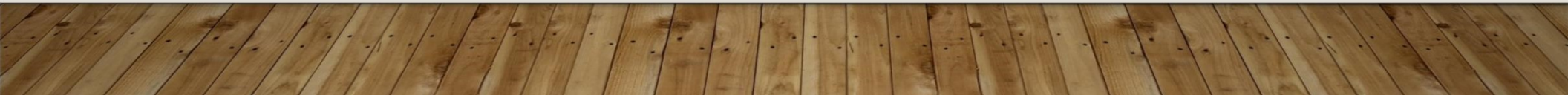


X EDUCATION - LEAD SCORING CASE STUDY

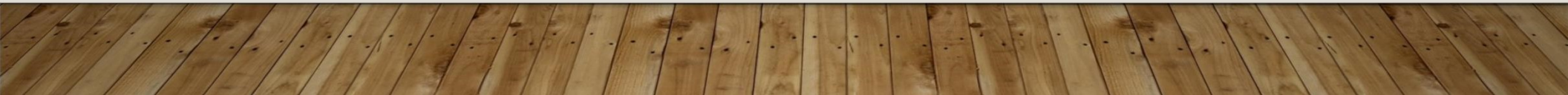
PRATIKSHA SHIRWALE, PRATIK
KALE & KARTHIKEYAN SANKAR



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.

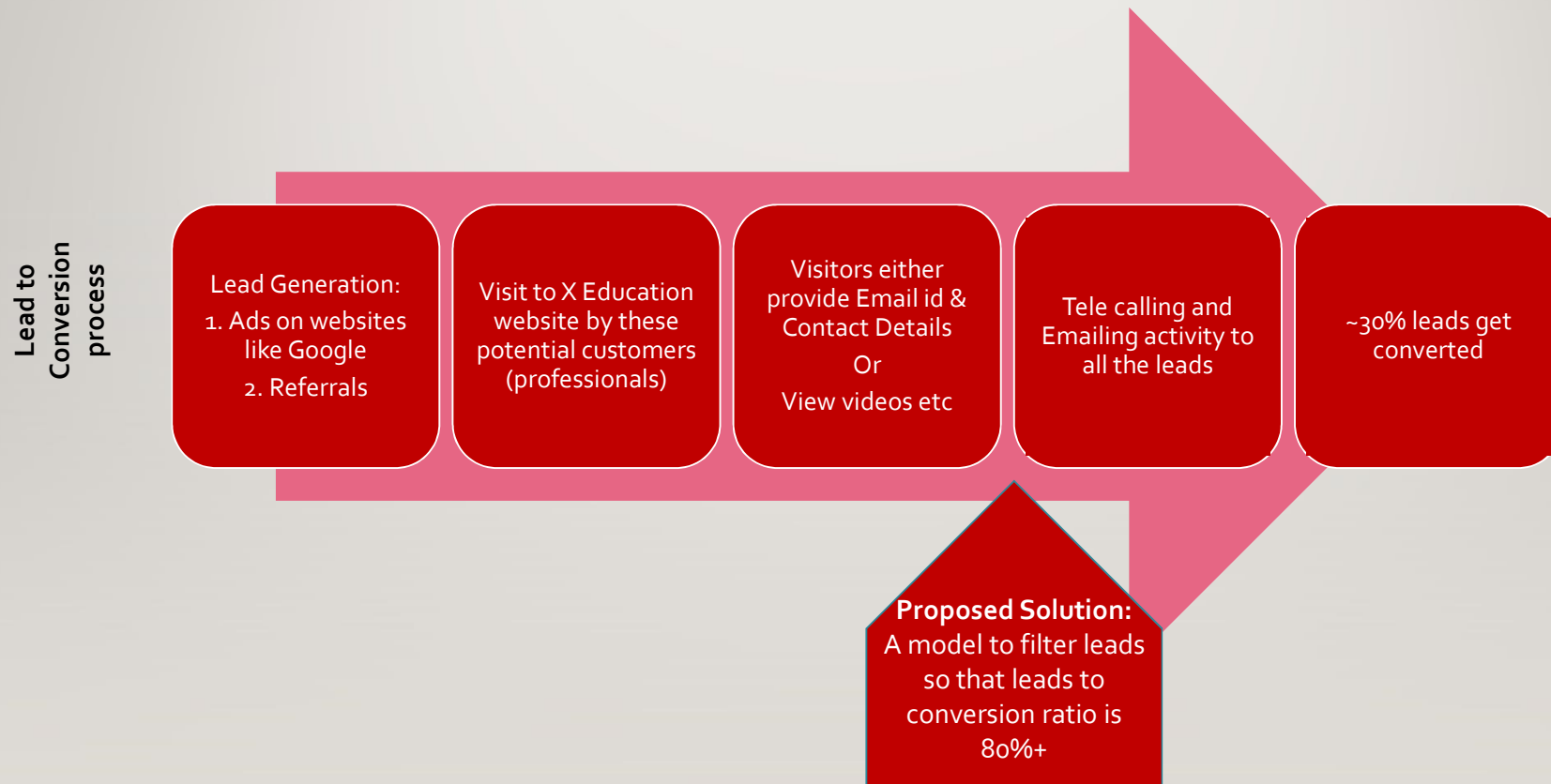
The company markets its courses on several websites and search engines like Google. Once these people land on the website, they might browse the courses or fill up a form for the course or watch some videos. When these people fill up a form providing their email address or phone number, they are classified 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. The typical lead conversion rate at X education is around 30%. Now, although X Education gets a lot of leads, its lead conversion rate is very poor. For example, if, say, they acquire 100 leads in a day, only about 30 of them are converted. 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. There are a lot of leads generated in the initial stage (top) but only a few of them come out as paying customers from the bottom. In the middle stage, you need to nurture the potential leads well (i.e. educating the leads about the product, constantly communicating etc.) in order to get a higher lead conversion.



GOALS OF THE CASE STUDY

1. Build a logistic regression model to assign a lead score between 0 and 100 to each of the leads which can be used by the company to target potential leads. A higher score would mean that the lead is hot, i.e. is most likely to convert whereas a lower score would mean that the lead is cold and will mostly not get converted.
2. There are some more problems presented by the company which your model should be able to adjust to if the company's requirement changes in the future so you will need to handle these as well. These problems are provided in a separate doc file. Please fill it based on the logistic regression model you got in the first step. Also, make sure you include this in your final PPT where you'll make recommendations.

Lead – Conversion Process



PROPOSED SOLUTION

Selection of Hot Leads

Leads Clustering

We cluster the leads into certain categories based on their tendency or probability to convert, thus, getting a smaller section of hot leads to focus more on.

Communicating with Hot Leads

Focus Communication

Since we would have a smaller set of leads to have communication with, we might make more impact with effective communication.

Conversion of Hot Leads

Increase conversion

Since we focussed on hot leads, which were more probable to convert, we would have a better conversion rate, and hence we can achieve the 80% target.

SOLUTION

For our Problem Solution, the crucial part is to accurately identify hot leads.

The more accurate we obtain the hot lead, the more chance we get of higher conversion ratio.

Since we have a target of 80% conversion rate, we would want to obtain a high accuracy in obtaining hot leads.

IMPLEMENTATION

Loading & Observing the
past data provided by
the Company

Univariate, Bivariate, and
Heatmap for numerical
and categorical columns

Performing pre-requisites
for RFE and Logistic
Regression

Data Gathering

Data
Cleaning

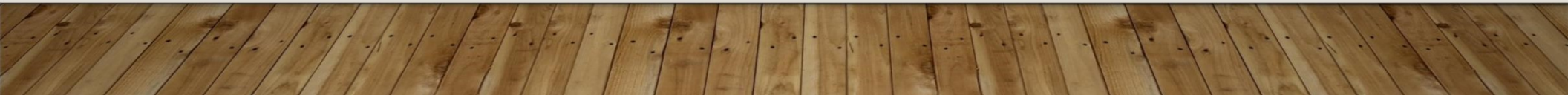
Performing
EDA

Data
Preparation

Model
Building

Duplicate removal, null value
treatment, unnecessary column
elimination, etc.

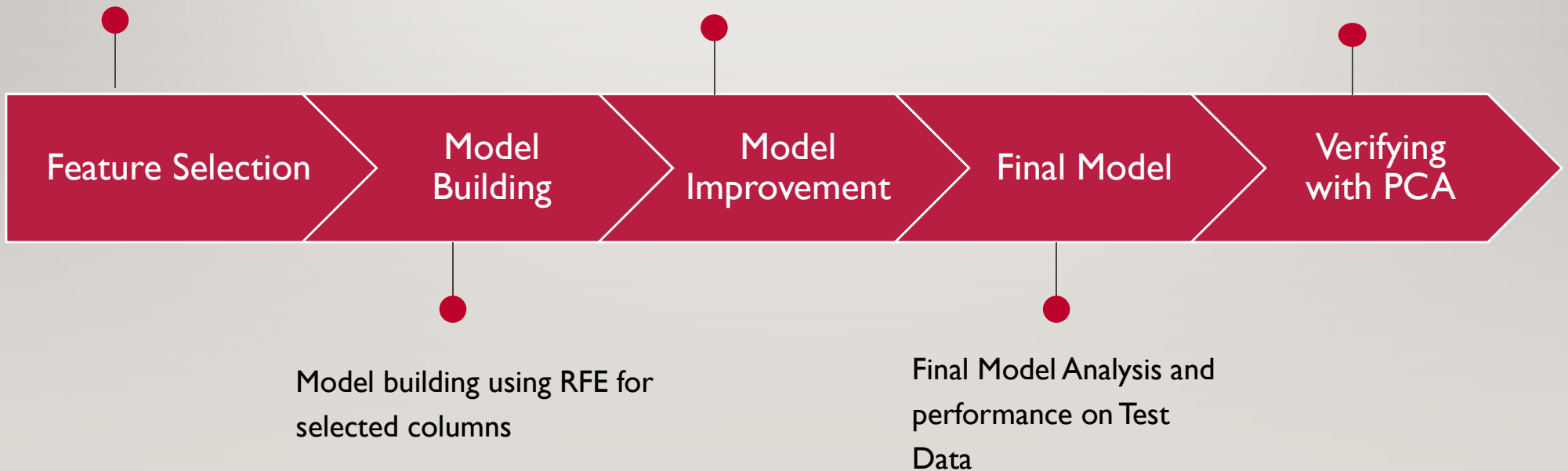
Outlier Treatment,
Feature-Standardization



Selection of top 25
features using RFE

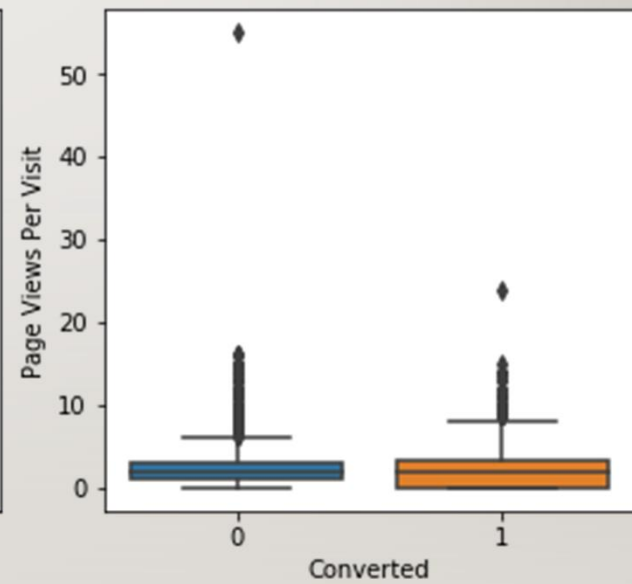
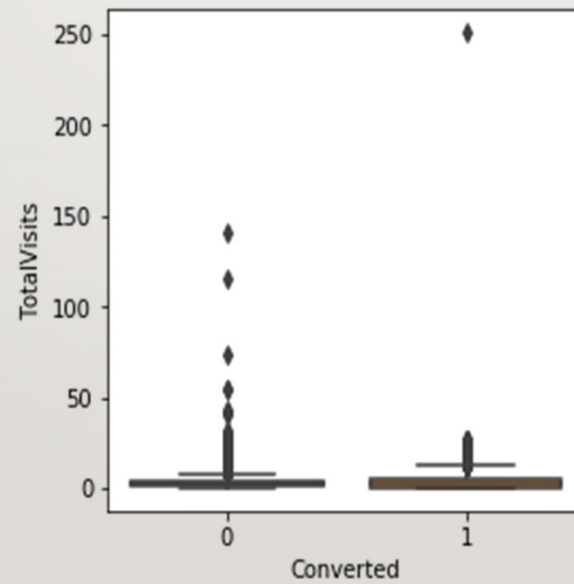
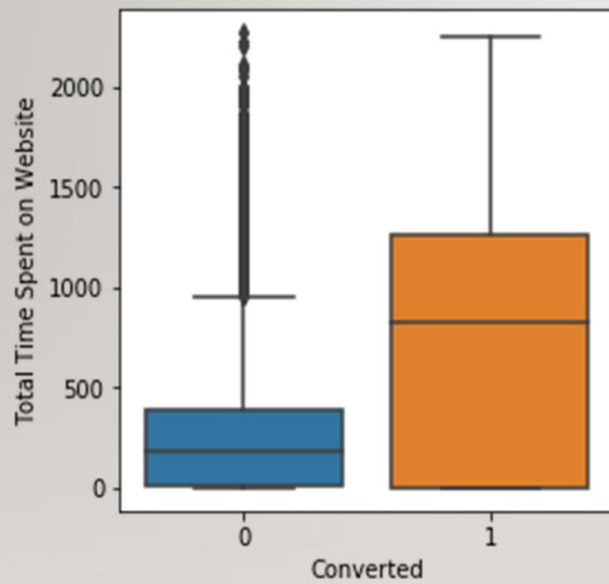
Reduction of columns
and Model re-building

Verifying our Final Model Accuracy
etc. with model built with PCA

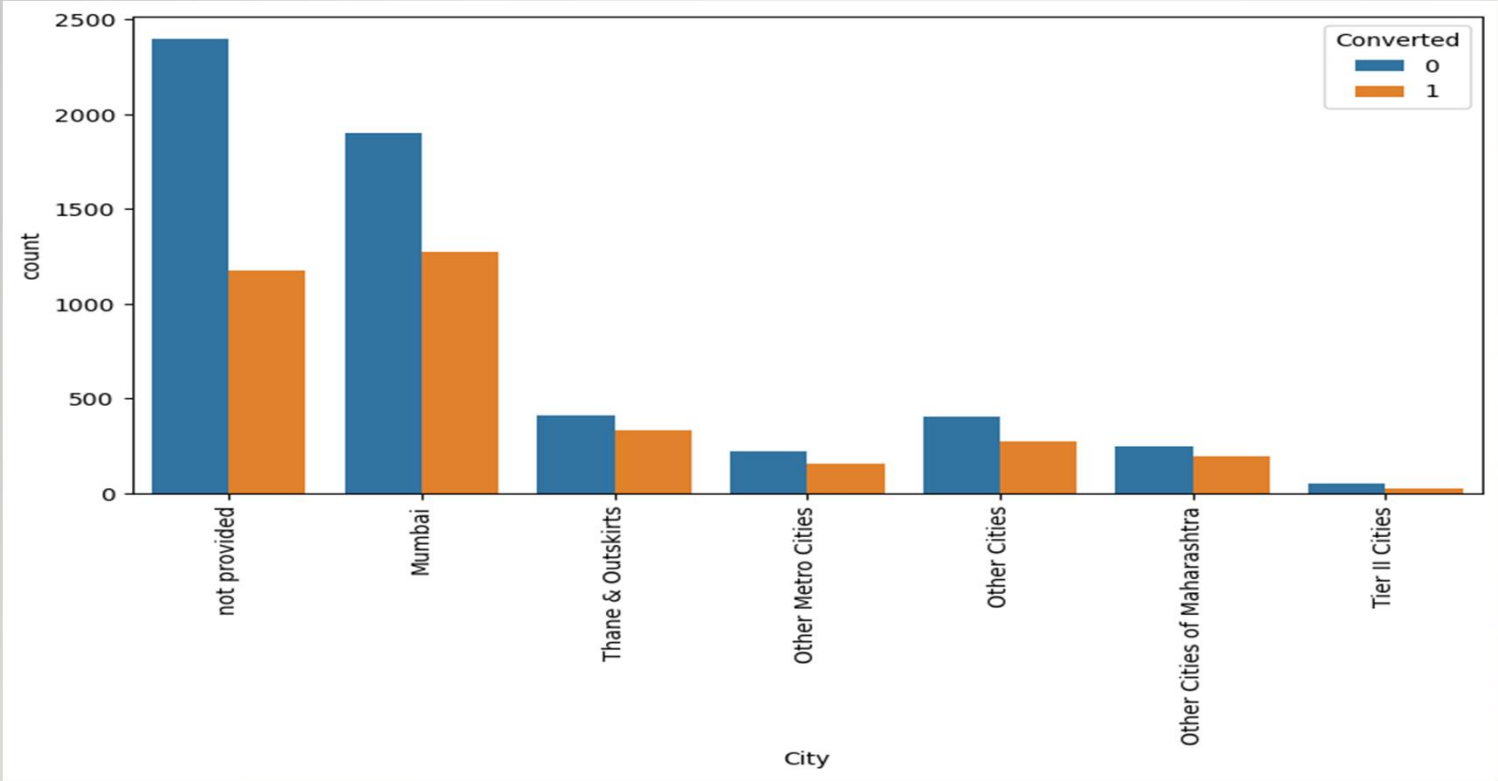


PLOTS

EDA plots depicting variation in numerical columns for those who Converted and those who didn't.



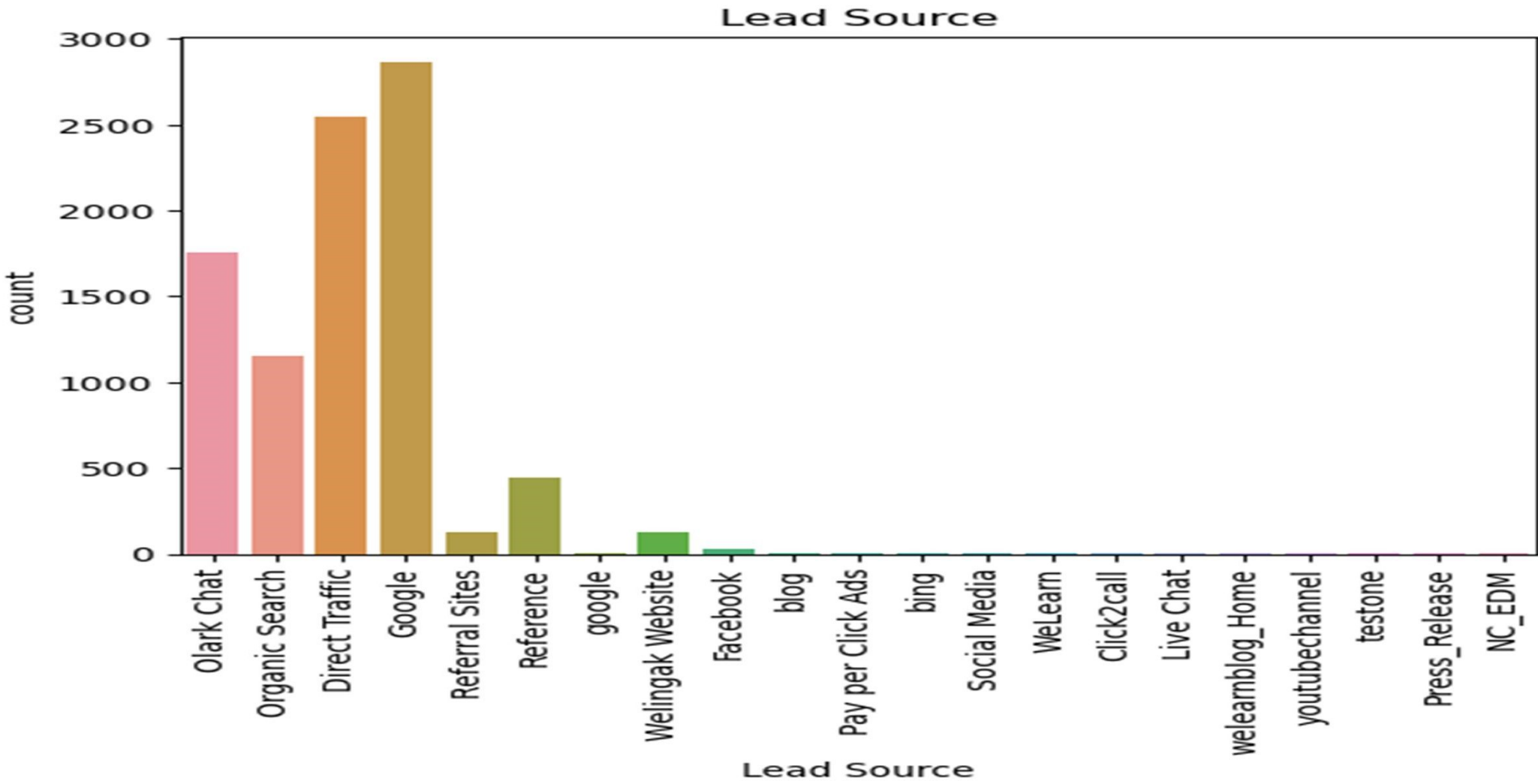
Spread of Cities

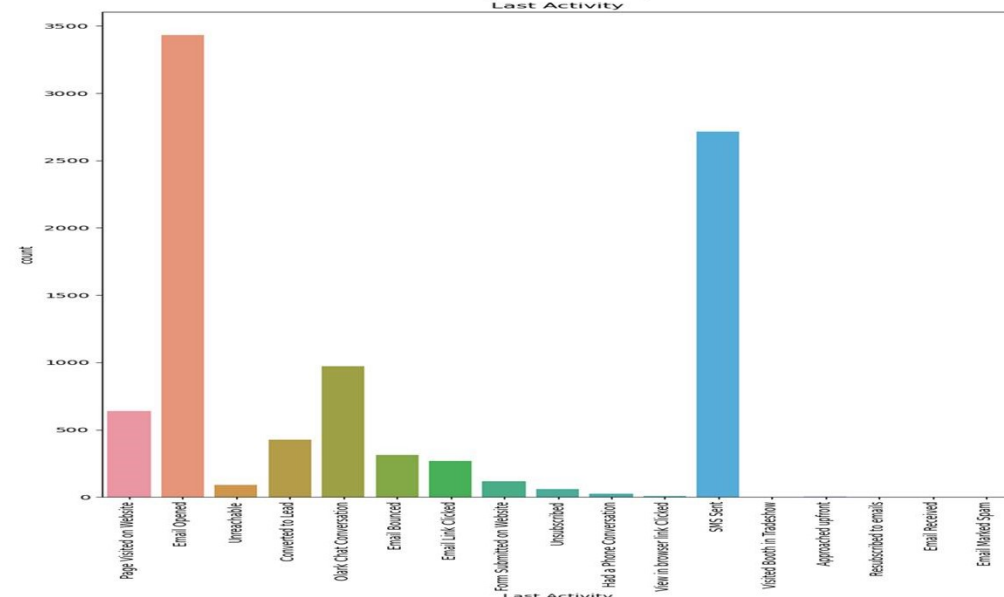
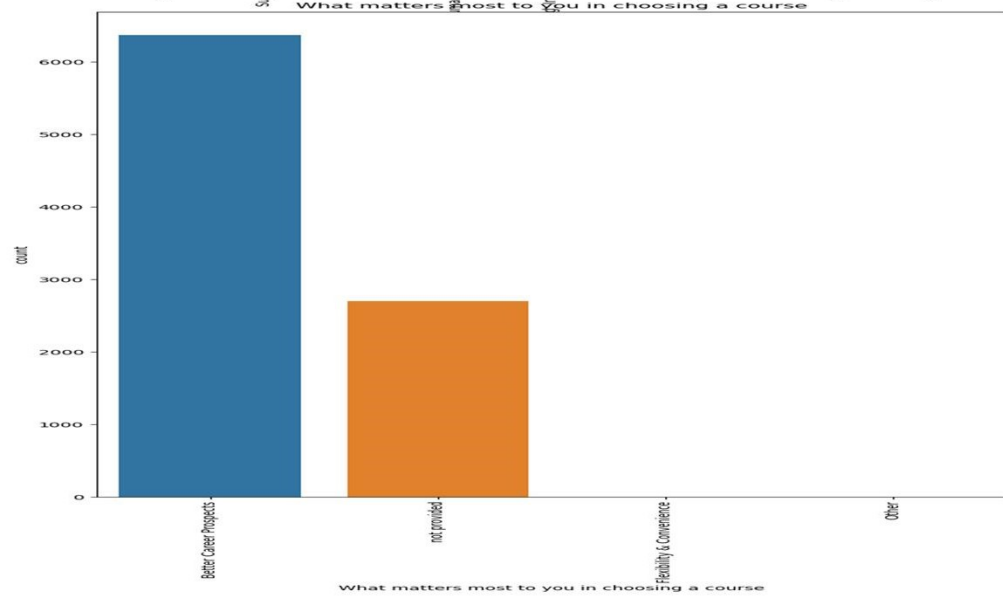
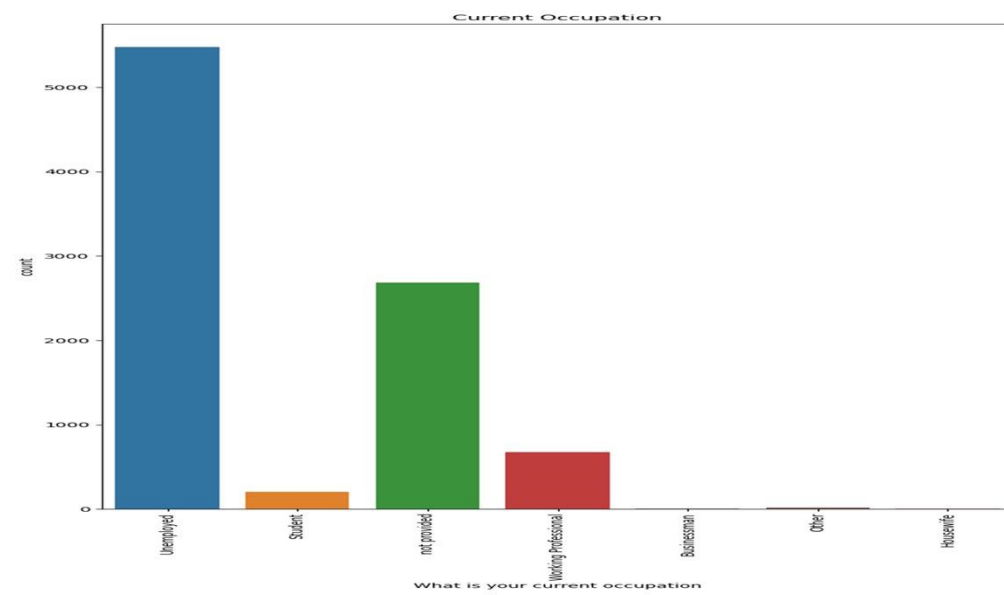
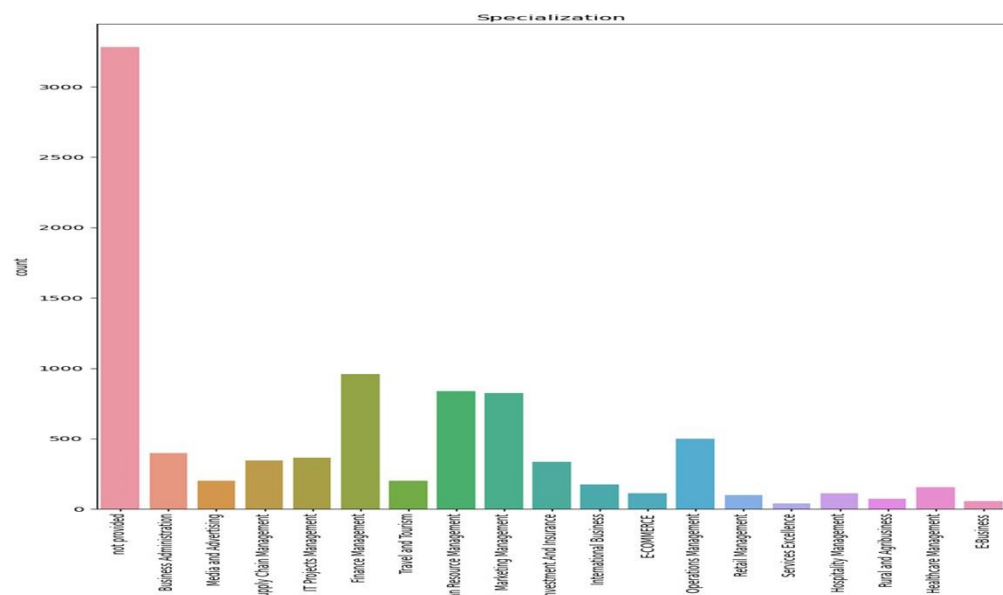


Observations

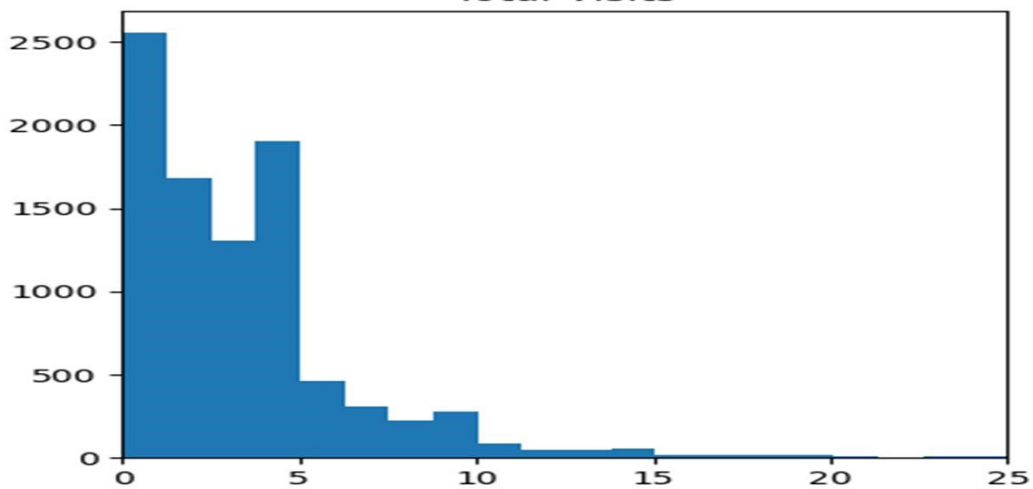


Lead Source

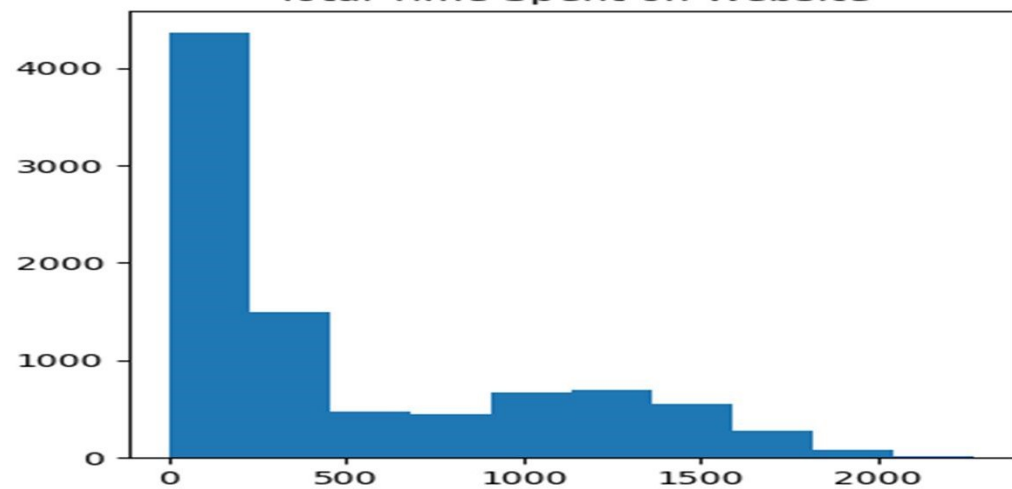




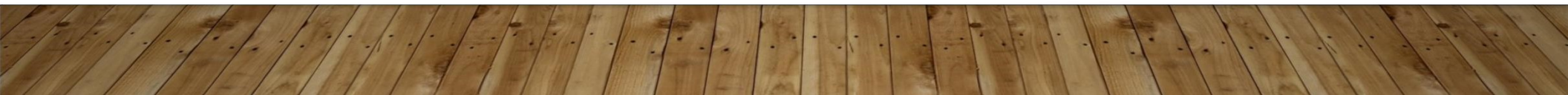
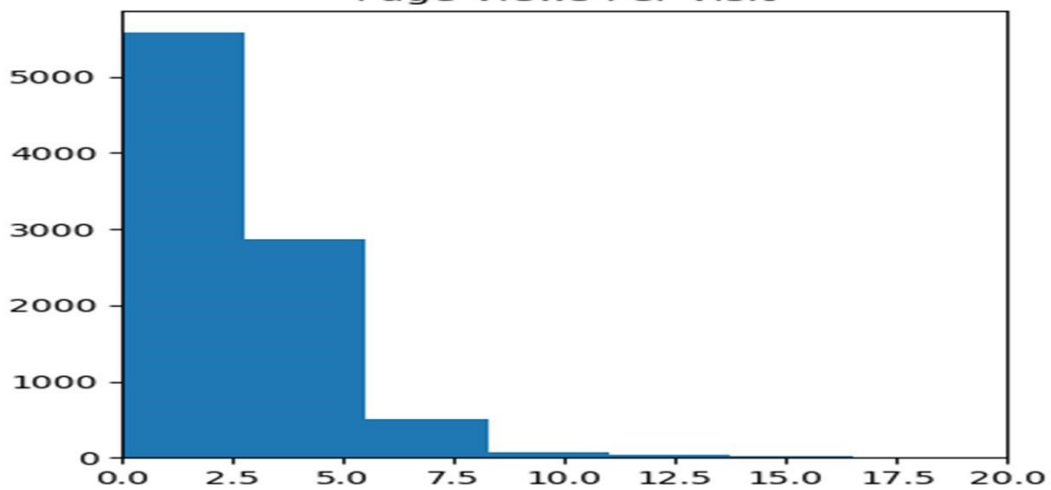
Total Visits

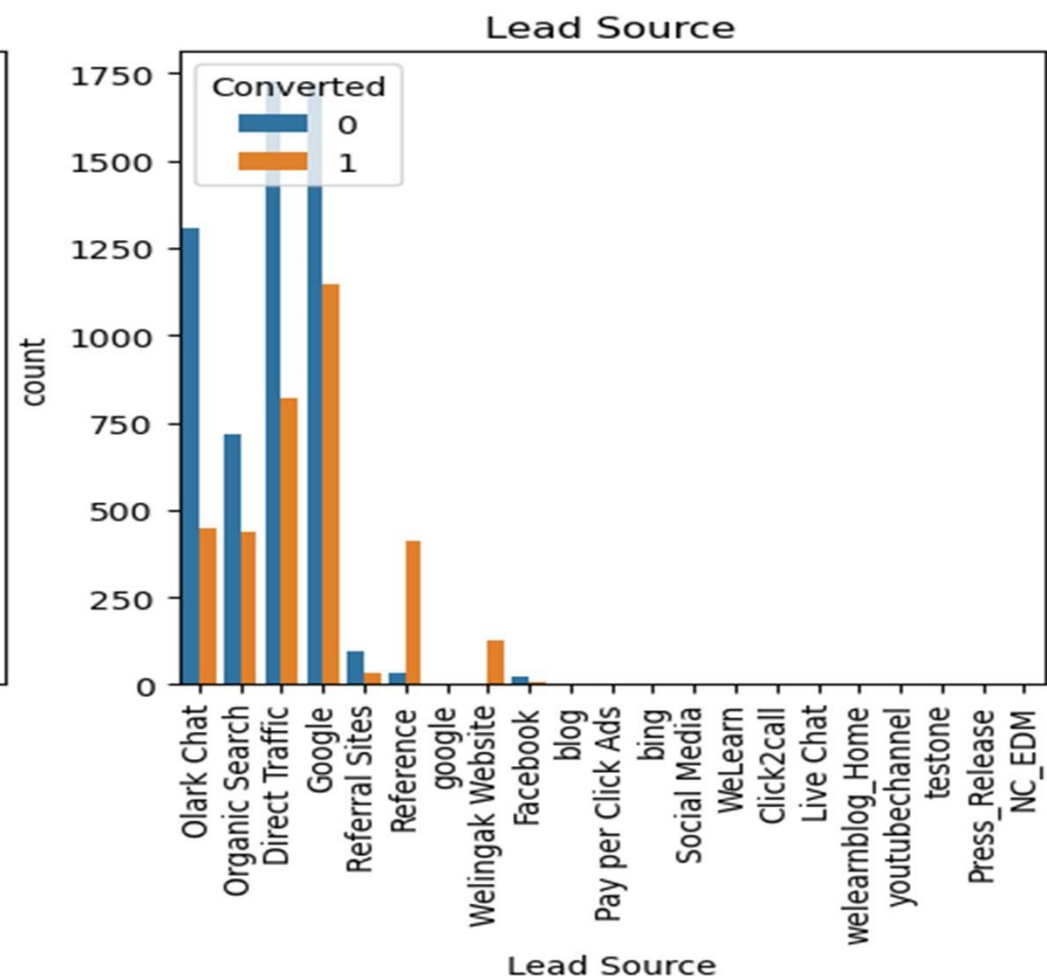
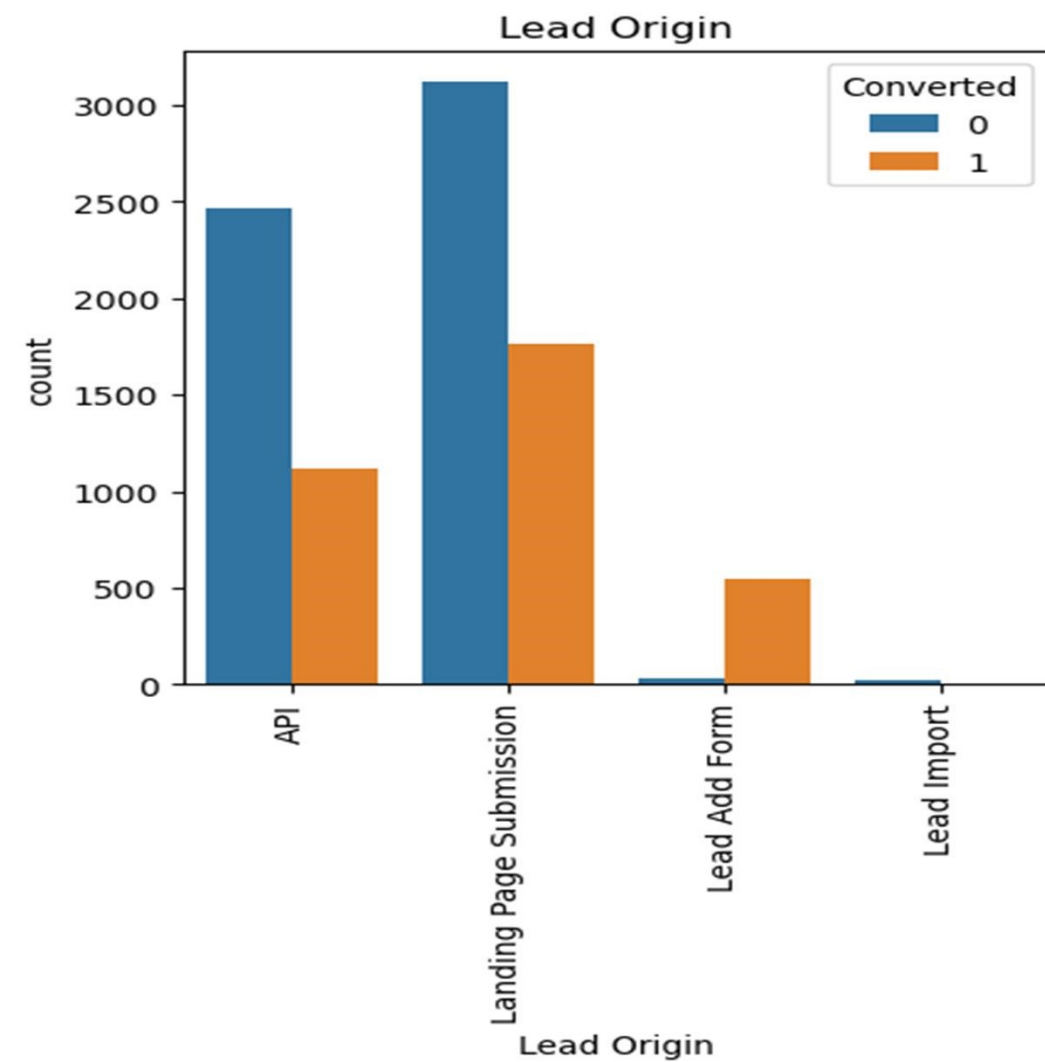


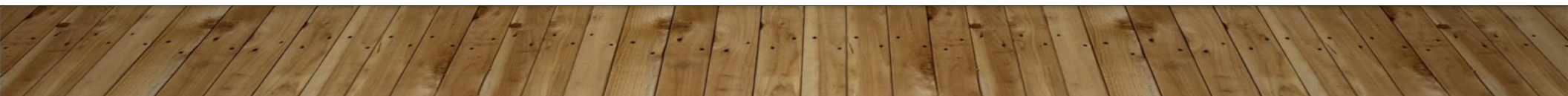
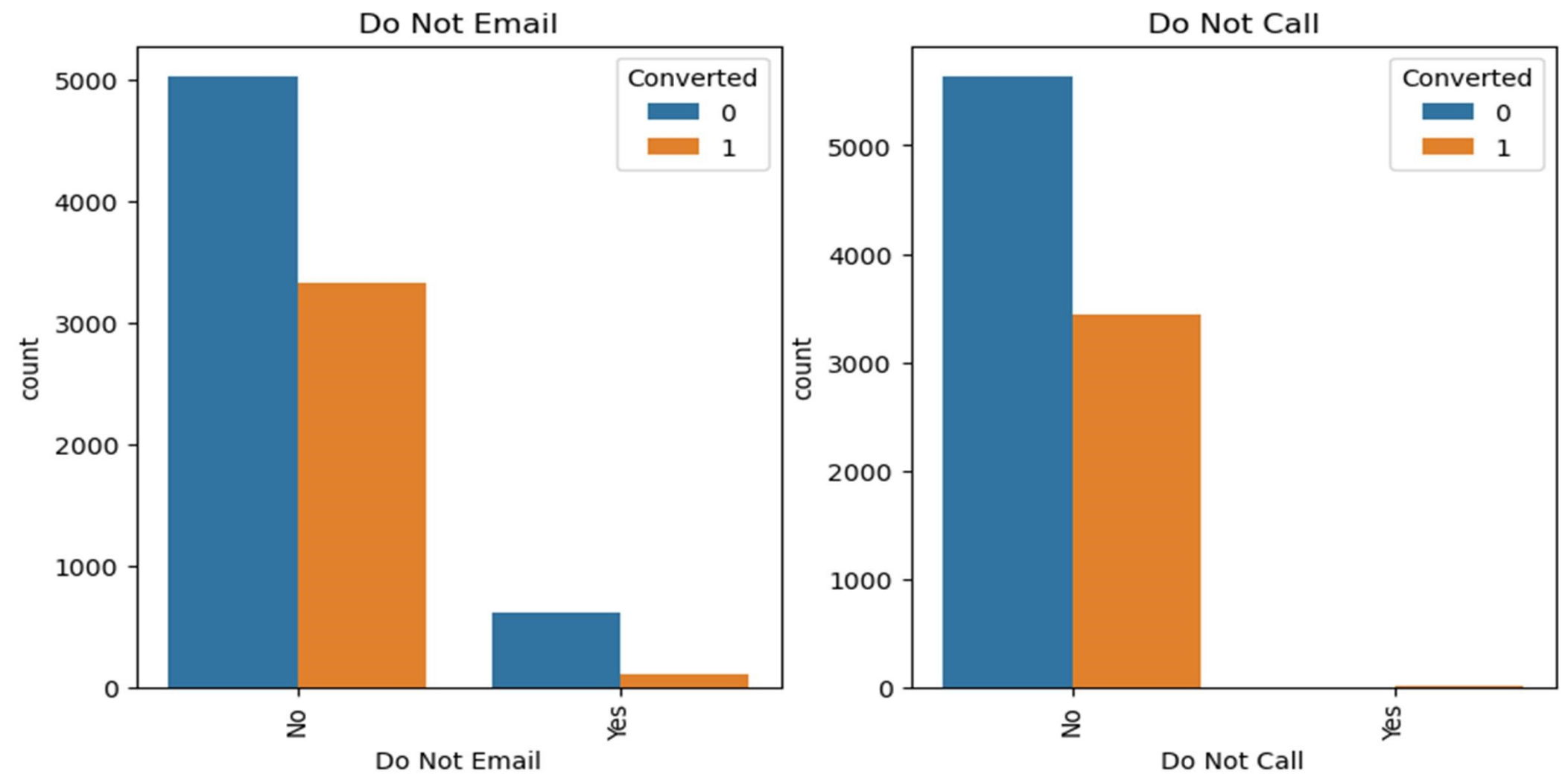
Total Time Spent on Website

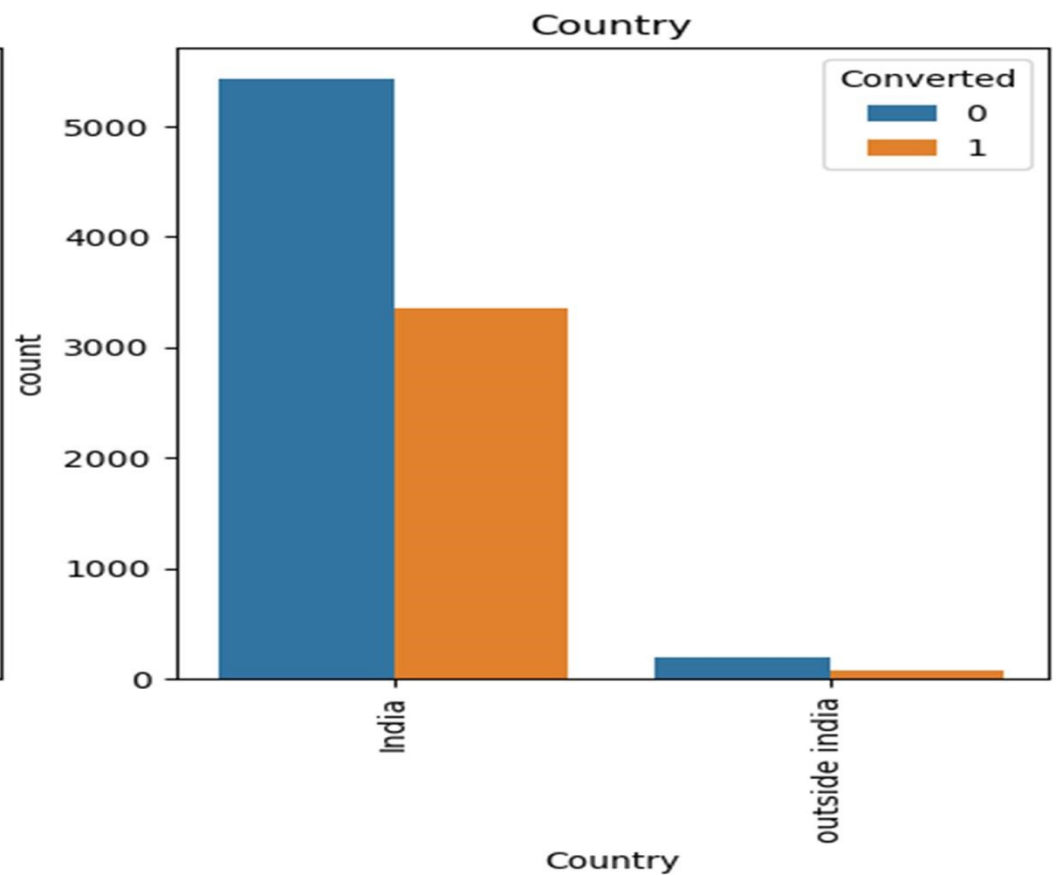
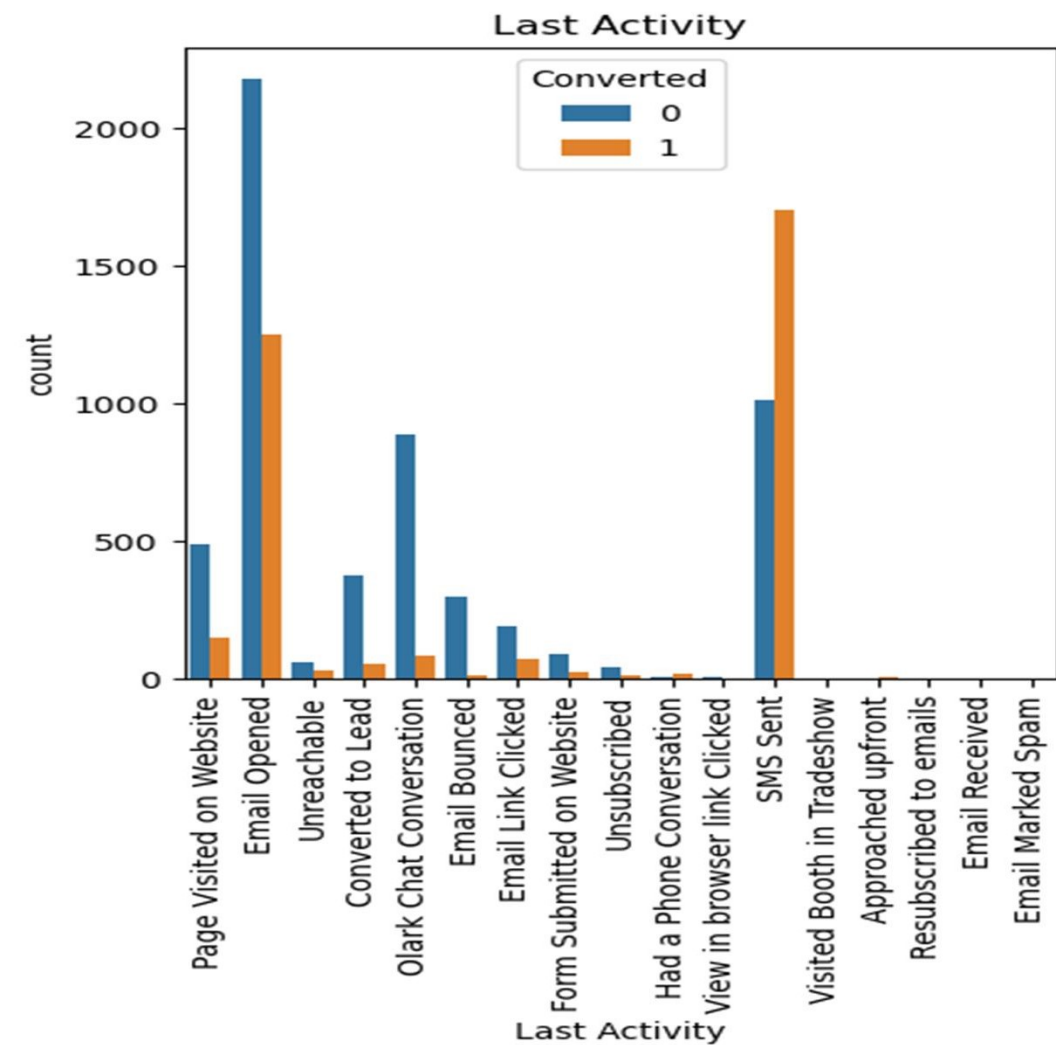


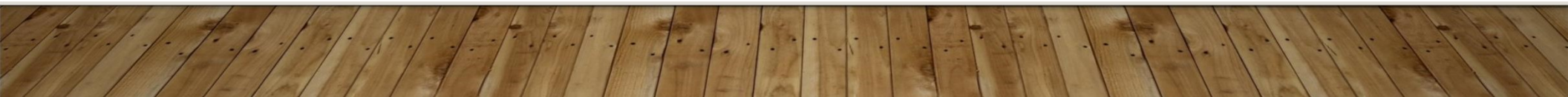
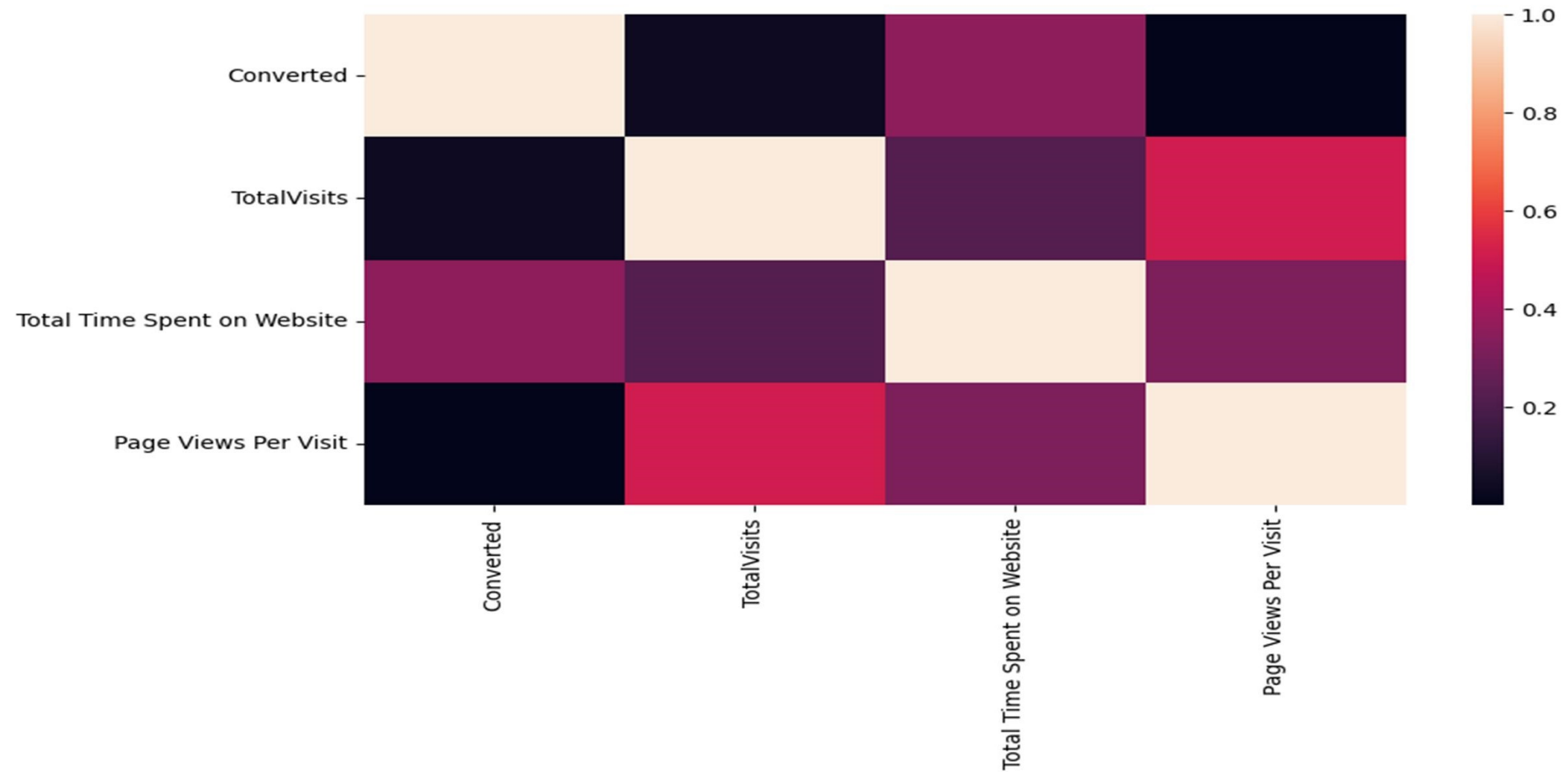
Page Views Per Visit

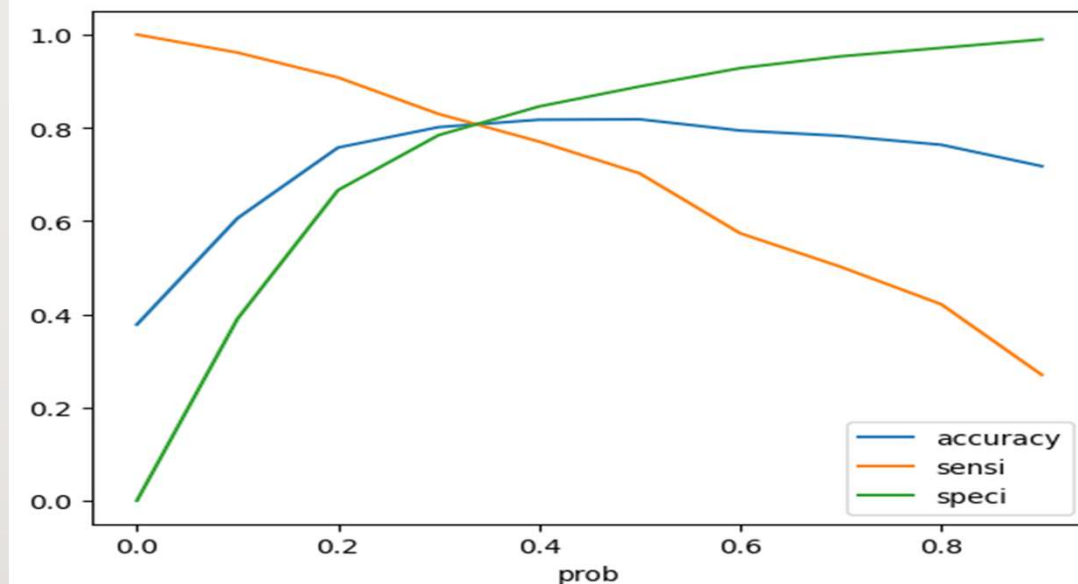
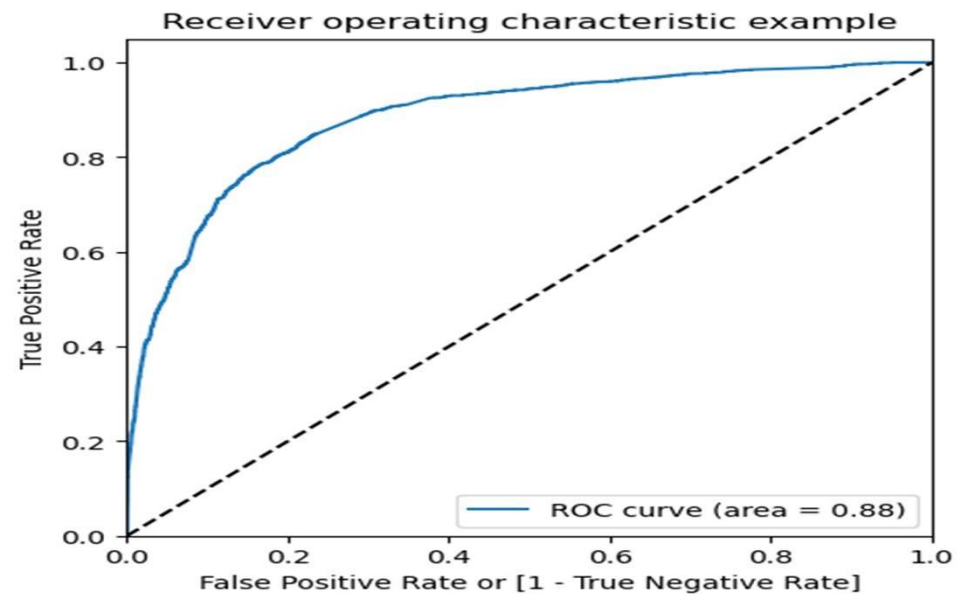




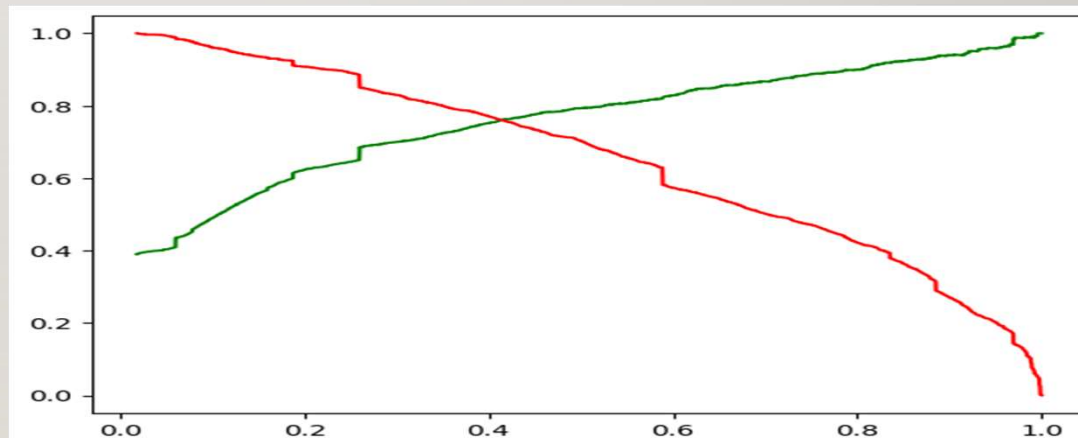








Linear Regression Final Model Parameters
Area under ROC = 0.88
Intermediate cut-off = 0.35
Final cut-off = 0.44



CONCLUSION

TOP VARIABLE CONTRIBUTING TO CONVERSION:

- LEAD SOURCE:
 - Total Visits
 - Total Time Spent on Website
- Lead Origin:
 - Lead Add Form
- Lead source:
 - Direct traffic
 - Google
 - Welingak website
 - Organic search
 - Referral Sites

Last Activity:

- Do Not Email_Yes
- Last Activity_Email Bounced
- Olark chat conversation

The Model seems to predict the Conversion Rate very well and we should be able to give the Company confidence in making good calls based on this model.

