# LEAD SCORE CASE STUDY

# BY

## SHUBHAM SINGH RANA
## ROOPA H

# Problem Statement

- X Education sells online courses to industry professionals.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%.

# Business Goal

- X Education needs help in selecting the most promising leads, i.e. the leads that are most likely to convert into paying customers.

- The company needs a model wherein you a lead score is assigned to each of the leads such that the customers with higher lead score have a higher conversion chance and the customers with lower lead score have a lower conversion chance.

- The CEO, in particular, has given a ballpark of the target lead conversion rate to be around 80%.

# Approaches

- Reading and understanding the dataset
- Checking the structures like shape, info and describe data
- Performing EDA
- Missing value check and Outlier check
- Data Pre-Processing, Splitting, Feature Scaling, Handling Multicollinearity within independent Features via Correlation
- Feature Selection and Model Building
- Here we are using ML model of Logistic Regression as it falls under Binary classification
- Optimal Probability Selection here we are using Precision – Recall, Trade-Off and Accuracy, Sensitivity & Specificity Trade-Off
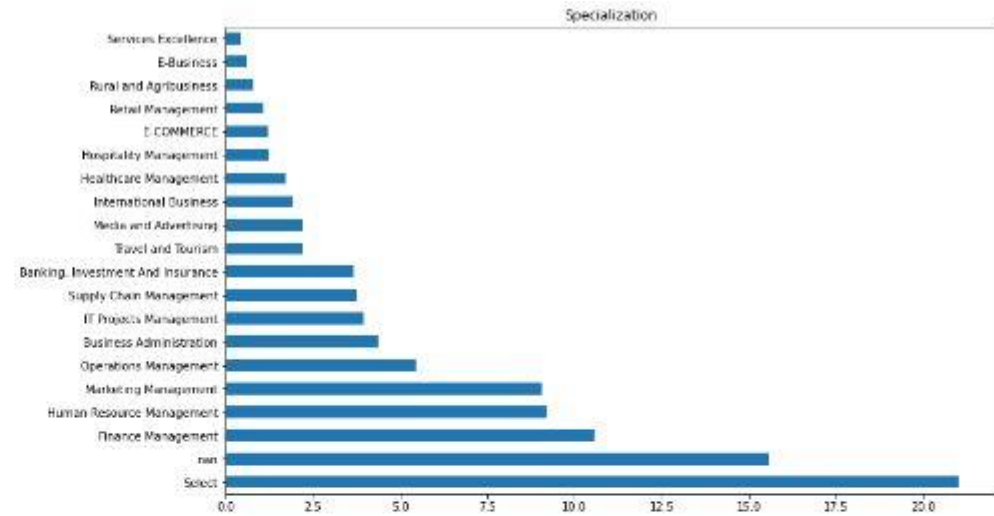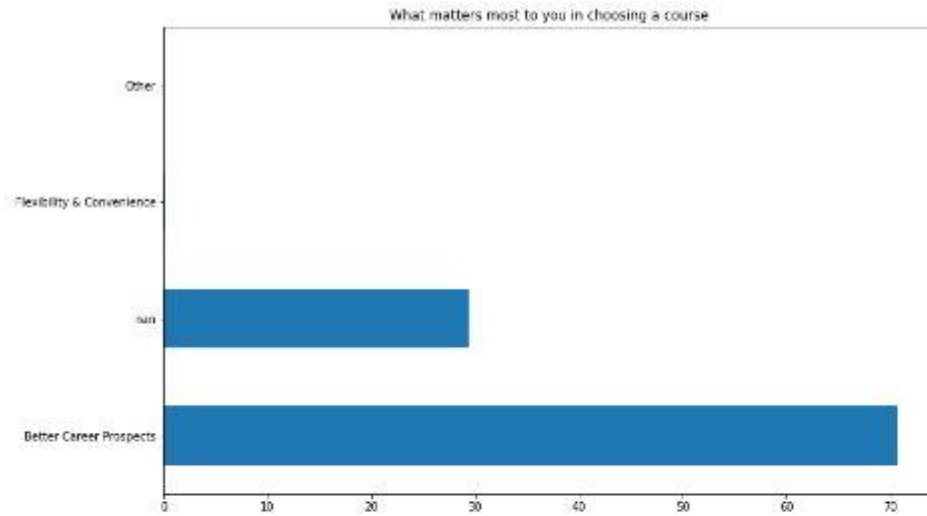
# Analysis

- We have found 9240 rows and 37 columns after reading the dataset
- We got many insights regarding features with Missing Values, Categorical features, Numerical features , Binary Features of Categorical & Numeric and their Data Imbalance after performing EDA.
- After EDA, we got to know that 10 Features have Missing Values more than 35% & 5 features have Missing Values less than 35%. For Features with less than 35% Missing Values, we have imputed their values accordingly.
- We dropped features on the basis of `More than 35% of Missing Values`, `High Multicollinearity` & `High Data Imbalance for Binary Categorical Features`.
- 7 Categorical Features had few categories which were having data less than 5%. Therefore, we merged them into a single category named `Other`.
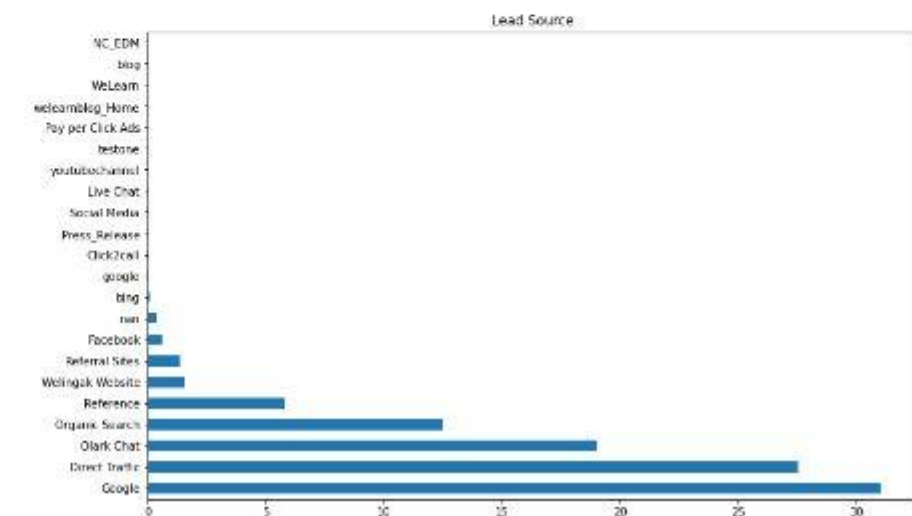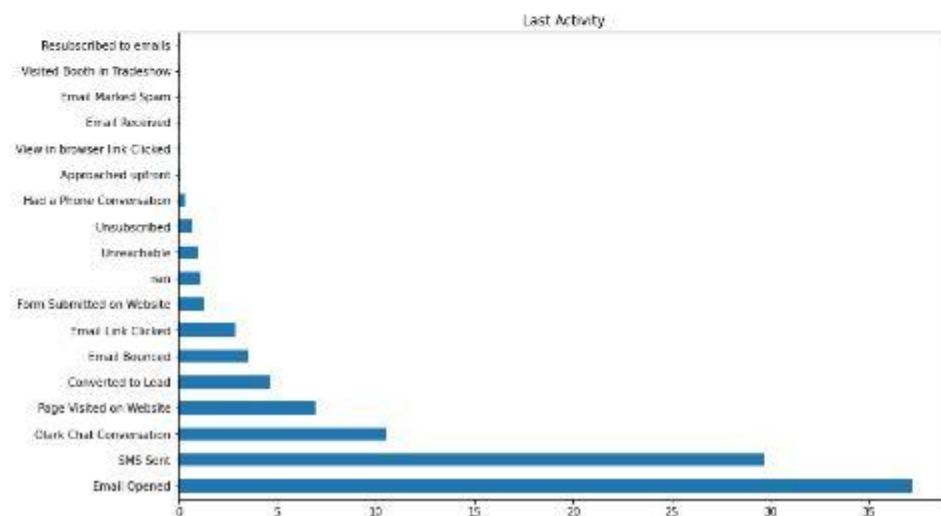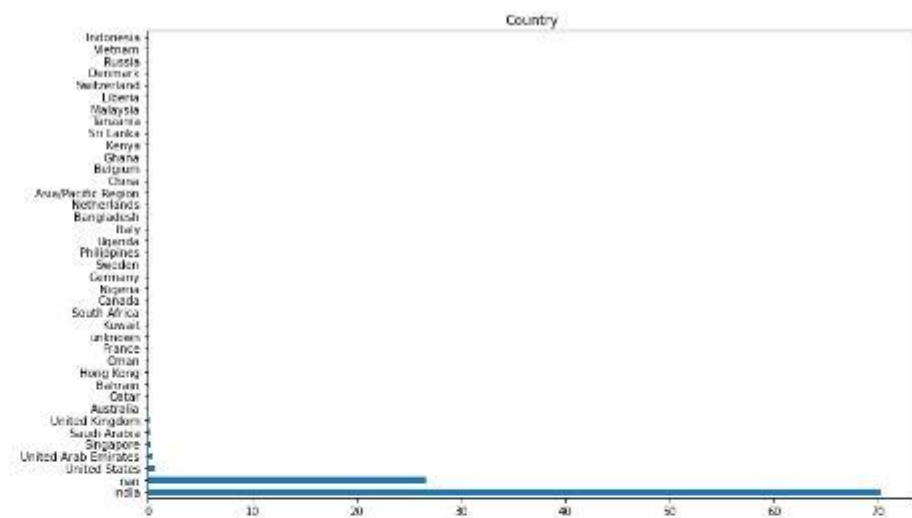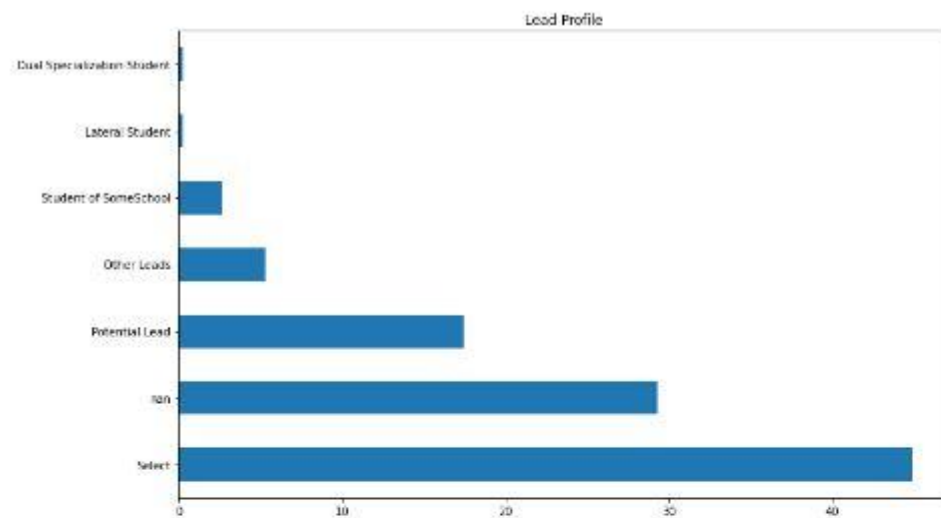
# Analysis

- 5 Binary Categorical Features had Very High Data Imbalance, therefore we dropped those features.

- High Multicollinearity was found in 4 features. Therefore, we dropped them as well.

- And after dropping we are left with 19 features in the dataset for modelling.

- While checking for the outliers we observed that there are outliers in Page Views Per Visit & Total Visits. Therefore, we dropped data which belonged to more than 99 Percentile of these features.
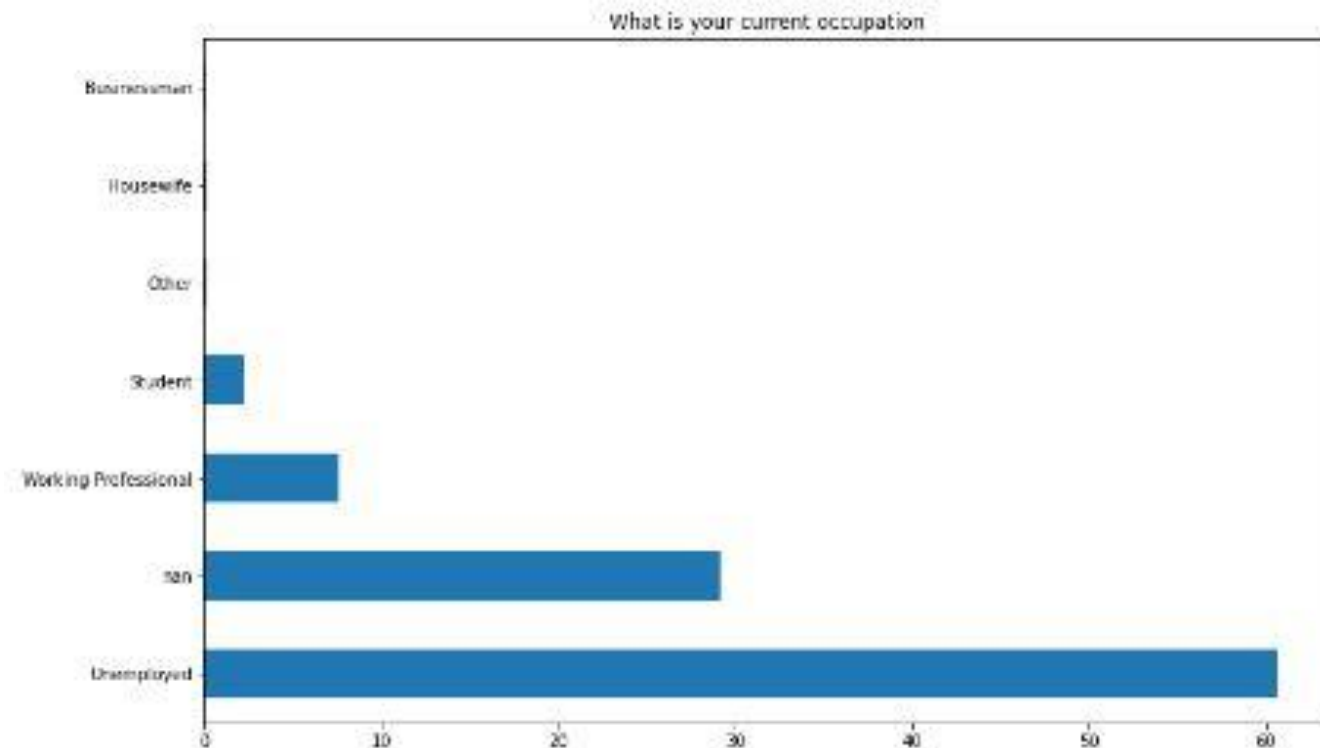
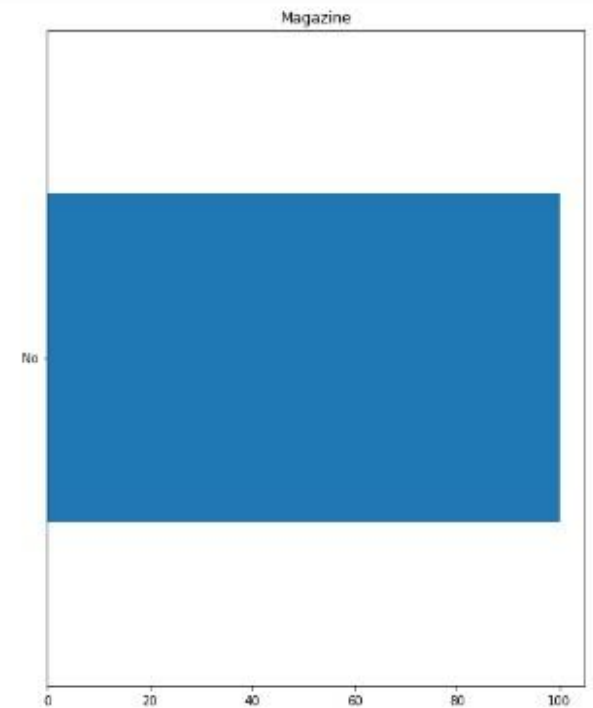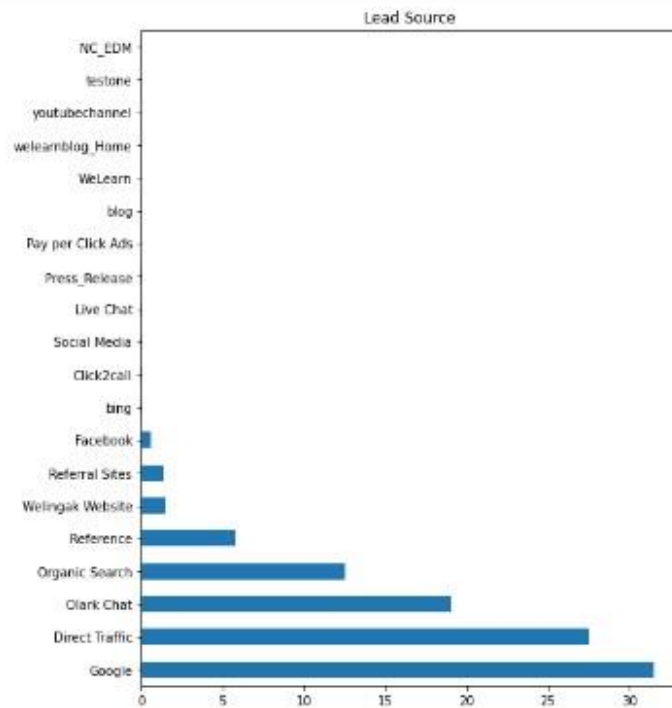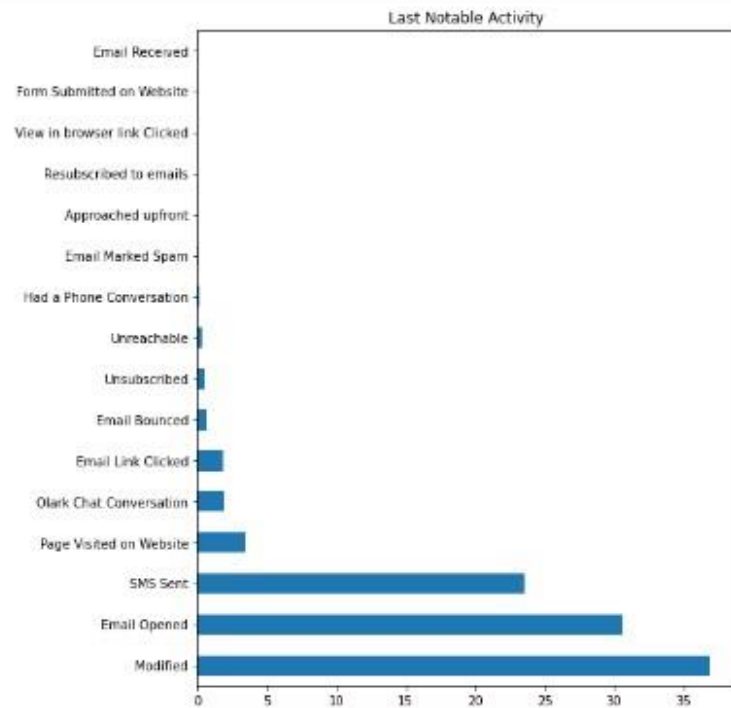- And all the plots are provided in the next slide

# PLOTS

# PLOTS

# PLOTS



What is your current occupation

**Points to be noted from above visualizations:**

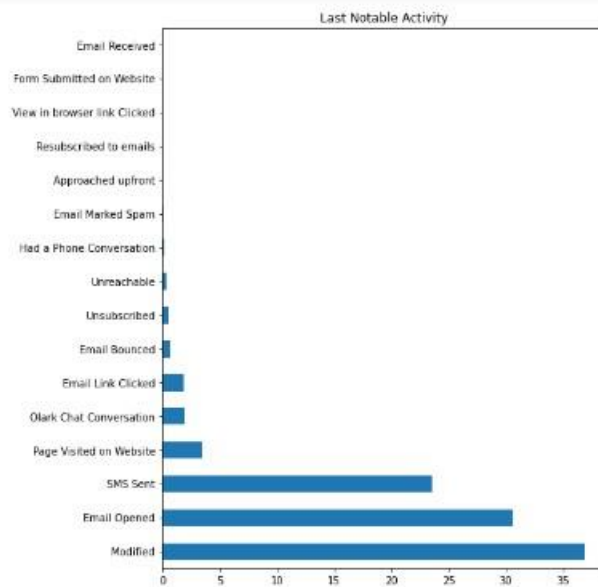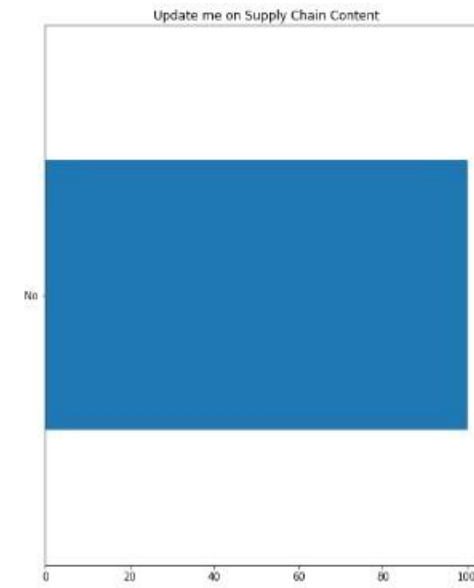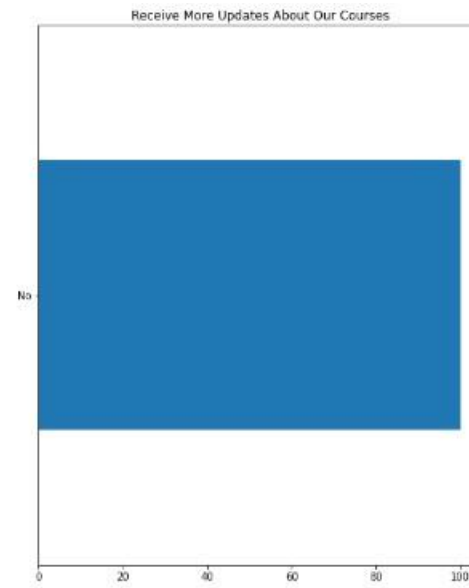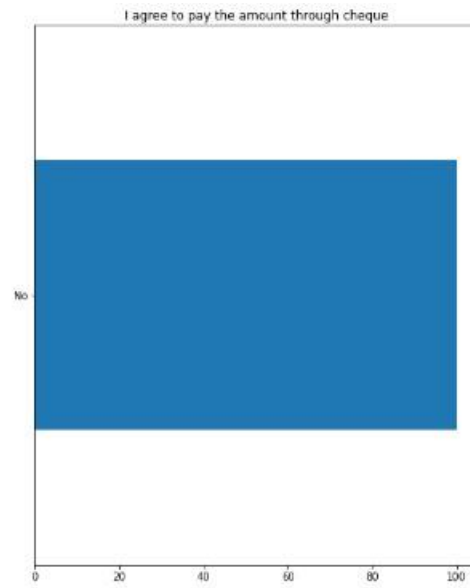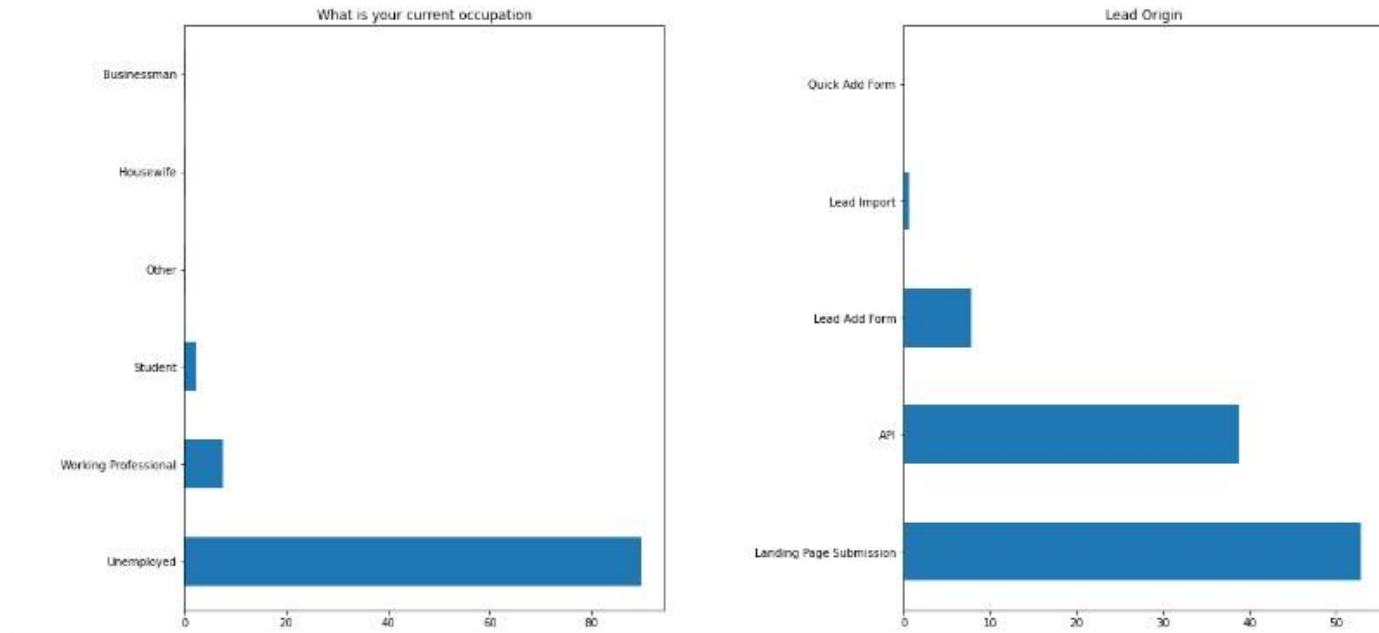- It can be seen that many features have Missing Values more than 50% again after combining `**Select**` & `**NaN**` Values. Let's again check the features having less than & more than 35% of Missing Values respectively.
- Also, in `**Lead Source**` feature, we can see that there are two different values named `Google` & `google`. Let's make them one as well.

# PLOTS

## I agree to pay the amount through cheque

No

## Receive More Updates About Our Courses

No

## Update me on Supply Chain Content

No

## Last Notable Activity

- Email Received
- Form Submitted on Website
- View in browser link Clicked
- Resubscribed to emails
- Approached upfront
- Email Marked Spam
- Had a Phone Conversation
- Unreachable
- Unsubscribed
- Email Bounced
- Email Link Clicked
- Olark Chat Conversation
- Page Visited on Website
- SMS Sent
- Email Opened
- Modified

## Lead Source

- NC_EDM
- testone
- youtubechannel
- welearnblog_Home
- WeLearn
- blog
- Pay per Click Ads
- Press_Release
- Live Chat
- Social Media
- Click2call
- bing
- Facebook
- Referral Sites
- Welingak Website
- Reference
- Organic Search
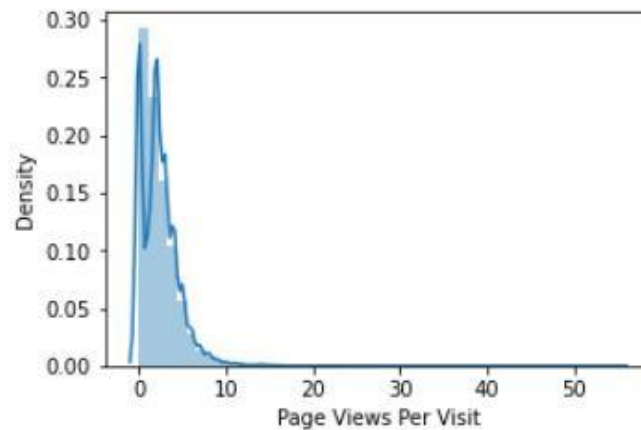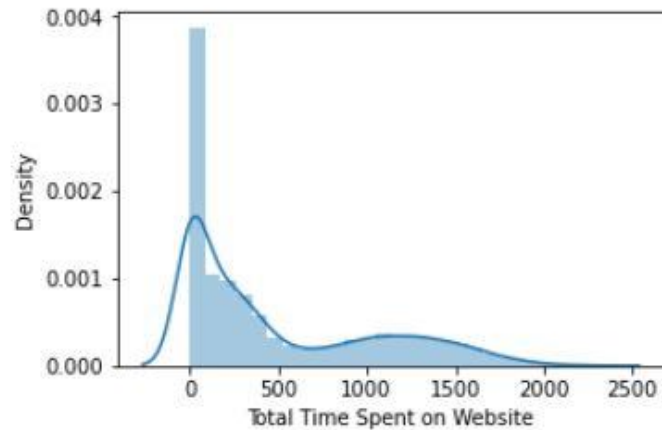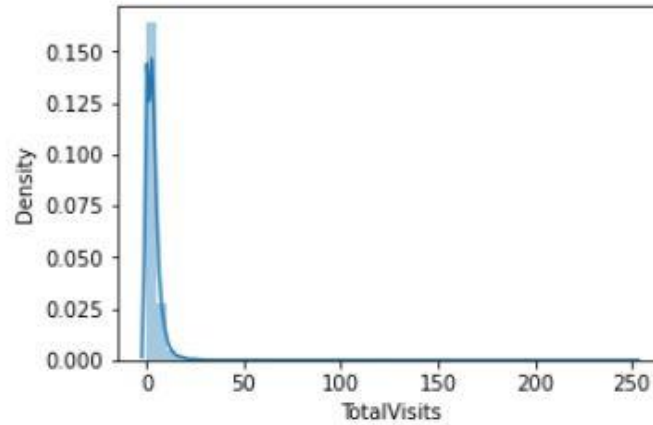- Olark Chat
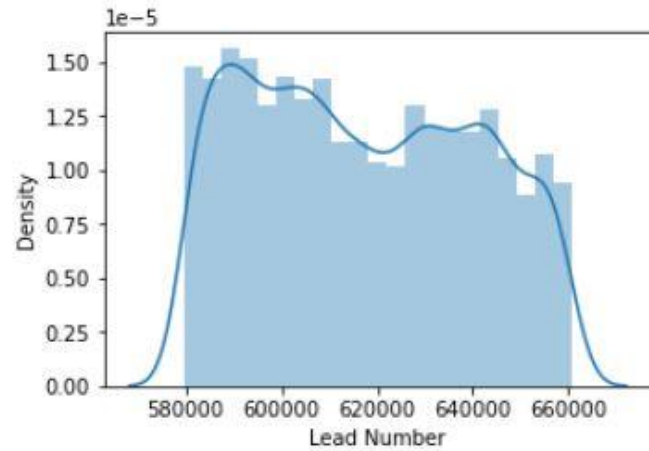- Direct Traffic
- Google

## Magazine

No

# PLOTS



Points to be noted from above Categorical Feature Analysis:(considering 10[th] 11[th] and 12[th] slides)_
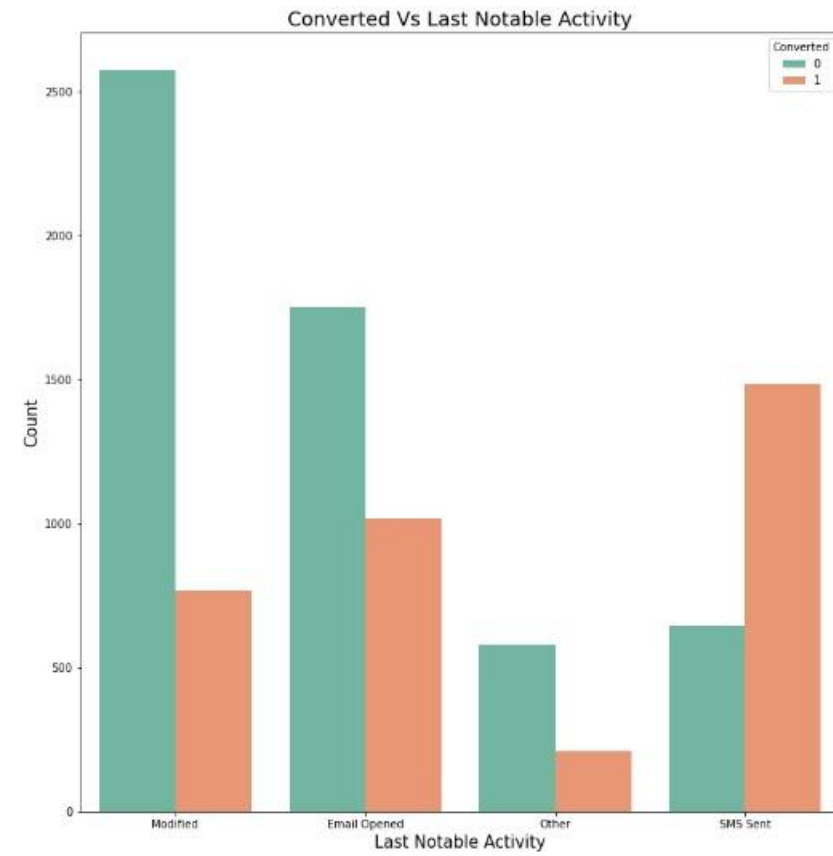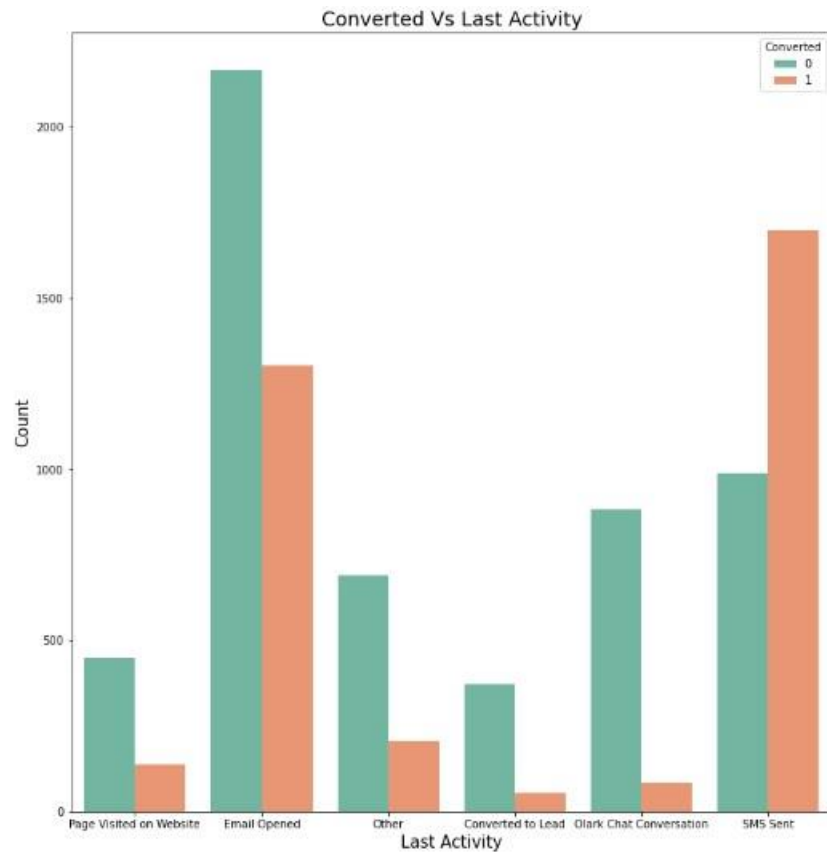
- `Lead Source`, `Country`, `Last Activity`, `What matters most to you in choosing a course`, `Last Notable Activity`, `What is your current Occupation` & `Lead Origin` have multiple categories with very less amount of data. Therefore, we can combine such categories into `Others`.
- Features such as `Receive More Updates About Our Courses`, `I agree to pay the amount through cheque`, `Get updates on DM Content`, `Update me on Supply Chain Content` & `Magazine` have only one category which will not add any value to make classification. Therefore, we will be dropping these features.
- Also `Prospect ID` Feature is having unique Id's for each customer/lead. Therefore, we will drop this as well.
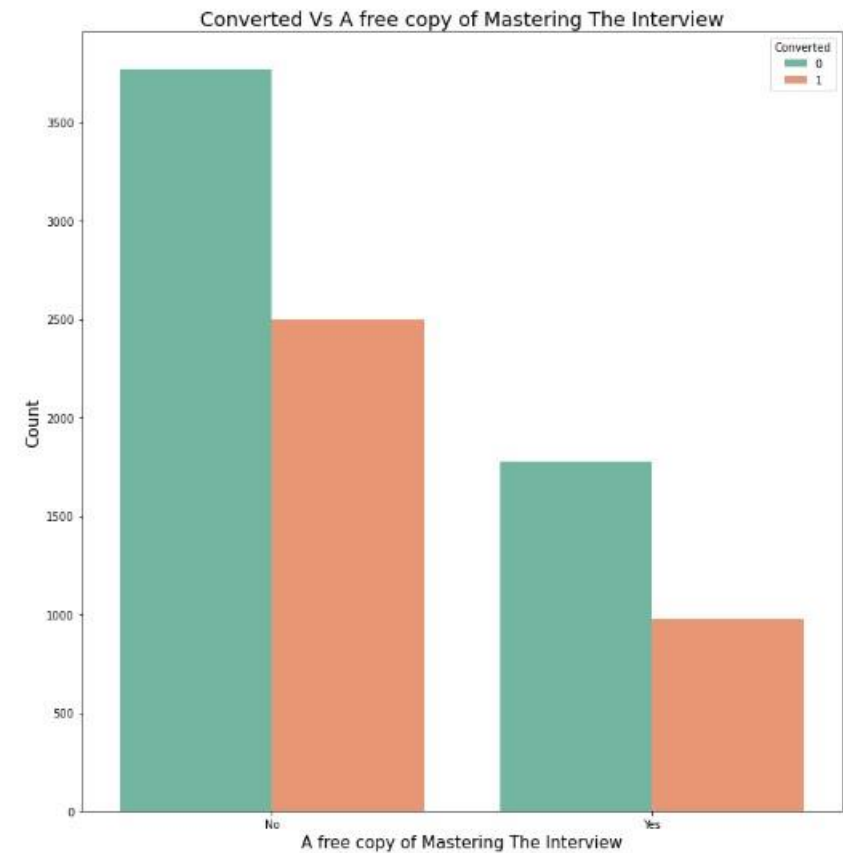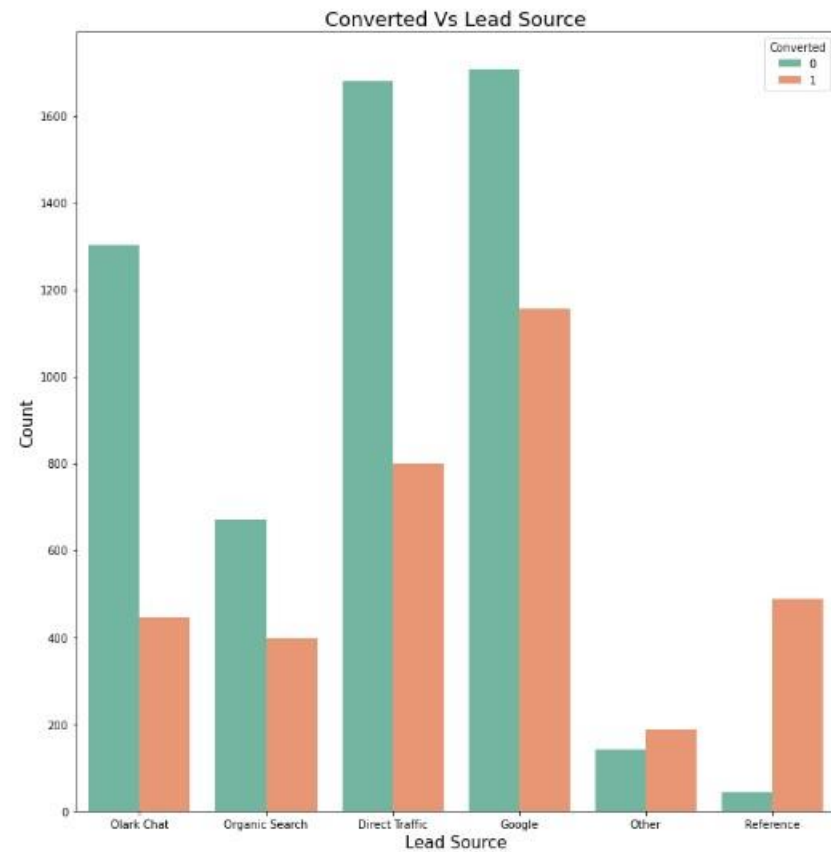
# PLOTS



From the given visualizations, we can observe that there is definitely a need for Feature Scaling which needs to be applied on Numeric Features such as `TotalVisits`, `Total Time Spent on Website` & `Page Views Per Visit` which are highly skewed.
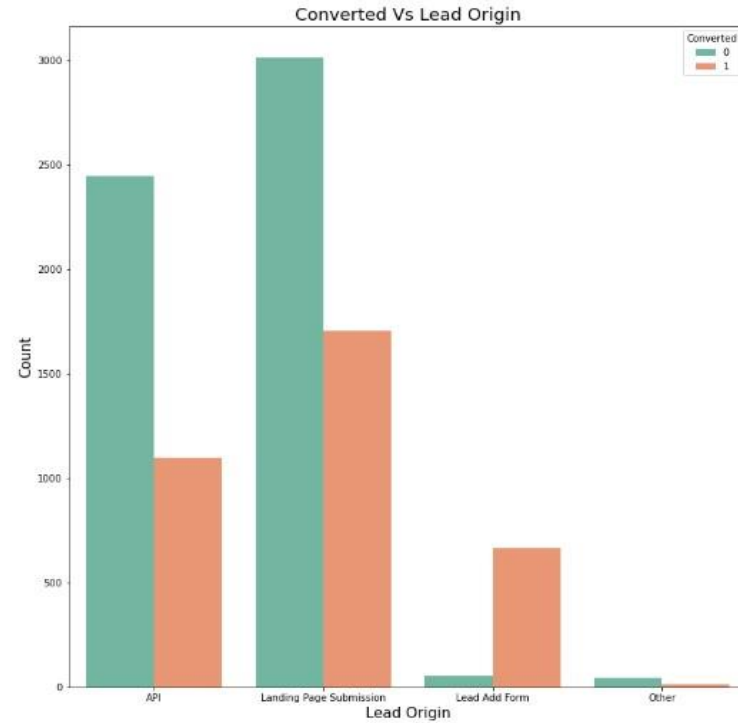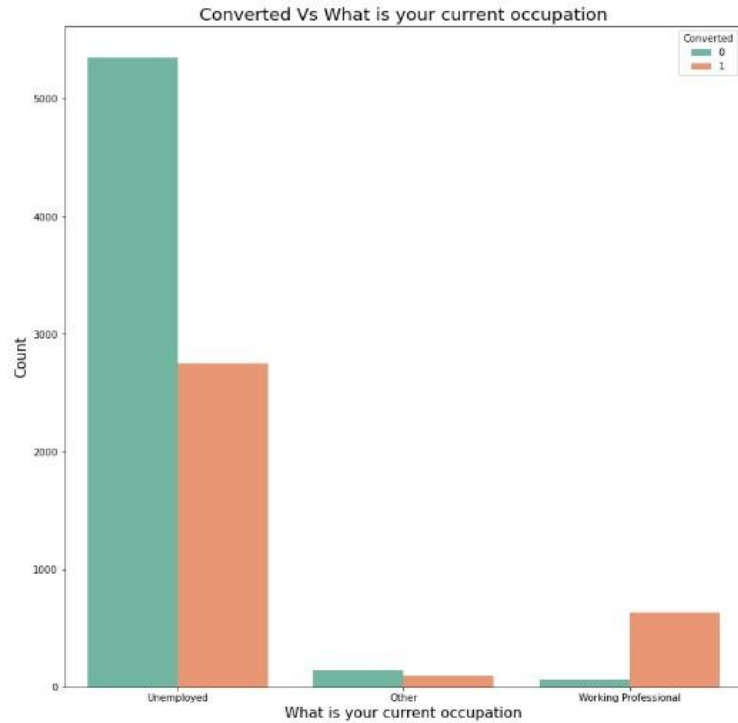
# PLOTS



Converted Vs Last Activity



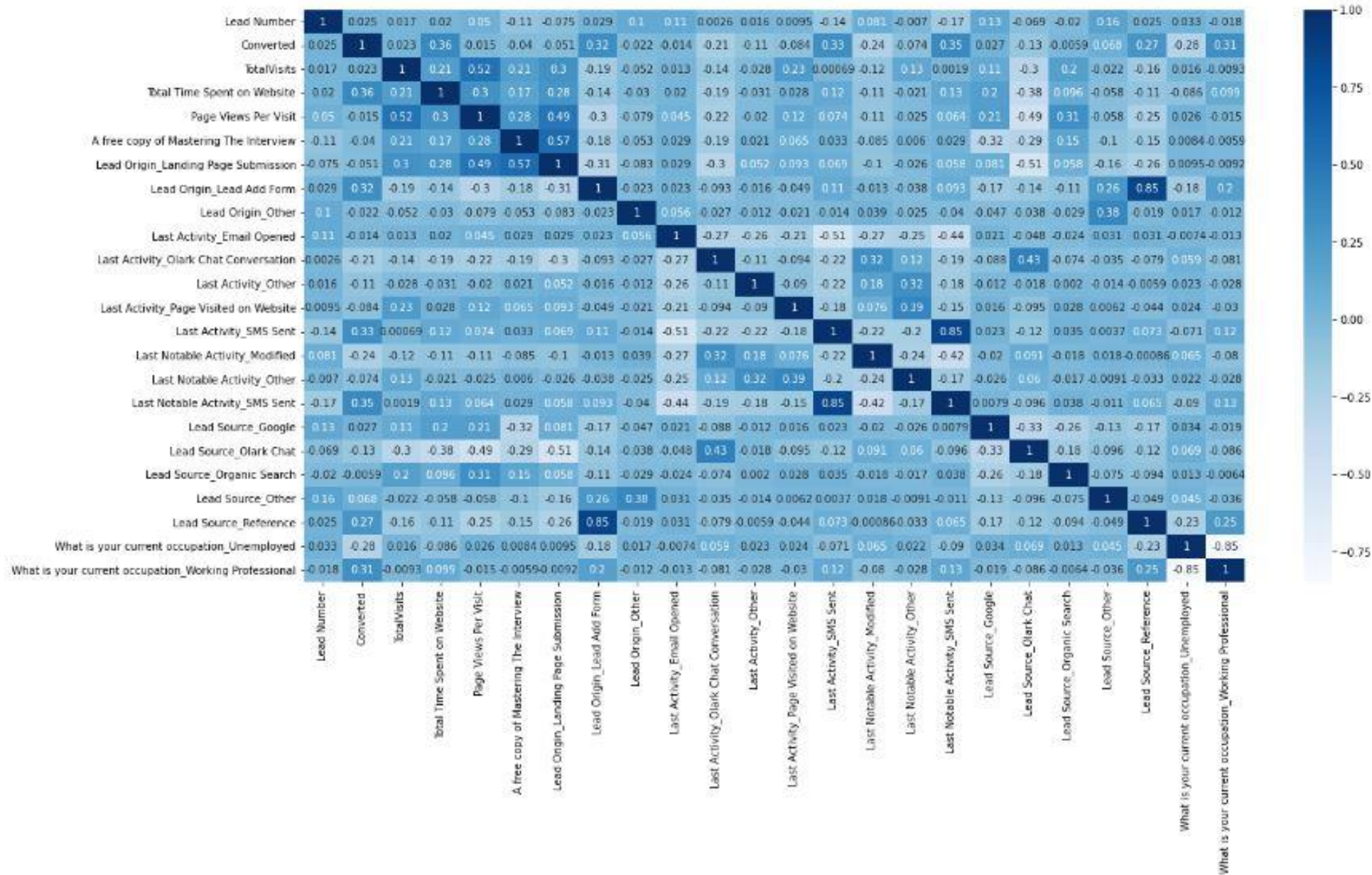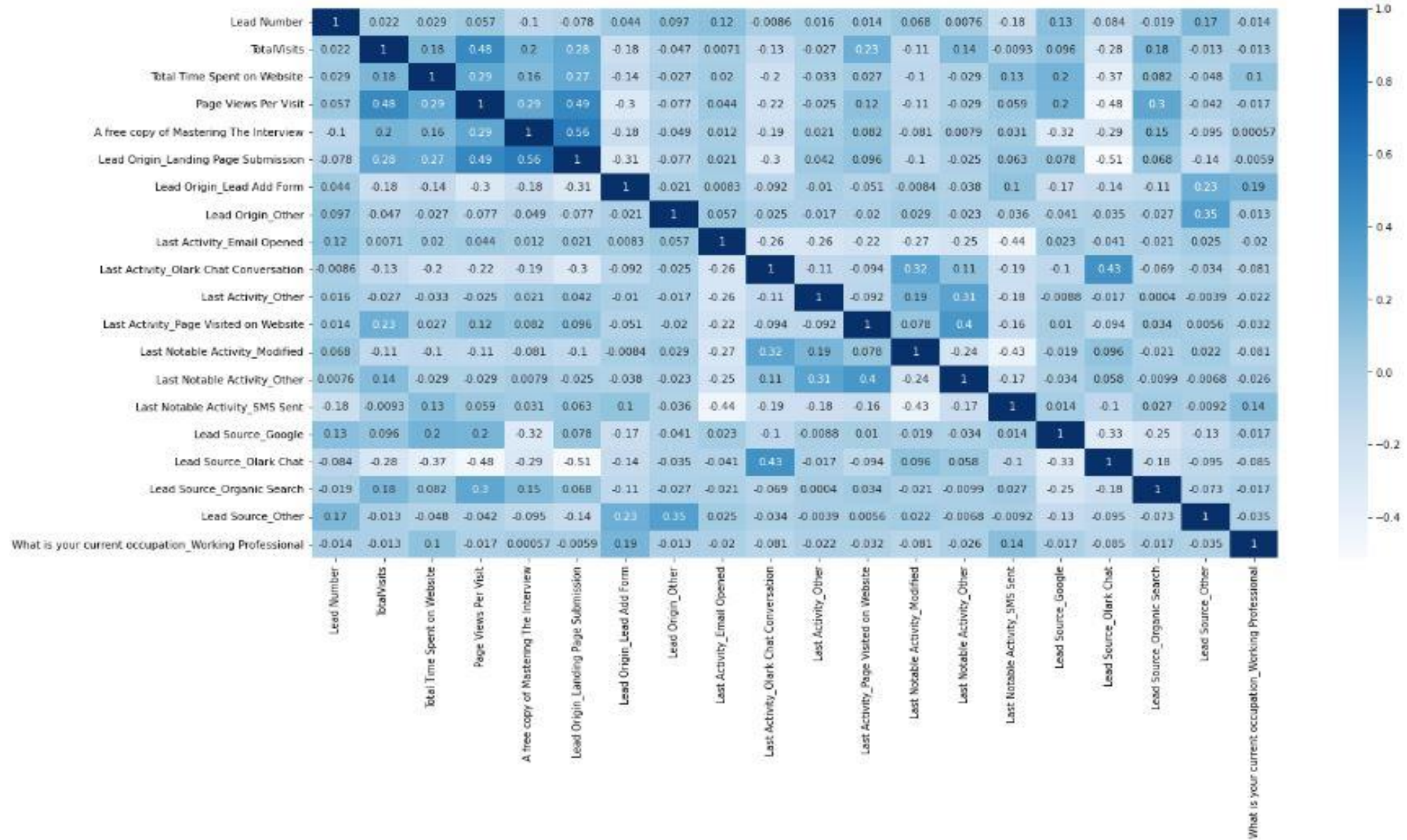Converted Vs Last Notable Activity

# PLOTS

# PLOTS



Following Points can be made from the analysis:(considering slide 14th , 15th and 16th )

- For `A free copy of Mastering The Interview`, it can be seen that the maximum conversion happened from the customers who do not want a free copy of Mastering the Interview.
- For `Last Notable Activity`, it can be seen that the maximum conversion happened from `SMS Sent`.
- For `What is your current occupation`, it can be seen that the maximum conversion happened from `Unemployed` whereas from `Working Professional` group there are more Successful Conversions than Non-Conversions.
- For `Last Activity`, it can be seen that the maximum conversion happened from `SMS Sent` whereas from `Email Opened` we have the highest amount of Non-Conversions.
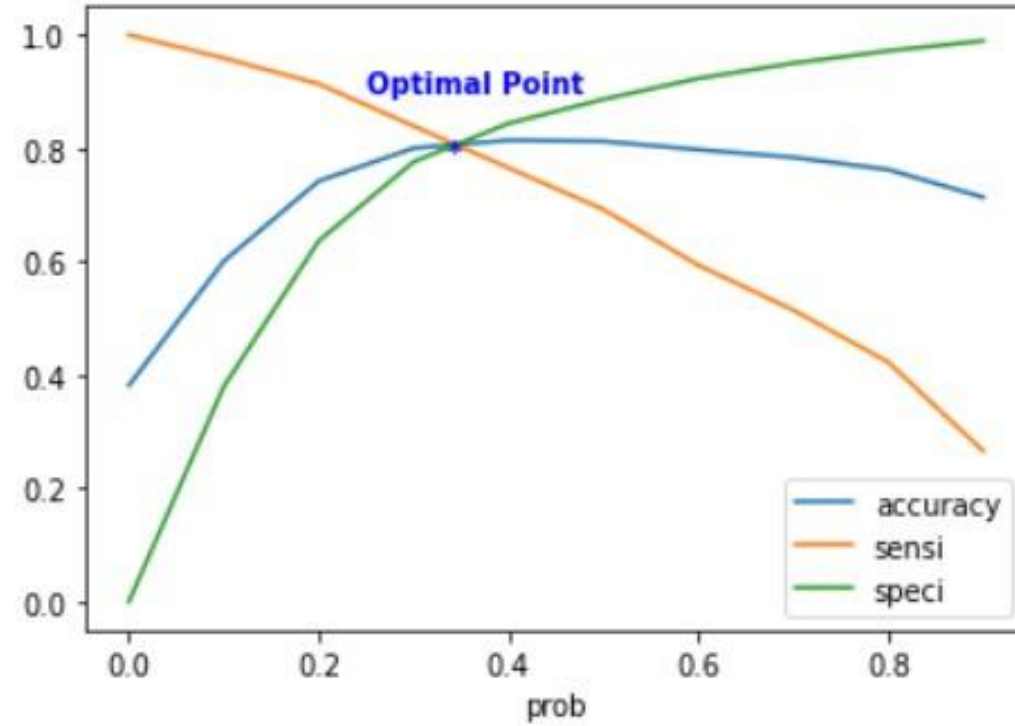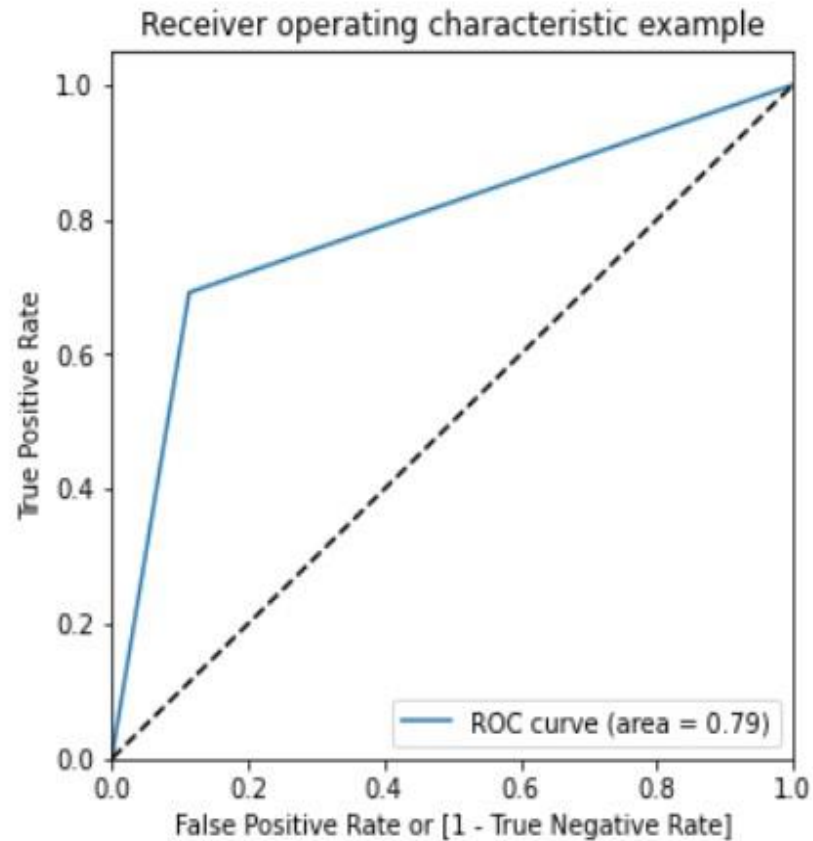
# PLOTS

# PLOTS

# PLOTS

From the above two plots(considering 17th and 18th slide):

- TotalVisits -> ['Page Views Per Visit']
- Lead Source_Reference -> ['Lead Origin_Lead Add Form']
- Last Activity_SMS Sent -> ['Last Notable Activity_SMS Sent']
- What is your current occupation_Unemployed -> ['What is your current occupation_Working Professional']
- All the above are highly correlated and we are dropping them
- By considering slide 18th visualization we found It can be seen that now no Independent Feature is highly inter-correlated.
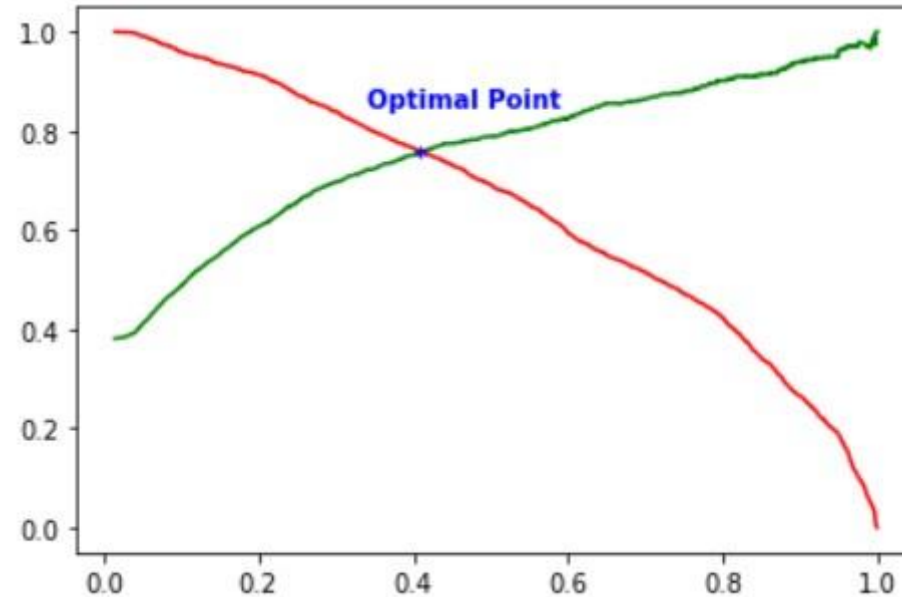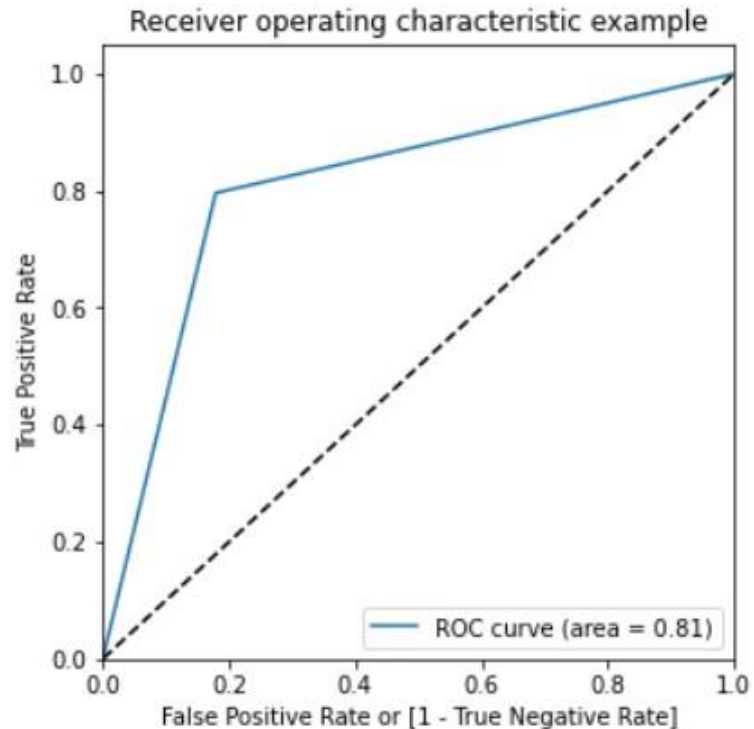
# PLOT OF ROC CURVE & OPTIMAL POINT



From the above visualization, we can observe that the Optimal Probabilistic CutOff is `0.335`.

Still, let's check from Precision-Recall Trade-Off Curve as well to find the optimal probability.

# PLOT OF ROC CURVE & OPTIMAL POINT



- In the above diagram, Red Line represents the Recall & Green represents the Precision.
- And the Optimal Probability is at `0.4`.
- Both Optimal Probabilities are close to each other. Therefore, let's take mean of these as the `Final Optimal Probability Value`
- Hence, the `Final Optimal Probability` is at `0.37`

# Conclusion

- After Lead Score we can see that the final prediction of conversions have a target of 80% (78.39%) conversion as per the X Education CEO's requirement

- Accuracy is 80.98%

- Recall is 78.39%

- Precision is 73.51%

- Specificity is 82.58%

- False Positive Rate is 17.42%

- Hence  it seems a Good Model for Lead Scoring