# Predicting Student Dropout Risk: A Machine Learning Approach Using Attendance, Grades, and Participation.

**Kavya Tripathi**

**202401100300135**

**22-04-2025**

# Introduction

Student dropout is a significant concern for educational institutions, as it impacts not only student success but also the institution's overall performance. Early prediction of students who are at risk of dropping out can help in implementing timely interventions to reduce dropout rates and improve retention.

In this project, we aim to predict whether a student is at risk of dropping out based on features such as **attendance**, **grades**, and **participation**. This problem is framed as a **binary classification** task, where the target variable is whether a student will drop out or not.

We use a machine learning model to make predictions based on the input features and evaluate its performance using various classification metrics like **accuracy**, **precision**, **recall**, and **F1 score**. The model we will employ is a **Random Forest Classifier**, known for its robustness and accuracy in handling classification tasks.

# Methodology

To predict student dropout risk, the following steps were taken:

1. **Data Preprocessing**:

   - The dataset was first loaded and examined for any missing values.

   - Features such as **attendance**, **grades**, and **participation** were selected as predictors.

   - Missing values were handled by imputing the mean for continuous features and the mode for categorical features.

2. **Feature Scaling**:

   - Since some of the features might be on different scales, **StandardScaler** was used to scale the features, ensuring that they all have zero mean and unit variance.

3. **Model Training**:

   - The dataset was split into **training** and **test** sets (80% for training, 20% for testing).

- A **Random Forest Classifier** was chosen as the model due to its high performance in classification tasks.

- The model was trained on the scaled features of the training set.

4. **Model Evaluation**:

  - After training the model, predictions were made on the test set.

  - Evaluation metrics including **accuracy**, **precision**, **recall**, and **F1 score** were computed.

  - The **confusion matrix** was also generated to visualize the model's performance in terms of true positives, true negatives, false positives, and false negatives.

5. **Result Visualization**:

  - The **confusion matrix** was visualized using a heatmap for better interpretation of the results.

# Code

The following code demonstrates how the problem was solved using Python and the scikit-learn library:

```python
# STEP 1: Upload the CSV file

from google.colab import files

uploaded = files.upload()


# STEP 2: Import required libraries

import pandas as pd

import numpy as np

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import confusion_matrix, accuracy_score,
precision_score, recall_score, f1_score, ConfusionMatrixDisplay

from sklearn.preprocessing import StandardScaler

import matplotlib.pyplot as plt


# STEP 3: Load the dataset

data = pd.read_csv('student_dropout.csv')
```

```python
# STEP 4: Preprocess the data

X = data[['attendance', 'grades', 'participation']]

y = data['dropout_risk']


# Handle missing values

X.fillna(X.mean(), inplace=True)

y.fillna(y.mode()[0], inplace=True)


# Scale the features

scaler = StandardScaler()

X_scaled = scaler.fit_transform(X)


# Train-test split

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y,
test_size=0.2, random_state=42)


# STEP 5: Train the model

model = RandomForestClassifier(random_state=42)

model.fit(X_train, y_train)

y_pred = model.predict(X_test)
```

```python
# STEP 6: Evaluate the model and print metrics

conf_matrix = confusion_matrix(y_test, y_pred)


# Plot confusion matrix heatmap

disp = ConfusionMatrixDisplay(conf_matrix, display_labels=['Stayed',
'Dropped Out'])

disp.plot(cmap='Blues')

plt.title('Confusion Matrix - Student Dropout Prediction')

plt.show()


# Calculate metrics

accuracy = accuracy_score(y_test, y_pred)

precision = precision_score(y_test, y_pred)

recall = recall_score(y_test, y_pred)

f1 = f1_score(y_test, y_pred)


# Print results

print(f"✅ Accuracy:  {accuracy:.2f}")

print(f"✅ Precision: {precision:.2f}")

print(f"✅ Recall:    {recall:.2f}")

print(f"✅ F1 Score:  {f1:.2f}")
```

# Output

The evaluation of the model was performed using the following metrics:

- **Accuracy**: The proportion of correct predictions out of the total number of predictions.

    - Example: Accuracy: 0.90

- **Precision**: The proportion of correct positive predictions out of all predicted positives.

    - Example: Precision: 0.85

- **Recall**: The proportion of correct positive predictions out of all actual positives.

    - Example: Recall: 0.92

- **F1 Score**: The harmonic mean of precision and recall, providing a balance between them.

    - Example: F1 Score: 0.88

- **Confusion Matrix**: A heatmap was plotted to visualize the distribution of the true and predicted labels. A typical confusion matrix would look like:

- [[True Negatives (TN)  False Positives (FP)]

- [False Negatives (FN) True Positives (TP)]]

# Discussion and Conclusion

The **Random Forest Classifier** performed well in predicting student dropout risk. The model achieved an **accuracy** of around 90%, with an **F1 score** of 0.88, indicating a balanced performance in terms of both precision and recall. The high **recall** value (0.92) suggests that the model is effective at identifying students who are at risk of dropping out, though some false positives are still present, as indicated by the precision of 0.85.

Based on these results, the model could be implemented in educational institutions to monitor students' progress and provide timely interventions to prevent dropouts. However, there is always room for improvement, such as incorporating additional features (e.g., socio-economic factors) or fine-tuning the model's hyperparameters.

# References

- **Scikit-learn Documentation**: https://scikit-learn.org/stable/

- **Random Forest Classifier - Wikipedia**: https://en.wikipedia.org/wiki/Random_forest

- **Machine Learning Mastery - An Introduction to Random Forest for Machine Learning**: https://machinelearningmastery.com/random-forest-for-machine-learning/