

LEAD SCORE CASE STUDY

GROUP MEMBERS:

AATISH KAUSHAL

HEMAVATHI A.B

PROBLEM STATEMENT

- X education sells online courses to industry professionals.
- X education gets a lot of leads, but its lead conversion rate is very poor.
- To make this process more efficient, the company wishes to identify the most potential leads, also known as 'Hot Leads'.
- If they successfully identify this set of leads, the lead conversion rate should go up as the sales team will now be focusing more on communicating with the potential leads rather than making calls to everyone.

Business Objective:

- X Education wants to identify the most promising leads
- With that they want to create a Model which identify hot leads
- Deployment of the model for the future use.

SOLUTION APPROACH

Data cleaning and manipulation:

- Identify and manage duplicate entries.
- Address NAN or missing values and option “SELECT”.
- Remove columns with significant missing data that are not essential for analysis.
- Impute missing values when necessary.
- Detect and manage outliers.

EDA:

Univariate Data Analysis: Analyze value counts and variable distributions.

Bivariate Data Analysis: Examine correlation coefficients and patterns between variables.

Feature Scaling & Dummy Variables and encoding of the data

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

Validation of the model

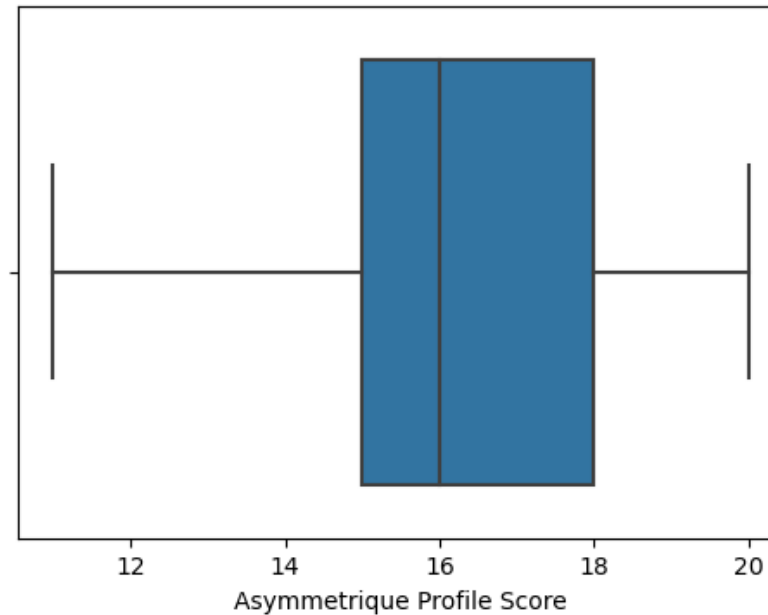
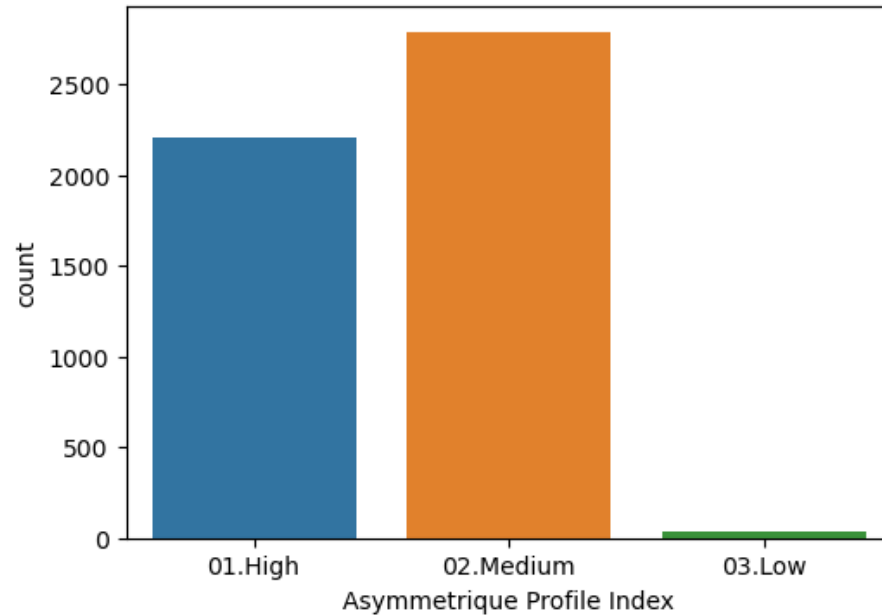
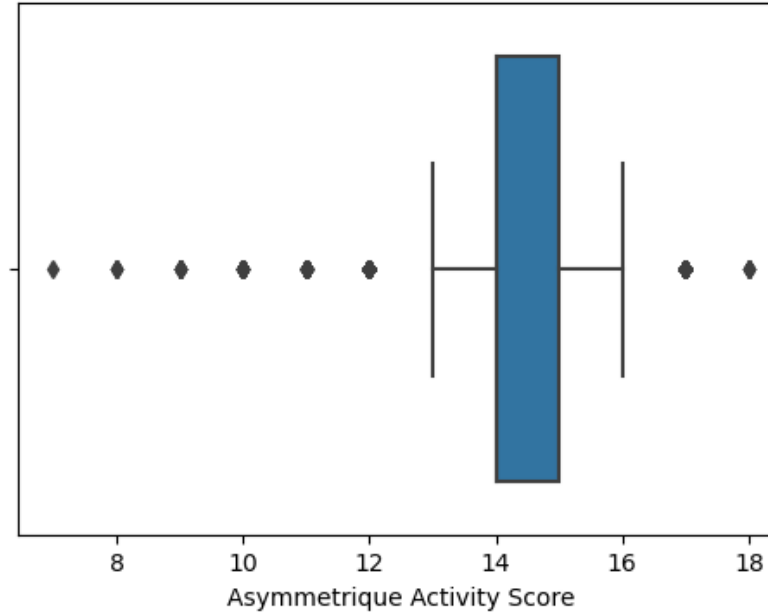
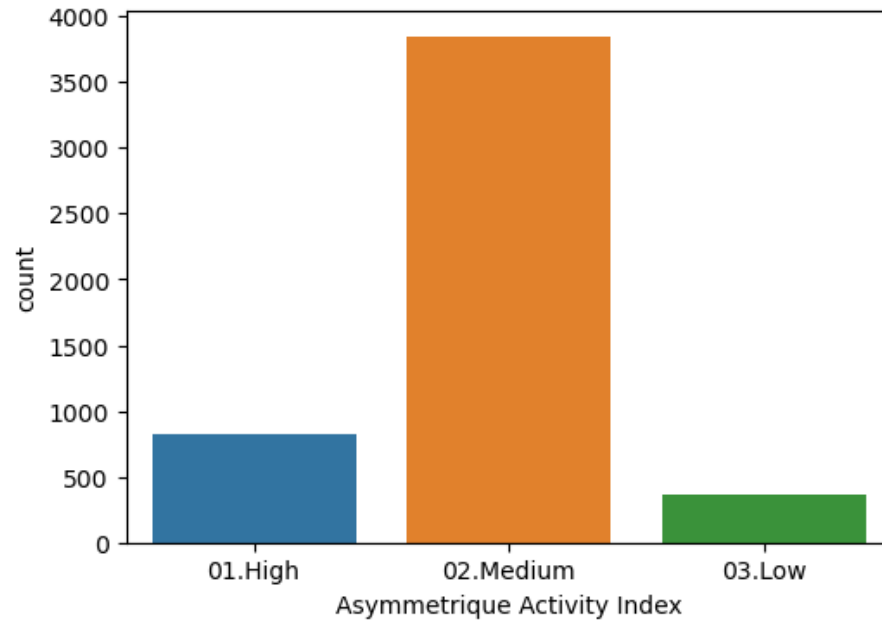
Model presentation

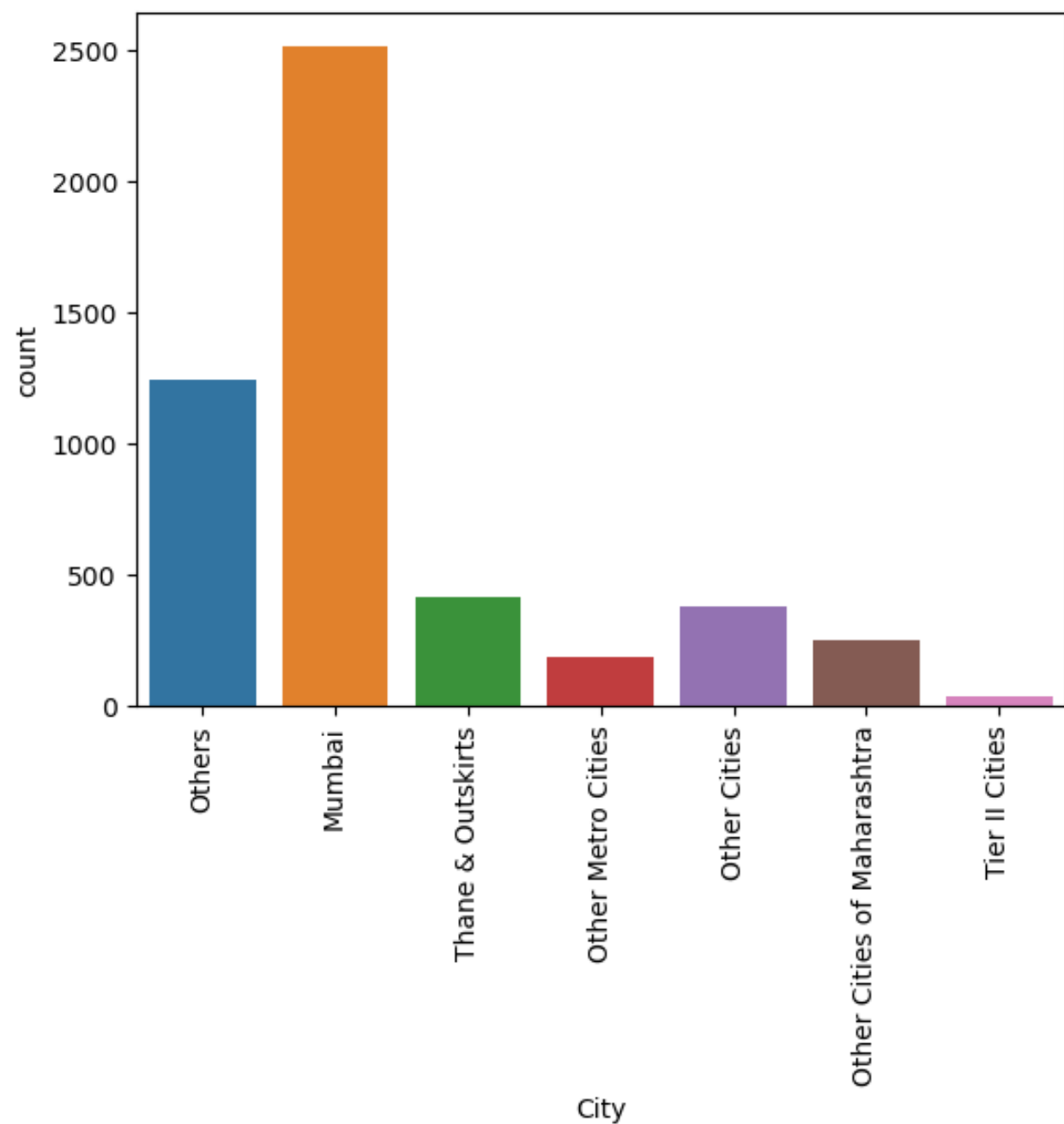
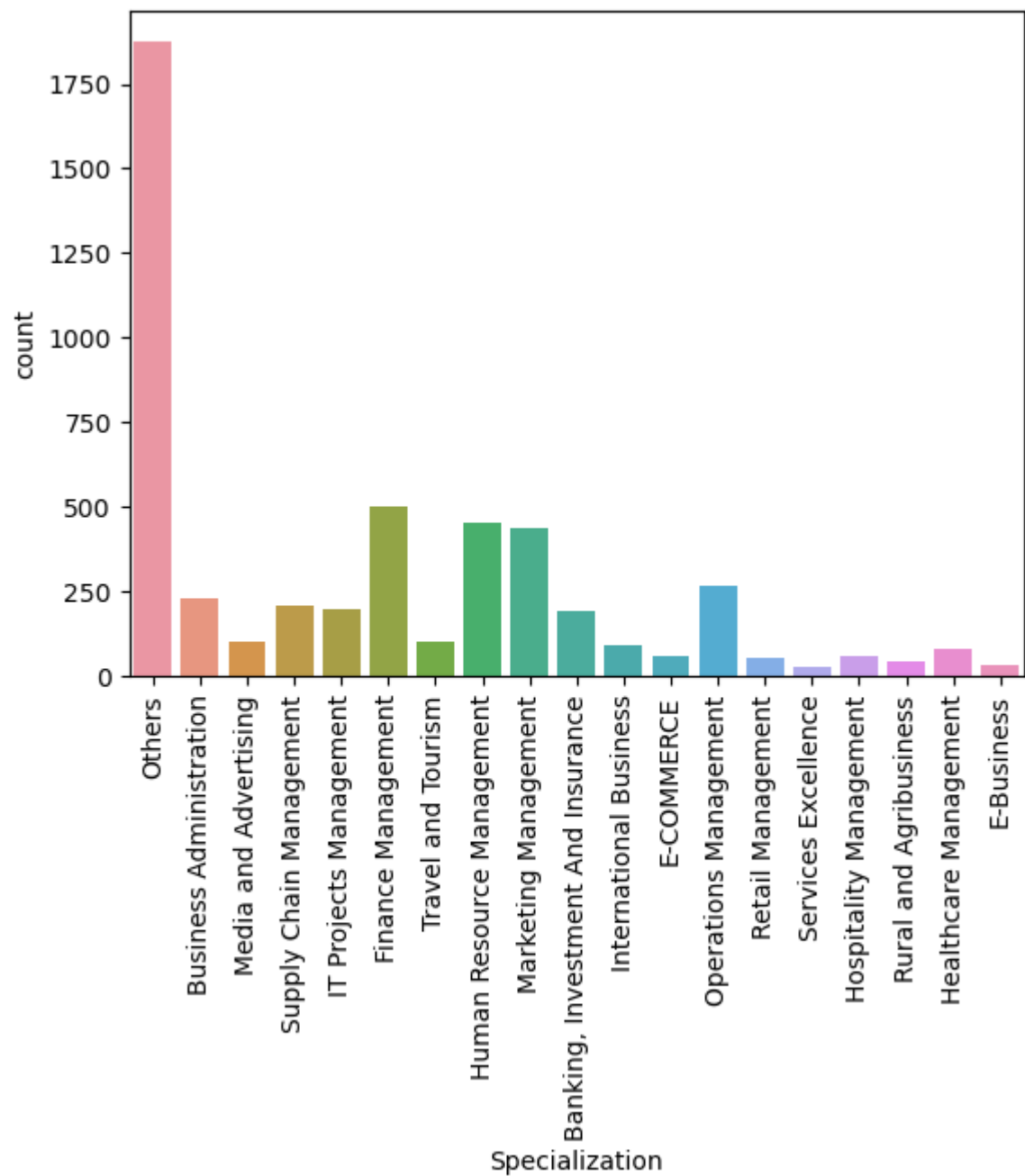
Conclusions and recommendations

DATA MANIPULATION

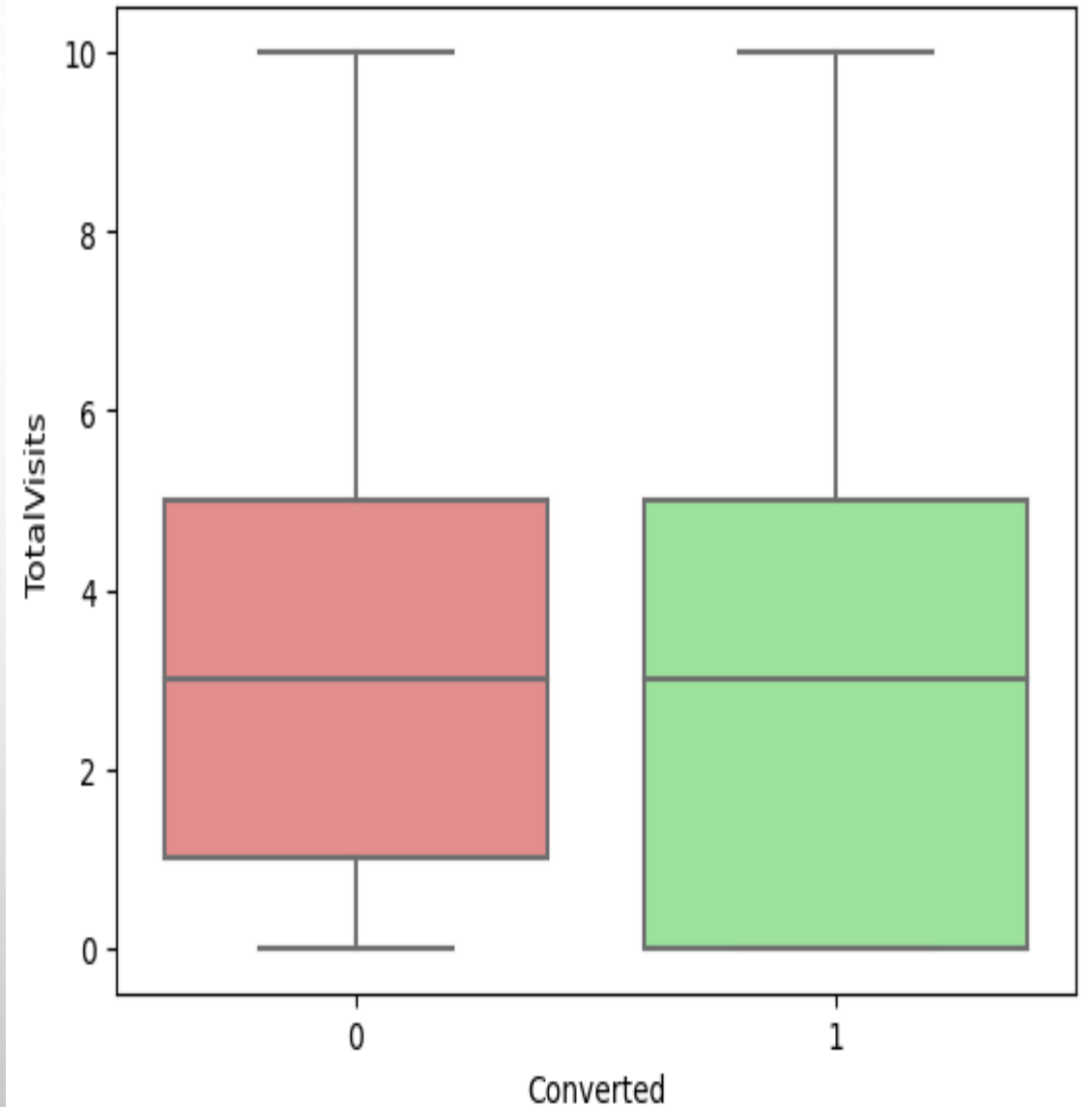
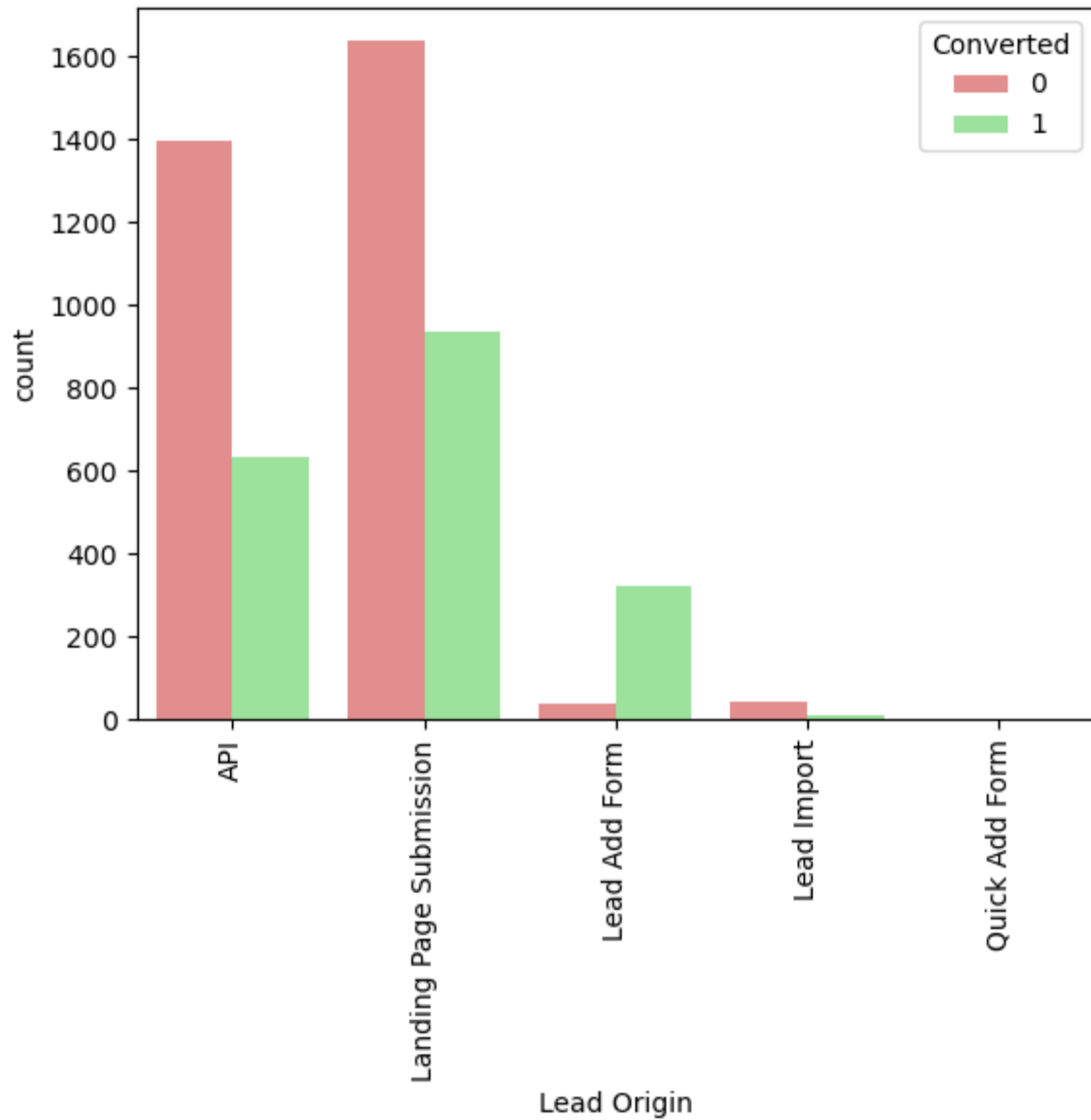
- TOTAL NUMBER OF ROWS = 9240, TOTAL NUMBER OF COLUMNS =37.
- MANY COLUMNS HAD SINGLE VALUE FEATURE “NO” LIKE INDICATING WHETHER THE CUSTOMER HAD SEEN THE AD IN ANY OF THE LISTED ITEMS FOR EXAMPLE “MAGAZINE” AND FOR , “RECEIVE MORE UPDATES ABOUT OUR COURSES”, “UPDATE ME ON SUPPLY CHAIN CONTENT”, “I AGREE TO PAY THE AMOUNT THROUGH CHEQUE” ETC.
- AN INDEX AND SCORE ASSIGNED TO EACH CUSTOMER BASED ON THEIR ACTIVITY AND THEIR PROFILE COLUMNS WERE DROPPED AS THEY SHOWED MORE THAN >45% NULL ROWS AND TOO MUCH OF VARIATIONS.
- CHOOSE MEDIAN AND MODE STATISTICS TO IMPUTE FOR MISSING ROWS AND REPLACED NAN AND SELECT WITH “OTHERS”
- IN THE “TAGS” COLUMN KEPT CONSIDERABLE LAST ACTIVITIES AS IS AND CLUBBED ALL OTHERS TO "OTHER_ACTIVITY"
- REMOVING THE “PROSPECT ID” AND “LEAD NUMBER” WHICH IS NOT NECESSARY FOR THE ANALYSIS.
- AFTER CHECKING FOR THE VALUE COUNTS FOR SOME OF THE OBJECT TYPE VARIABLES, WE FOUND SOME OF THE FEATURES WHICH HAD NOT ENOUGH VARIANCE, IS DROPPED, THE FEATURES ARE: “DO NOT CALL”, “WHAT MATTERS MOST TO YOU IN CHOOSING COURSE”, “SEARCH”, “NEWSPAPER ARTICLE”, “X EDUCATION FORUMS”, “NEWSPAPER”, “DIGITAL ADVERTISEMENT” ETC.
- DROPPING THE COLUMNS LEAD ORIGIN', 'LEAD SOURCE', 'LAST ACTIVITY', 'SPECIALIZATION', 'WHAT IS YOUR CURRENT OCCUPATION', 'TAGS', 'LEAD QUALITY', 'HOW DID YOU HEAR ABOUT X EDUCATION', 'LEAD PROFILE', 'CITY', 'LAST NOTABLE ACTIVITY

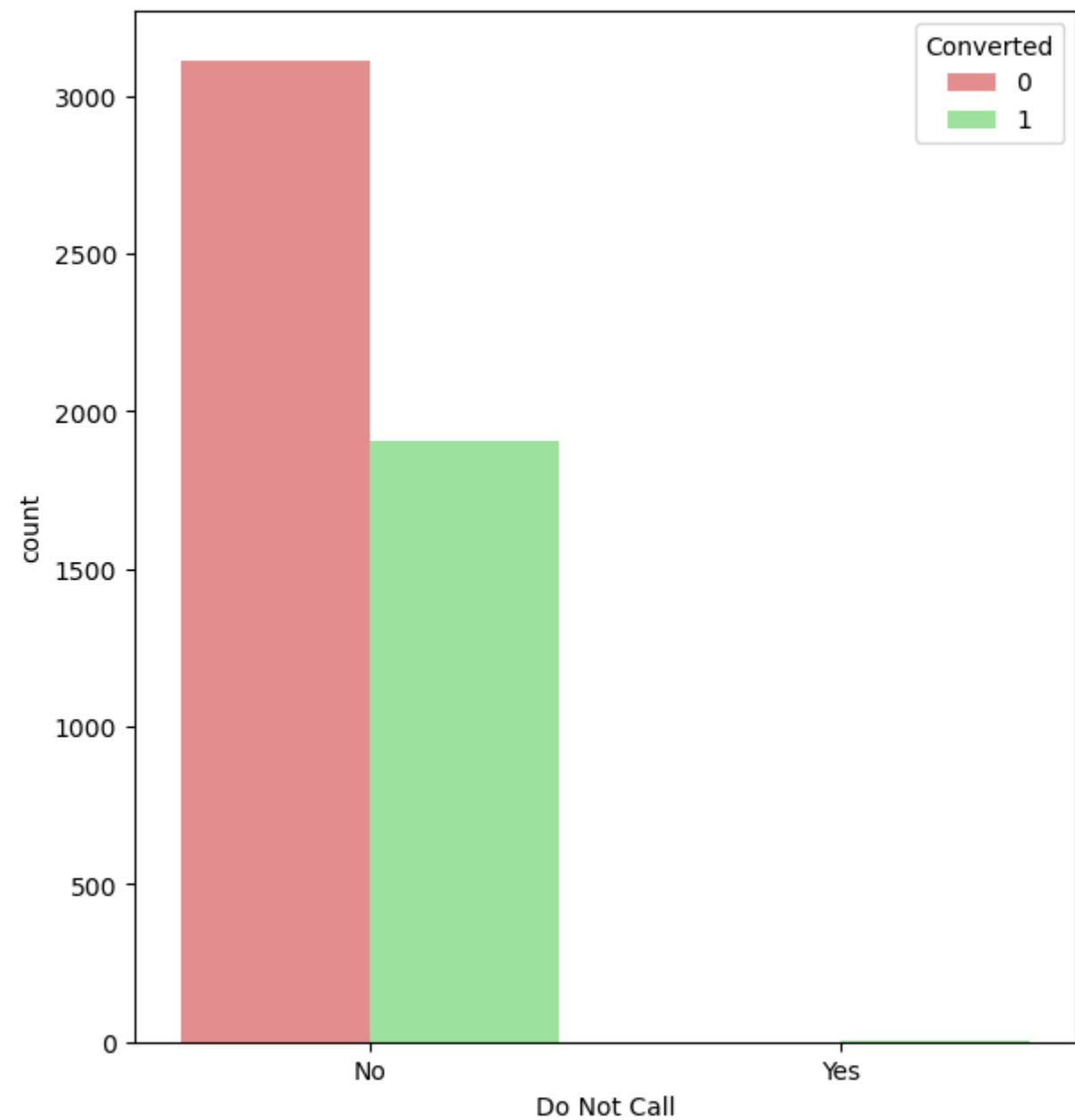
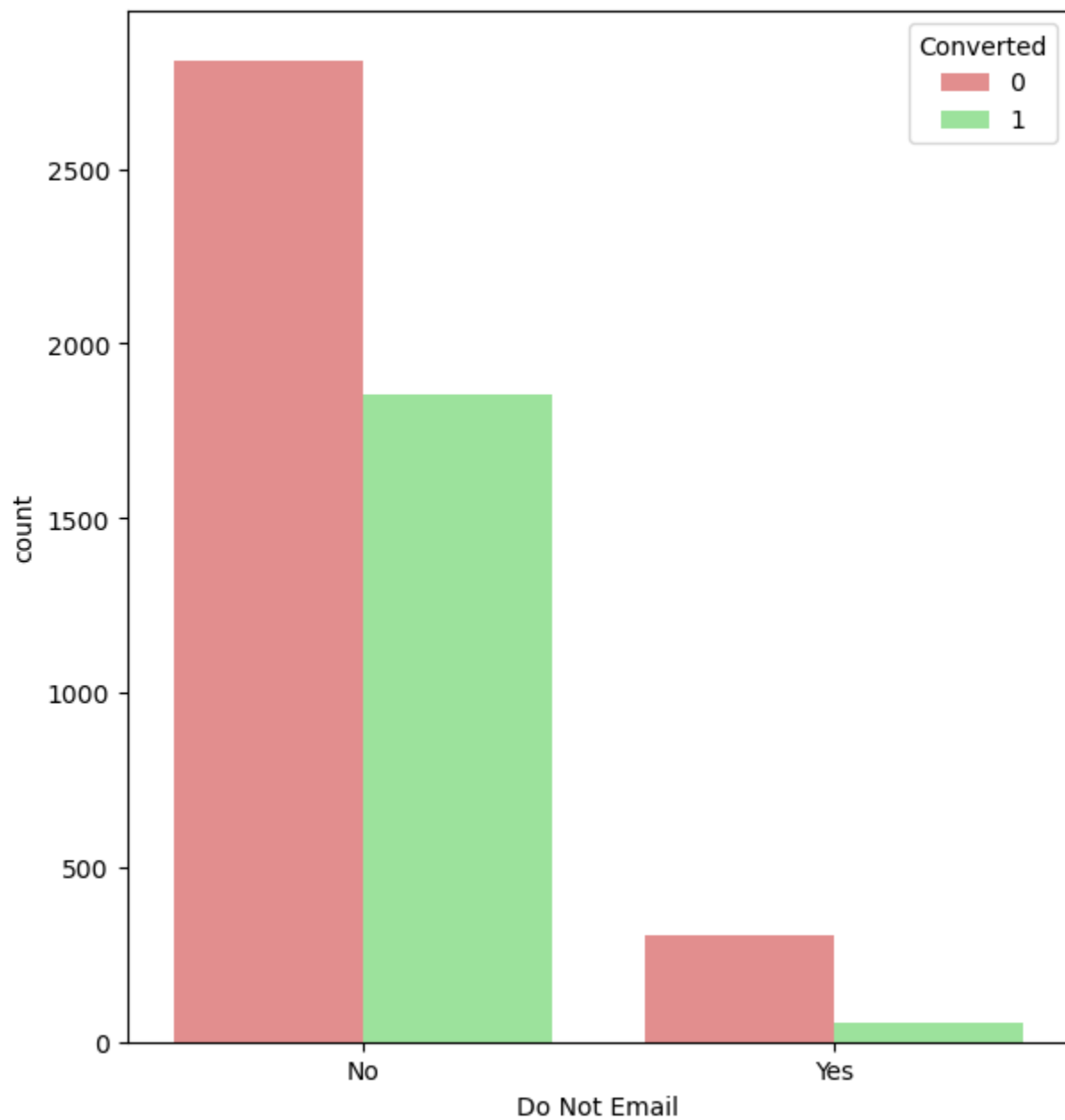
EDA

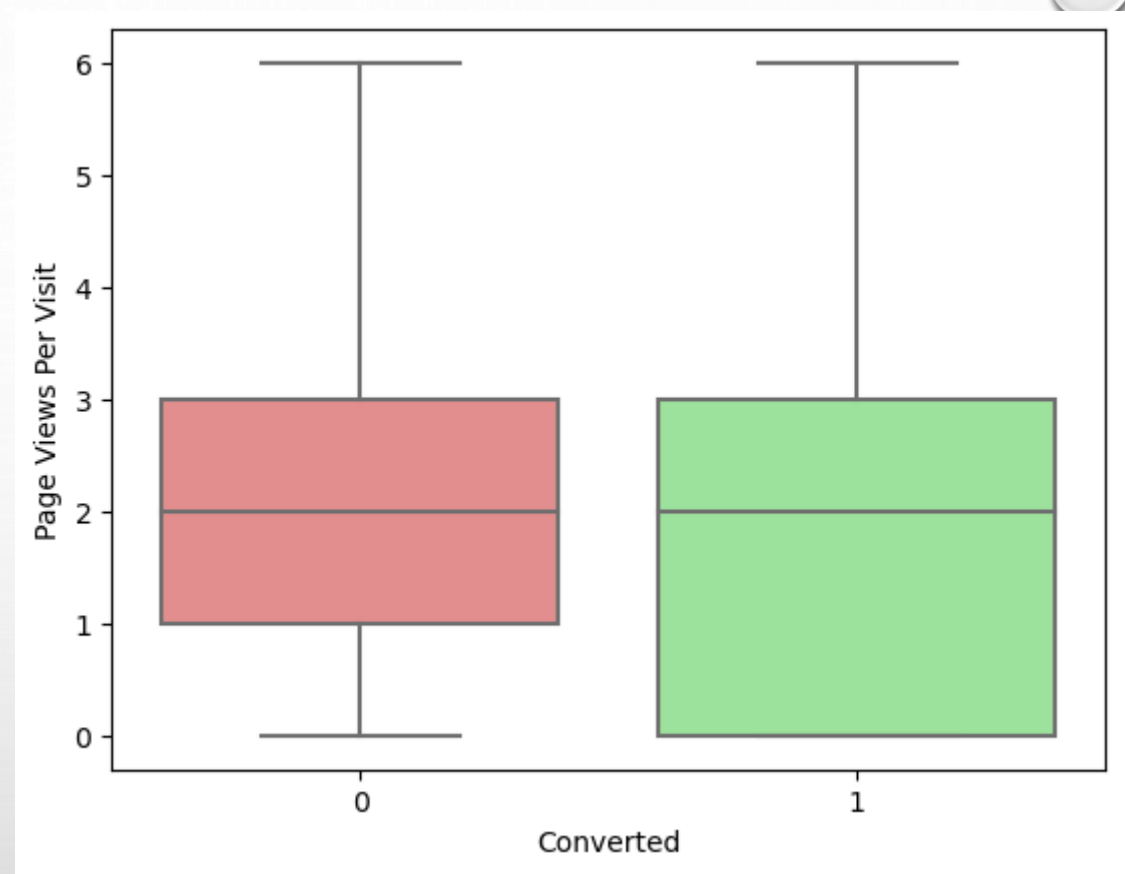
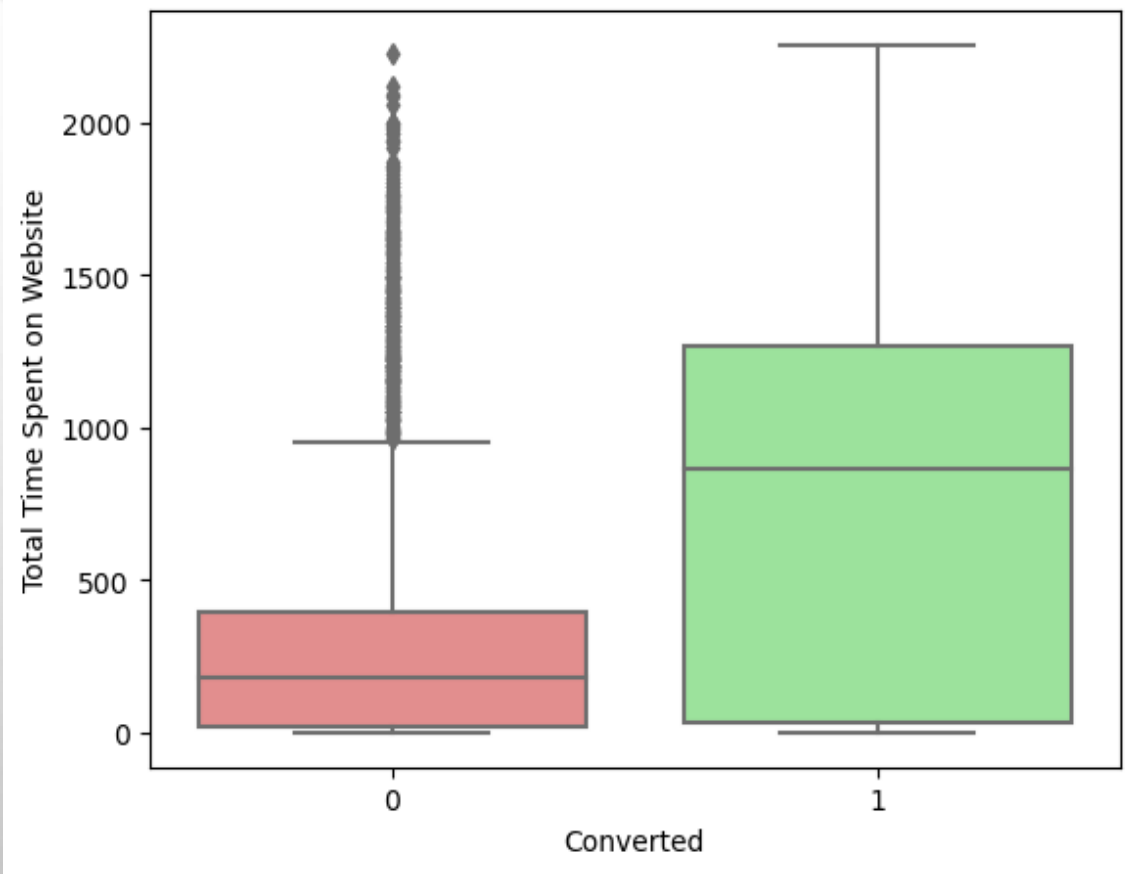




CATEGORICAL VARIABLE RELATION







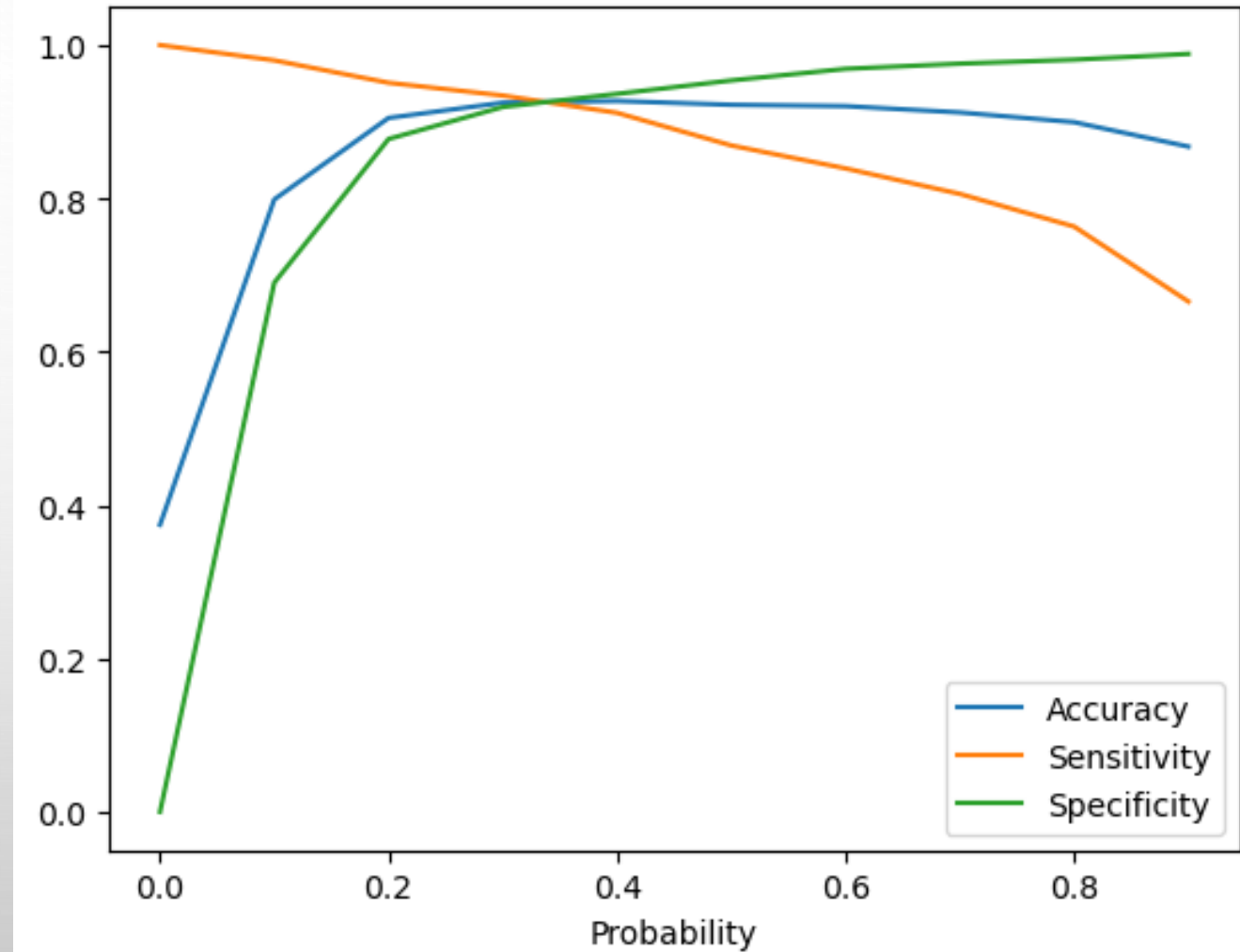
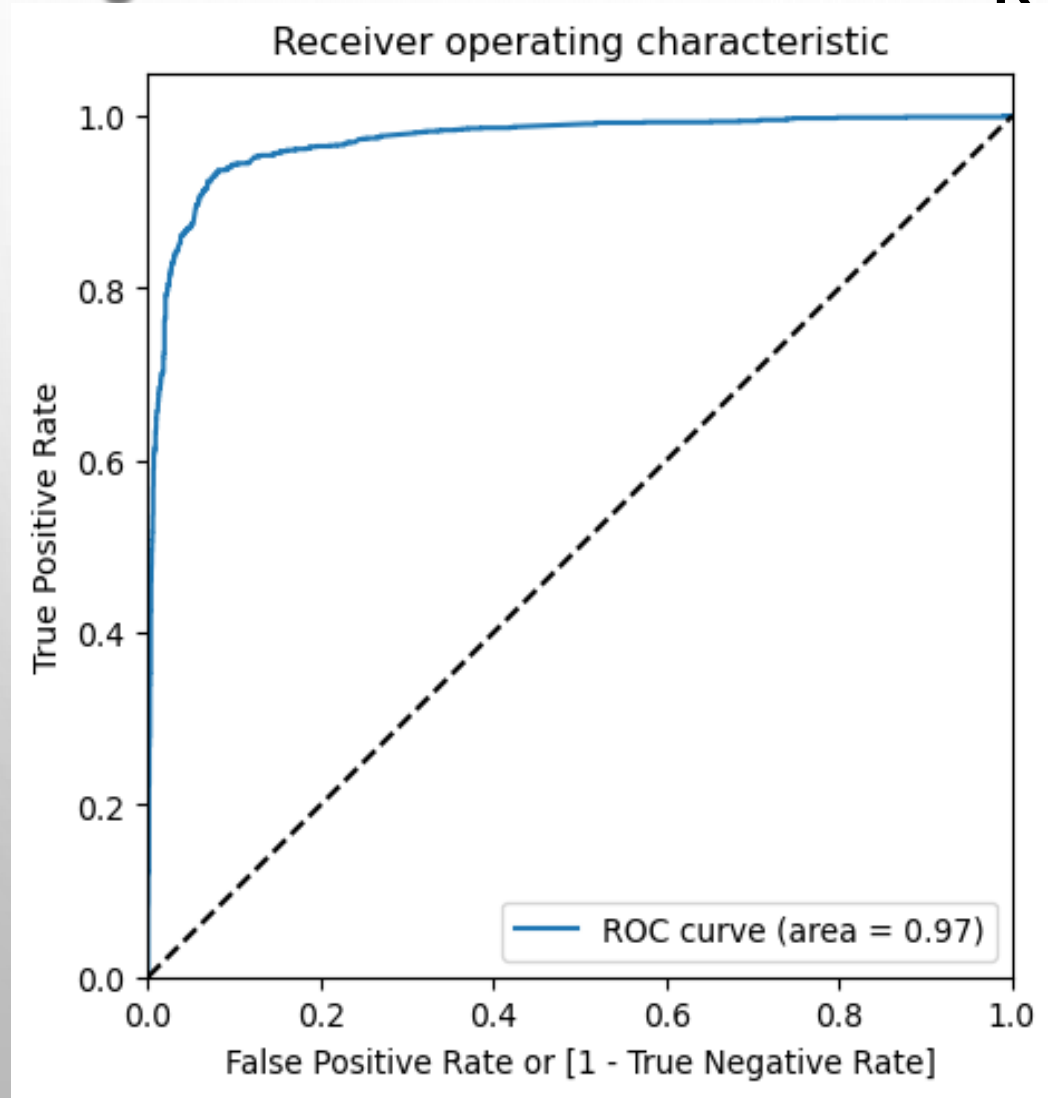
DATA CONVERSION

- Numerical variables are normalised using `StandardScaler()`.
- Dummy variables are created for object type variables.
- Total rows for analysis: 5022.
- Total columns for analysis: 86.

MODEL BUILDING

- SPLITTING THE DATA INTO TRAINING AND TESTING SETS
- THE FIRST BASIC STEP FOR REGRESSION IS PERFORMING A TRAIN-TEST SPLIT, WE HAVE CHOSEN 80:20 RATIO.
- USE RFE FOR FEATURE SELECTION
- RUNNING RFE WITH 20 VARIABLES AS OUTPUT
- BUILDING MODEL BY REMOVING THE VARIABLE WHOSE P- VALUE IS GREATER THAN 0.05 AND VIF VALUE IS GREATER THAN 5
- PREDICTIONS ON TEST DATA SET
- OVERALL ACCURACY 92%

ROC CURVE



CONCLUSION

THE KEY FACTORS INFLUENCING POTENTIAL BUYERS, IN DESCENDING ORDER OF IMPORTANCE, ARE:

- TOTAL TIME SPENT ON THE WEBSITE
- TOTAL NUMBER OF VISITS
- LEAD SOURCE, SPECIFICALLY:
 - GOOGLE
 - DIRECT TRAFFIC
 - ORGANIC SEARCH
 - WELINGAK WEBSITE
- LAST ACTIVITY, SPECIFICALLY:
 - SMS
 - OLARK CHAT CONVERSATION
- LEAD ORIGIN BEING LEAD ADD FORMAT
- CURRENT OCCUPATION AS A WORKING PROFESSIONAL BY FOCUSING ON THESE FACTORS, X EDUCATION CAN SIGNIFICANTLY INCREASE THEIR CHANCES OF CONVERTING POTENTIAL BUYERS INTO CUSTOMERS, THEREBY BOOSTING THEIR COURSE SALES.