# Lead Score Case Study

**Created By:**

Nayanshi Sahu

Sameer Kumar

Sakshi Bhasin

## Problem Statement

An X Education need help to select the most promising leads, i.e. the leads that are most likely to convert into paying customers. The company requires us to build a model wherein you need to assign a lead score 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%.

## Goals of Case Study

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.

## Steps Followed:

- Importing Libraries and Data
- Inspecting the Dataframe
- Data Cleaning
- Data Preparation
- Test-Train Split
- Feature Scaling
- Model Building using Stats Model & RFE
- Predicting a Train model
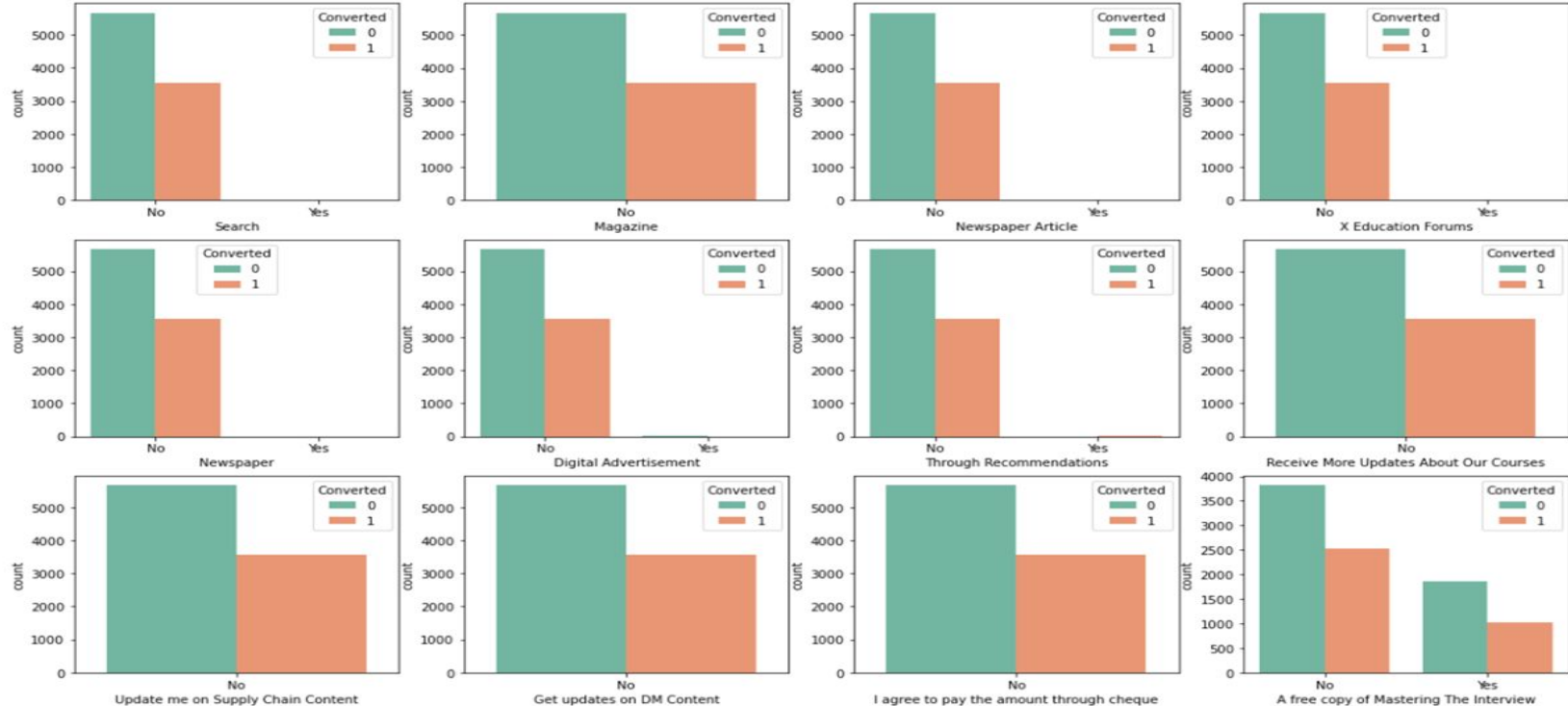- Predictions on  the test set
- Conclusion

**Library Used:**

- Numpy
- Pandas
- Matplotlib
- Seaborn
- Sklearn

**Data Inspection:**

- Checked data Dimension size and description
- Checked for Duplicity
- Deleted Unique keys columns

# Data Cleaning:

- Dropping Unique Keys columns and columns with more than 35% null value
- Dropped few column whose data are highly imbalanced shown in the below graph

Examined every Categorical column and if found any null values, we have replaced it with Mode in most of the cases, In some cases we replace null values with text like not provided.

Some columns that looked screwed and had redundant information we have dropped them.

We have noticed Below 2 point from above exercise:
1. Maximum leads are generated having last activity as Email opened but conversion rate is not too good.
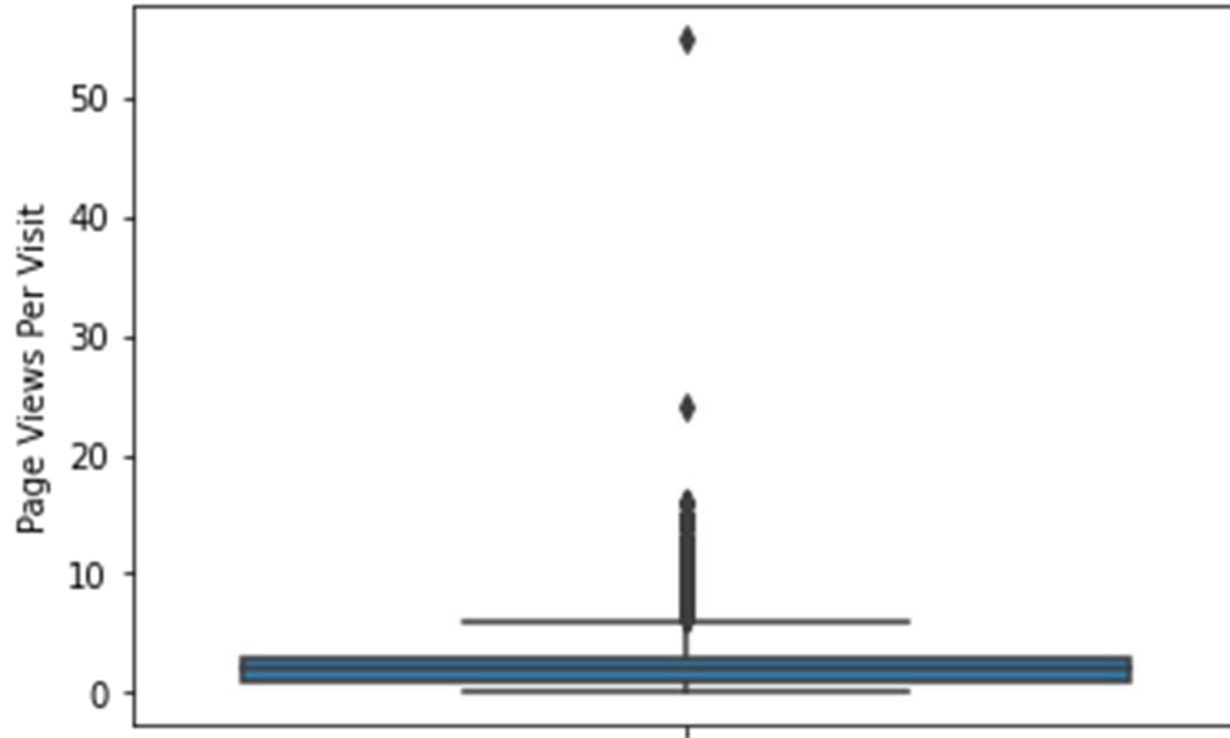2. SMS sent as last activity has high conversion rate.

## Correlation:

•Numerical Columns examin by using correlation Matrix as shown in the right side image.

•By Correlation we have figured out that mostly People who has spend more time on website has higher chance to get converted.
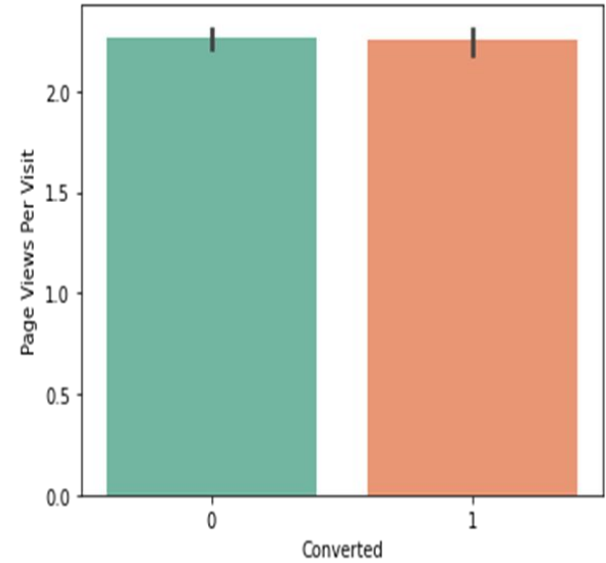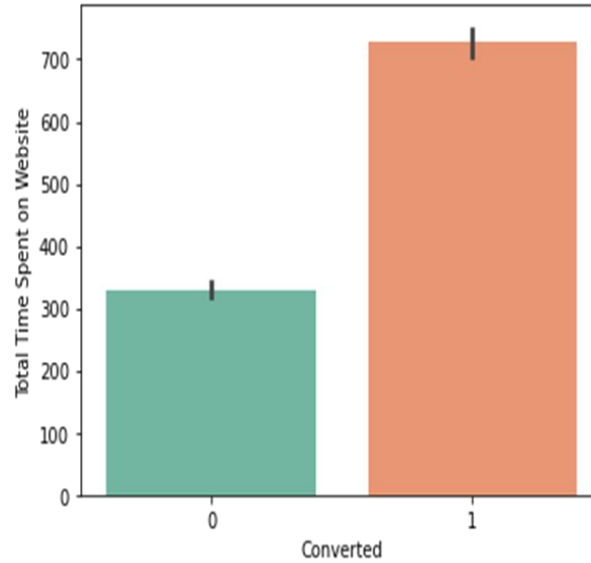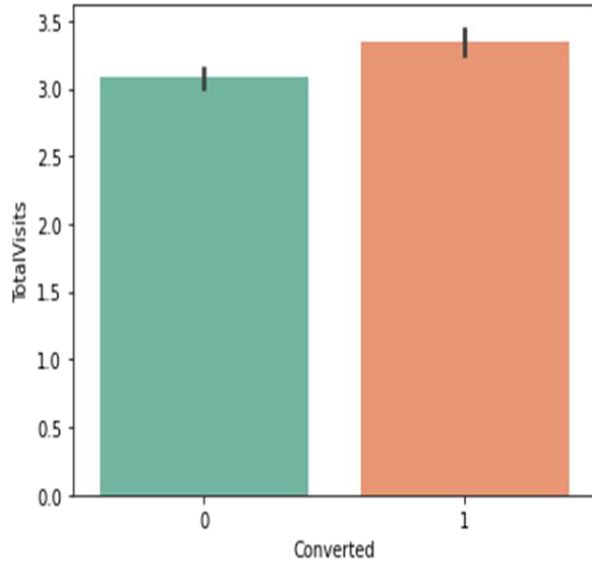
# Outlier:

Outlier found only in "Page Views per Visit" Column and we have treated them by replacing all the values which are higher than one percentile with one percentile.
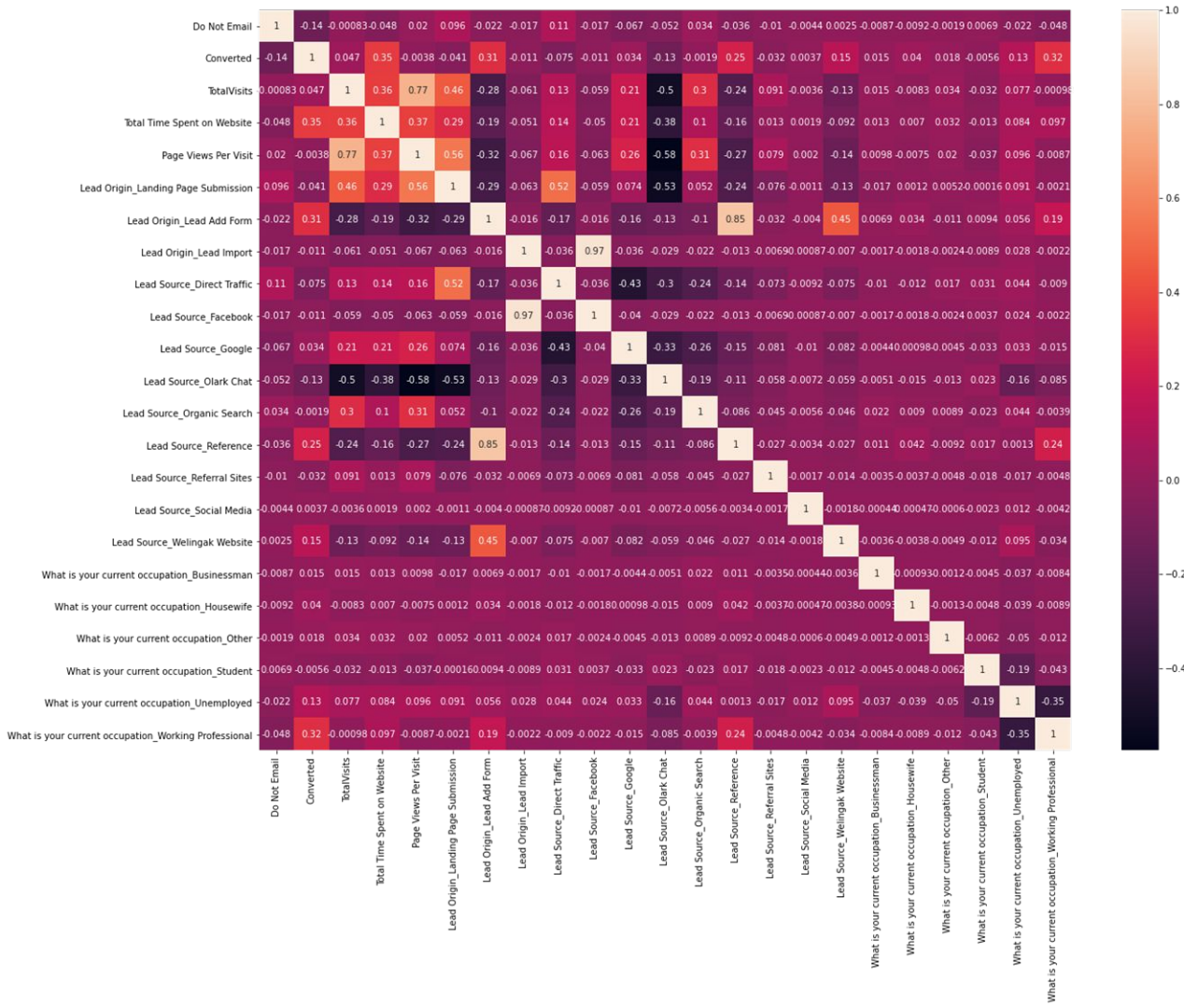
# Checking Numerical value conversion rate:



The conversion rate is high for all of the above three Columns.

**Data Preparation:**

- Converting some binary variables (Yes/No) to 0/1
- Created Dummy Variables
- Test Train Split

# Feature Scaling:

- Using Correlation matrix after creating Dummy Variables and dropped highly correlated matrix

**Model Building:**

•Used Stats Model & RFE

•By optimizing p value we have rebuilt our model 5 times.

•After 5$^{th}$ time rebuilding the model the all features VIF So we need to further

optimize our model and we can proceed with making predictions using this model only

•Predicting a Train Model and calculated Metrics -Accuracy, Sensitivity, Specificity, False Positive Rate, Positive Predictive Value and Negative Predictive Value.

**Confusion Matrix:**

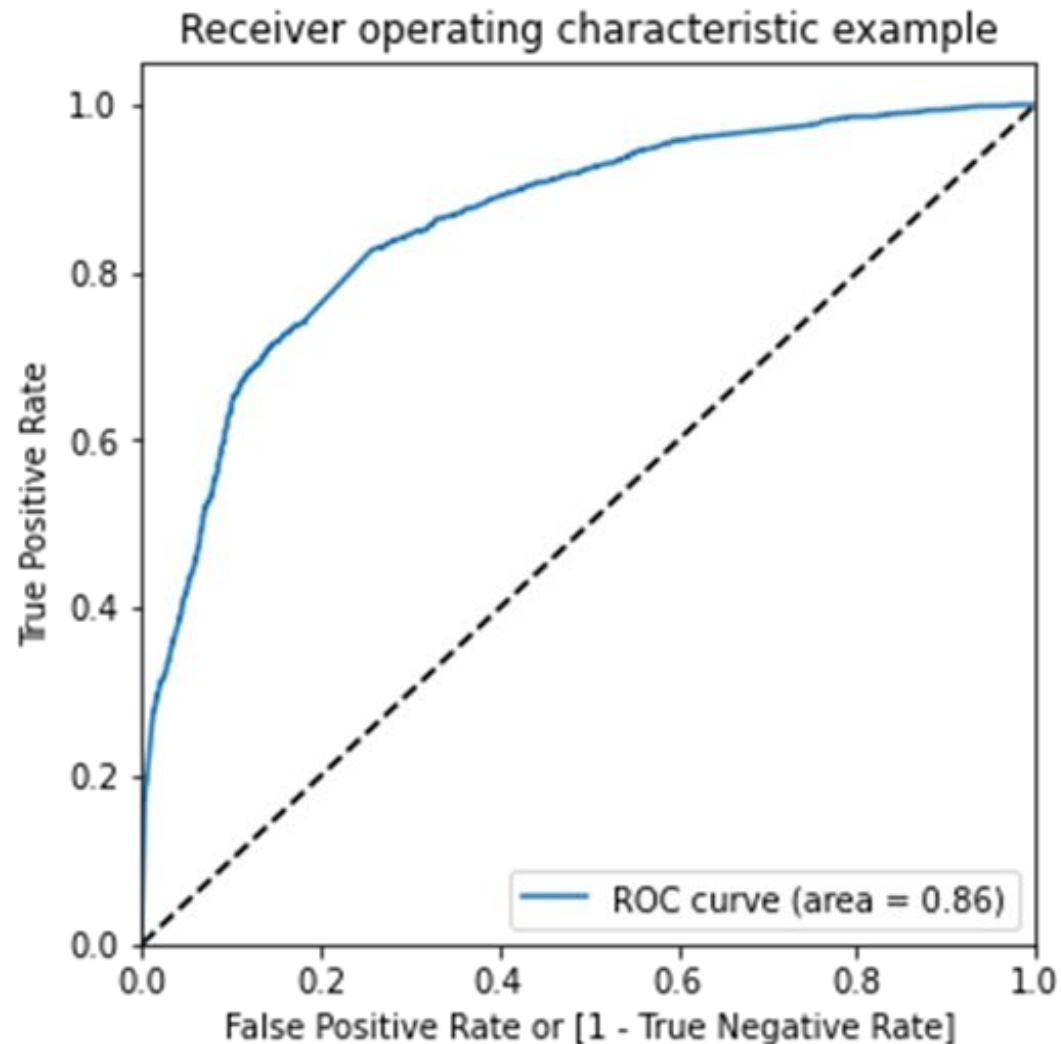| | |
|------|------|
| 2905 | 1048 |
| 414 | 2005 |

**Accuracy: 77.05%**
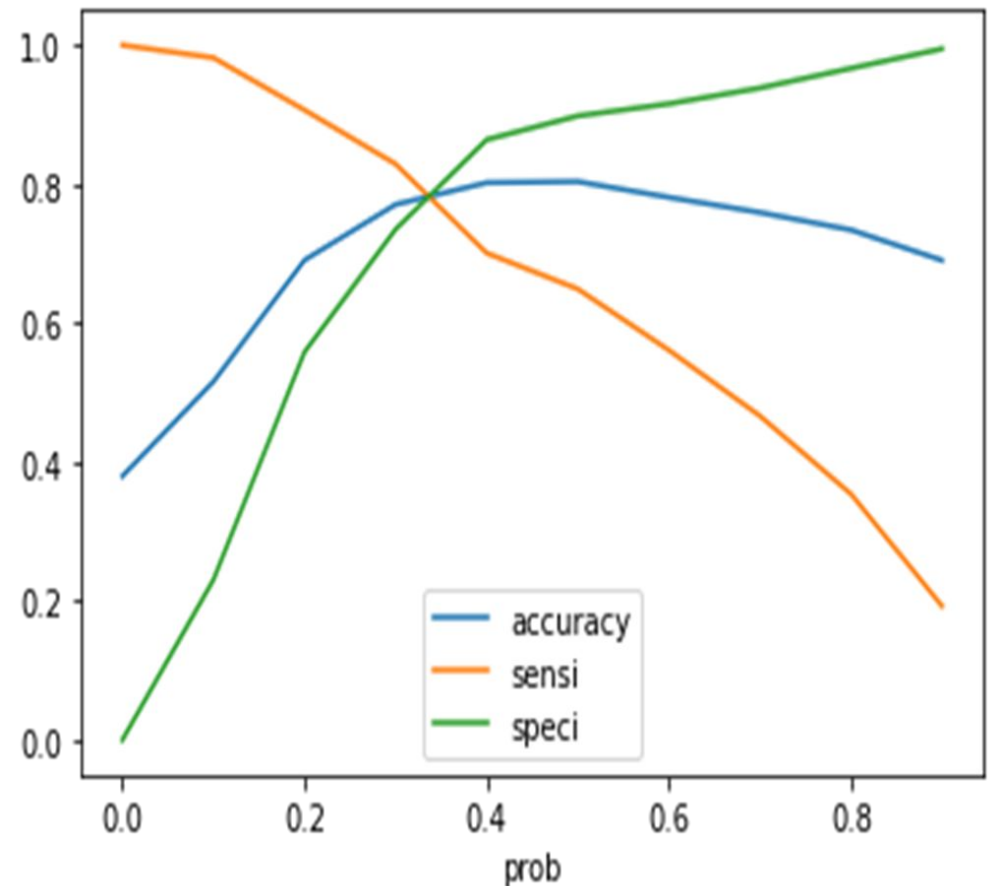**Sensitivity: 82.89%**
**Specificity: 73.9%**
**Precision: 65.67%**
**Recall: 82.88**

# PLOTTING ROC CURVE

The ROC Curve should be a value close to 1. We are getting a good value of 0.86 indicating a good predictive model
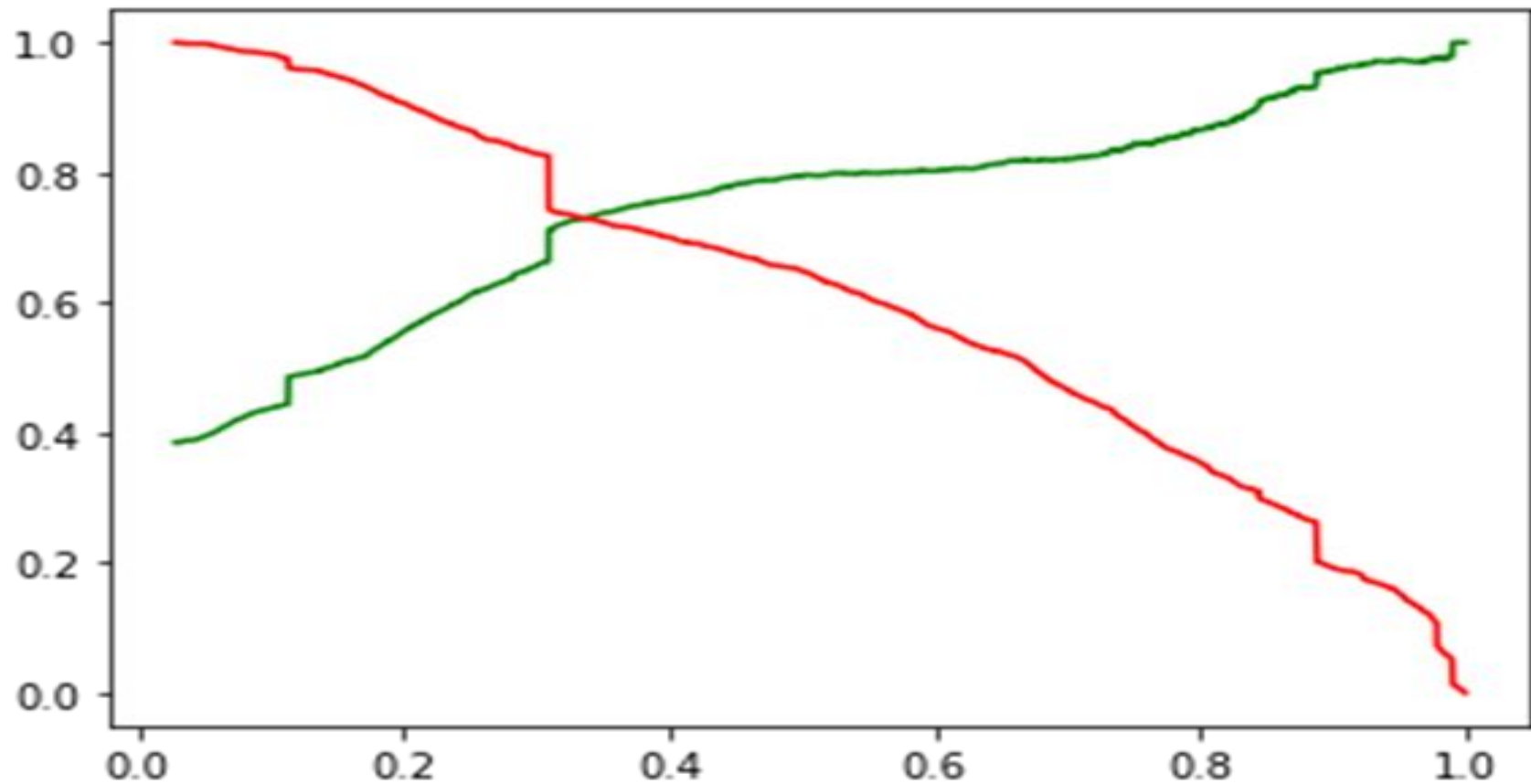
# Accuracy sensitivity and specificity Plot



From the curve above, 0.3 is the optimum point to take it as a cutoff probability.

# Precision and Recall Trade-off

## Prediction on the Test Set

- The final prediction of conversions have a target rate of 83% (same as predictions made on training data set)
- Overall Metrics - Accuracy, Confusion Metrices, Sensitivity, Specificity, False Positive Rate, Positive Predictive Value, Negative Predictive Value on final prediction on test set
- After running the model on the Test Data these are the figures we obtain:
  - Accuracy : 77.52%
  - Sensitivity :83.01%
  - Specificity : 74.13%

**Conclusion:**

- While we have checked both Sensitivity-Specificity as well as Precision and Recall Metrics, we have considered the optimal cut off based on Sensitivity and Specificity for calculating the final prediction.
- Accuracy, Sensitivity and Specificity values of test set are around 77%, 83% and 74% which are approximately closer to the respective values calculated using trained set.
- Also the lead score calculated in the trained set of data shows the conversion rate on the final predicted model is around 80%.
- Hence overall this model seems to be good.

**Important features responsible for good conversion rate or the ones' which contributes more towards the probability of a lead getting converted are :**

- Lead Origin_Lead Add Form
- What is your current occupation_Working Professional
- Total Time Spent on Website

# Thank You