significant factors associated with dropout risk to inform targeted intervention strategies.
3. To provide recommendations for implementing data-driven early warning systems in educational institutions.
4. To contribute to the understanding of dropout patterns and risk factors in the contemporary educational landscape.

1.4. Significant research

This research has significant implications for educational practice and policy:

1. **Early Intervention:** By identifying students at risk of dropping out before they show obvious signs, schools can implement timely interventions that may prevent dropout.
2. **Resource Allocation:** Understanding the most significant dropout risk factors can help schools allocate limited resources more effectively toward the most impactful interventions.
3. **Policy Development:** Findings can inform evidence-based policies aimed at reducing dropout rates at the school, district, and state levels.
4. **Educational Equity:** Identifying patterns in dropout risk can highlight systemic inequities and help target support to vulnerable student populations.
5. **Methodological Contribution:** The study demonstrates the application of advanced machine learning techniques to educational data, providing a framework for similar analyses in other educational contexts.

CHAPTER II: LITERATURE REVIEW

The phenomenon of high school dropout has been extensively studied across various disciplines including education, sociology, psychology, and more recently, data science. Early research on dropout prediction by Rumberger (1983) identified demographic and family background factors as key predictors of dropout risk. Subsequent work by Finn (1989) proposed theoretical models explaining the process of disengagement that leads to dropout, highlighting both academic and behavioral factors.

More recent studies have focused on the identification of early warning indicators. Balfanz et al. (2007) demonstrated that attendance, behavior, and course performance in middle school could predict over 50% of eventual dropouts. Bowers et al. (2013) conducted a comprehensive review of dropout predictors, finding that academic performance, especially GPA, was among the strongest predictors across multiple studies.

The application of machine learning to educational data has grown rapidly in recent years. Baker and Inventado (2014) reviewed educational data mining techniques, highlighting their potential for identifying at-risk students. Researchers have explored various algorithms, with ensemble methods like Random Forest and Gradient Boosting often showing superior performance (Adejo and Connolly, 2018).

Several studies have addressed the challenge of imbalanced datasets in dropout prediction. Thammasiri et al. (2014) compared several techniques for handling class imbalance, finding that oversampling methods improved model performance. Meanwhile, the interpretability of machine learning models has been addressed by Xing et al. (2019), who emphasized the importance of feature importance analysis in educational contexts.

The present study builds upon this rich literature by applying contemporary machine learning techniques to a comprehensive dataset of US high school students. Unlike many previous studies that focused on limited geographical areas or small sample sizes, this research leverages a diverse national dataset. Additionally, the study's emphasis on both prediction and interpretation addresses a critical gap in the literature, as many previous works have prioritized predictive accuracy over actionable insights.

CHAPTER IV: METHODOLOGY

4.1. Data Collection

This study utilized a comprehensive dataset of US high school students containing 5,400 records collected from various states across America. The dataset, named "US_Highschool_Student_data.csv," includes a wide range of student information, including demographic details, academic performance metrics, family background, and behavioral indicators.

The dataset includes variables such as:

1. Demographic information: Sex, Age, Name, State, Address
2. Family structure: Family size (Famsize), Parent's cohabitation status (Pstatus)
3. Parental education: Mother's education (Medu), Father's education (Fedu)
4. Parental occupation: Mother's job (Mjob), Father's job (Fjob)
5. Academic performance: Math scores, Reading scores, Writing scores
6. Behavioral indicators: Attendance rate, Suspensions, Expulsions
7. Support systems: Teacher support, Counseling, Social worker visits
8. Parental involvement
9. GPA (Grade Point Average)

The data represents students from various states including California, Nevada, Utah, Oregon, and Arizona, with both urban and rural settings represented. The dataset provides a cross-sectional snapshot of student characteristics, allowing for analysis of factors associated with dropout risk.

4.2. Data Preprocessing

Several preprocessing steps were performed to prepare the data for analysis and modeling:

1. Removal of unnecessary columns: Name and Address fields were removed as they were not relevant for predicting dropout risk.

2. Missing value handling: The dataset contained missing values in several columns including Math_Score (8.1%), Reading_Score (8.0%), Writing_Score (8.1%), Attendance_Rate (1.3%), Expulsions (5.9%), and GPA (12.9%). These missing values were imputed using the median value for each column to maintain the overall distribution of the data.
3. Categorical data standardization: Categorical variables such as Teacher_Support were standardized to ensure consistency. For instance, numerical ratings (1-3) were converted to "Low," ratings (4-5) to "Medium," and rating 6 to "High" to match the existing categorical labels.
4. Feature encoding: Categorical variables were encoded using Label Encoding to convert them into a format suitable for machine learning algorithms. This included variables like Sex, State, Family size, Parent's status, and educational levels.
5. Target variable creation: A binary "Dropout" variable was created based on three criteria:
 6. GPA less than 2.0
 7. Attendance rate below 60%
 8. Any record of expulsions
9. Students meeting any of these criteria were labeled as dropouts (1), while others were labeled as non-dropouts (0). This resulted in 531 students (9.8%) classified as dropouts and 4,869 (90.2%) as non-dropouts.
10. Correlation analysis: A correlation matrix was generated to understand the relationships between different variables and the target dropout variable. This helped identify potentially important predictors and redundant features.
11. Class imbalance assessment: The significant imbalance between dropout (9.8%) and non-dropout (90.2%) classes was noted for consideration during model training and evaluation.

4.3. Prediction Models

4.3.1. Random Forest

The first predictive model implemented was Random Forest, an ensemble learning method that operates by constructing multiple decision trees during training and outputting the class that is the mode of the classes from individual trees. Random Forest was chosen for its ability to handle non-linear relationships, robustness to overfitting, and capability to provide feature importance scores.

The Random Forest classifier was configured with the following parameters:

1. Number of estimators (trees): 100
2. Random state: 42 (for reproducibility)
3. Default values were used for other parameters

The dataset was split into training (80%) and testing (20%) sets using stratified sampling to maintain the same class distribution in both sets. The model was trained on the training data and evaluated on the test data.

4.3.2. XGBoost

The second model implemented was XGBoost (Extreme Gradient Boosting), an optimized distributed gradient boosting library. XGBoost was selected due to its performance advantages in structured data problems, efficiency, and ability to handle imbalanced datasets effectively.

The XGBoost classifier was configured with the following parameters:

1. Number of estimators: 100
2. Learning rate: 0.1
3. Maximum depth: 5
4. Random state: 42 (for reproducibility)

Based on the feature importance analysis from the Random Forest model, less important features were excluded from the XGBoost model to improve computational efficiency and

reduce noise. The features removed included "State," "Fjob," "Mjob," "Age," "Guardian," "Counseling," "Sex," "Teacher_Support," "Famsize," and "Pstatus."

4.4. Hyperparameter Turning

During the analysis, a concerning pattern was observed: the `Social_Worker_Visits` variable showed an extremely high correlation (0.82) with the target `Dropout` variable, suggesting potential data leakage or a circular definition. After investigation, this feature was removed from the final models to ensure that predictions were based on genuinely predictive features rather than potentially problematic data relationships.

The models were then retrained with the refined feature set. While formal grid search or random search hyperparameter tuning was not explicitly shown in the notebook, the parameters were manually adjusted based on initial performance results. The final models used the parameters specified in the previous section, which provided a good balance between computational efficiency and predictive performance.

4.5. Evaluation Metrics

Several metrics were used to evaluate the performance of the predictive models:

1. **Accuracy:** The proportion of correct predictions (both true positives and true negatives) among the total number of cases examined.
2. **Classification Report:** Providing precision, recall, and F1-score for each class:
3. **Precision:** The ratio of correctly predicted positive observations to the total predicted positives
4. **Recall:** The ratio of correctly predicted positive observations to all actual positives
5. **F1-score:** The weighted average of Precision and Recall
6. **Feature Importance:** Analysis of which features contributed most significantly to the model's predictions, providing insights into the factors most associated with dropout risk.

Given the class imbalance in the dataset, particular attention was paid to the performance metrics for the minority class (dropout), as high overall accuracy could be misleading if the model simply predicted the majority class (non-dropout) in most cases.

CHAPTER V: RESULTS AND DISCUSSION

5.1. Exploratory Data Analysis

The exploratory data analysis revealed several important insights about the dataset and the relationships between various student characteristics and dropout risk:

1. **Dropout Distribution:** Of the 5,400 students in the dataset, 531 (9.8%) were classified as dropouts based on our criteria ($\text{GPA} < 2.0$, $\text{attendance} < 60\%$, or $\text{expulsions} > 0$), while 4,869 (90.2%) were non-dropouts. This imbalance is typical in educational datasets focused on dropout prediction.
2. **Correlation Analysis:** The correlation matrix revealed significant relationships between several variables and dropout status
3. Social Worker Visits showed an extremely high positive correlation (0.823) with dropout, which raised concerns about potential data leakage or circular definitions in the dataset.
4. Expulsions had a strong positive correlation (0.574) with dropout, which is expected given that expulsions were part of our dropout classification criteria.
5. GPA (-0.398), test scores (Math: -0.457, Reading: -0.449, Writing: -0.444), and attendance rates (-0.538) all showed moderate negative correlations with dropout, indicating that higher academic performance and attendance are associated with lower dropout risk.
6. Parental education levels (Mother's education: -0.567, Father's education: -0.567) showed strong negative correlations with dropout risk, suggesting the significant impact of family educational background.
7. **Geographic and Demographic Patterns:** The analysis indicated some variations in dropout rates across different states and between urban and rural settings, though these were not among the strongest predictors. Similarly, while there were slight differences between male and female dropout rates, sex was not a strong predictor.
8. **Family Structure Impact:** Family size and parent's cohabitation status showed weak correlations with dropout, suggesting that these structural factors may be less influential than other variables like parental education and involvement.

5.2. Results of Prediction Model

Both machine learning models demonstrated strong performance in predicting student dropout, though with some important nuances:

1. Random Forest Model :

- Achieved an overall accuracy of approximately 94%
- For the majority class (non-dropout), the model achieved high precision (0.95) and recall (0.98)
- For the minority class (dropout), the model achieved high precision (0.85) but moderate recall (0.65), indicating that while the model's dropout predictions were usually correct, it missed some actual dropout cases
- The F1-score for dropout prediction was 0.73, reflecting the balanced performance between precision and recall

2. XGBoost Model:

- After removing potentially problematic features (including Social_Worker_Visits), the XGBoost model achieved comparable overall accuracy to the Random Forest model
- The model showed slightly improved recall for the dropout class (0.68) compared to the Random Forest model
- The F1-score for dropout prediction was 0.75, slightly better than the Random Forest model
- **Feature Importance Analysis:**
 - Both models identified similar key predictors of dropout risk
 - The top predictors from the Random Forest model, in order of importance, were:
 - Attendance Rate
 - GPA

- Parental Education (Mother's and Father's)
- Math, Reading, and Writing Scores
- Parental Involvement
- Suspensions
- The XGBoost model, after feature selection, showed a similar pattern of feature importance but with a slightly stronger emphasis on academic performance metrics

3. **Model Refinement:**

- After removing the Social_Worker_Visits variable due to concerns about data leakage, the model performance decreased slightly but remained strong, confirming that the models were not overly dependent on this potentially problematic feature
- The refined models still achieved accuracy above 90%, demonstrating robust predictive capability based on legitimate predictors

5.3. Feature importance Analysis

The analysis of feature importance provided valuable insights into the factors that most strongly predict dropout risk:

1. **Attendance Rate:** Emerged as the most important predictor in both models, highlighting the critical role of consistent school attendance in student success. This aligns with previous research suggesting that attendance patterns can be early warning signs of disengagement.
2. **Academic Performance:** GPA and standardized test scores (Math, Reading, Writing) were consistently among the top predictors. This supports the notion that academic struggle is a key risk factor for dropout and suggests that academic support interventions may be particularly important.
3. **Parental Education:** Both mother's and father's education levels were strong predictors, underscoring the influence of family educational background on student

outcomes. This suggests that additional support may be needed for first-generation students whose parents have lower educational attainment.

4. **Parental Involvement:** Ranked as a significant predictor, indicating that student success is influenced by the level of engagement from parents in their education. This points to the potential value of programs that strengthen school-family connections.
5. **Behavioral Indicators:** Suspensions showed moderate predictive importance, suggesting that disciplinary issues can be warning signs of dropout risk. However, their lower ranking compared to academic factors indicates that behavioral issues alone may not be the primary drivers of dropout.
6. **Demographic and Geographic Factors:** Variables such as sex, state, and urban/rural setting had relatively low importance in the models. This suggests that while demographic factors may have some influence, they are less predictive than academic performance, attendance, and family background factors.

CHAPTER VI: CONCLUSION AND RECOMMENDATION

6.1. Conclusion

This study has successfully developed machine learning models capable of predicting high school student dropout risk with high accuracy. The research has identified key factors associated with dropout risk and demonstrated the potential of data-driven approaches to inform early intervention strategies.

Key findings from the research include:

1. Attendance rates, academic performance metrics (GPA and standardized test scores), and parental education levels emerged as the strongest predictors of dropout risk. This highlights the multifaceted nature of dropout risk, encompassing both in-school factors and family background characteristics.
2. Both Random Forest and XGBoost models achieved high accuracy (>90%) in predicting dropout risk, with the XGBoost model showing slightly better performance in identifying students at risk of dropping out. This demonstrates the viability of machine learning approaches for developing early warning systems.
3. The high predictive power of attendance rates suggests that monitoring and addressing attendance issues early could be one of the most effective strategies for preventing dropout. This is a particularly actionable insight as attendance is relatively easy to track in real-time.
4. The strong influence of parental education levels and involvement indicates that family factors play a significant role in student outcomes. This suggests that effective dropout prevention may need to include family engagement and support components.
5. The relatively lower importance of demographic and geographic factors suggests that dropout risk transcends these categories, and interventions should focus primarily on addressing academic struggles, attendance issues, and family support regardless of student demographics.

6.2. Recommendation

Based on the findings of this research, several recommendations can be made for educational institutions, policymakers, and future research:

For Educational Institutions:

1. Implement automated early warning systems that monitor attendance patterns and academic performance to identify at-risk students before they show obvious signs of disengagement.
2. Develop targeted intervention strategies that address the specific factors most strongly associated with dropout risk, particularly attendance issues and academic struggles.
3. Create support programs specifically designed for students whose parents have lower educational attainment, such as additional academic guidance and mentorship opportunities.

4. Strengthen family engagement initiatives that recognize the important role of parental involvement in student success and provide resources to help parents support their children's education.
5. Use data-driven approaches to evaluate the effectiveness of intervention programs and continuously refine strategies based on outcomes.

For Policymakers:

1. Allocate resources for dropout prevention based on evidence-based risk factors rather than assumptions or traditions. This research provides a data-driven foundation for determining which factors most merit attention and funding.
2. Support professional development for educators in identifying and addressing early warning signs of dropout risk, particularly in the areas identified as most predictive in this research.
3. Develop policies that reduce barriers to parental involvement in education, particularly for parents with lower educational attainment or demanding work schedules.
4. Fund research on effective interventions for the specific risk factors identified in this study, particularly around attendance improvement and academic support programs.

For Future Research:

1. Investigate the effectiveness of interventions targeted at the specific risk factors identified in this research to determine which approaches most effectively reduce dropout rates.
2. Explore the relationships between different risk factors to better understand how they interact and potentially compound dropout risk. For example, how attendance, academic performance, and parental factors might interact.
3. Conduct longitudinal studies that track students over time to better understand the causal relationships between risk factors and dropout decisions.
4. Refine predictive models with additional data including more detailed information on students' social-emotional factors, school climate measures, and community characteristics.

REFERENCES

- Adejo, O., & Connolly, T. (2018). Predicting student academic performance using multi-model heterogeneous ensemble approach. *Journal of Applied Research in Higher Education*, 10(1), 61-75.
- Alexander, K. L., Entwisle, D. R., & Kabbani, N. S. (2001). The dropout process in life course perspective: Early risk factors at home and school. *Teachers College Record*, 103(5), 760-822.
- Baker, R. S., & Inventado, P. S. (2014). Educational data mining and learning analytics. In J. A. Larusson & B. White (Eds.), *Learning Analytics: From Research to Practice* (pp. 61-75). Springer.
- Balfanz, R., Herzog, L., & Mac Iver, D. J. (2007). Preventing student disengagement and keeping students on the graduation path in urban middle-grades schools: Early identification and effective interventions. *Educational Psychologist*, 42(4), 223-235.
- Bowers, A. J., Sprott, R., & Taff, S. A. (2013). Do we know who will drop out?: A review of the predictors of dropping out of high school: Precision, sensitivity, and specificity. *The High School Journal*, 96(2), 77-100.
- Finn, J. D. (1989). Withdrawing from school. *Review of Educational Research*, 59(2), 117-142.
- Márquez-Vera, C., Cano, A., Romero, C., & Ventura, S. (2016). Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data. *Applied Intelligence*, 44(3), 699-712.
- Rumberger, R. W. (1983). Dropping out of high school: The influence of race, sex, and family background. *American Educational Research Journal*, 20(2), 199-220.
- Rumberger, R. W., & Lim, S. A. (2008). Why students drop out of school: A review of 25 years of research. *California Dropout Research Project Report #15*.
- Sara, N. B., Halland, R., Igel, C., & Alstrup, S. (2015). High-school dropout prediction using machine learning: A Danish large-scale study. In *Proceedings of the 7th European Symposium on Computational Intelligence and Mathematics* (pp. 319-325).

Thammasiri, D., Delen, D., Meesad, P., & Kasap, N. (2014). A critical assessment of imbalanced class distribution problem: The case of predicting freshmen student attrition. *Expert Systems with Applications*, 41(2), 321-330.

Xing, W., Chen, X., Stein, J., & Marcinkowski, M. (2019). Temporal predicting of dropouts in MOOCs: Reaching the low hanging fruit through stacking generalization. *Computers in Human Behavior*, 96, 29-39.