

# Leading Score Case Study

**Presented by:**  
**Praveen Thota**  
**Pravin Shinde**

# 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%.

X Education has appointed you to help them select the most promising leads, i.e., the leads that are most likely to convert into paying customers. The company requires you 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 the Case Study

There are quite a few goals for this 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.
- 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.

## Approach

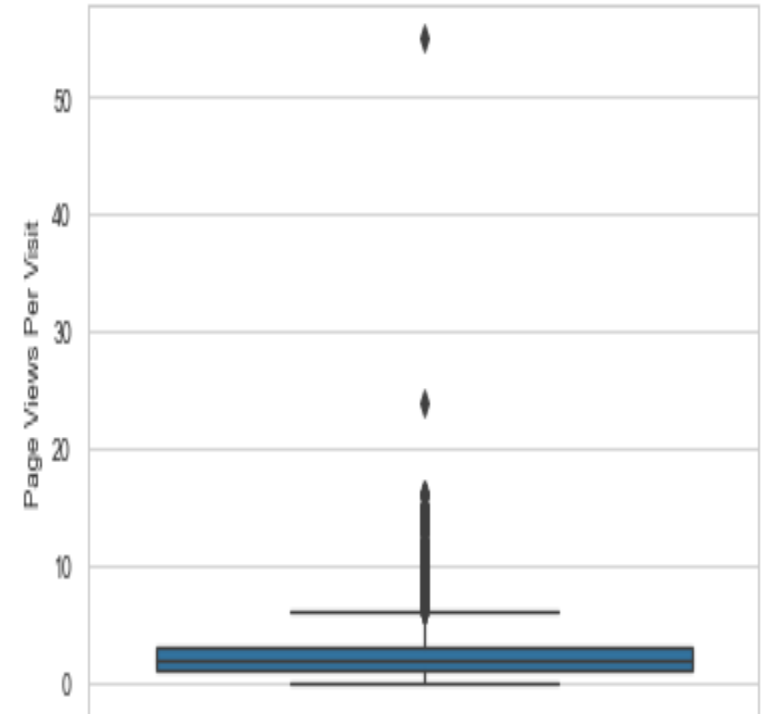
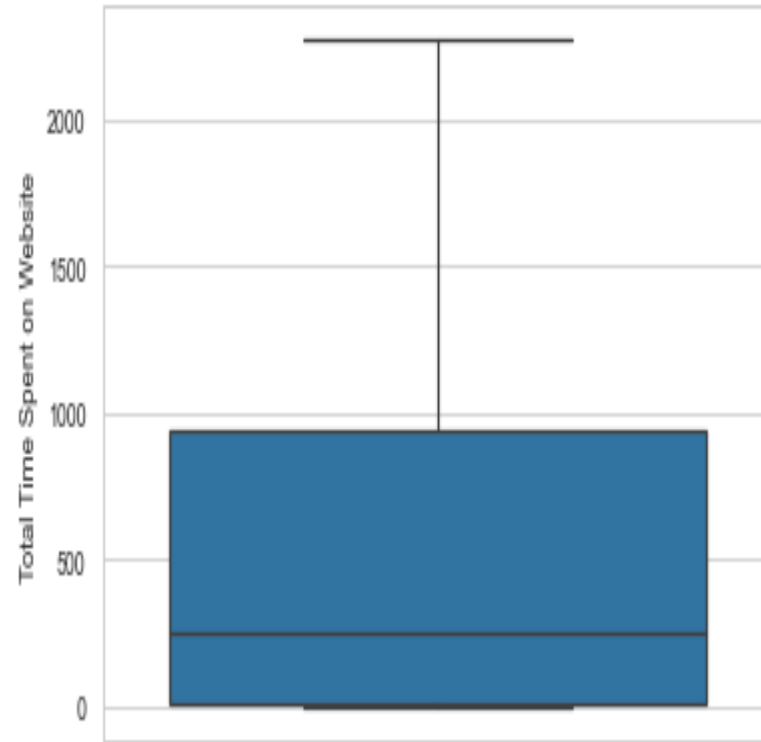
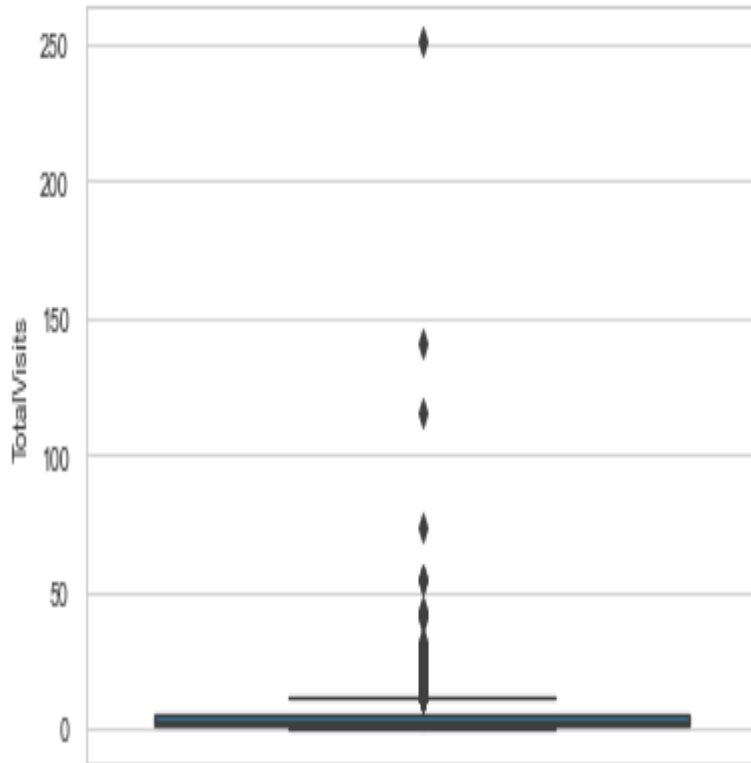
- Loading Libraries and Inspecting Data Frame
- Reading & Understanding the Data
- Data Cleaning and Data Preparation
- EDA
- Feature Scaling
- Train Test split data
- Prepare the data for modelling
- Build model
- Model Evaluation
- Accuracy, specificity, sensitivity
- Precision, Recall and F1 score
- Making predictions on the test set
- Conclusion

# Data Sourcing, Cleaning & Preparation

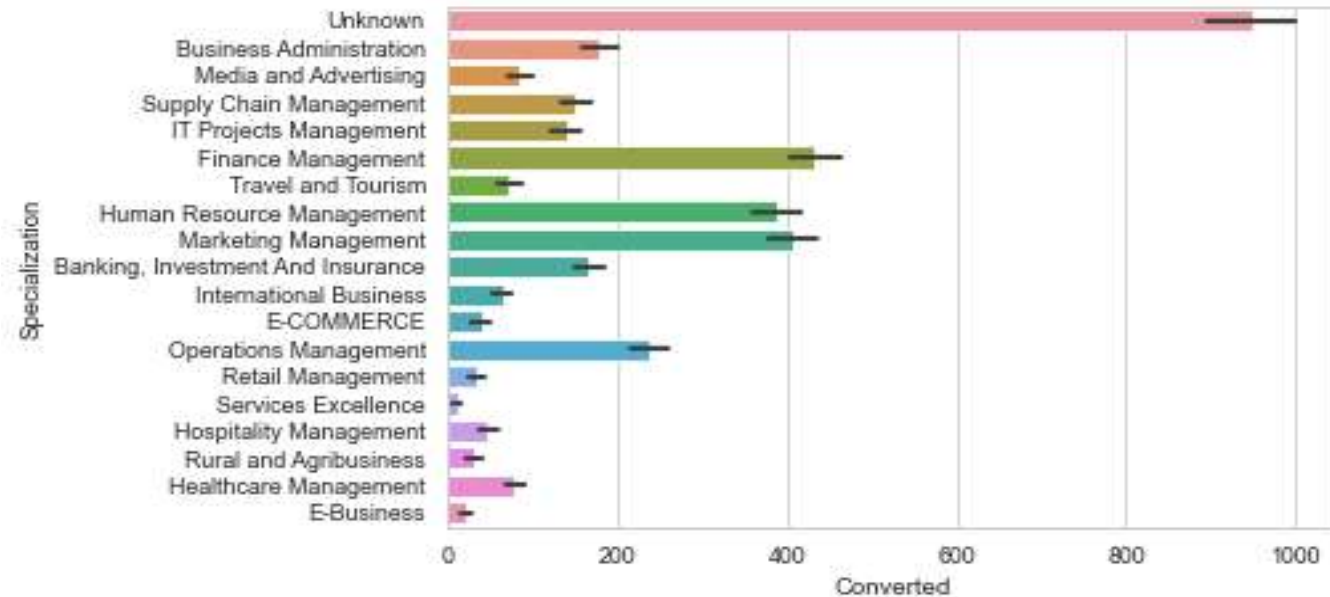
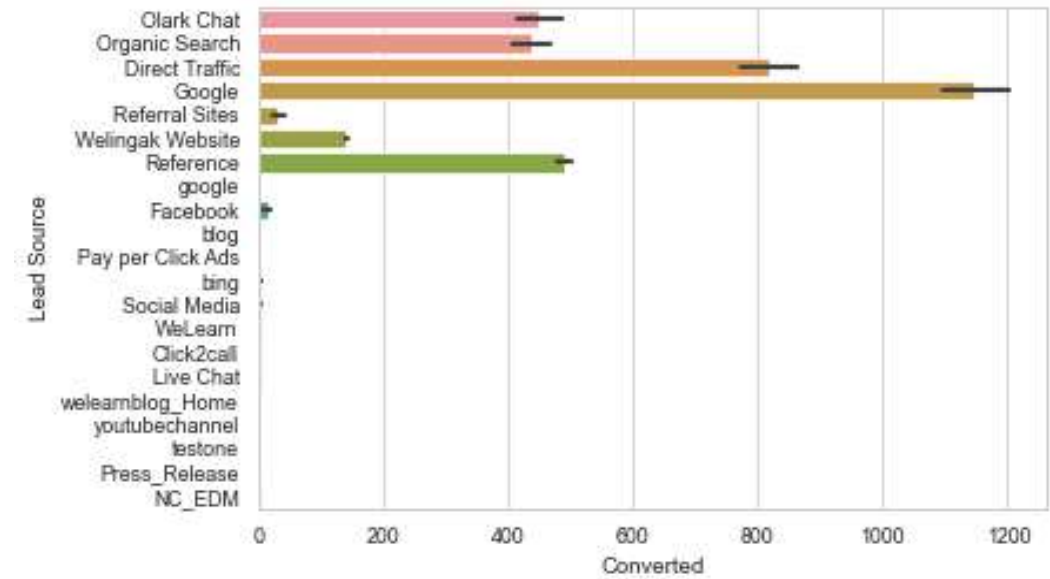
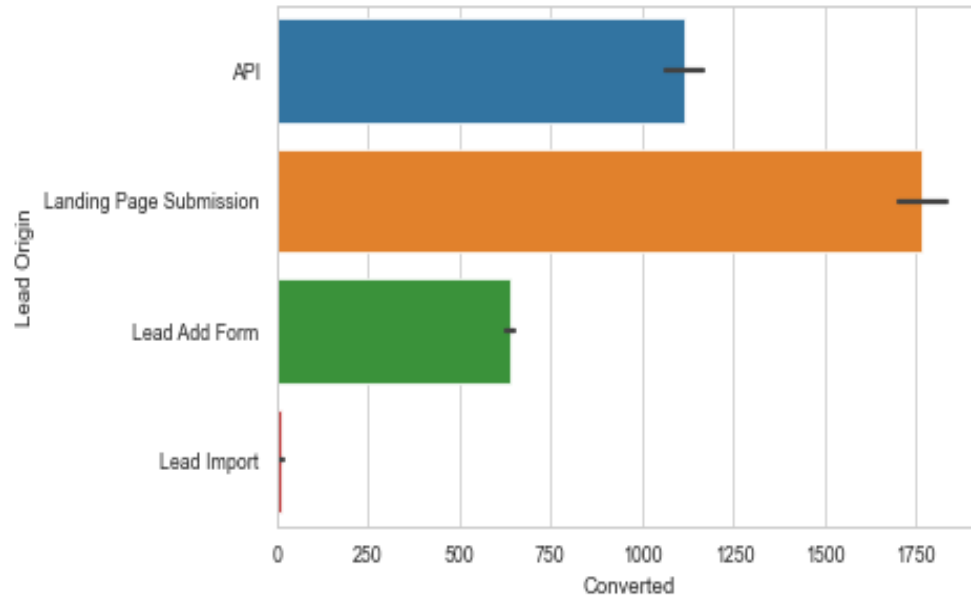
- Importing libraries
- Read the data from the csv file
- Inspecting the data frame
- Handling null values & Removing cols with high null values data
- Assigning an Unique Category to NULL/SELECT values (Unknown values)
- Dropping the columns which have only unique values (Magazine etc.,)
- Deleting the columns 'Asymmetrique Activity Score' & 'Asymmetrique Profile Score' as they will be represented by their corresponding index columns as per correlation
- Imputing null values with mode() and median()
- Removing rows with minor null values in a particular column (Lead Source)
- Outlier treatment
- Removing redundant columns 'Prospect ID' and 'Lead Number'
- Removing outlier values based on the Interquartile distance for some of the continuous variables (TotalVisits, Page Views Per Visit)

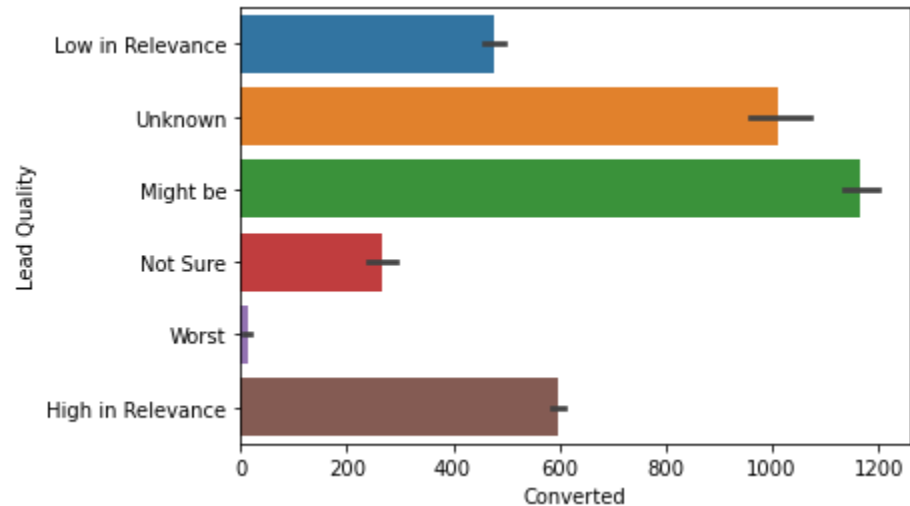
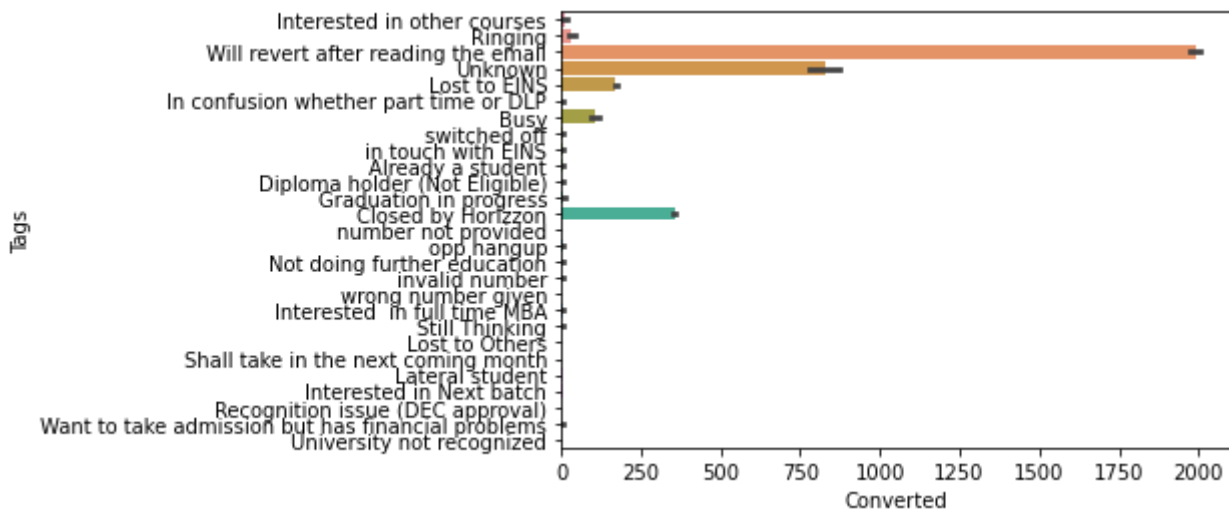
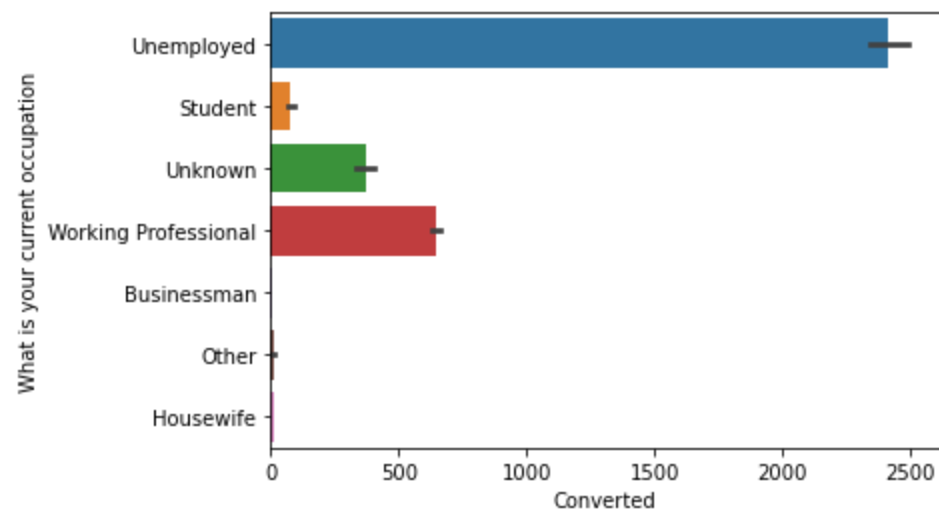
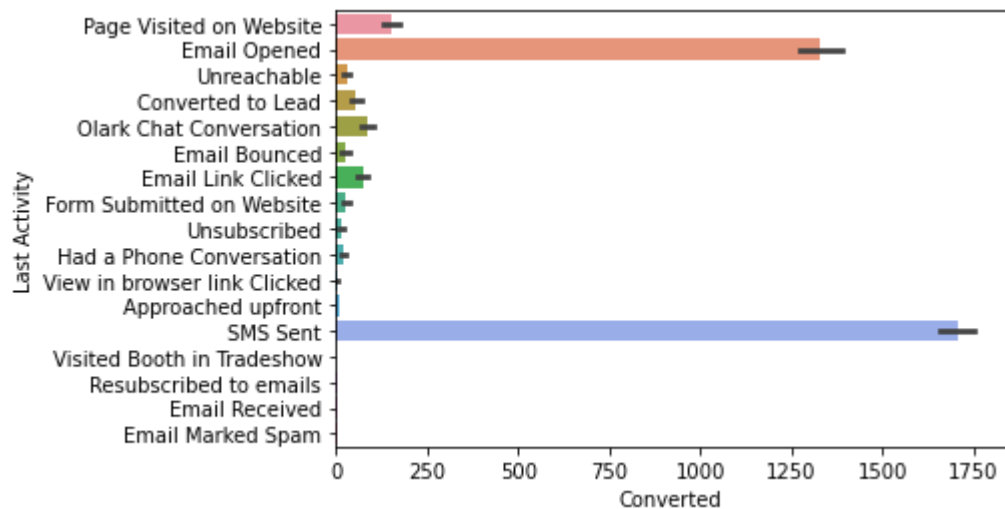
# Exploratory Data Analysis

- Box Plot (Checking Outliers):

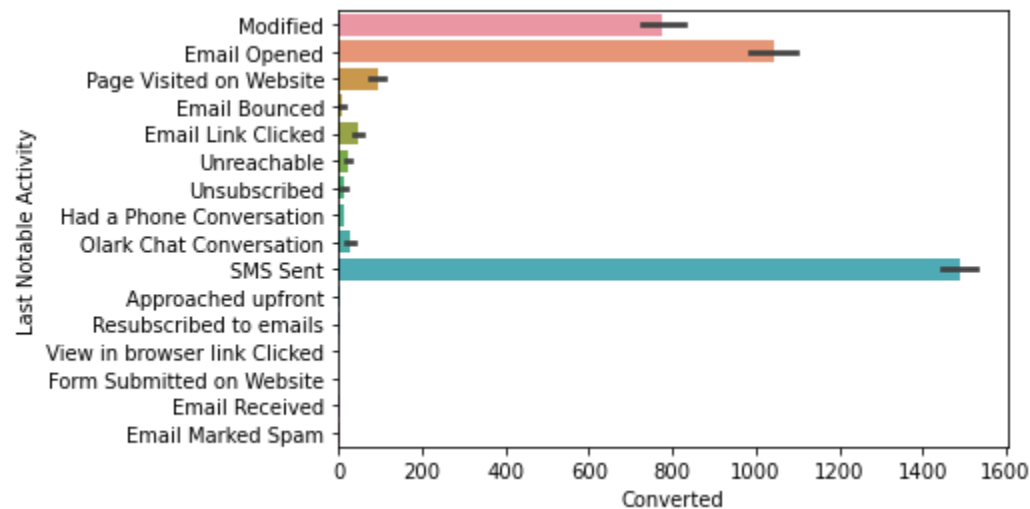
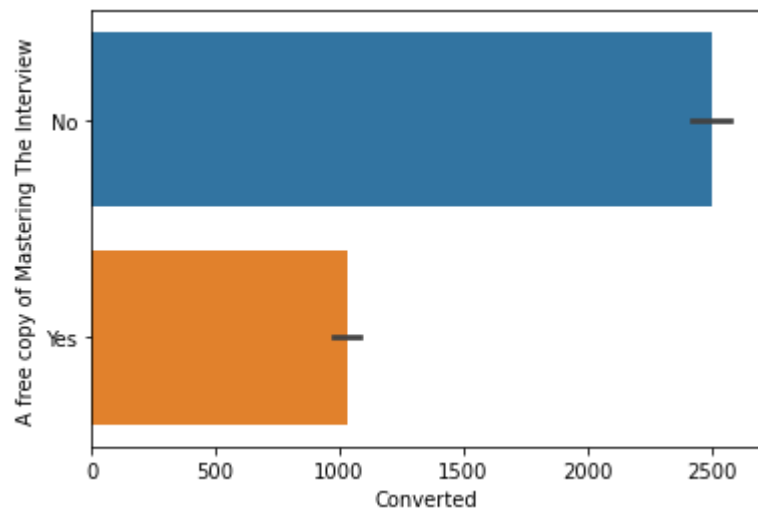
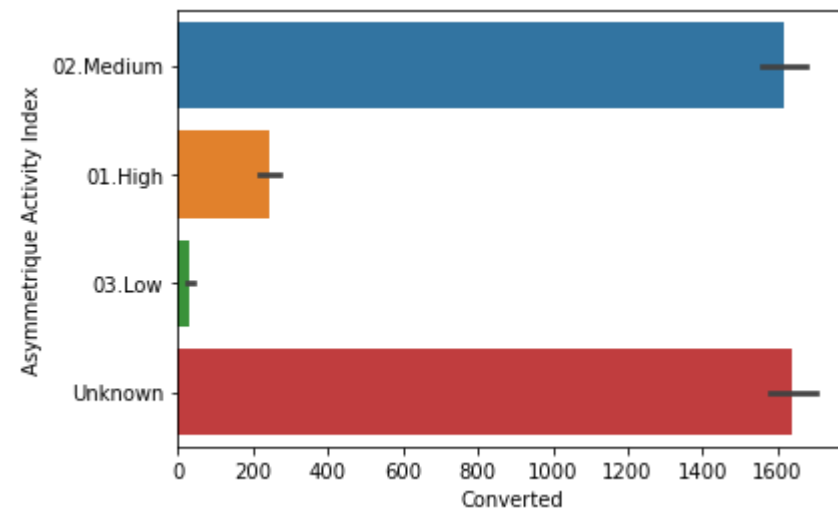
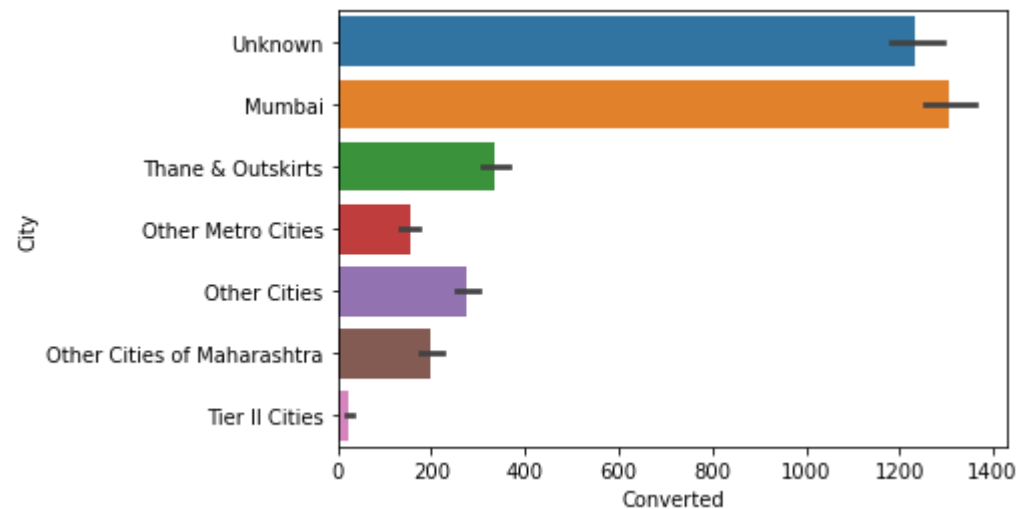


# Bar Plots

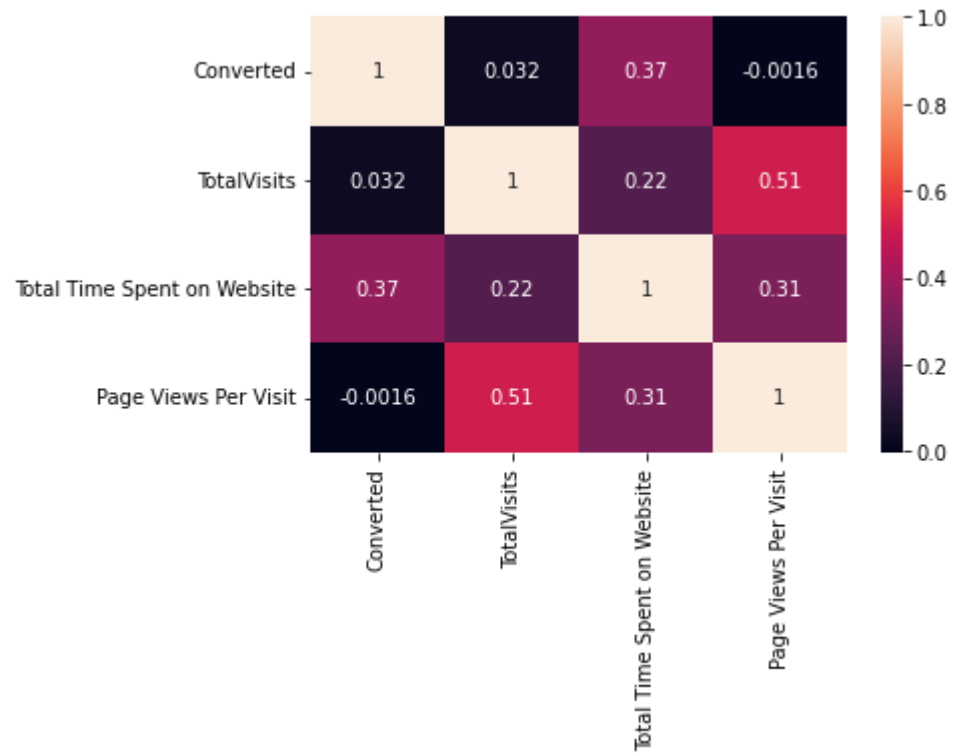


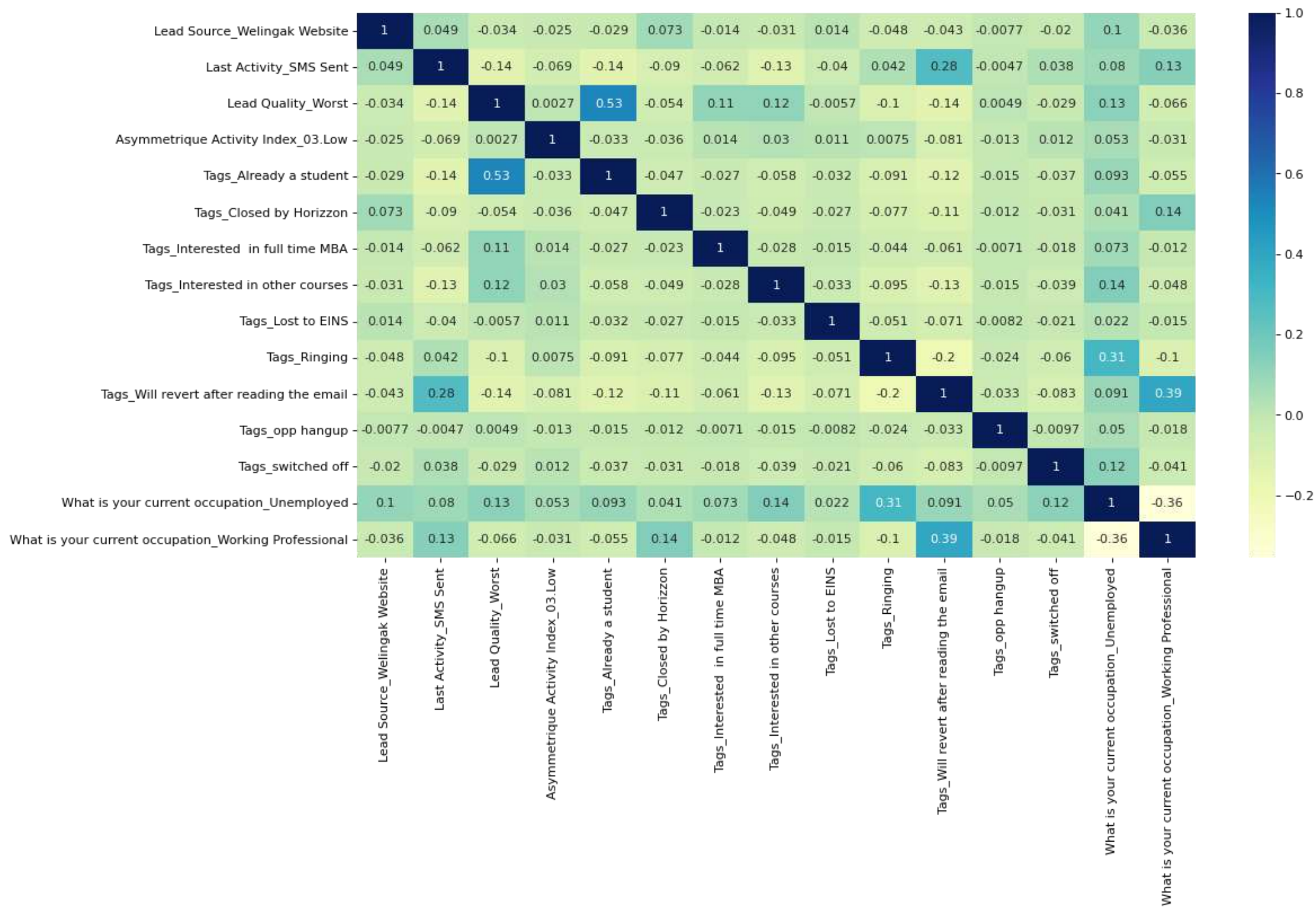






# Correlation among numeric variables





# Data Preparation for Logistic Regression

- Converting some binary variables (Yes/No) to 0/1
- Create dummy variables for categorical variables
- *Creating dummy variables for the remaining categorical variables and dropping the level called 'Unknown' which represents null/select values*
- Dropping the repeated features from the data frame

## Split Train and Test set & Feature scaling

- Split data set into train and test data set
- Scaling of features for standardization (numerical data)

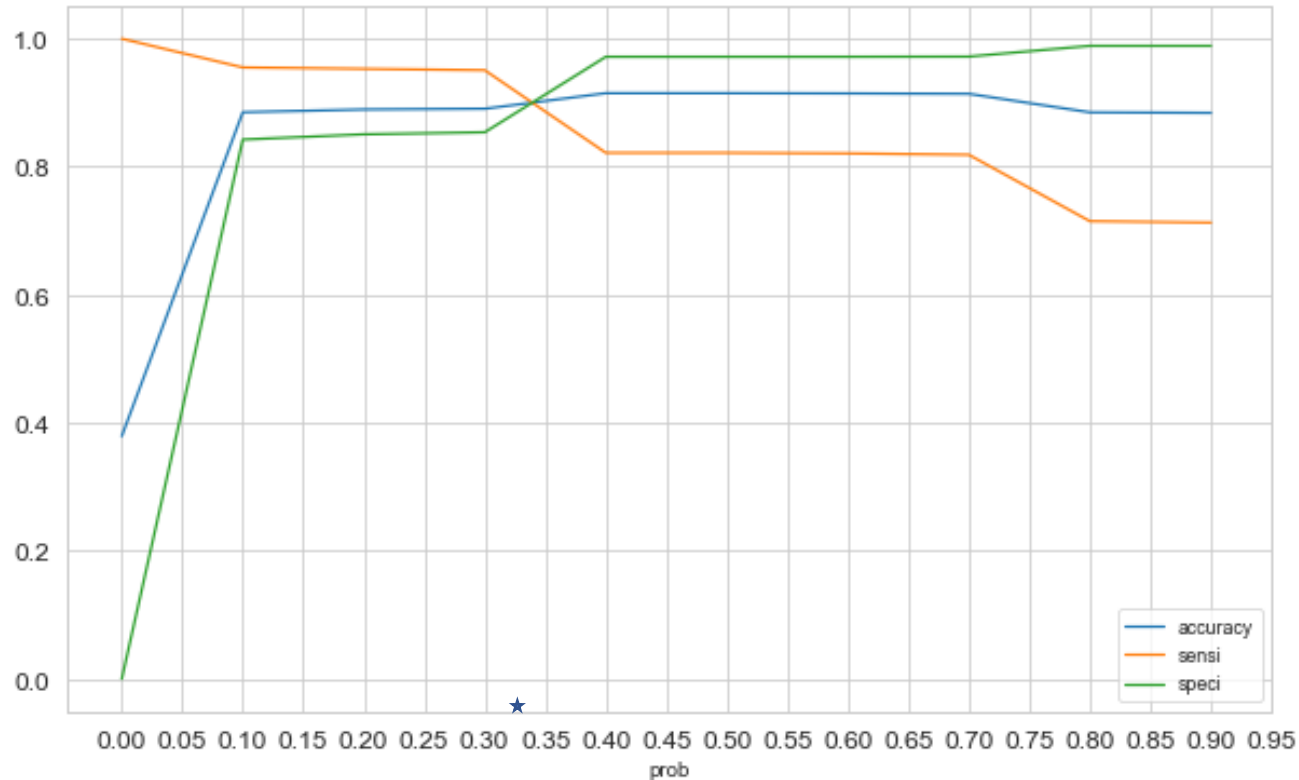
## Model Building

- Feature selection using RFE on the trained data set
- Find out optimal model using Logistic Regression
- Calculate p-values and VIF values for each trained model
- Find out correlation between numeric variables
- Plot the ROC curve and calculate GINI
- Find out the optimal cutoff probability for the data set from default
- Find out the precision and recall trade off (in our case for 80%)
- Calculate accuracy, sensitivity, specificity, precision, recall, f1 score, Positive predictive value, negative predictive value from confusion matrix etc.

## Plotting the ROC Curve

- It shows the tradeoff between sensitivity and specificity (any increase in sensitivity will be accompanied by a decrease in specificity).
- The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the test.
- The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test.

## Model Evaluation (Accuracy, Sensitivity, Specificity on Trained data set)

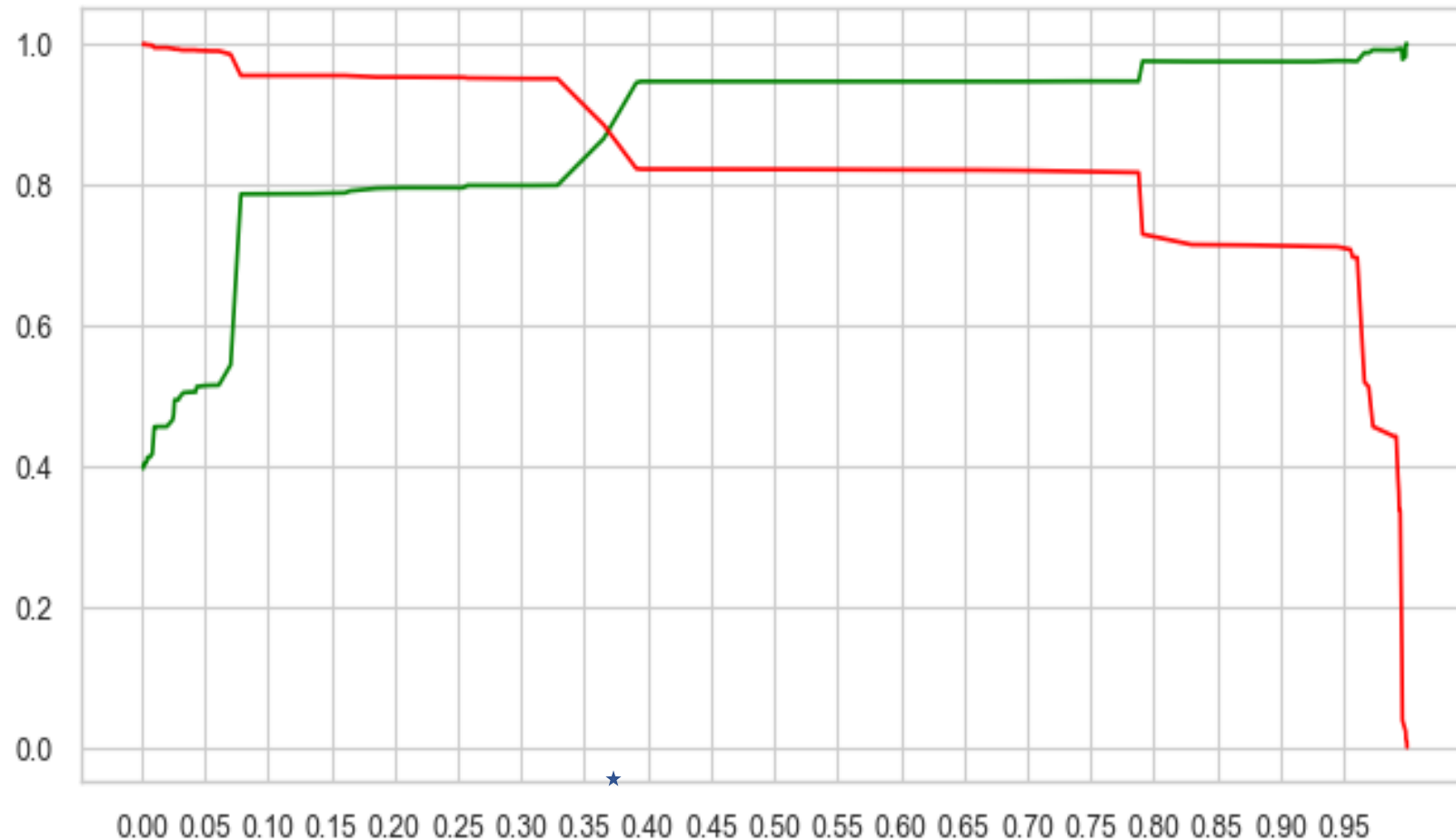


Accuracy – 0.90  
Sensitivity – 0.88  
Specificity – 0.91

**From the curve above, 0.33 is the optimum point to take it as a cutoff probability.**



## Model Evaluation (Precision and Recall on Trained data set)



Precision – 0.86  
Recall – 0.88  
F1 score – 0.875

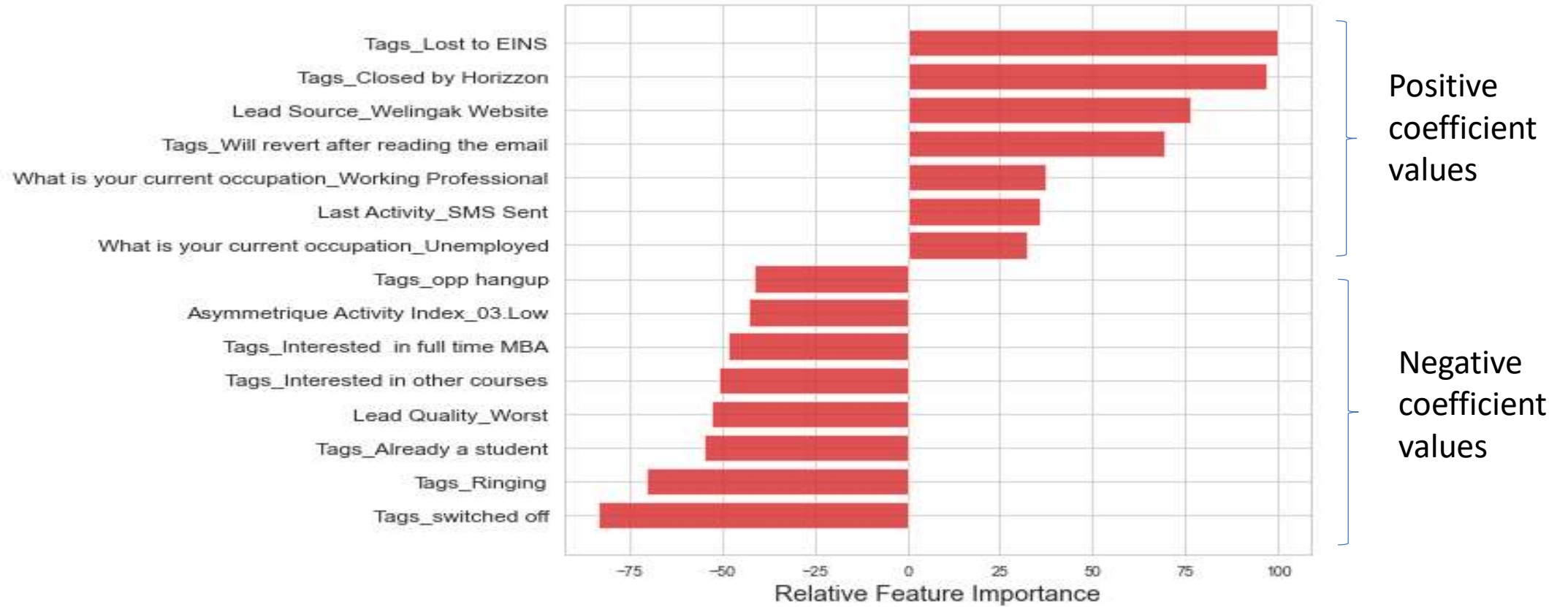
From the precision-recall graph above, we get the optimal threshold value as close to .37. However our business requirement here is to have Lead Conversion Rate around 80%.

This is achieved with our earlier threshold value of around 0.33. So we will stick to this value.

## Model Evaluation (Accuracy, Sensitivity, Specificity on Test data set)

- Accuracy – 0.90
- Sensitivity – 0.89
- Specificity – 0.90
- Precision – 0.86
- Recall – 0.89

# Feature Importance



## Summary

- Optimal cutoff probability is that prob where we get balanced sensitivity and specificity.
- After trying several models, we finally chose a model with the following characteristics.
  - All variables have p-value  $< 0.05$ .
  - All the features have very low VIF values, meaning, there is `hardly any multicollinearity` among the features. This is also evident from the heat map.
  - The overall accuracy of 0.914 at a probability threshold of 0.33 on the test dataset is also very acceptable.
- While we have checked both sensitivity & specificity as well as precision & recall metrics , we have considered the optimal cut off based on sensitivity & specificity for calculating the final prediction.
- Accuracy, sensitivity & specificity values of test set are approx. closer to values calculated using trained data set.
- Depending on the business requirement, we can increase or decrease the probability threshold which ultimately results in either decrease or increase in sensitivity and increase or decrease in specificity of the model.
- High sensitivity will ensure that almost all leads who are likely to Convert are correctly predicted, while high specificity will ensure that leads that are on the probability of getting Converted are not selected.
- Hence, overall model seems to be good.