# Student Dropout Prediction System – Project Report

**Project Title:** Student Dropout Prediction System

**Objective:**
The primary objective of this project is to predict whether a student is likely to drop out based on historical academic, personal, and demographic data. This system aims to help educational institutions identify at-risk students and take timely intervention measures, thereby improving student retention and performance.

**Introduction:**
Student dropout is a significant challenge for educational institutions. Understanding the factors that contribute to dropout is essential for improving student success and institutional efficiency. The Student Dropout Prediction System uses **Machine Learning techniques** to analyze patterns in student data and predict potential dropouts. This allows institutions to focus resources on students who need support and take preventive actions.

**Dataset Description:**
The dataset used in this project consists of multiple features, including:

- **Academic performance:** Grades, test scores, attendance percentage.

- **Demographic information:** Age, gender, socioeconomic status, family background.

- **Behavioral data:** Participation in extracurricular activities, engagement in online learning platforms.

- **Previous academic history:** Previous year performance, course repetition.

The dataset contains both **numerical and categorical features**, with a mix of complete and missing data entries, reflecting real-world conditions.

**Data Preprocessing:**
Data preprocessing is a critical step in ensuring the model performs accurately. The following steps were undertaken:

1. **Handling Missing Values:** Missing numeric values were filled using the mean of the column, while missing categorical values were filled using the mode.

2. **Removing Duplicates:** Duplicate records were identified and removed to ensure data integrity.

3. **Feature Selection:** Features with low relevance to the target variable (dropout) were removed based on correlation analysis and domain knowledge.

4. **Encoding Categorical Variables:** Categorical data such as gender, parental education, and activity participation were converted into numerical format using **One-Hot Encoding**, making them suitable for ML algorithms.

5. **Data Normalization:** Numerical features were normalized to ensure uniformity and improve model training efficiency.

**Exploratory Data Analysis (EDA):**
EDA was conducted to understand the dataset and identify key patterns:

- **Distribution Analysis:** Histograms and boxplots were created to visualize numeric feature distributions, such as grades and attendance.

- **Correlation Analysis:** A correlation heatmap was used to identify features strongly associated with dropout likelihood.

- **Categorical Analysis:** Count plots and bar charts visualized categorical feature distributions, such as gender or parental education level.

The insights gained from EDA helped refine feature selection and model development.

**Model Selection and Training:**
Multiple machine learning algorithms were evaluated to predict student dropout, including:

1. **Logistic Regression:** Used for binary classification of dropout (Yes/No).

2. **Random Forest Classifier:** Ensemble method providing higher accuracy and robustness against overfitting.

3. **Decision Trees:** Simple interpretable model for understanding feature importance.

The dataset was split into **training and testing sets** (80%-20%). Models were trained using the training set, and hyperparameters were optimized using cross-validation techniques to improve performance.

**Model Evaluation:**
The models were evaluated using standard metrics:

- **Accuracy:** Percentage of correctly predicted dropouts.

- **Precision and Recall:** Evaluated the ability to correctly identify students at risk.

- **F1 Score:** Harmonic mean of precision and recall to balance false positives and false negatives.

- **Confusion Matrix:** Visual representation of prediction results for true positives, true negatives, false positives, and false negatives.

The Random Forest Classifier achieved the **highest accuracy**, around 90%, and demonstrated better recall in identifying students likely to drop out.

**Results and Interpretation:**
The model successfully identified key factors contributing to student dropout, including low attendance, poor academic performance, lack of extracurricular engagement, and socioeconomic challenges. The predictive system can flag at-risk students, enabling institutions to provide counseling, tutoring, or other support programs proactively.

**Tools and Technologies Used:**

- **Programming Language:** Python

- **Libraries:** Pandas, NumPy, Matplotlib, Seaborn, Scikit-learn

- **Environment:** Jupyter Notebook

- **Visualization Tools:** Matplotlib and Seaborn for data analysis and pattern recognition

**Conclusion:**
The Student Dropout Prediction System demonstrates how Machine Learning can be effectively applied in education to reduce student dropout rates. By analyzing historical data and predicting at-risk students, the system allows institutions to take timely intervention measures. Through this project, I learned the **end-to-end workflow of a Machine Learning project**, from data preprocessing and exploratory analysis to model building, evaluation, and interpretation. This experience strengthened my skills in Python programming, data handling, visualization, and practical application of Machine Learning algorithms.