

A wooden abacus with dark brown beads on a dark background. The abacus is made of light-colored wood and has several rows of beads. The beads are dark brown and are arranged in a grid-like pattern. The background is a dark, solid color.

# LEAD SCORING CASE STUDY

ANSUMAN PATNAIK – 8763073935

MUKARRAM ALI – 8555008016

ABHISHEK PARAYIL - 9446534873

# PROBLEM STATEMENT

- X Education sells online courses to industry professionals and Leads are generated from various sources are captured.
- To make this process more efficient, the X Education company wants to identify the potential leads so that the rate of conversion can be higher.
- If the company identifies potential leads, the employees will focus more on communicating with potential leads rather than writing emails or calling every user.



# SOLUTION APPROACH

## 1. Understanding the data

## 2. Clean the data

1. Handle the duplicate data, missing values and drop columns irrelevant to our analysis.

## 3. Model Building Preparation

1. Univariate and Bivariate Data Analysis.

## 4. Model Building

- 4.1. Scaling and Dummy variables
- 4.2. Split the data into Train and Test

## 5. Model Evaluation

1. Creation of a Confusion matrix and finding the data's overall accuracy, specificity and sensitivity.

## 6. Prediction on the Test set

1. Predictions were performed on the Test data and then find the data's overall accuracy, specificity and sensitivity.

## 7. Conclusion





# UNDERSTAND THE DATA

Leads.csv contains all the leads generated through various sources. This file contains the following:

- ▶ The file contains 9240 rows and 37 columns
- ▶ Out of 37 columns, 7 are numerical columns and the rest 30 are categorical columns.

Leads Data Dictionary.csv file describes the meaning of all the variables involved in the Leads dataset.



# DATA CLEANING AND PREPARATION

The columns were analysed and the columns where only one unique value is present have been dropped. Below are the columns:

- Magazine
- Receive More Updates About Our Courses
- Get updates on DM Content
- Update me on Supply Chain Content
- I agree to pay the amount through a cheque

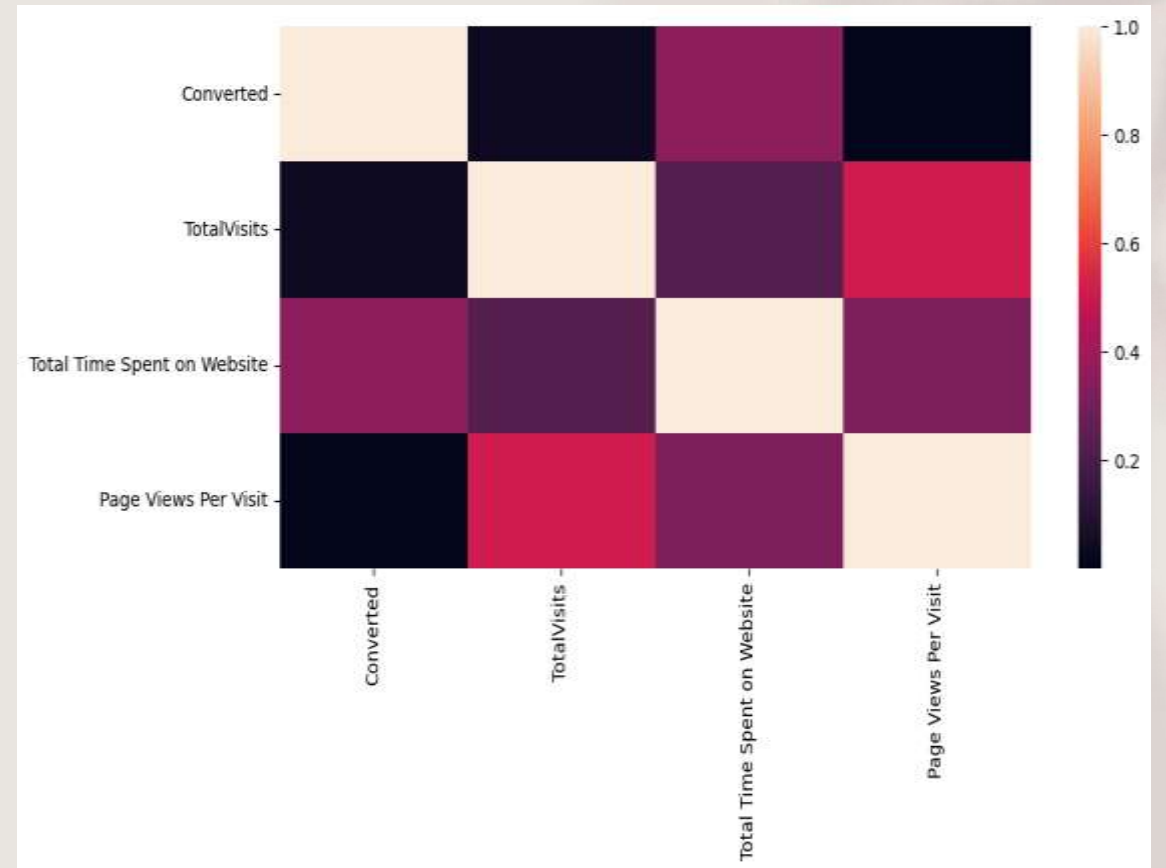
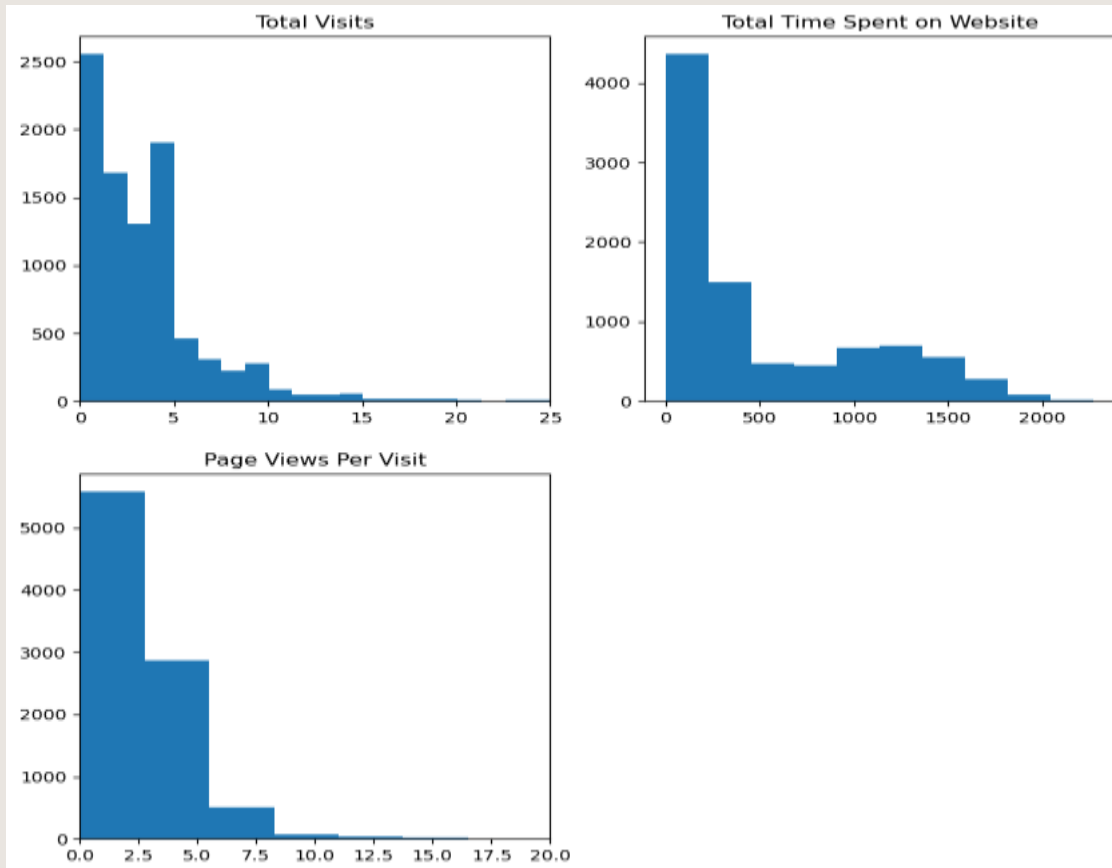
Below are the columns which are having more than 35% of the null values and the columns that were dropped:

- How did you hear about X Education
- What is your current occupation
- What matters most to you in choosing a course
- Tags
- Lead Quality
- Lead Profile
- City
- Asymmetric Activity Index
- Asymmetric Profile Index
- Asymmetric Activity Score
- Asymmetric Profile Score

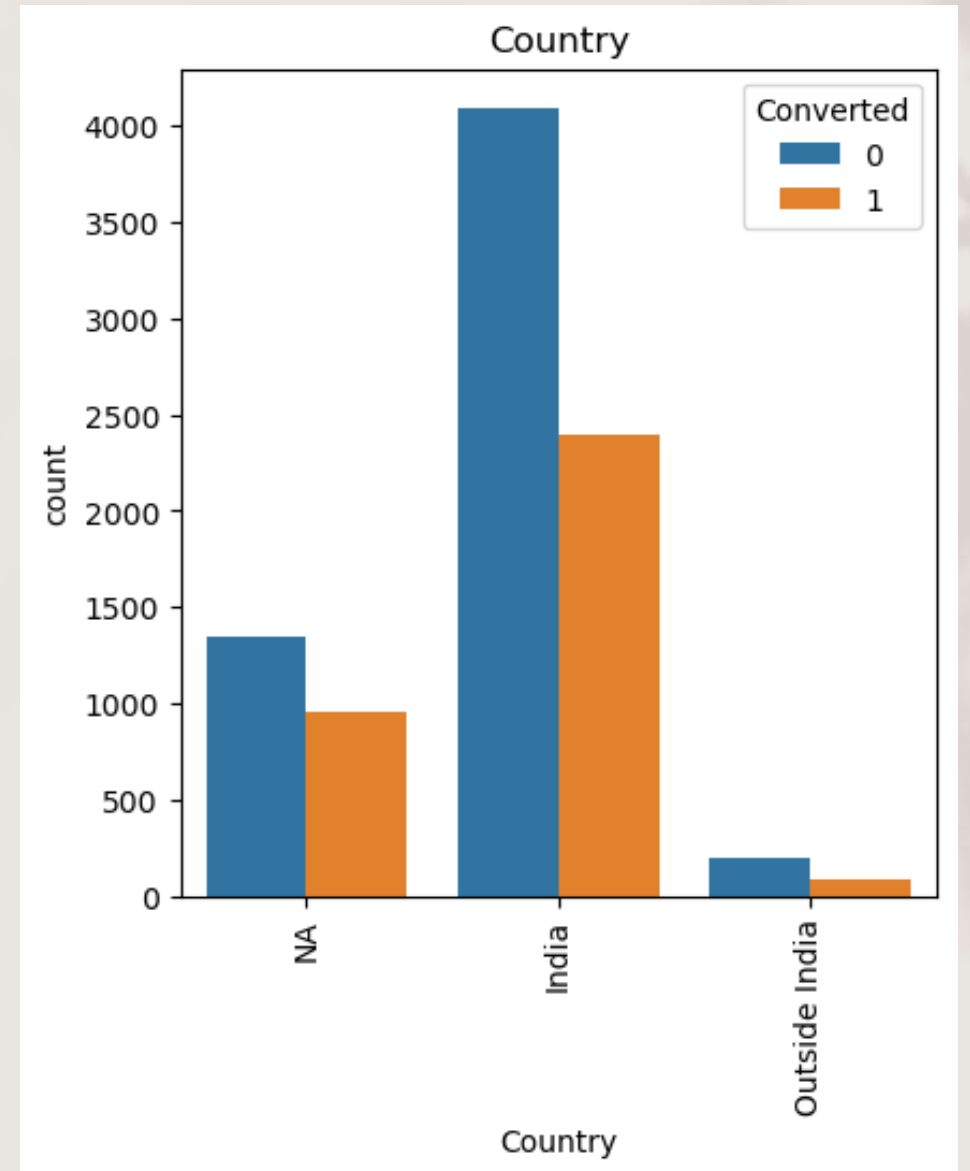
Following are the columns where Select value has been replaced with NA and later removed these data from the analysis:

- Specialization
- How did you hear about X education
- Lead Profile
- Country

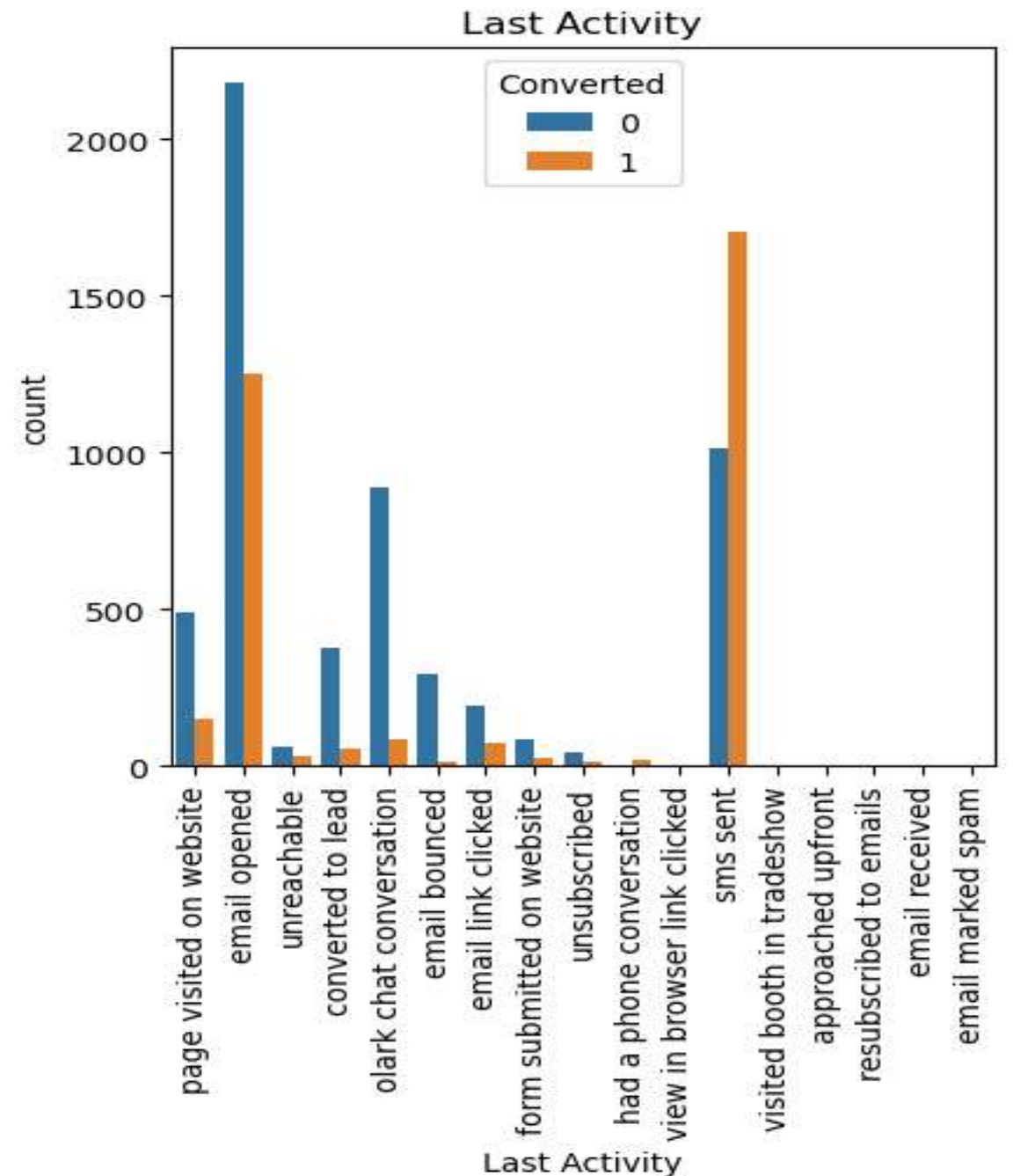
# EXPLORATORY DATA ANALYSIS



# EXPLORATORY DATA ANALYSIS

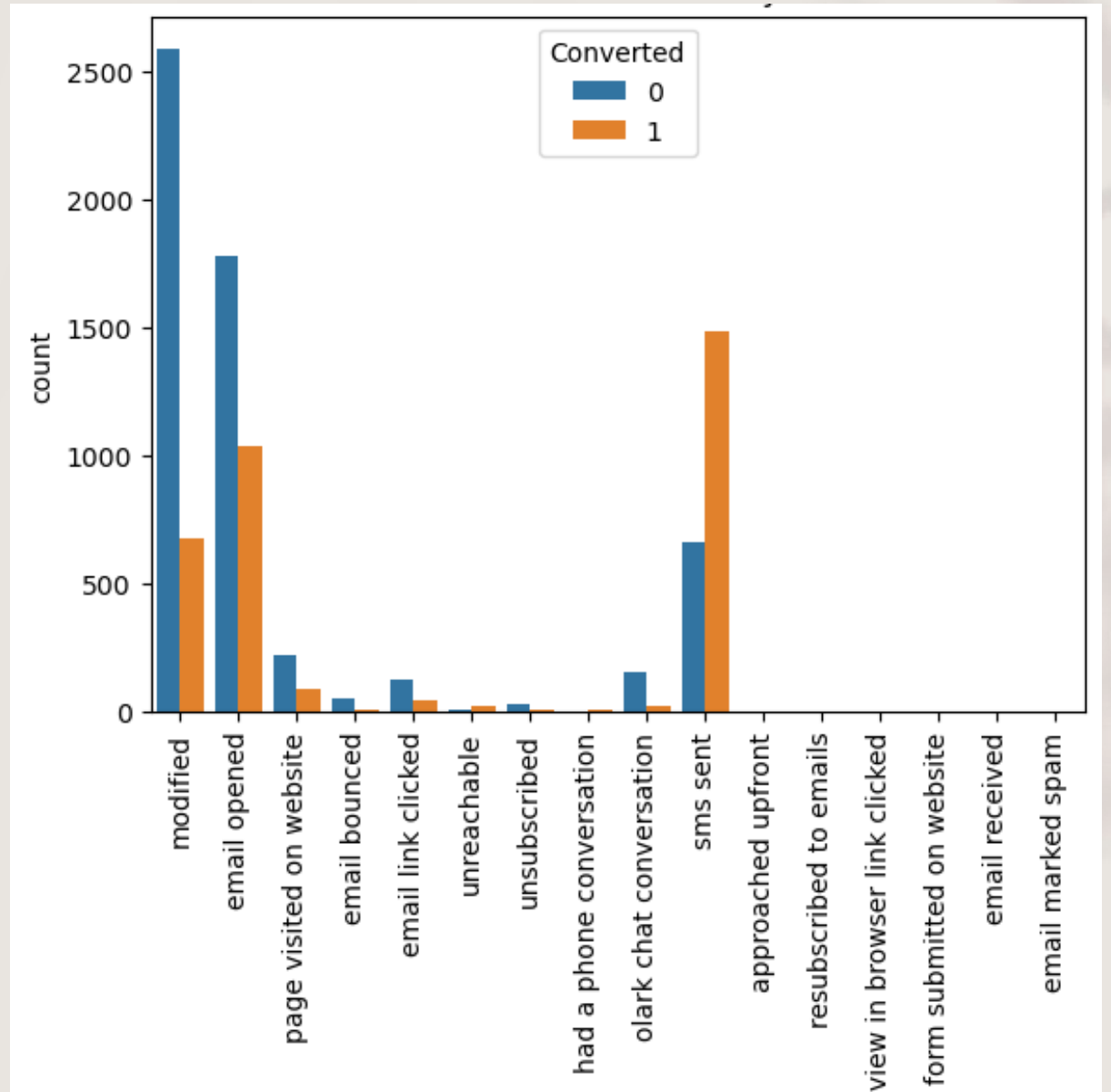


# EXPLORATORY DATA ANALYSIS

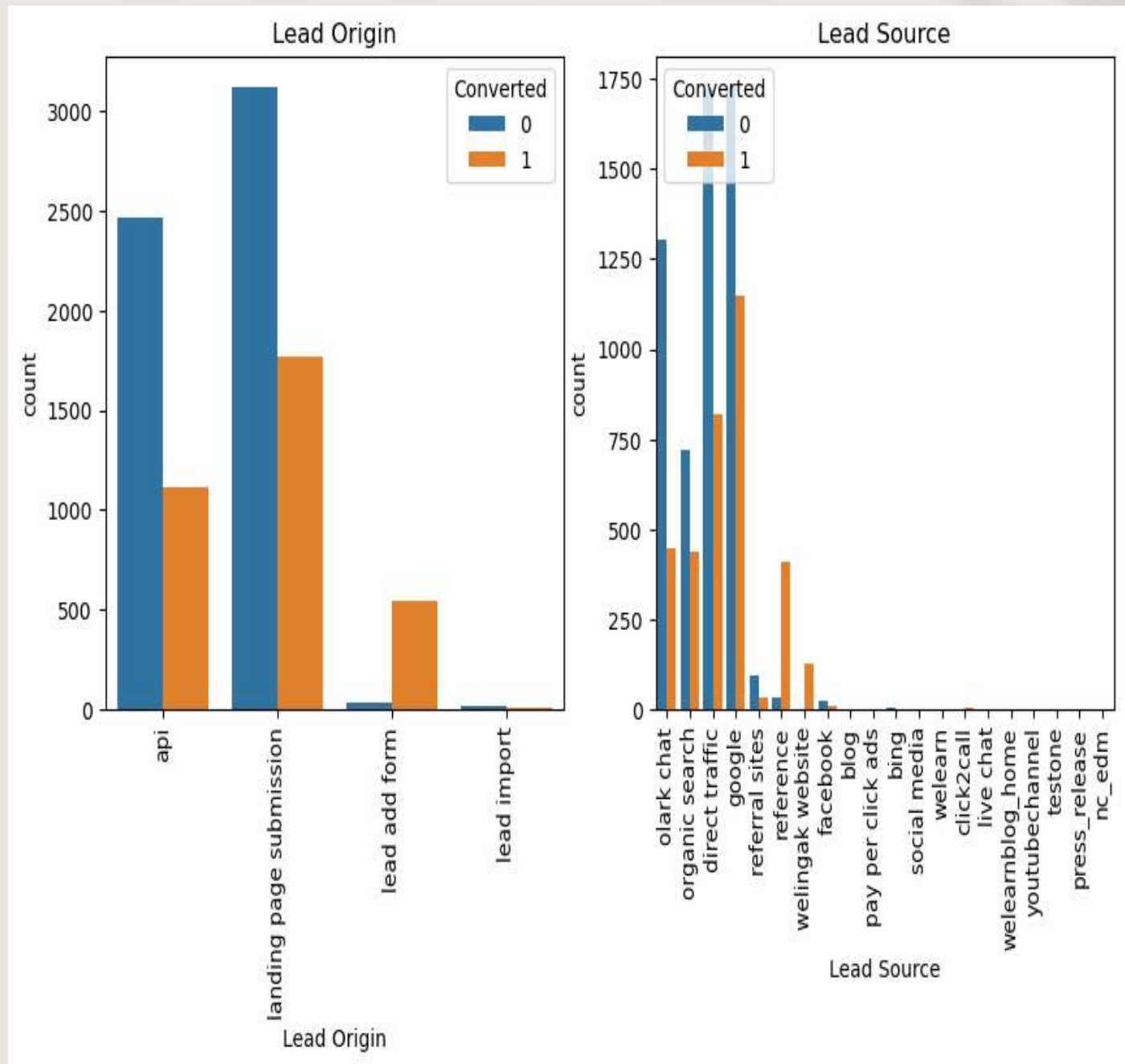




# EXPLORATORY DATA ANALYSIS



# EXPLORATORY DATA ANALYSIS



# MODEL BUILDING

To build the model we first created dummy variables for all the categorical variables using `get_dummies` method and all the NA variables were removed. For numerical variables we have used `MinMaxScaler`.

Then we have split the data into Train and Test where 70% of the data was for Train and 30% of the data was for Test.

Then we have performed RFE to attain top 15 relevant variables. Depending on the VIF and the P values i.e.,  $VIF < 5$  and  $P \text{ value} < 0.05$ , we have also dropped some columns manually using `drop`.

# MODEL EVALUATION

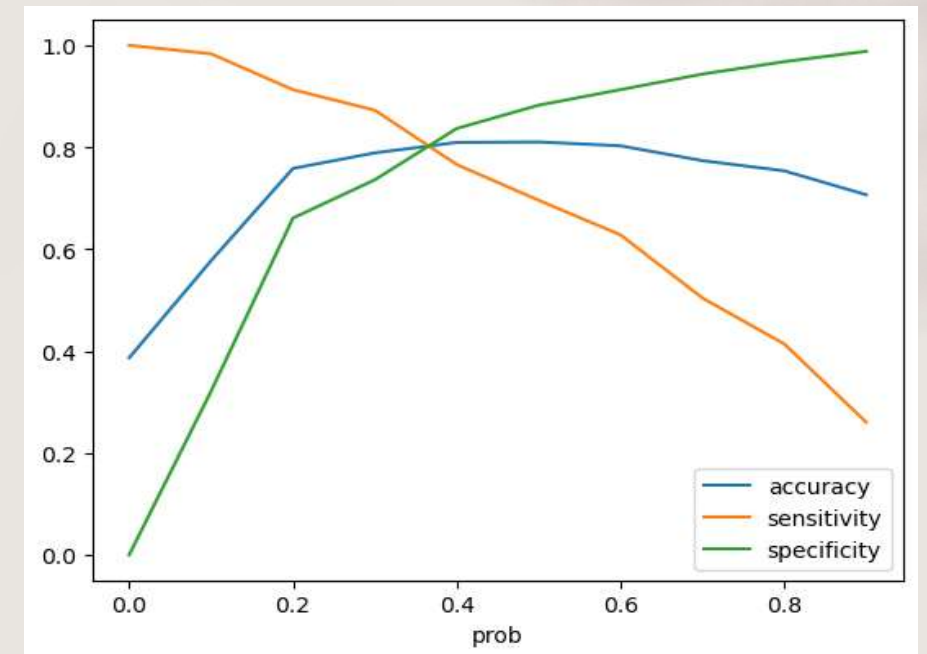
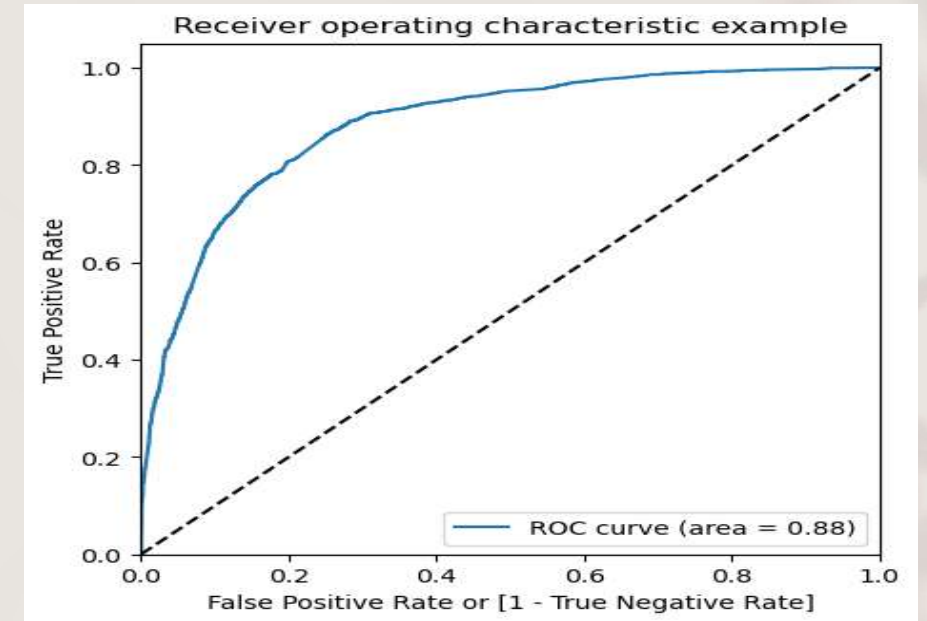
A prediction was also performed on the Test data with cut value as 0.35 and found the overall accuracy, sensitivity and specificity which turns around 80% on an average.

This method was also used to recheck with an optimum cut off value of 0.4 and the same turns around 75% on the test data frame.

A confusion matrix was created and with an optimum cut off value using ROC curve we found the overall accuracy, sensitivity and specificity which turns around 80% on an average.

# ROC CURVE

From the second chart it is viewed that optimal cut off is at 0.35 and the overall accuracy, specificity and sensitivity comes around to be around 80%





# CONCLUSION

- From the second chart it is viewed that optimal cut off is at 0.35 and the overall accuracy, specificity and sensitivity comes around to be around 80%
- Total number of visits
- When their current occupation is working professional
- The total time spend on the Website
- When the Lead Source is :
  - Direct Traffic
  - Google
  - Organic Search
  - Wellingak Website
  - Olark Chat
  - Reference
- X Education can keep the above points to get more conversion rate for their courses.

THANK YOU

