**Course Name: Artificial Intelligence**
**Course Code: CSE422**
**<u>Lab Project Report</u>**


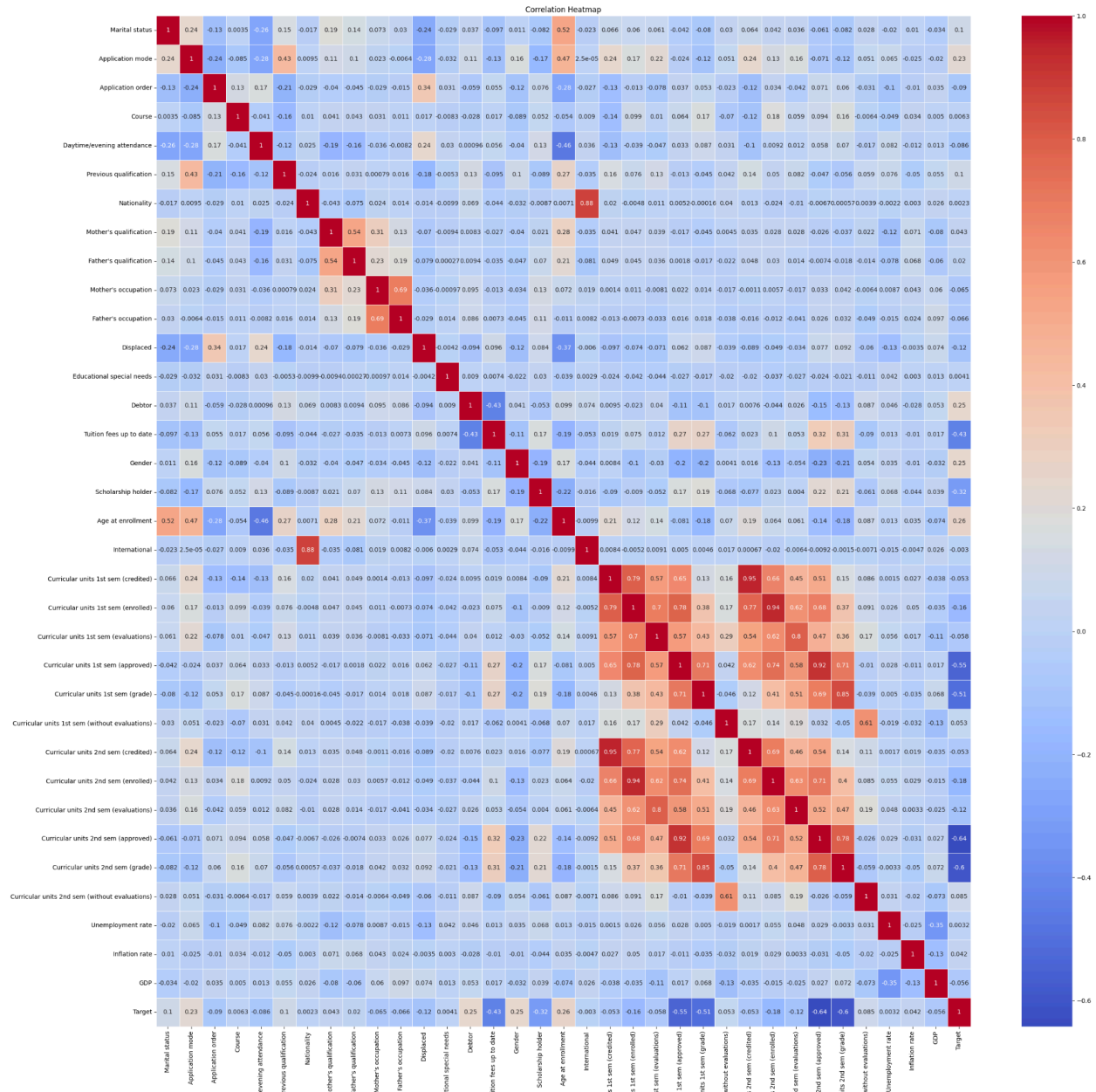**Project Topic: Student Dropout Prediction**

# Table of contents

1. **Introduction:**

   The project aims to predict whether a student will drop out or graduate with the help of a number of features. This could help educational institutions to get an estimation of which students are likely to drop out and provide them with additional support to minimize the chances of them dropping out. It could boost the graduation rate of an institution by providing data-driven insights.
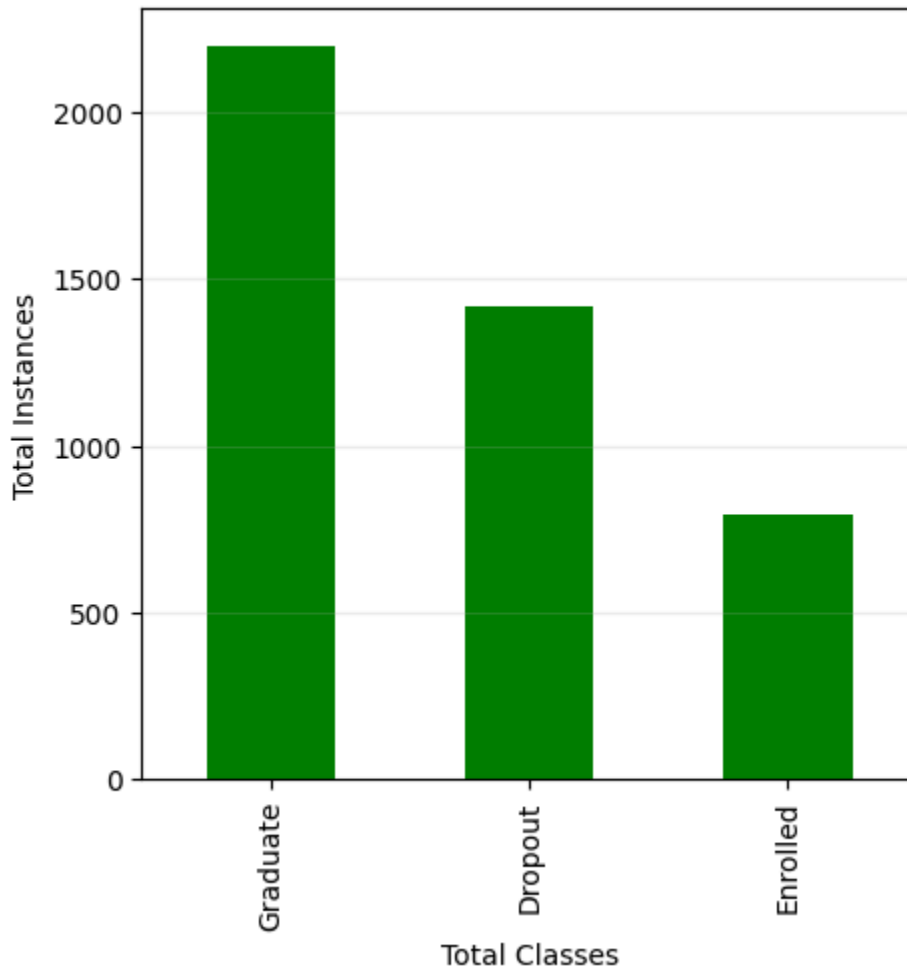
2. **Dataset Description:**

   - **Source:** [Predict students' dropout and academic success](#)

   - **Number of features:** 35

   - **Problem type:** Classification, as the output will be a Categorical Value, which is 0 for graduate and 1 for dropout.

   - **Number of Columns:** 4424

   - **Feature Types**: Quantitative and Categorical both

● Correlation of all features(Using Heatmap):



Correlation Heatmap

● Imbalance dataset: Here, for the output feature, all the unique values do not have an equal number of instances.



3. **Dataset Pre-processing**

- **Modified the dataset:** We had to modify the data a bit in order to make the dataset unprocessed; this included nullifying some of the values. This is the link to the updated dataset.

- **Null-values:** There are some columns which included null values; We dropped those columns

- **Column name fix:** The "Nacionality" column name was misspelled, so fixed the spelling.

- **Mapping Categorical columns with Numerical Values:** The Columns "Country Names" & "Target" had categorical values so we replaced them with some arbitrary numerical values.

- **Dropping Unnecessary Columns:** There were some columns which we felt were a bit redundant and did not add much to the training data; therefore, we decided to drop them.

- **Dropping "Enrolled" Rows from Target Column:** The number of instances between all the three features were too imbalanced and the "Enrolled" instance was not giving us any conclusive decision that the student will drop out or not; So, we dropped the rows with "Enrolled" instances.

- **Dropping Correlated Features:** In the heat map, some features were highly correlating; so we chose to drop some of them to reduce redundancy.

4. **Feature Scaling**

- We scaled the columns with quantitative values from the processed dataset using the MiniMax Scaler algorithm.

## 5. Dataset Splitting

- 70% of the data were used as Train set and 30% of the data were used for Test set
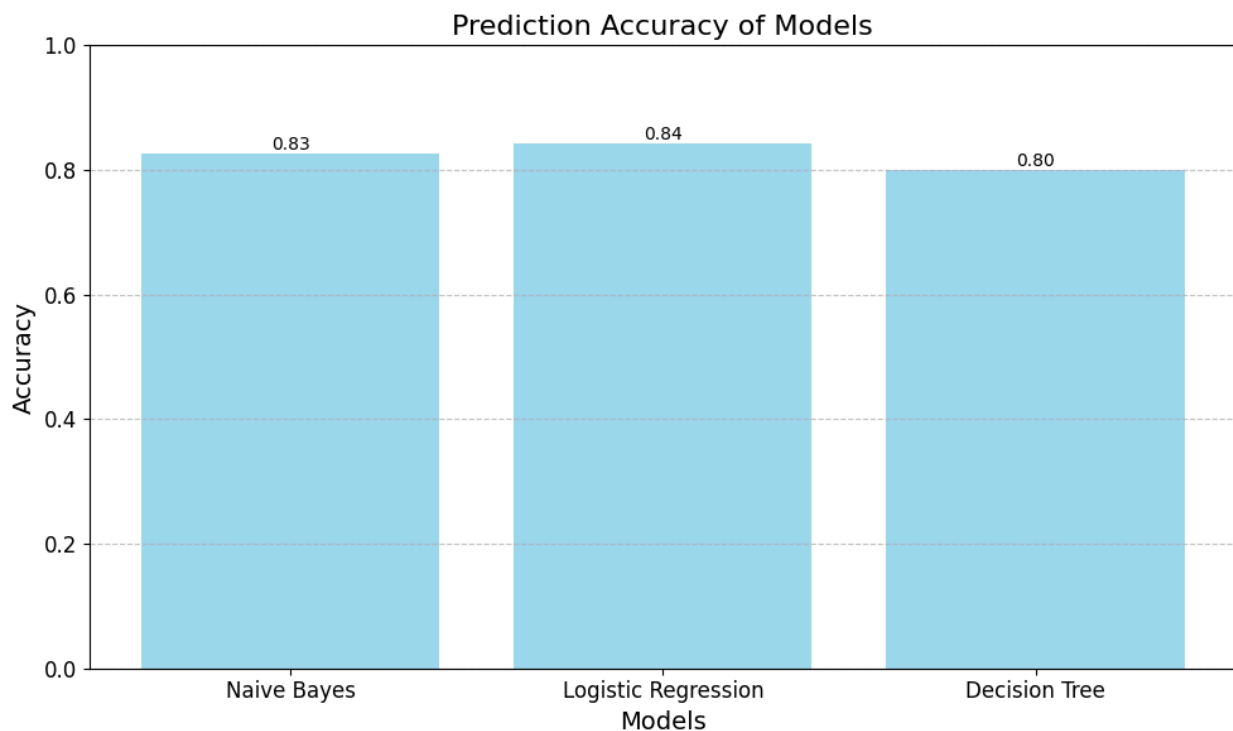
## 6. Model Training and Testing

- **Naive Bayes:** We chose naive bayes for its simplicity and for the fact that it performs well with both large high-dimensional datasets as well as small datasets. It assumes conditional independence between features making the calculation simpler, therefore being able to generate responses faster.

- **Logistic Regression:** Logistic Regression uses sigmoid function to map the predicted values to [0, 1] of binary classification problems. It works great with binary classification problems and it is quite easy to implement and very efficient for training. Also, it provides good accuracy and performs well for fairly simple datasets and when the dataset is linearly separable; hence, we chose to use this model.

- **Decision Tree:** We used decision tree for its insensitivity to outliers as it splits data based on feature thresholds rather than depending on distance calculations which gives a better probability to find the optimal class of the target.
  We used "class_weight = balanced" as a parameter so that it doesn't give a biased class because of target feature imbalance.

## 7. Comparison Analysis

- **Bar Chart showcasing prediction accuracy of models**

  We can notice that Logistic Regression gives us the most accurate prediction with 84% correct prediction, while Decision Tree gives us the worst with 80%; Naive Bayes is in the middle with 83%.



- **Precision and recall comparison**

| Models | Precision | Recall |
|---|---|---|
| **Naive Bayes** | 0.83 | 0.83 |
| **Logistic Regression** | 0.84 | 0.84 |
| **Decision Tree** | 0.80 | 0.80 |

- **Confusion Matrix of each model**

    **Naive Bayes:**

    ```
    [[488  64]
     [ 93 264]]
    ```

    **Logistic Regression:**

    ```
    [[510  42]
     [101 256]]
    ```

    **Decision Tree:**

    ```
    [[459  93]
     [ 88 269]]
    ```

8. **Conclusion**

   This project was designed to predict the chances of a student dropping out. By training the model, we tried finding an appropriate answer to the query. Although the models were inaccurate for about ~20% of the data, it can be used as a guide to predicting the likelihood of the student dropping out. Educational institutions could use something of this nature and enhance their graduation rate and help the students out greatly.