

*WELCOME*

# ***INTRODUCTION***

## **Problem Statement:**

- An education company named X Education sells online courses to industry professionals. On any given day, many professionals who are interested in the courses land on their website and browse for courses.
- The company markets its courses on several websites and search engines like Google. Once these people land on the website, they might browse the courses or fill up a form for the course or watch some videos. When these people fill up a form providing their email address or phone number, they are classified to be a lead. Moreover, the company also gets leads through past referrals. Once these leads are acquired, employees from the sales team start making calls, writing emails, etc. Through this process, some of the leads get converted while most do not.

## **Business Objectives:**

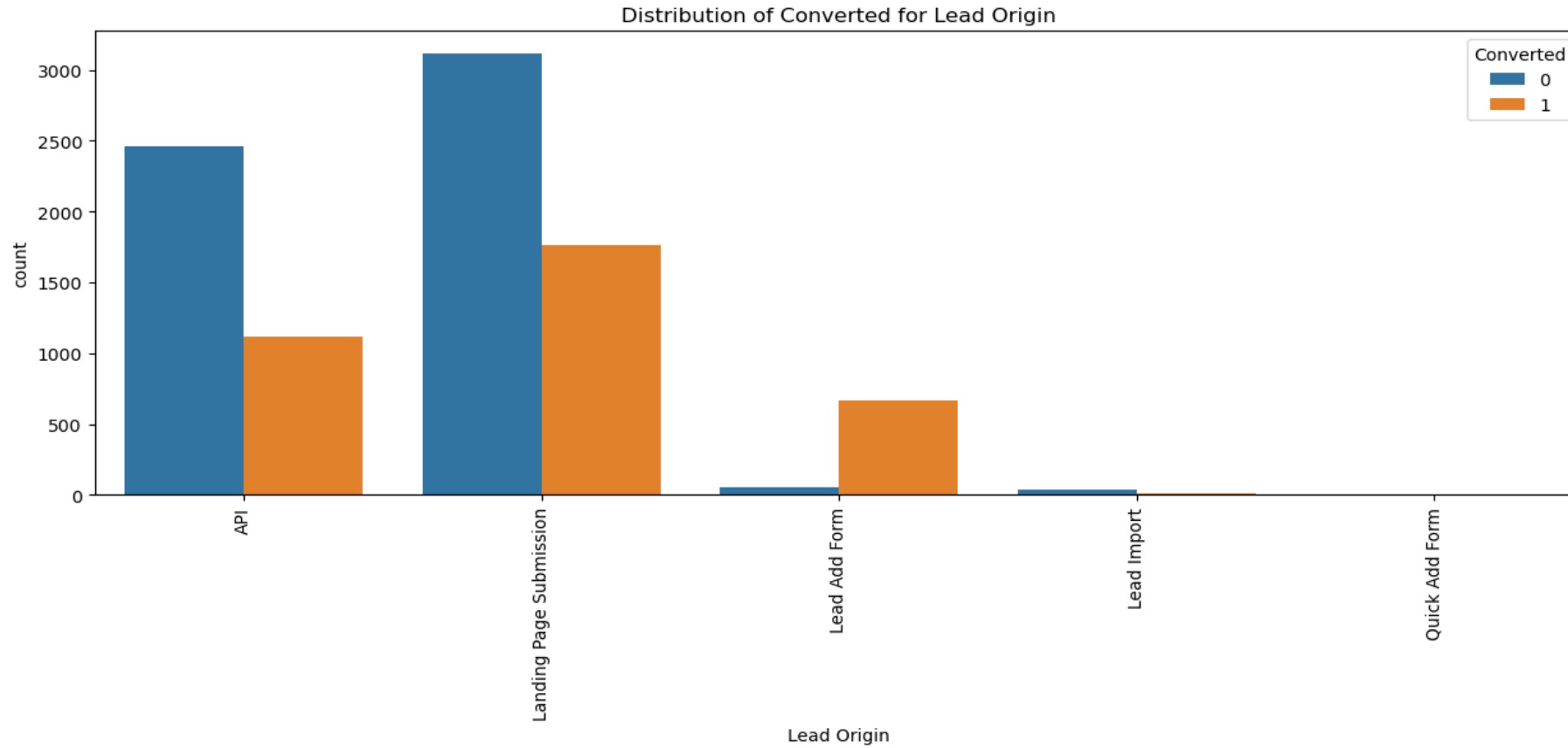
- X education wants to know most promising leads.
- For that they want to build a Model which identifies the hot leads.
- Deployment of the model for the future use.
- The CEO, in particular, has given a ballpark of the target lead conversion rate to be around 80%.

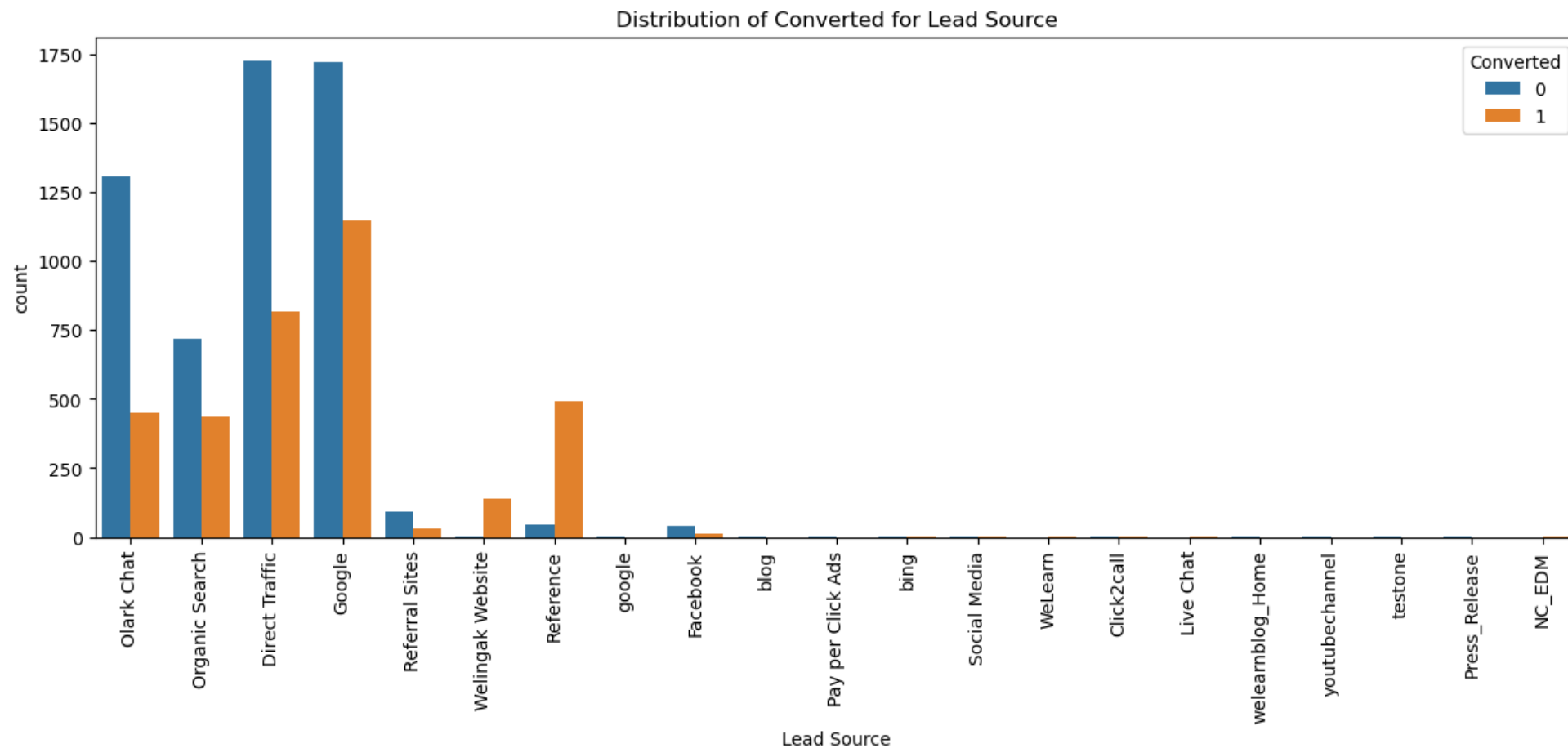
# ***Solution Methodology***

- Data cleaning and data manipulation.
  1. Check and handle duplicate data.
  2. Check and handle NA values and missing values.
  3. Drop columns.
  4. Imputation of the values.
  5. Check and handle outliers in data.
- EDA
- Feature Scaling & Dummy Variables and encoding of the data.
- Classification technique: logistic regression used for the model making and prediction.
- Validation of the model.
- Model presentation.

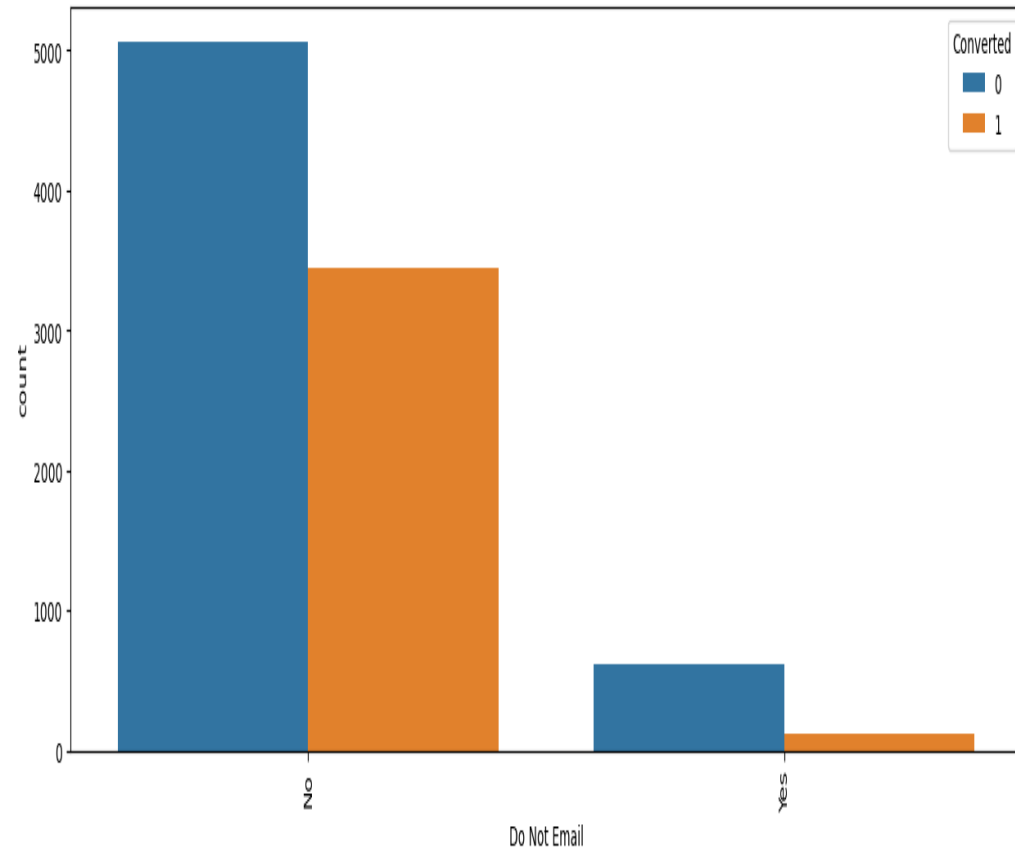
# ***Data cleaning and data manipulation***

- Raw data is having 9240 rows and 37 columns.
- Raw data is having 39% conversion rate.
- Out of 37 columns 16 columns having missing values.
- There are lots of columns which are having 'select' as value probably due to columns are not marked as mandatory.
- There are no duplicates in the data.
- 'Prospect ID', 'Lead Number' are dropped as this are unique value and not required for our analysis.
- Columns(Lead\_Quality,Asymmetrique\_Activity\_Index,Asymmetrique\_Profile\_Score,Asymmetrique\_Activity\_Score,Asymmetrique\_Profile\_Index) which are having more the 40% missing values are dropped.

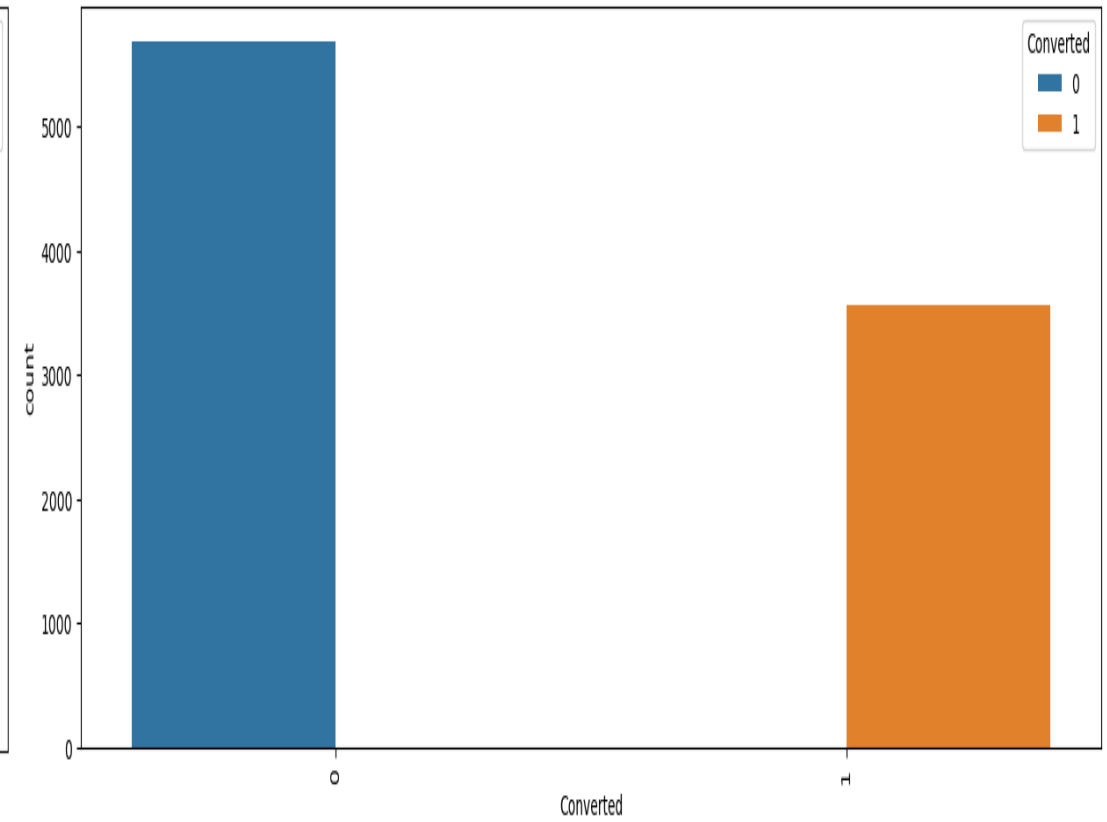




Distribution of Converted for Do Not Email

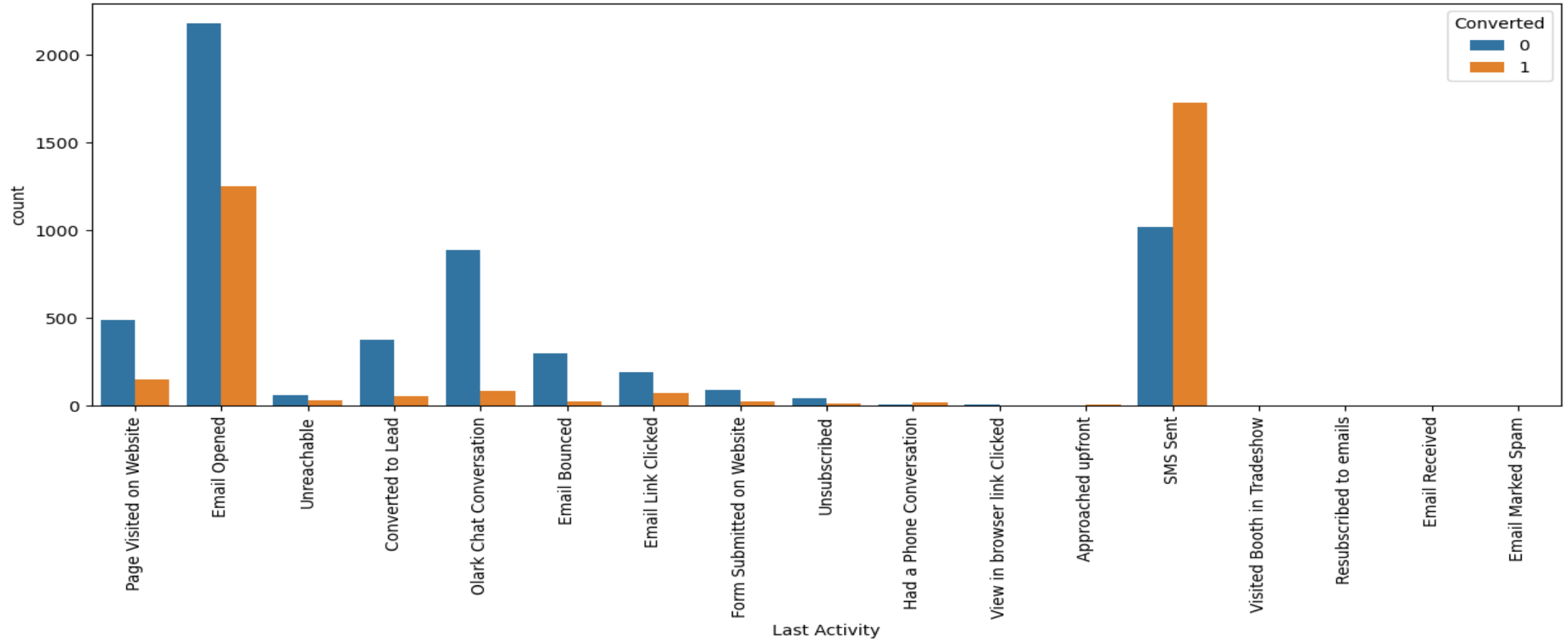


Distribution of Converted for Converted





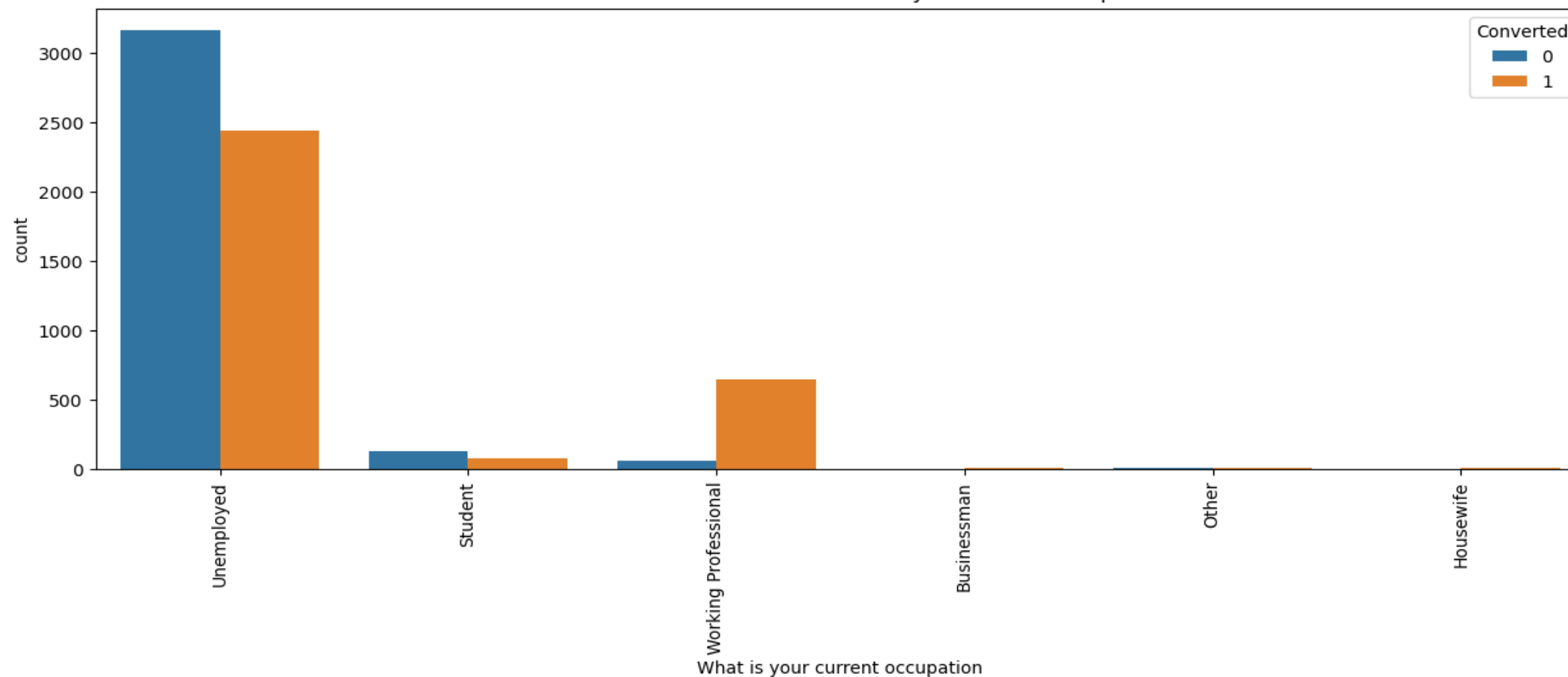
Distribution of Converted for Last Activity

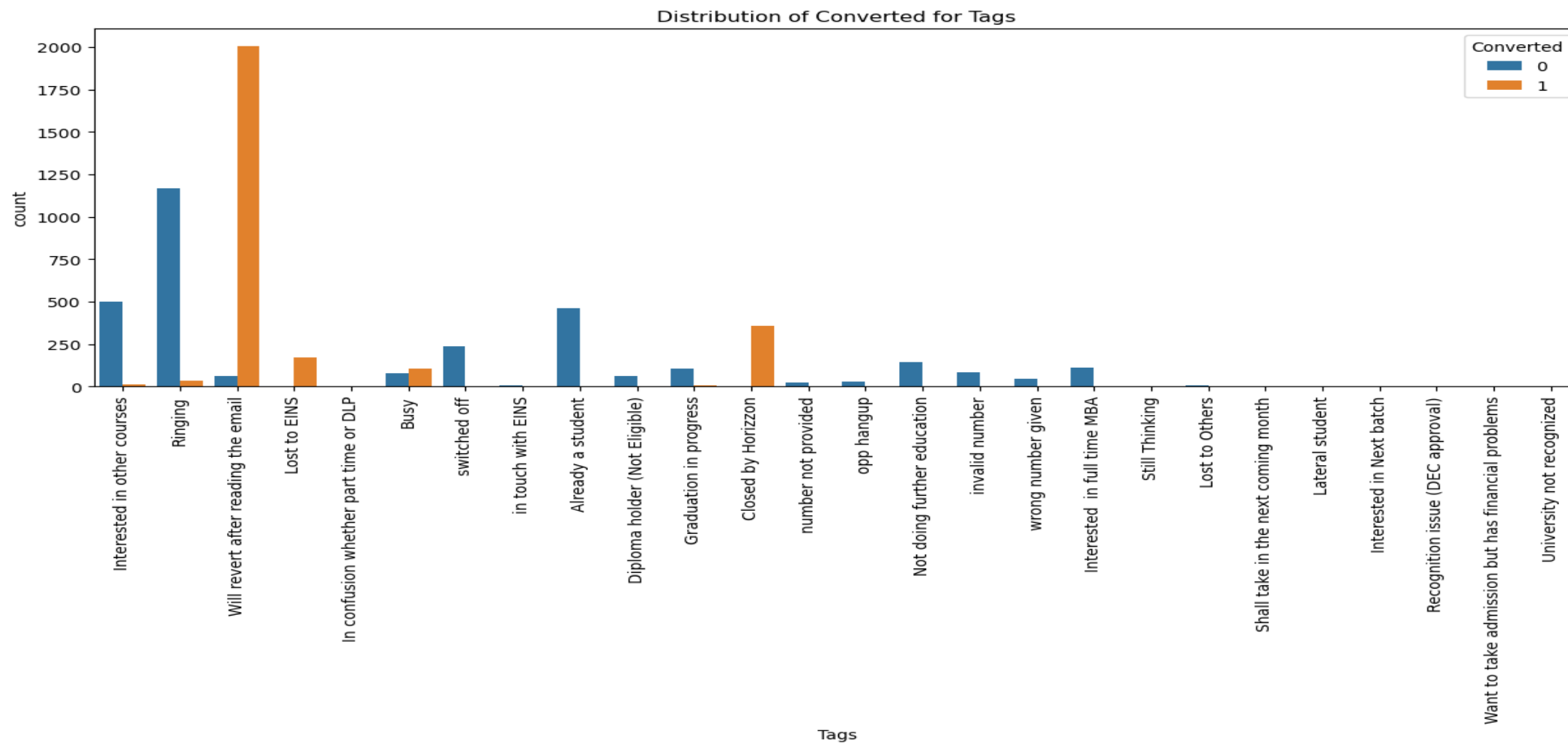




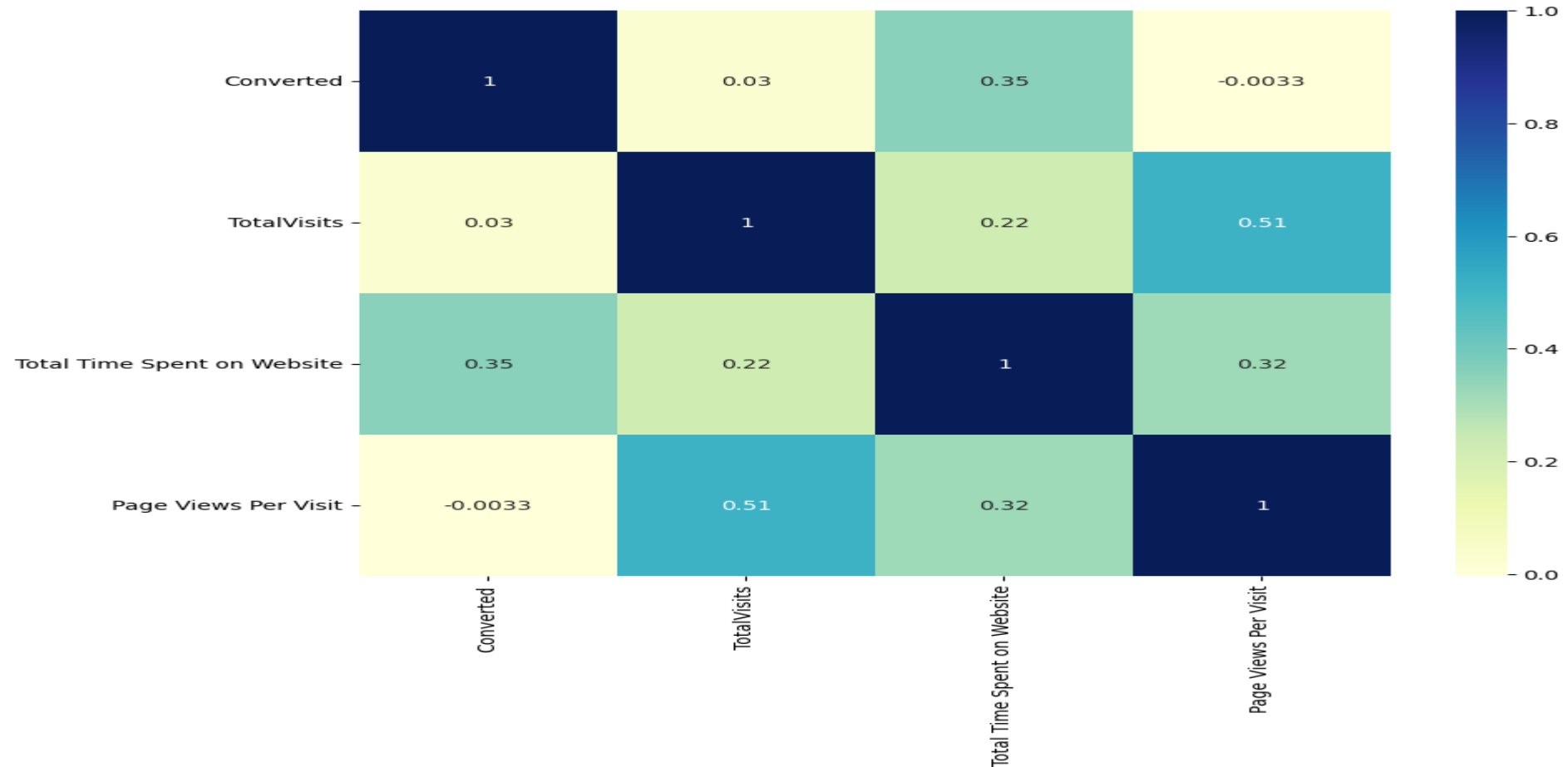


Distribution of Converted for What is your current occupation

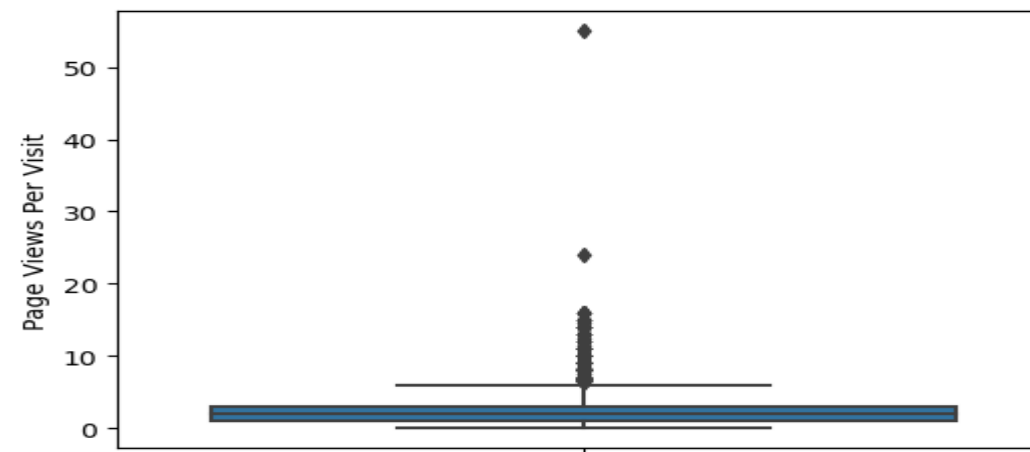
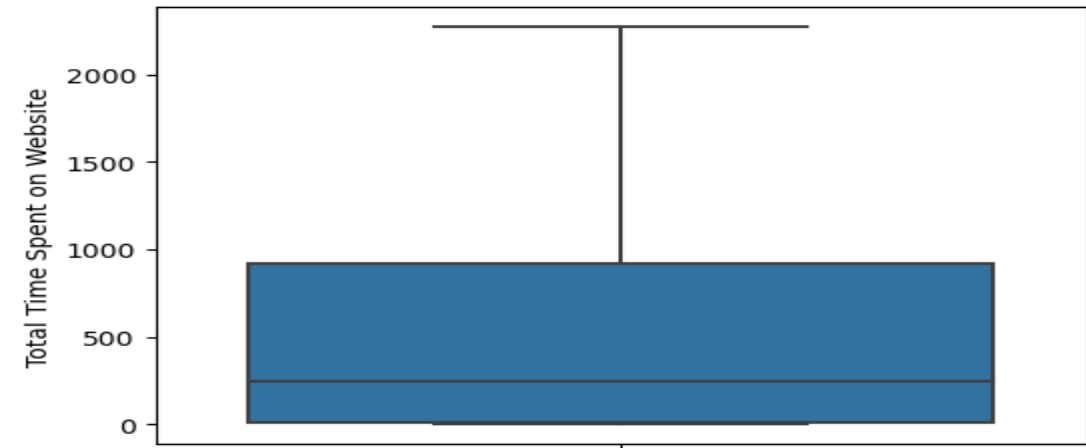
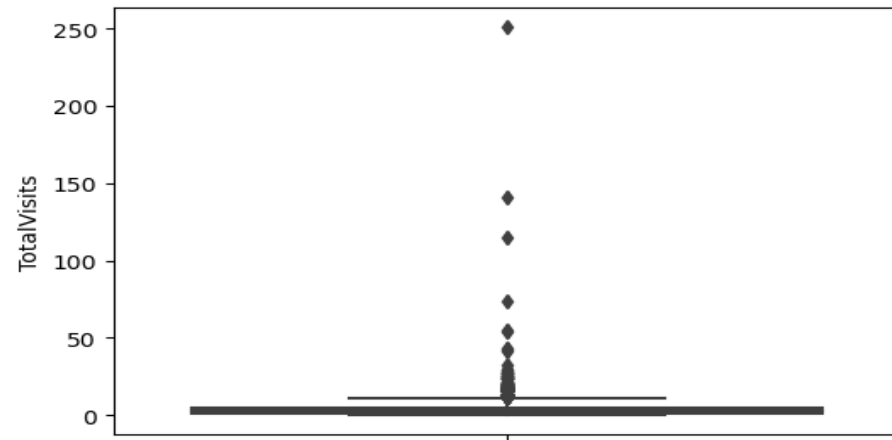




# Correlation Matrix For Numerical Variables



# *Numerical Outlier Checking*



## ***Feature Scaling & Dummy Variables and encoding of the data.***

- All numerical outlier variables are normalized.
- Dummy Variables for binary columns and categorical columns are created.
- Total no of columns 50 and total no of rows 8953 are created for Logistic regression model building and analysis

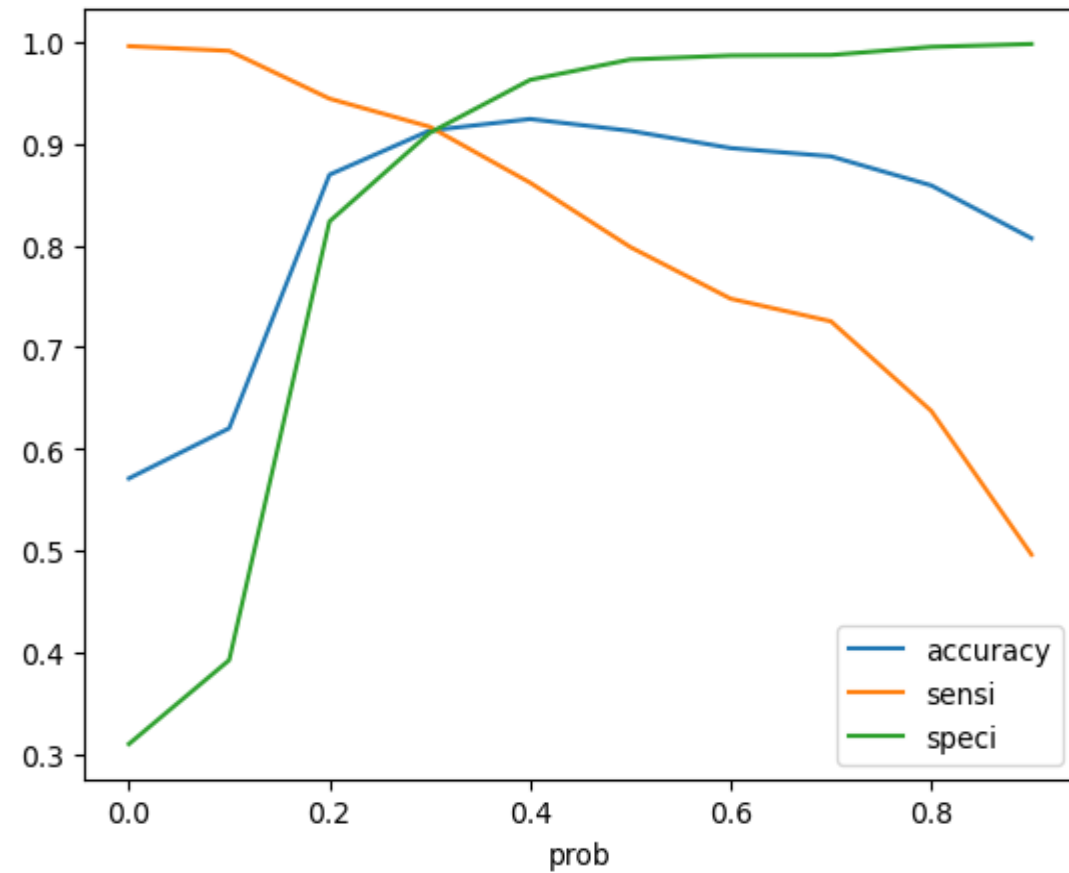
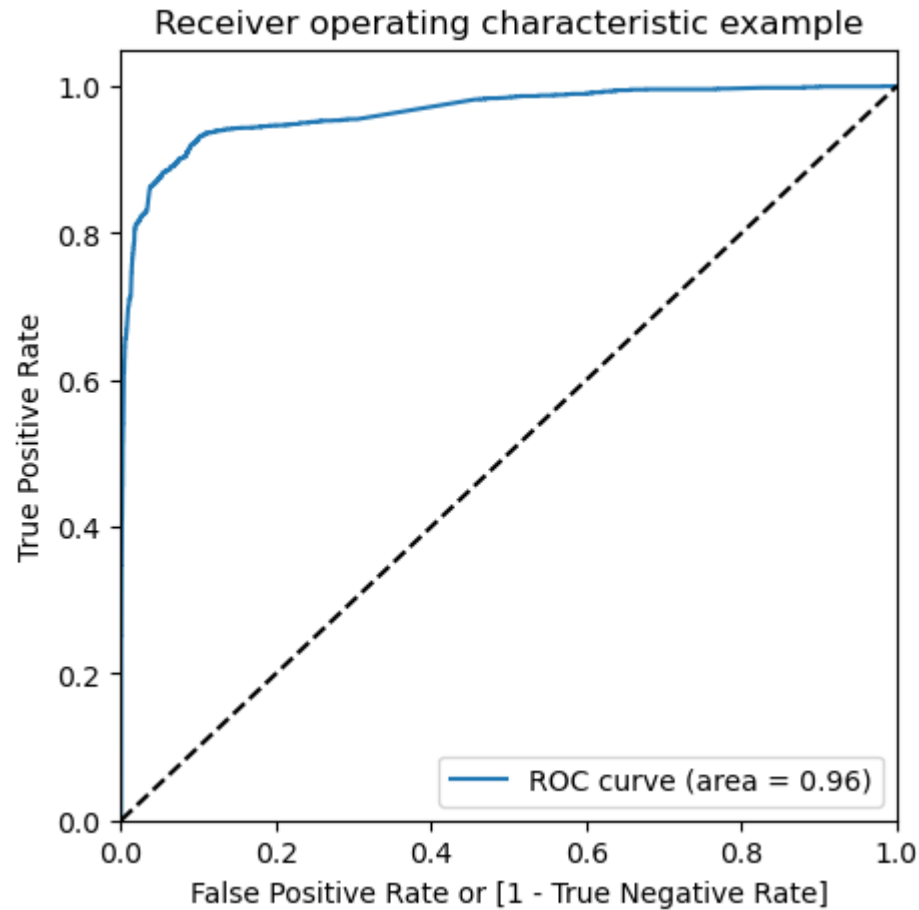
## ***Classification technique: logistic regression used for the model making and prediction***

- Splitting the data in train and test for 70% and 30% ratio.
- Scaling and fitting the train data set using SKLEARN MinMaxScaler.
- Top 15 important feature selection using RFE.
- Automation for Stats model and VIF checking.
- Elimination of features which are having VIF >5 and P value >0.05.
- Making Prediction on train data set.
- Achieved accuracy of 91.28% on trained data.
- 91.65% Specificity : 91.06%
- Optimal cut off value found as 0.3

# *Validation of the model*

- Achieved overall accuracy of 92.33%.
- Achieved Sensitivity : 91.49% Specificity : 92.84% on test data.
- Precision 88.50% achieved.
- Recall 91.48% achieved.
- Point where Accuracy, Sensitivity, Specificity intersects each other is the point for choosing optimal cutoff point.
- Area of ROC curve found 0.9 which is very nearly to 1.

# ROC Curve





## *Summary*

- Our final model performance was attained using accuracy, sensitivity, specificity and ROC-AUC scores. The model demonstrated a good ability to differentiate between converting and non\_converting leads, with sensitivity being particularly important given the business goal of identifying hot leads.
- We also observed that the Total Time Spent on Website, Lead Origin\_Landing Page Submission and Last Activity\_SMS Sent were the top predictors of conversion suggesting that engagement metrics are key drivers of customer interest.
- One key learning from this assignment was the importance of thorough data cleaning and preparation, especially handling missing values and redundant categories in categorical variables. Proper feature engineering significantly impacts model performance by reducing noise and focusing on the most relevant information.