

Taiyō  
Machine Learning  
Assignment  
Natural Language  
Understanding

**By Remita Austin**

# PROBLEM STATEMENT

Use World Bank Projects dataset.

Option 2. Binary Classifier. Using the status variable build a binary classifier to predict the probability whether a project will be “closed” or “canceled/distressed”.



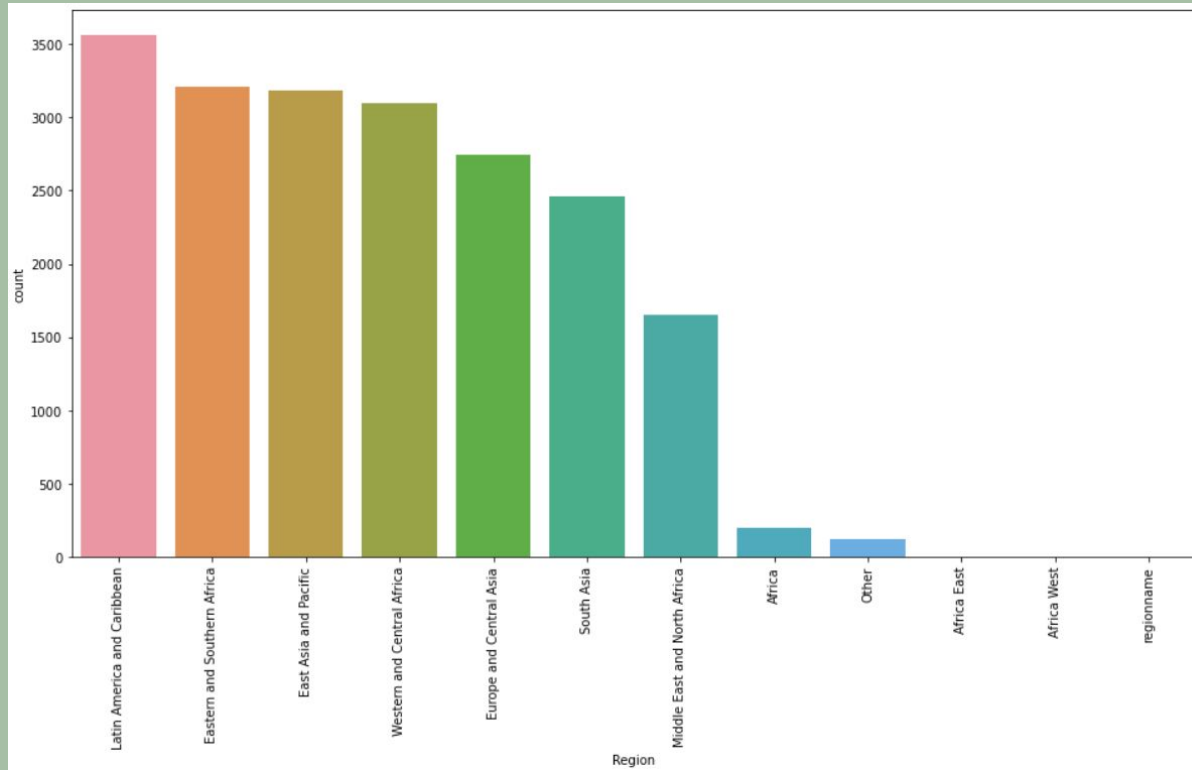
# IMPLEMENTATION

## 1. Data Pre-processing

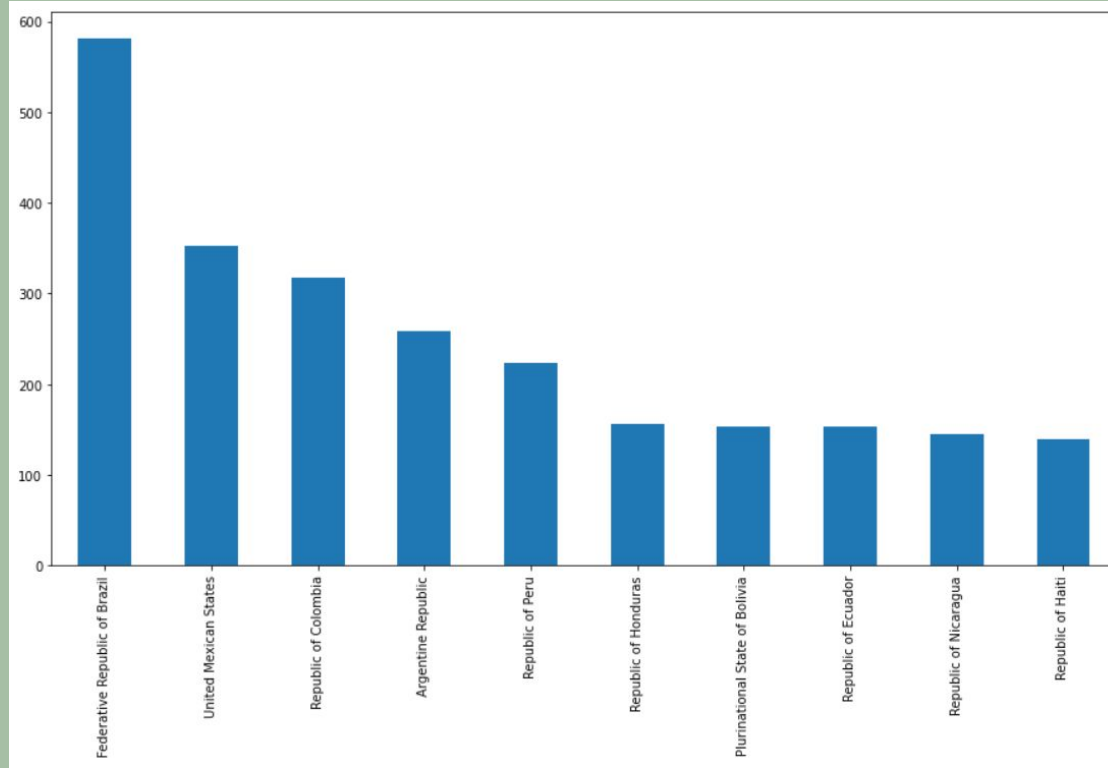
The Projects dataset was imported from World Bank Projects and loaded. We checked if the dataset contained missing/null values. The dataset was cleaned by removing NA values, and dropping irrelevant fields and records with Project Status=Pipeline since this is not required for the given problem statement.

## 2. Exploratory Analysis

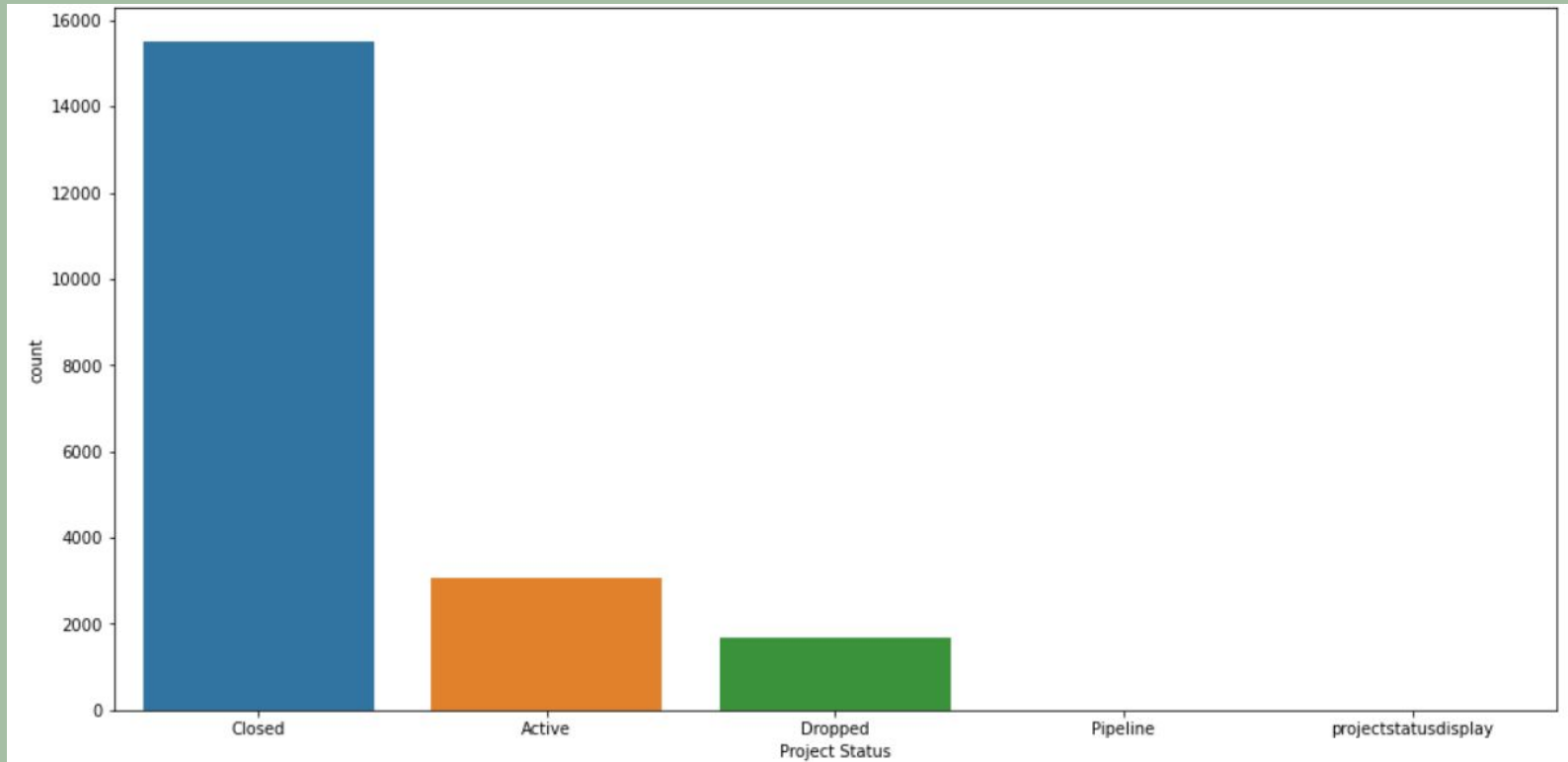
We need to observe and analyze the data to see what we are working with. For this, the matplotlib and seaborn libraries were used and the data was visualized and inferences were observed.



Most of the projects have been implemented in Latin America and Caribbean region.



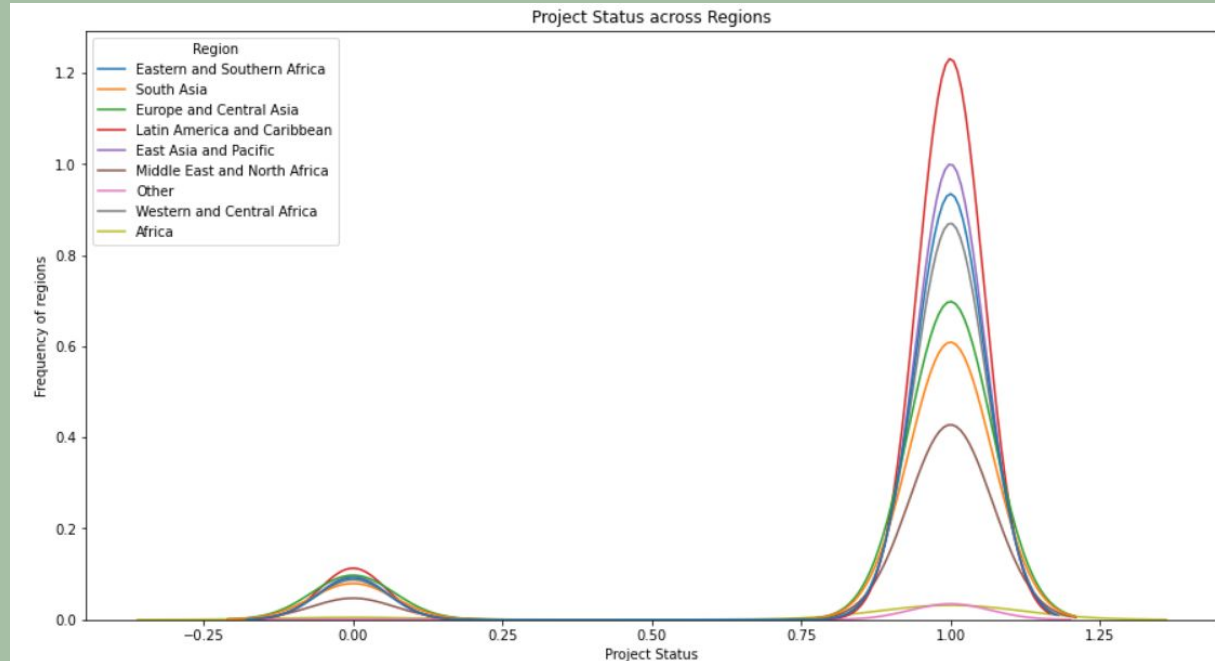
Here the top ten countries of Latin America and Caribbean region with the most number of projects are displayed.



Most of the projects have their status as Closed. Few projects are Active and even fewer projects have been Dropped.

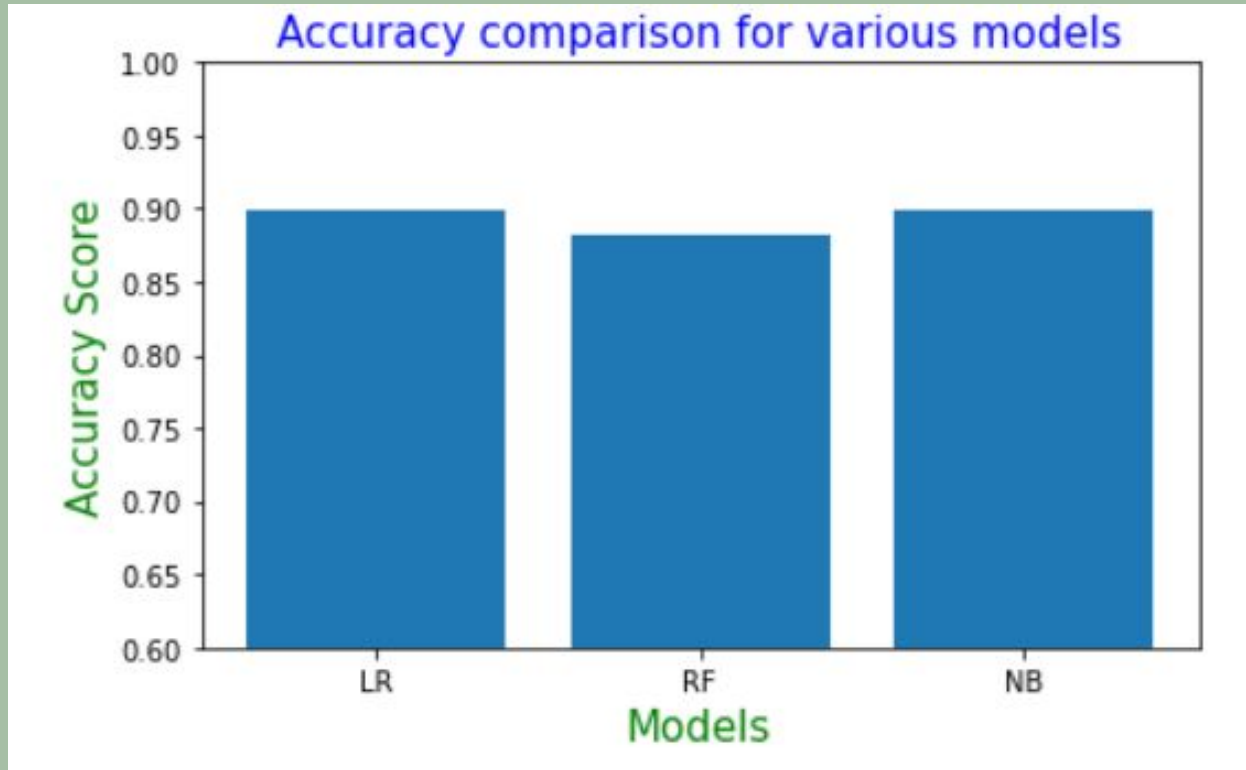
### 3. Split the data wrt Project Status

The dataset was split such that active projects consisted of those whose Project Status=Active and closed/dropped projects consisted of those whose Project Status!=Active. Then the relevant features were selected and dummy variables were given for categorical data for ease of prediction.









From the barplot, it is clear that LR and NB shows highest accuracy score, and RF shows least accuracy score.

# CONCLUSION

The World Bank Projects data was analyzed, and the data was cleaned and modified in order to train the model using Logistic Regression, Random Forest, and Naive Bayes algorithms. We then computed the accuracy scores of these models and visualized them using a barplot. The LR and NB models gave the highest accuracy. Then the project status of sample input data was predicted using the LR model.