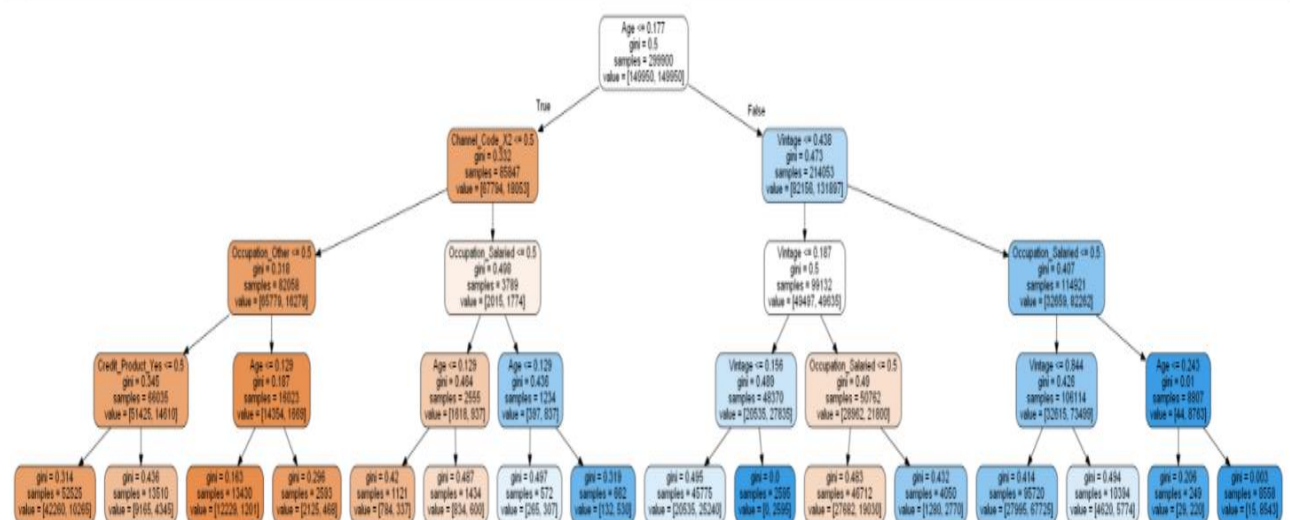**The Steps I have taken for classifying or finding the probabilities of customers who would turn on for credit card purchase from the Happy Bank are:**

1. Firstly, I have viewed the datasets given in the train and test set to work on it
2. I have then checked for missing values in the dataset i,.e, Missing value analysis. In this step I have filled the missing values of categorical columns with mode. There are not many missing values. The percentage of missing values are very less can say as negligible amount. There are null values in only one column in the dataset and it comprises to 11% of the data.
3. After missing values imputation, I started cleaning the data before I proceeded to exploratory data analysis. In data cleaning I found the whether the datatypes are correct to the data or not. Next have gone through outlier analysis using IQR method.
4. In EDA, I have found the categorical and numerical features with their statistics interpretation using describe method and pandas profiling
5. Next, I have plotted some graphs in data visualization to find the insights from the data.
6. Next comes the feature engineering and feature selection out the new data formed from the raw data. Here I have done some normalization and one hot encoded (created dummies) the data.
7. Split the data into train and validation sets for training purpose and evaluating it. Also tried different ways here like balanced the data for some models and used unencoded values for some models according the possibilities present in the ML models.
8. Build the models using various ML algorithms like Logistic, Decision tree, Random Forest, XG Boost, Gradient Boosting, Cat Boost, Stacking Classifier, Microsoft AutoML.
9. Evaluated the models and performed hyperparameter tuning to improve the model performance.
10. Final model was decided based on roc_auc_score of various models.
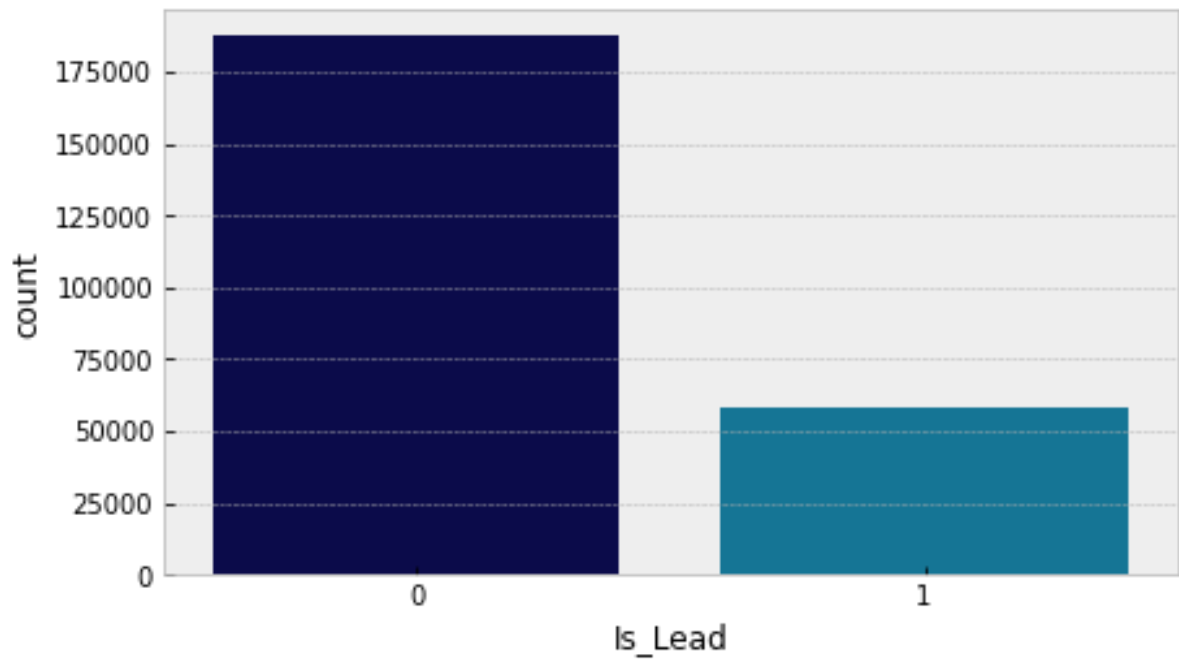11. Made predictions on the test set.
12. Submission file was created.

## RESULTS:

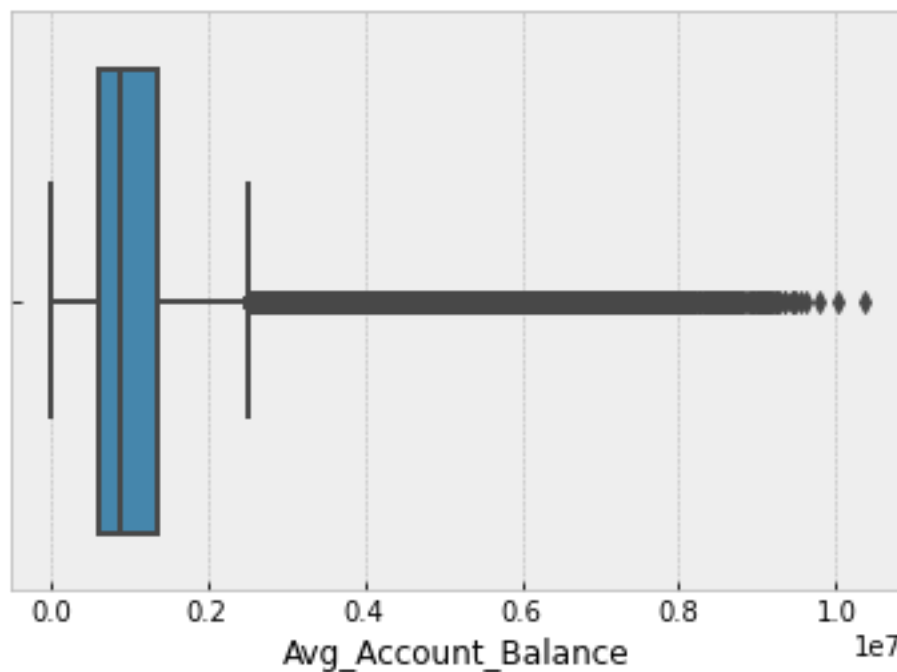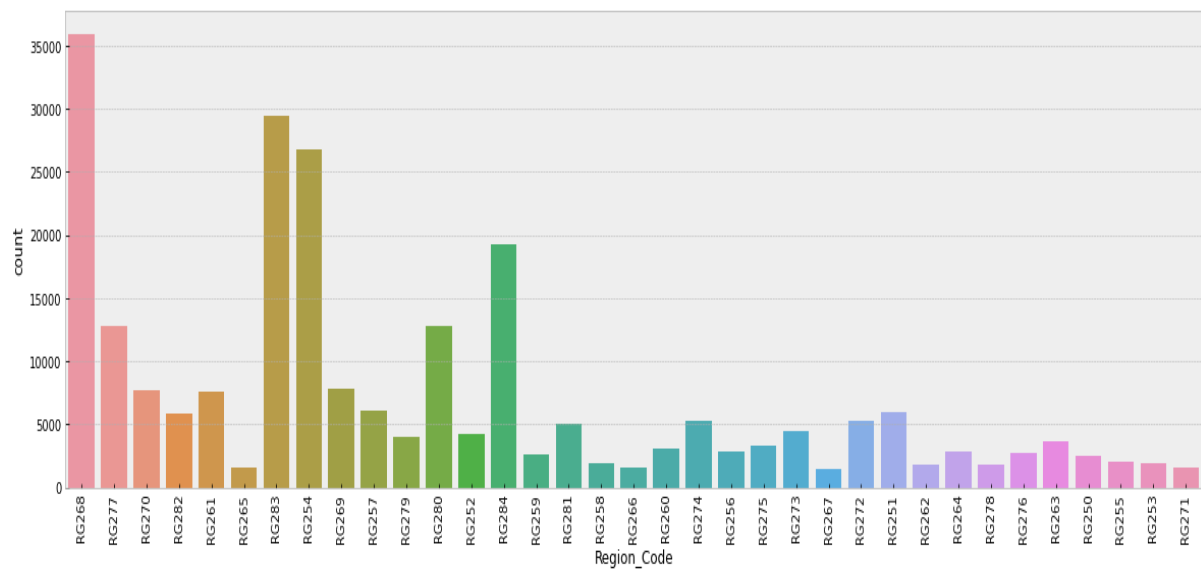| | Models | Roc_auc_score |
|---|---|---|
| 0 | Logistic Regression | 0.732933 |
| 1 | Decision Tree | 0.741120 |
| 2 | Random Forest | 0.725390 |
| 3 | XG Boost | 0.776558 |
| 4 | Gradient Boost | 0.772198 |
| 5 | Light GBM(only Label encoding) | 0.786487 |
| 6 | Light GBM | 0.786934 |
| 7 | Catboost | 0.787897 |
| 8 | AutoML2 | 0.788189 |

## SOME DATA VISUALIZATION FROM THE DATASET:



- This is Decision tree graph obtained from decision tree modelling

- This is the target imbalance present in the dataset which can we can overcome using oversampling , undersampling and some boosting models which can deal with such type of data.



- From this it can be observed that this is right tail distribution where there are many customers with balance in their accounts is greater than average balance. Also can say as there are more high income customers.

- These are the various region codes present in the dataset which are in the encoded format. It can also be observed that there are more customer or high number in RG268 code region.