

Lead Scoring Case Study

By

T Sai Kamal Chand

Sampath

Chandan Singh

Process Summary

- **Data Cleaning and Preparation**
- **EDA**
- **Dummy Variables**
- **Dataset split into Train and Test**
- **Feature Scaling**
- **Model Building**
- **Model Evaluation and Prediction on Test dataset**
- **Precision – Recall**
- **Conclusion and Findings**

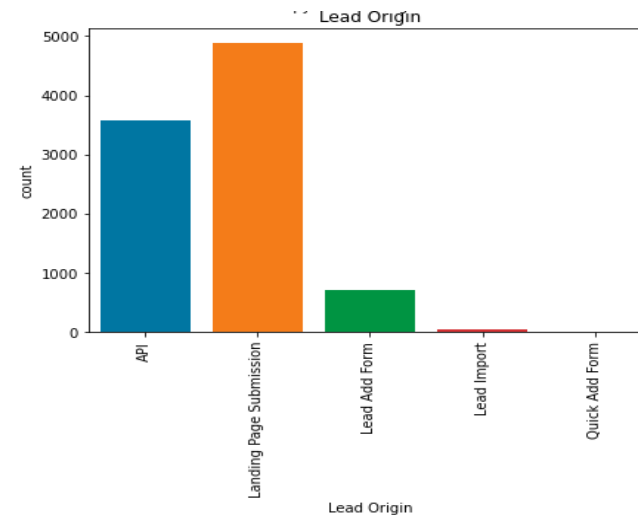
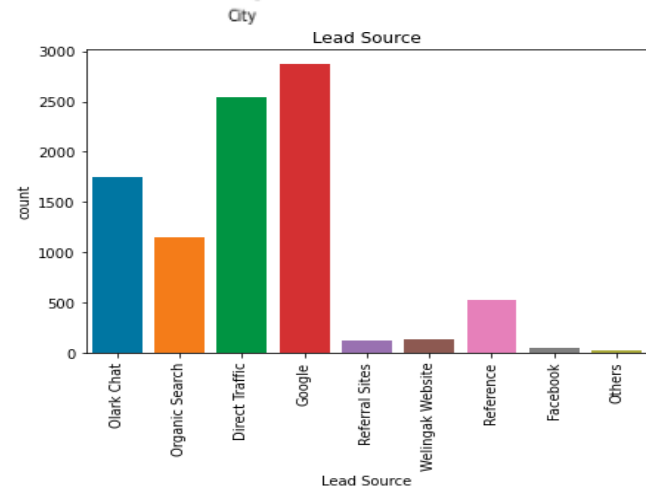
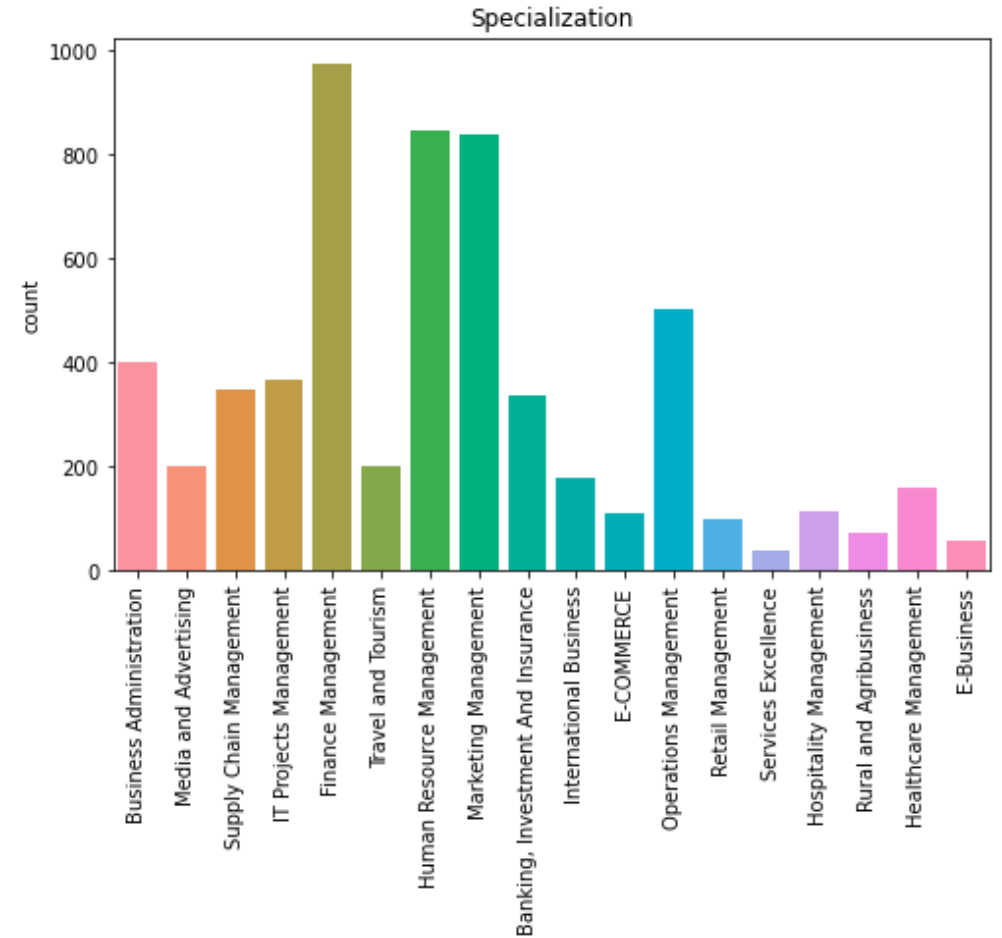
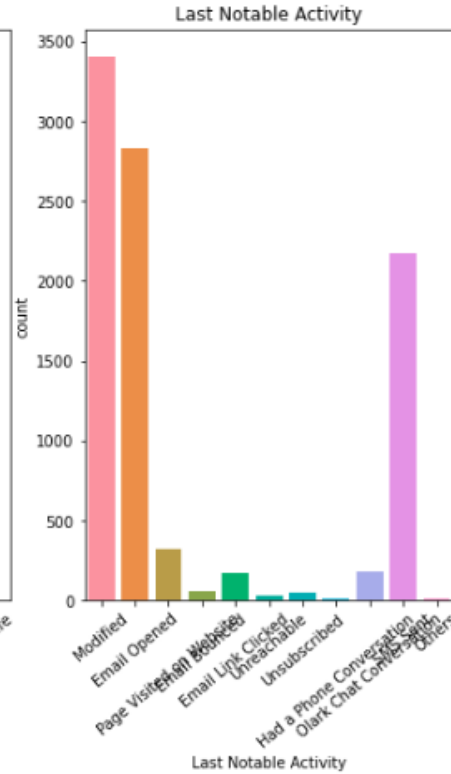
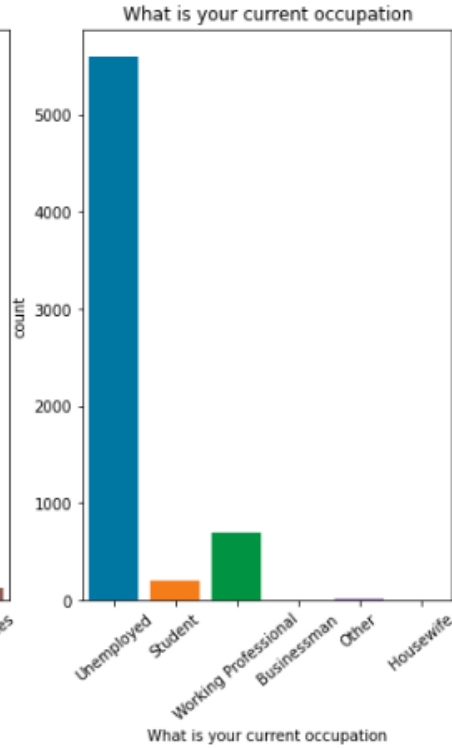
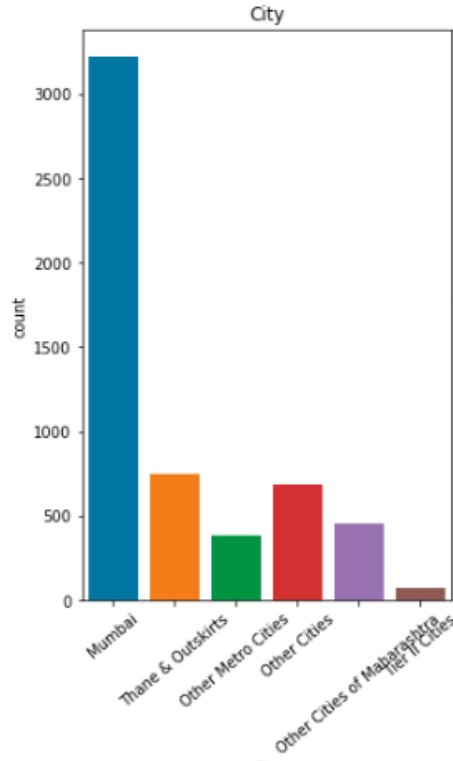
Data Cleaning and Preparation

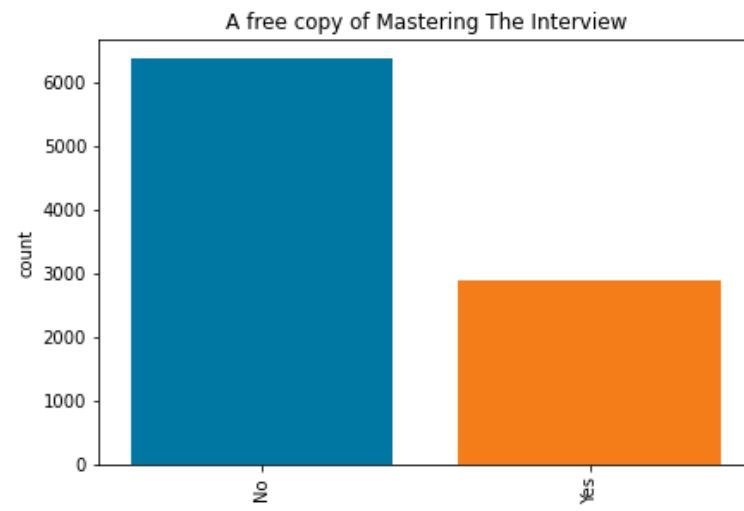
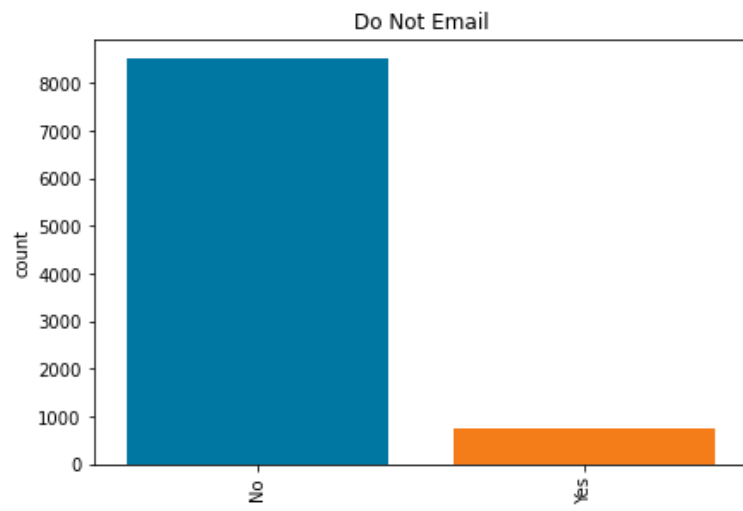
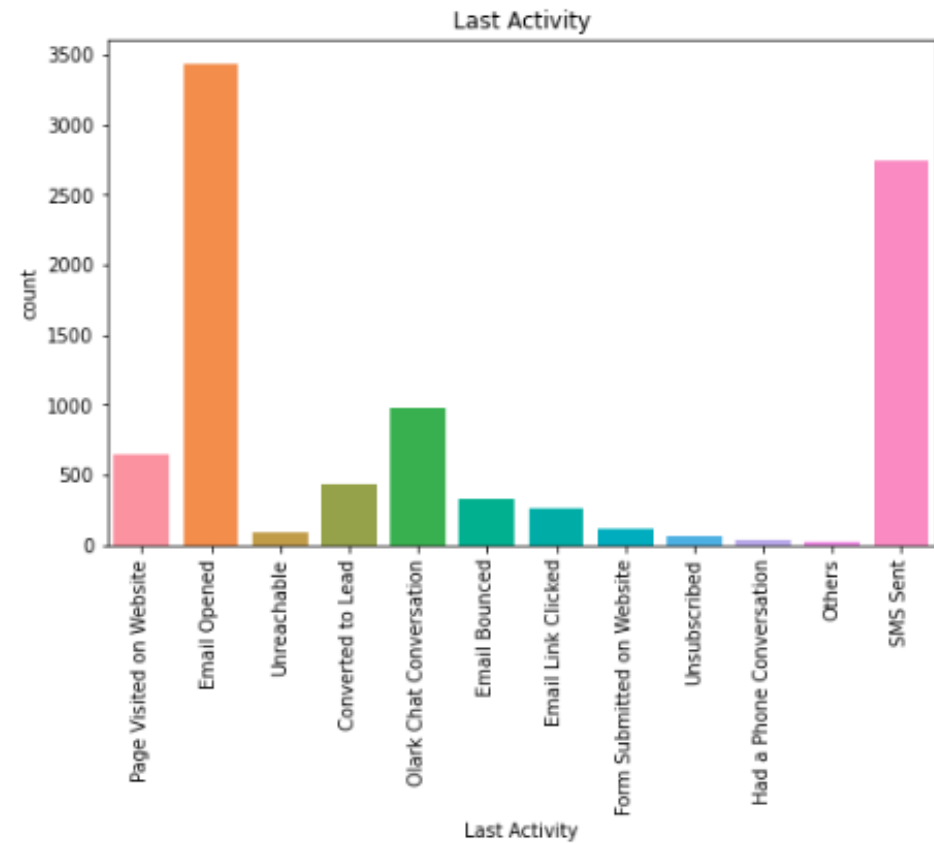
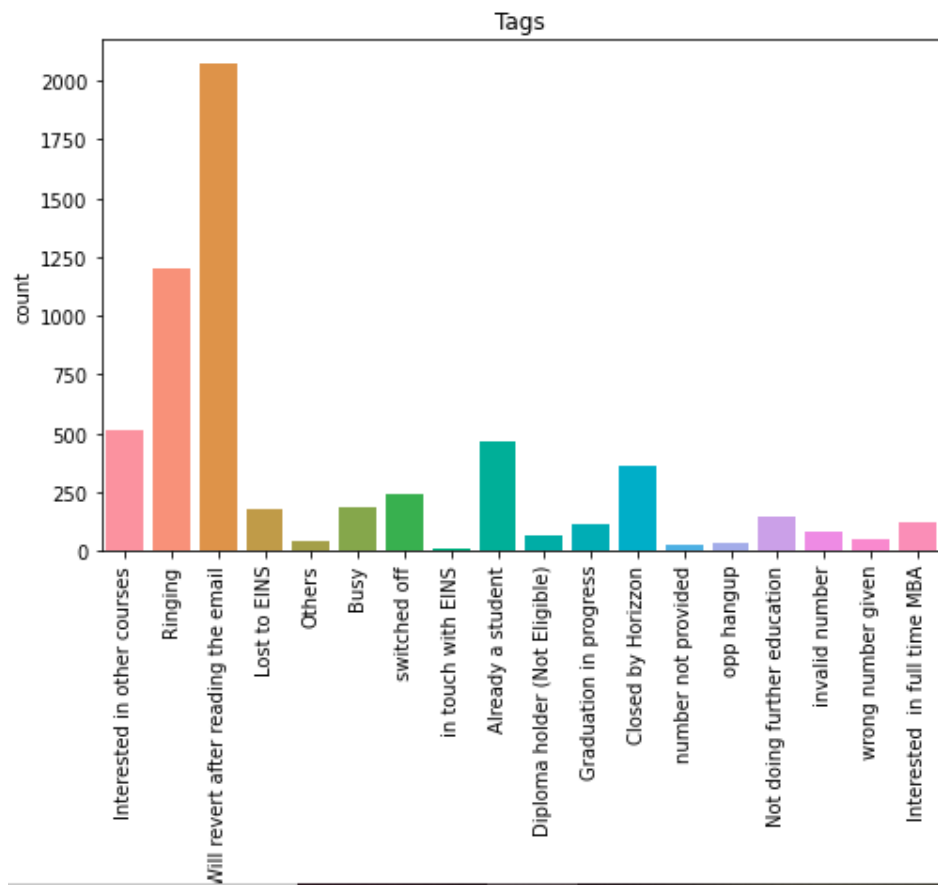
- There are 9240 rows and 37 columns.
- Dropped the columns with NULL values that are more than 40% and some redundant variables with unique values of 1 and other variables where only one value is predominantly seen.
- Changed the values with “Select” to NAs.

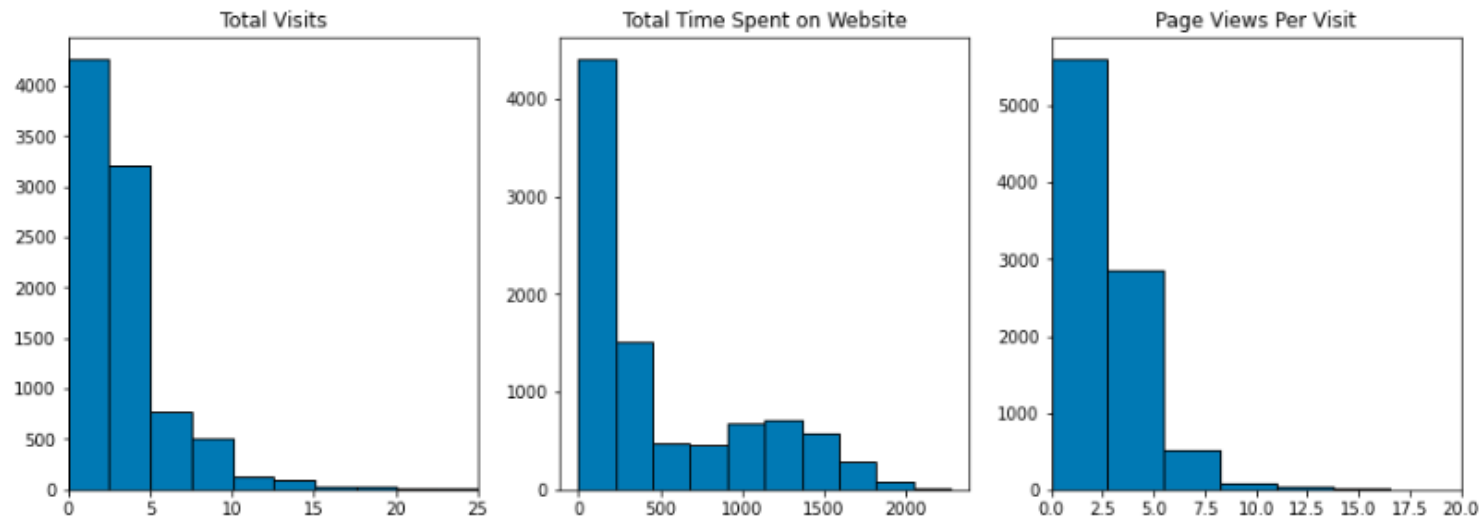
Data Visualization - EDA

- Categorical and Numerical Variables
 - Univariate Analysis
 - Bivariate Analysis

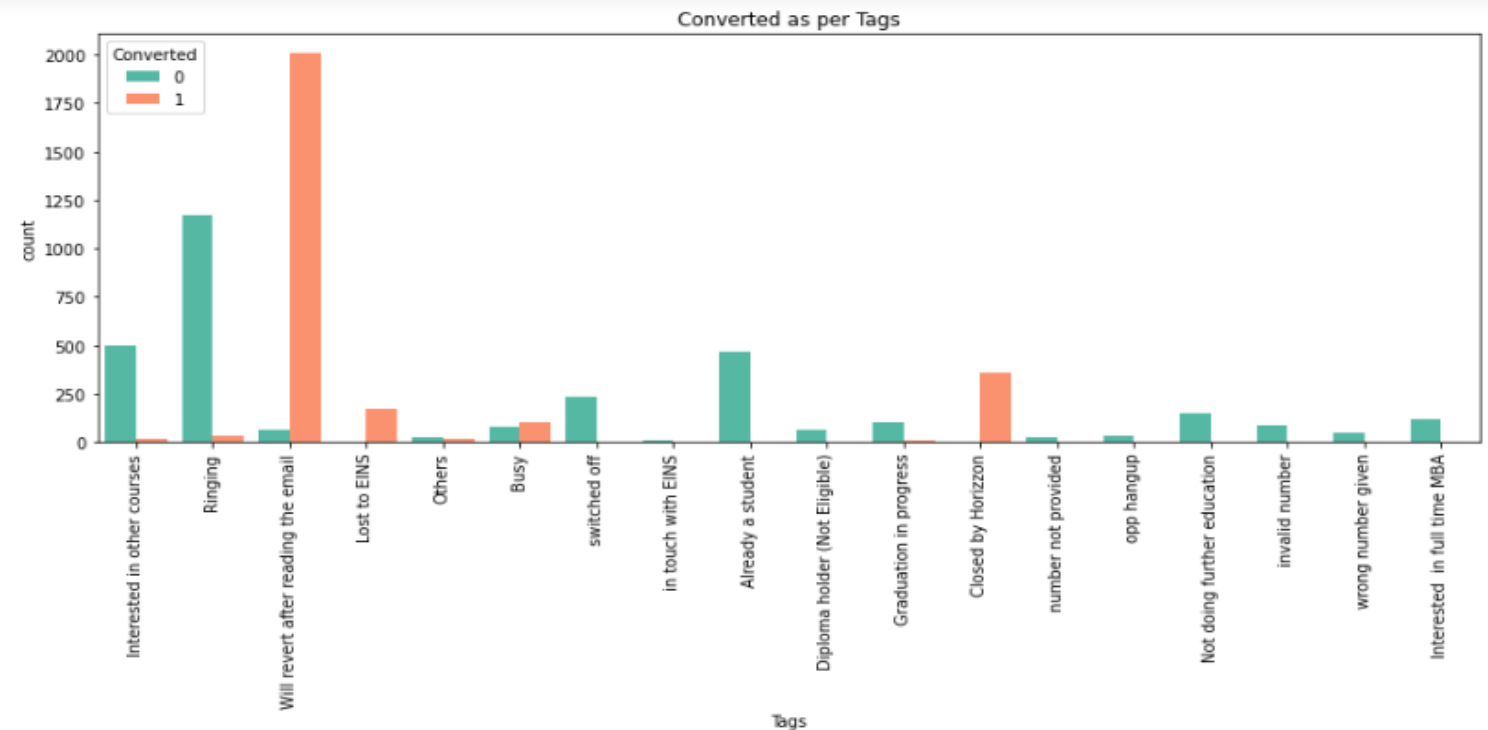
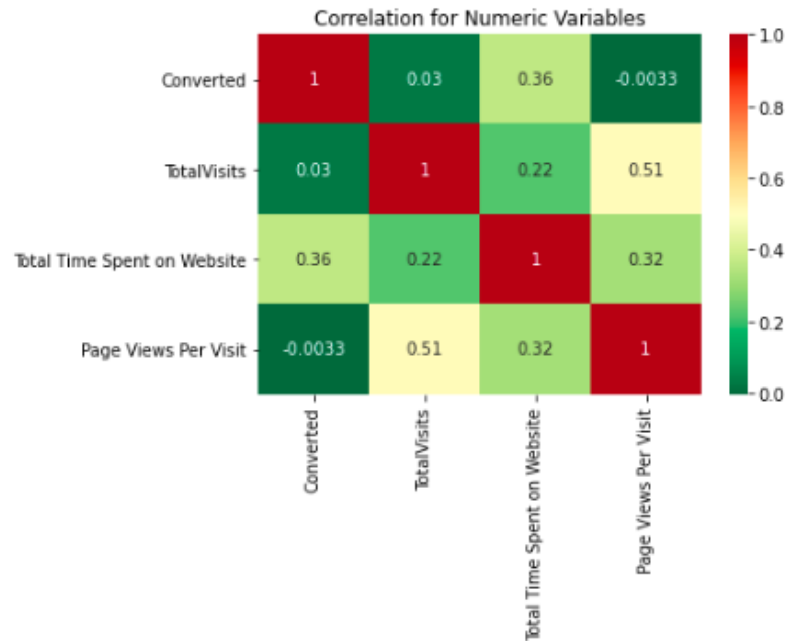
Univariate Analysis

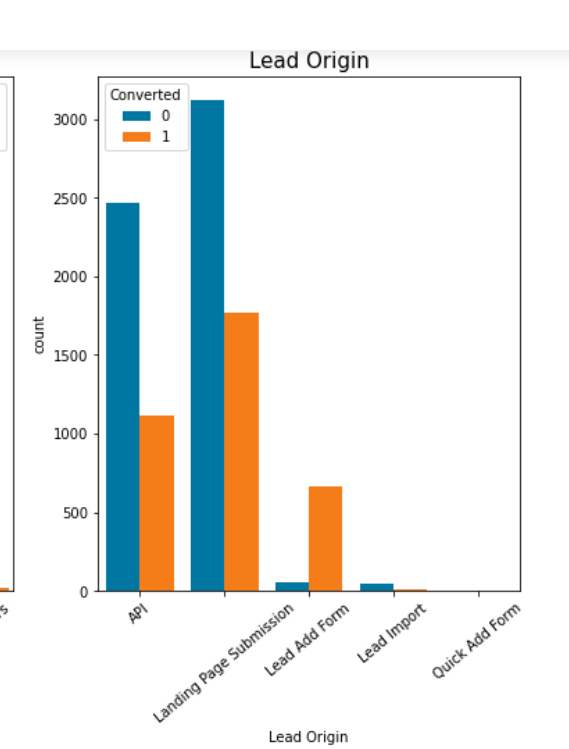
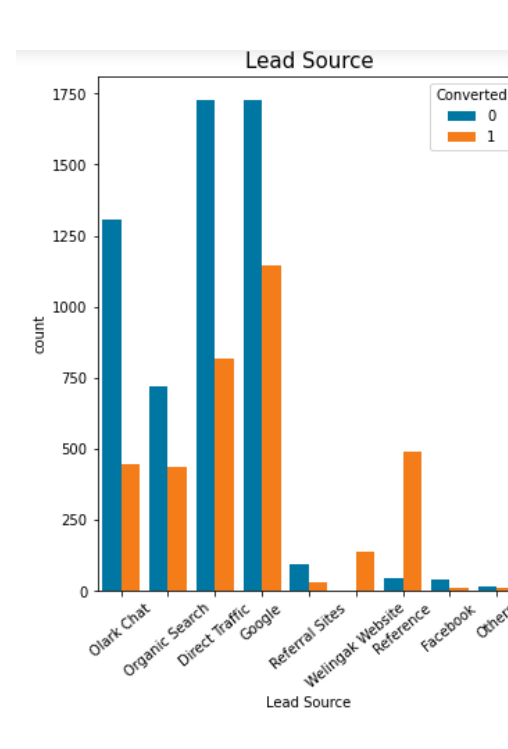
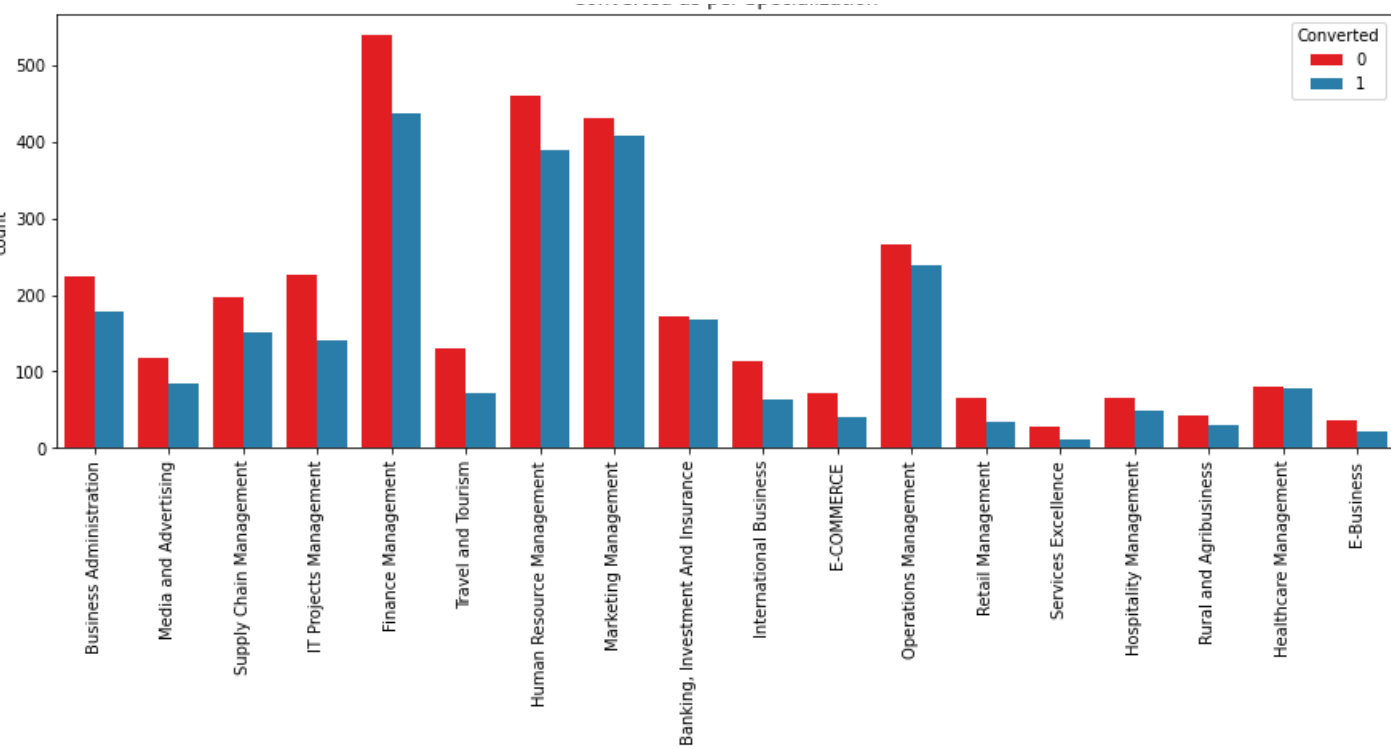
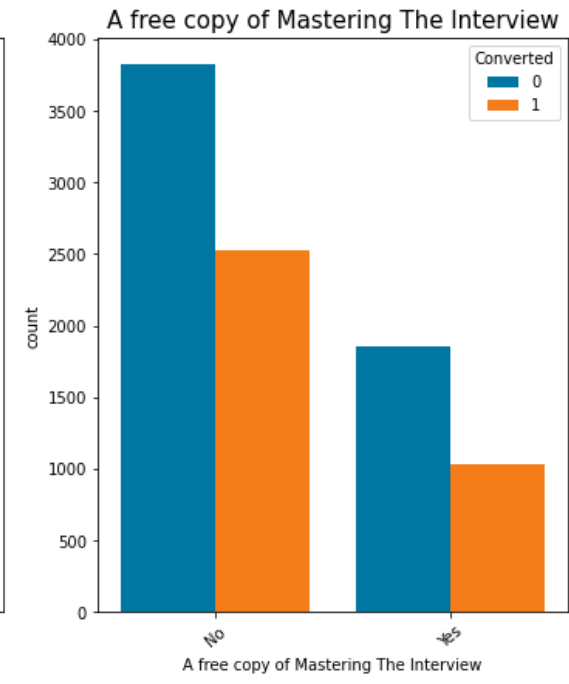
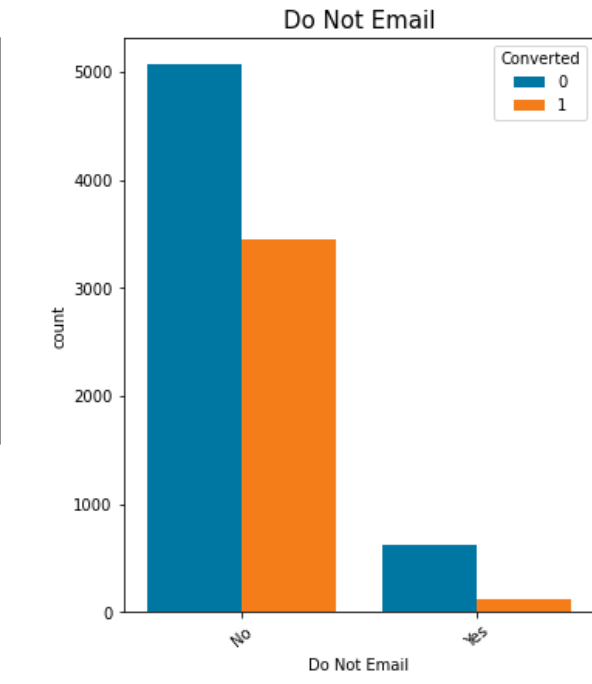
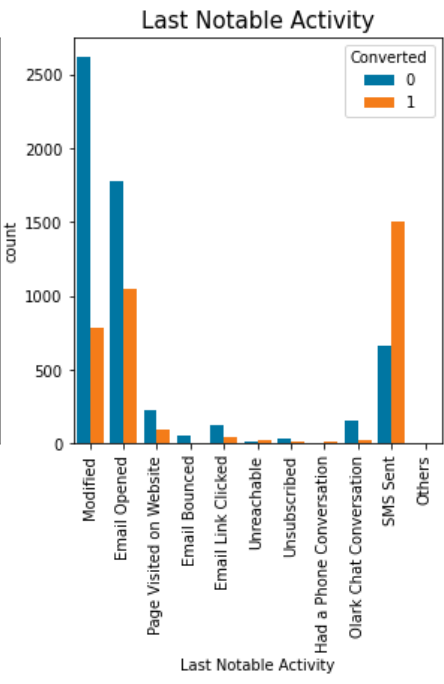
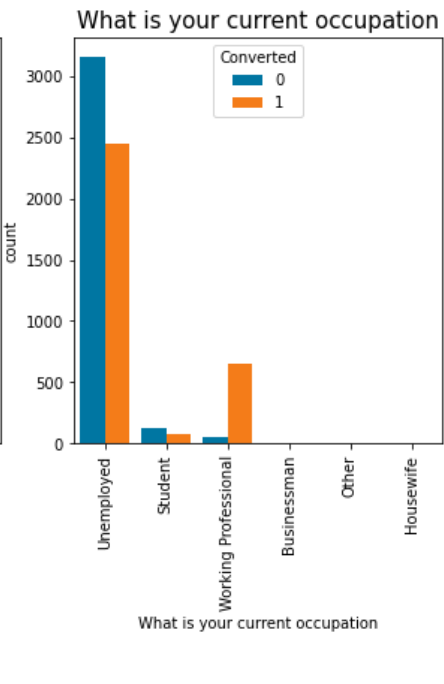
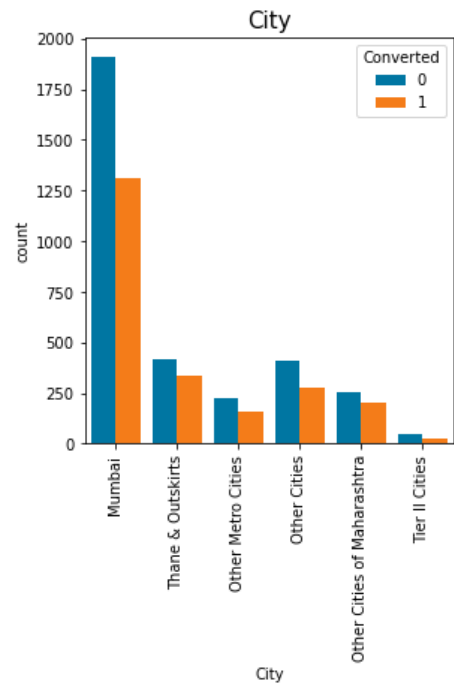






Bivariate Analysis and Correlation Matrix





- **Dummy Variables**

- Created dummy variables for all categorical variables.

- **Dataset split into Train and Test**

- The dataset was split into 70% train and 30% test.
- Missing value treatment was done for X train dataset.

- **Feature Scaling**

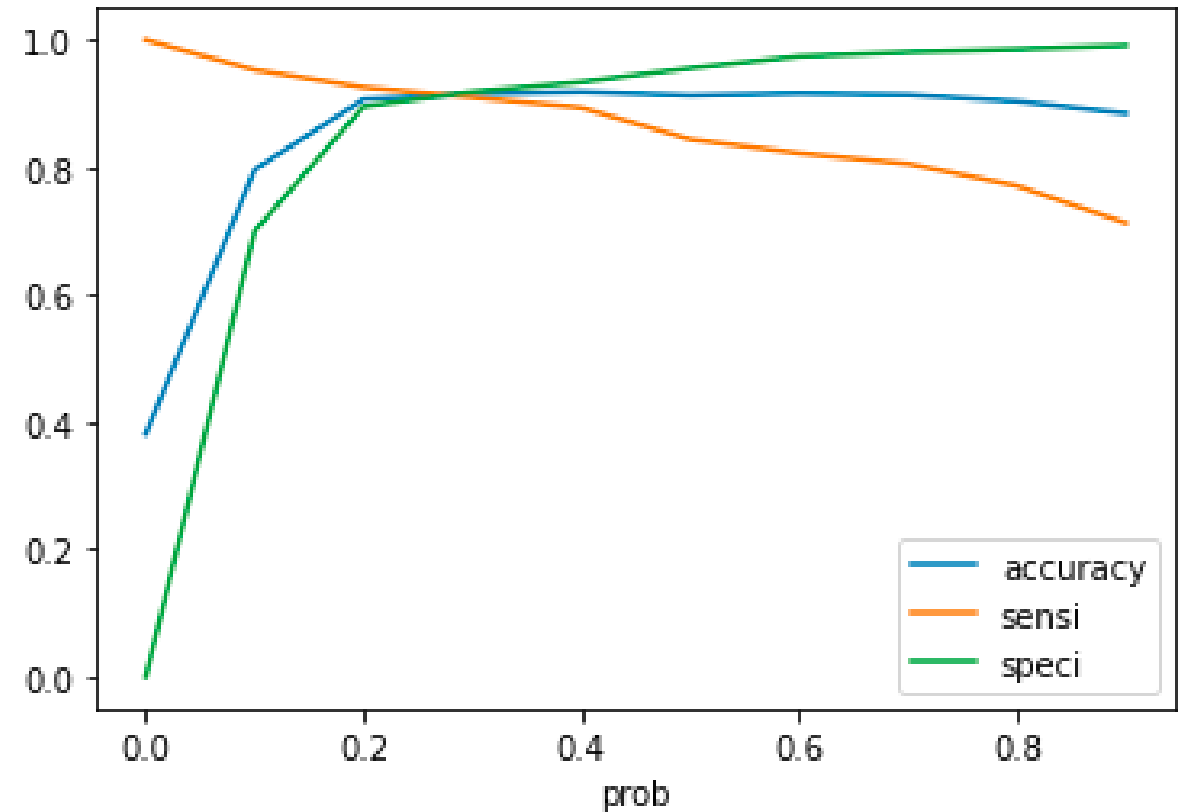
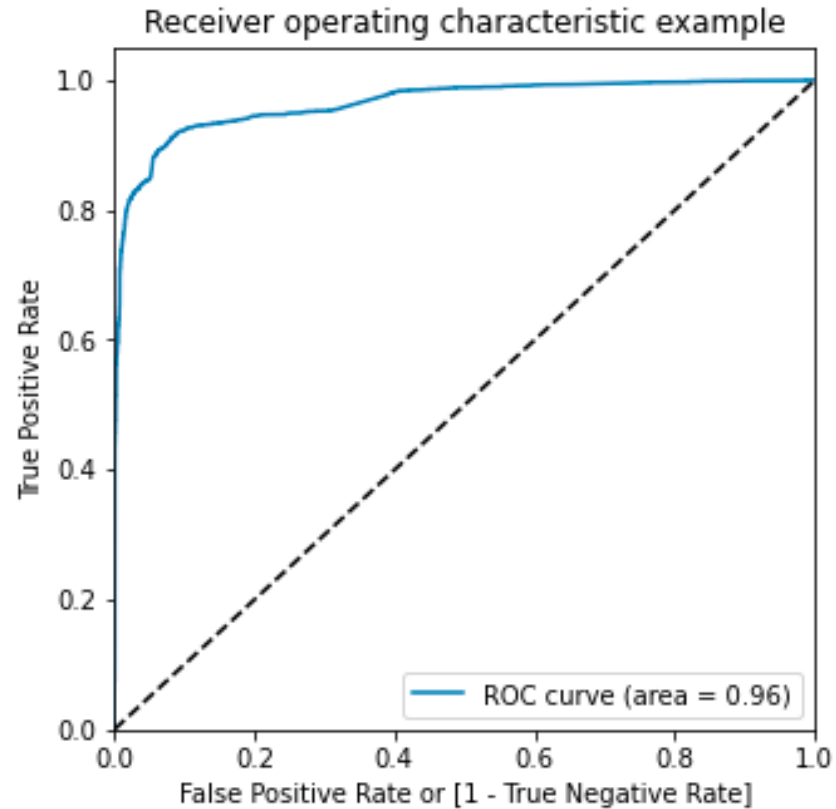
- Using MinMaxScaler, all the numeric variables were re-scaled.

- **Model Building**

- Using Recursive Feature Elimination (RFE) for selecting 15 top important features.
- Built model with the variables which have $VIF < 5$ and $p\text{-value} < 0.05$.
- Checked the optimal probabilities and the cutoff points on the train set.

- **Model Evaluation and Prediction on Test dataset**
- Area under ROC Curve is 0.96.
- For the train dataset, with a 0.3 cutoff, Accuracy, Sensitivity, and Specificity 91%.
- Predicted the test dataset with Accuracy (91%), Sensitivity (92.5%), and Specificity (91%).
- **Precision – Recall**
- With cutoff 0.41,
 - Accuracy - 91.9%, sensitivity - 89.2% and specificity - 93.5% on the train dataset.
 - Accuracy - 91.6%, sensitivity - 90.4% and specificity - 92.4% on the test dataset.
 - And, precision is 88%, and recall is 90% on the test dataset.

ROC Curve and Optimal Probabilities



- **Conclusion and Findings**

- The lead score on the test dataset - 90.41%.

- Features with high probability of a lead getting converted are:

- Tags_Closed by Horizzon
- Tags_Lost to EINS
- Lead Source_Welingak Website
- Tags_Will revert after reading the email.
- Last Activity_SMS Sent

- Negatively affected variables to avoid. Some of them are

- Tags_Ringing
- Tags_switched off
- Tags_invalid number