# Lead Score Case Study

SUBMITTED BY:
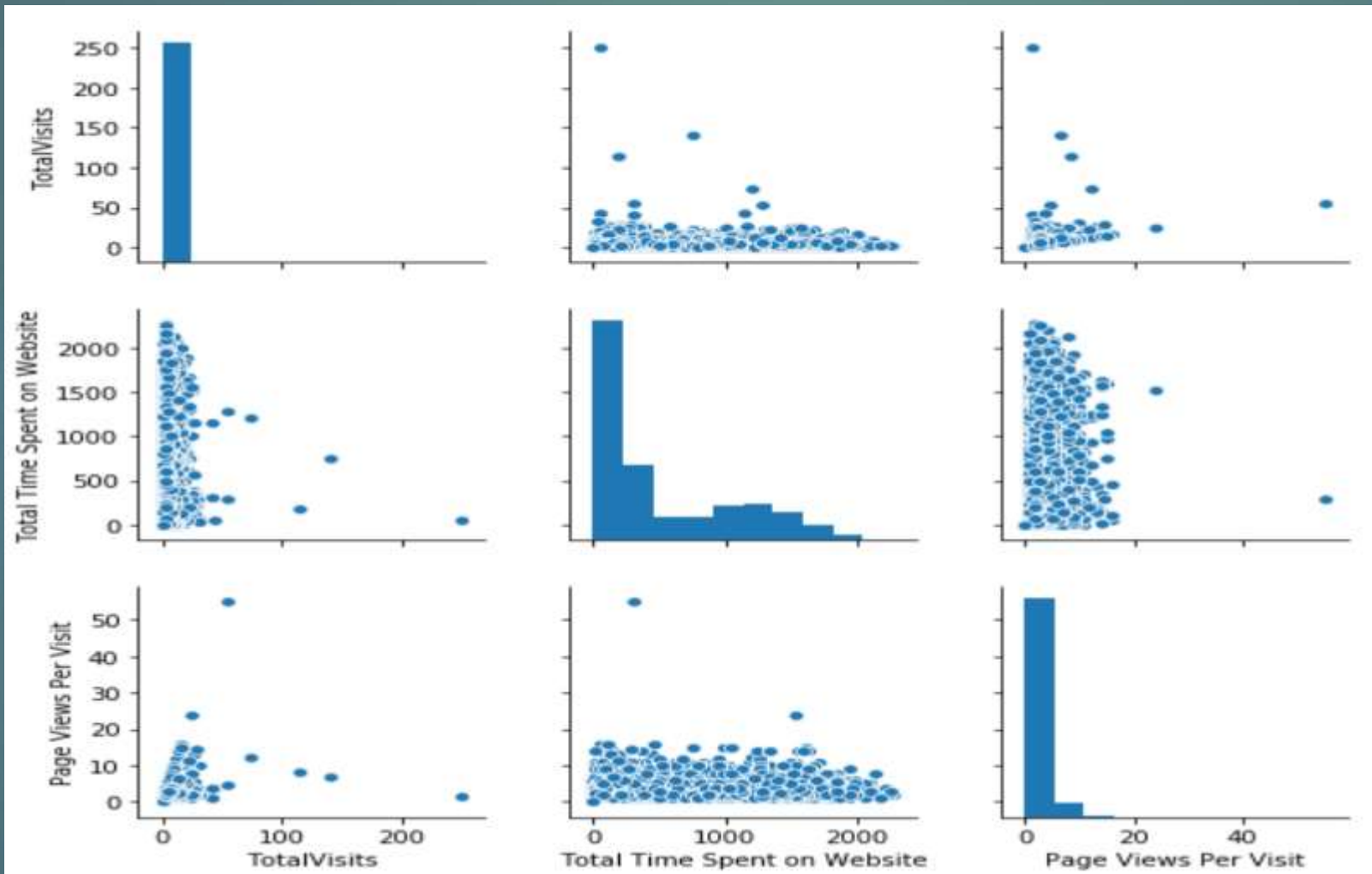
PRAVEEN MAURYA

SOURAV KARMAKAR

# Problem Statement

► An 'X' Education Company, offers online courses to industry professionals, gets many people landing on their website from the advertisements posted on various platforms and search engines like Google. The interested people provide their contact details to the company.

► These people are called as leads and are followed by the sales team to convert them, i.e. sell them the course. The conversion to lead ratio is currently 30%. The company wants to improve this ratio by following up only the potential leads that would most probably be interested in buying the course and filter out other.

► As a Data Scientist, we have to build a Machine Learning model using Logistic Regression that helps the company in achieving this and reducing the manpower in irrelevant follow ups by the sales team.

# Reading, understanding and cleansing data

▶ The given data set has 37 columns and 9240 rows with two unique columns, 'Lead Number' and 'Prospect ID'. It has few insignificant columns, 'Tags', 'Asymmetrique Activity Index','Asymmetrique Profile Index','Asymmetrique Profile Score','Asymmetrique Activity Score'.

▶ It has many non-responsive values with name 'Select' that we replaced with null values. It has 3 columns with more than 50% null values.

▶ After removing the insignificant variables, variables with high null values and records with more than 5 null values we are left with 28 variables and 9191 rows. 99.47% of rows have been retained.

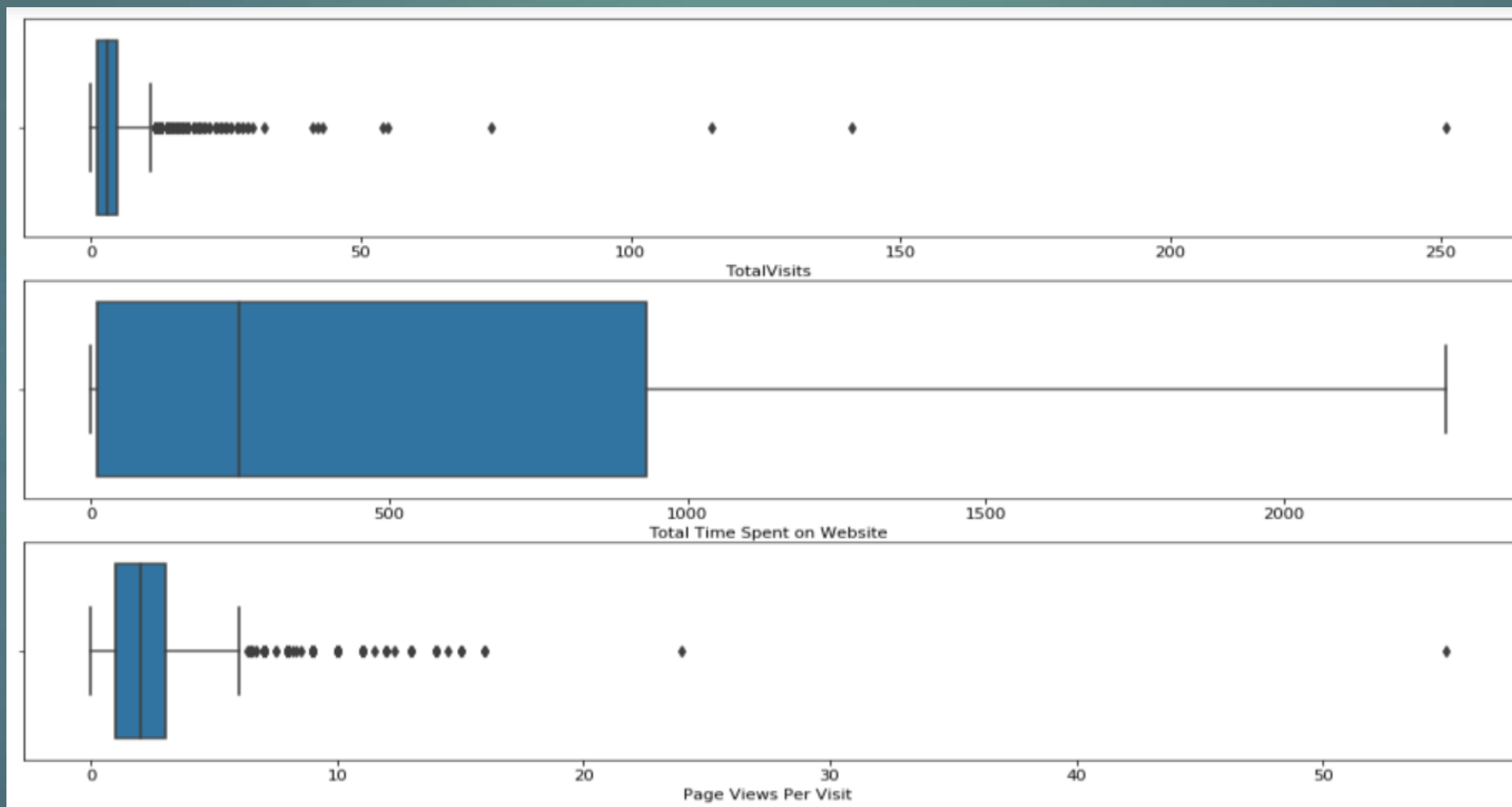▶ Imputed the null values in remaining variables with mean/median/mode.

# Visualizing data

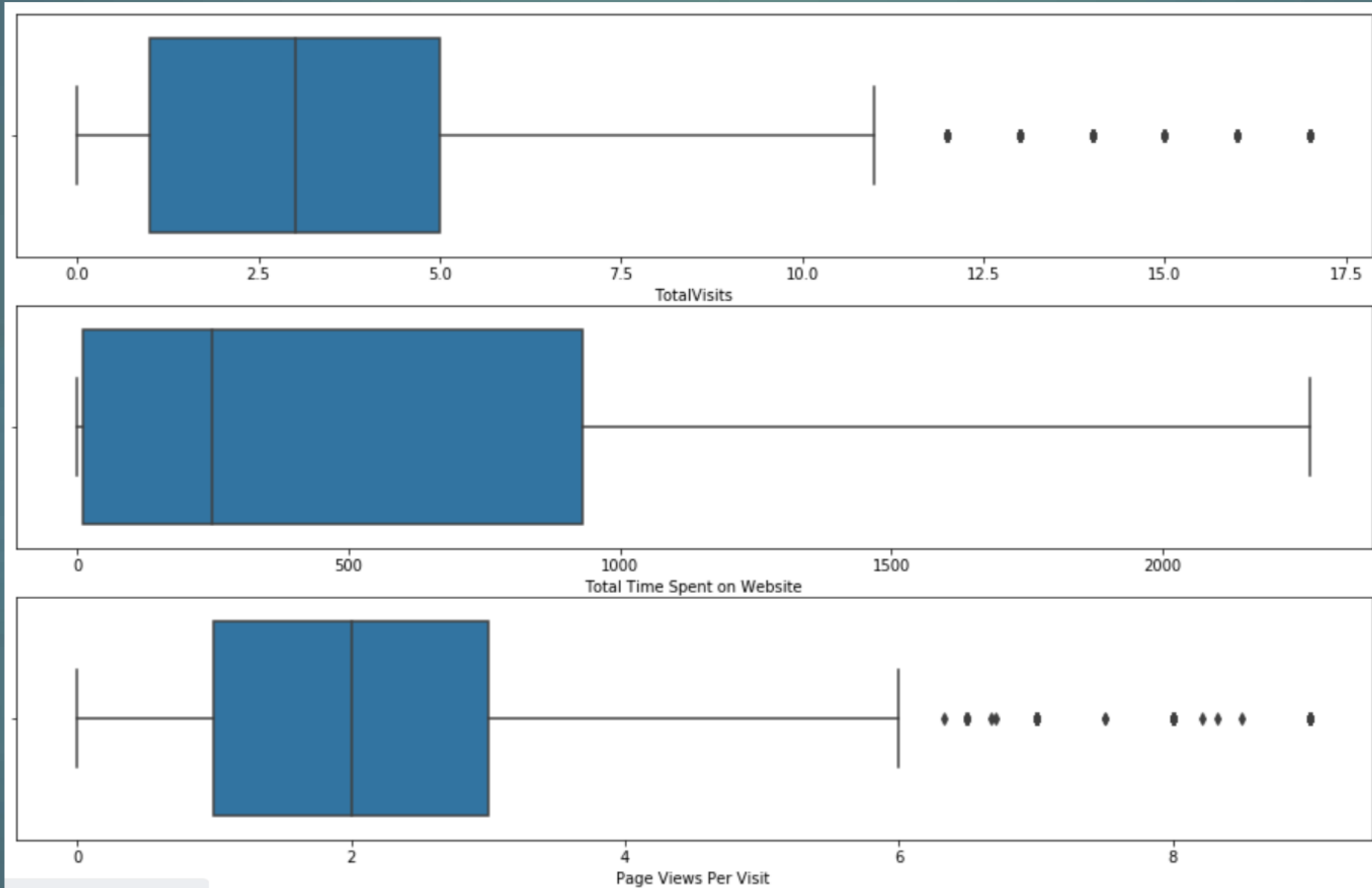► From the below pairplots we can that the variable are not collinear.

# Outlier Analysis

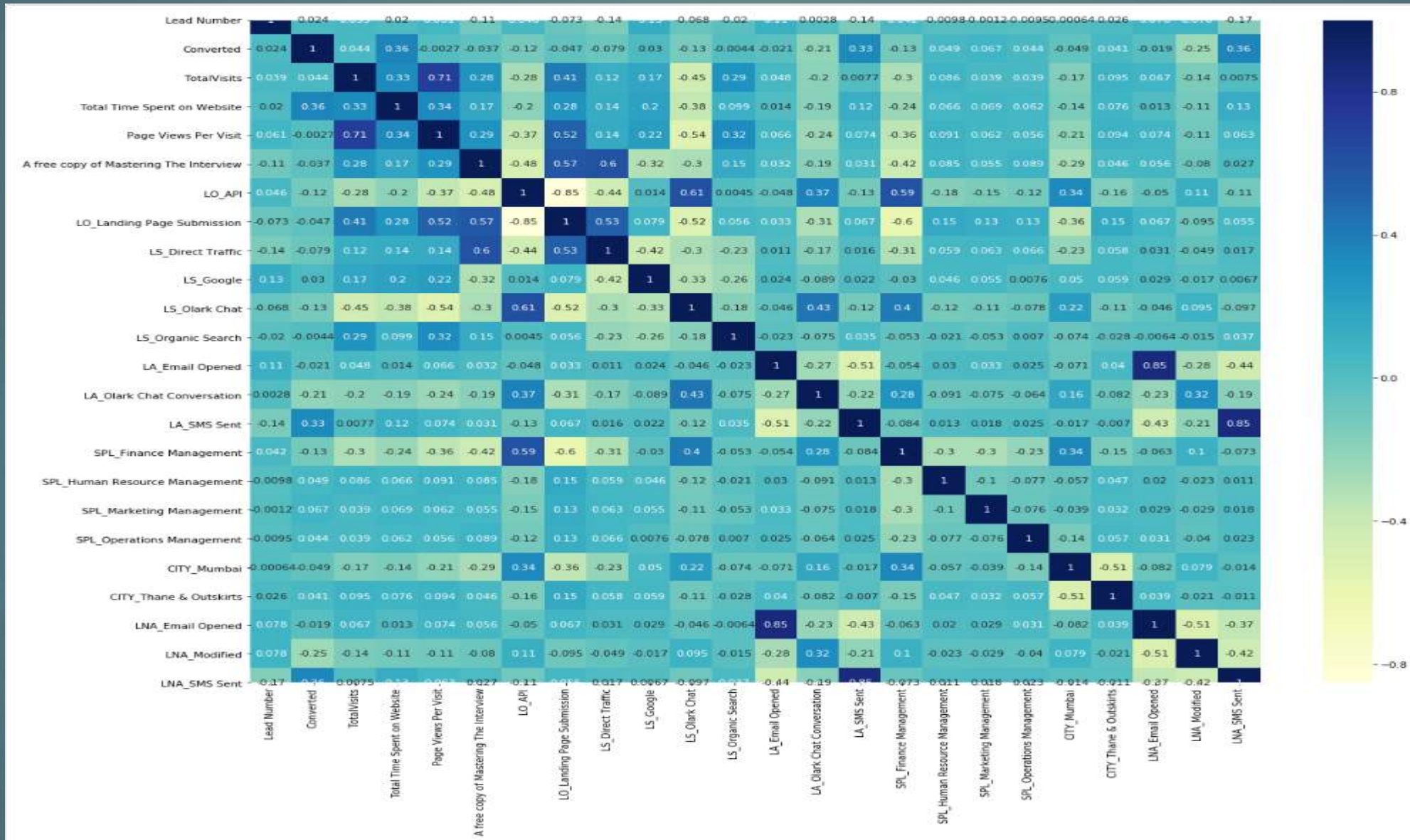▶ Outliers found in 'Total Visits' and 'Page Views Per Visit'

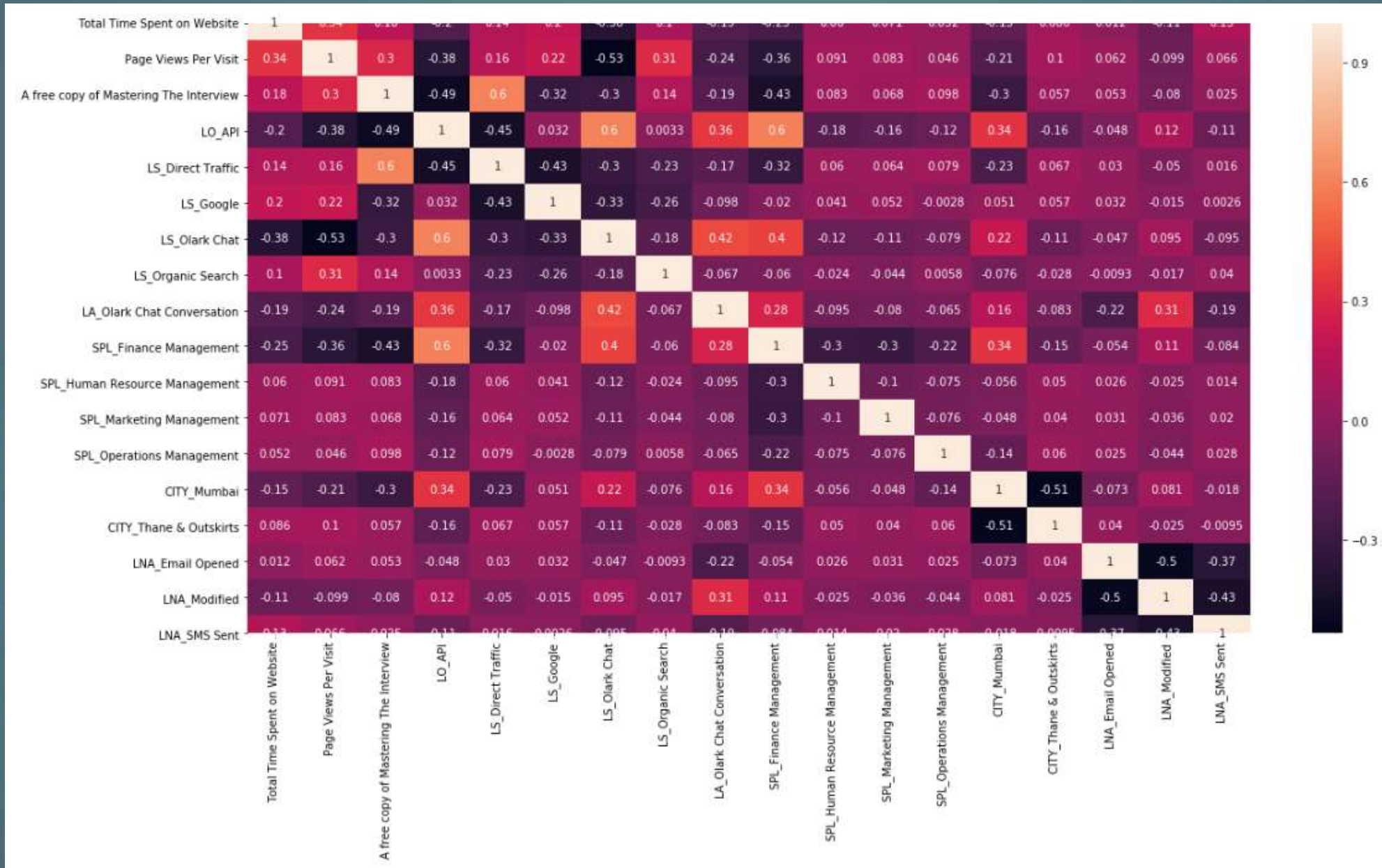► Treated outliers by capping the data with an upper limit of 99%

# Data Preparation

- Removed 16 variables with very highly skewed values

- Mapped Yes/No values to binary 1/0

- Performed bucketing to keep most significant values in categorical variables and treated remaining as 'Others'

- Created dummy variables for 6 columns and assigned them prefixes

- Dropped columns for which dummy variables have been created and concatenated the dummy columns with the original dataframe and called it 'Lead'

- We have 24 columns and 9191 rows in Lead dataframe

- Created X=Predictor variables and y=Response variable(Converted)

- Split this data into train and test with a ratio of 70 : 30.

- Scaled the numerical variable using StandardScaler in X_train

- The lead conversion rate is 38.34%

► Plotted correlation matrix and dropped highly correlated variables from both train and test datasets

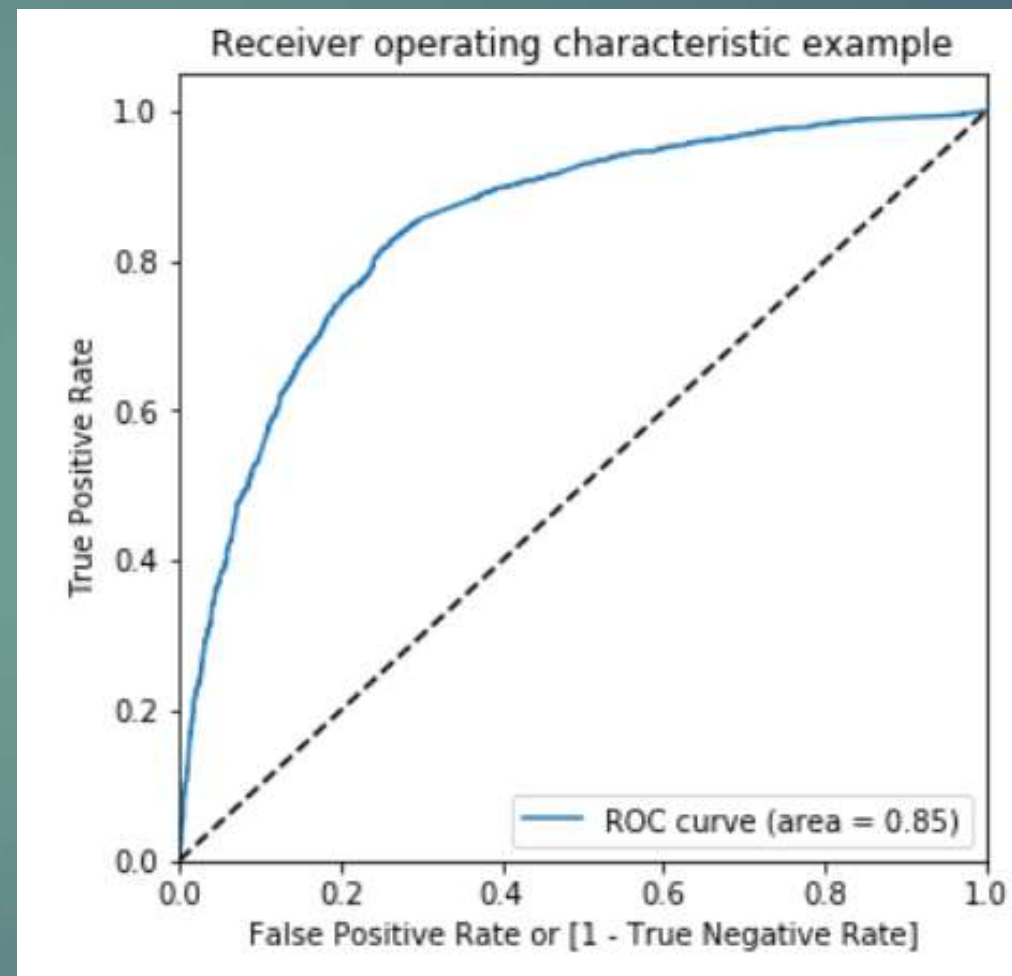► Correlation matrix after dropping highly correlated variables.

# Building Model

▶ Using Recursive Feature Elimination(RFE) selected top 13 most significant variables and created of list of these variables

▶ Added constant and build the model with 13 variables

▶ Variables 'SPL_Operations Management' and 'LNA_Email Opened' were statistically insignificant as their p-values were high. The VIF values were low

▶ Predicted using the above model and defined initial cut-off of 0.5 conversion probability. Lead Score(Conversion Probability > 0.5) ➜ 1

▶ Sensitivity = 66%, Accuracy = 78% ,based on confusion matrix created for above predicted values

▶ After dropping all insignificant variables over iteration we were left with 11 significant variables

▶ There is no significant difference in accuracy and sensitivity at this stage for the model
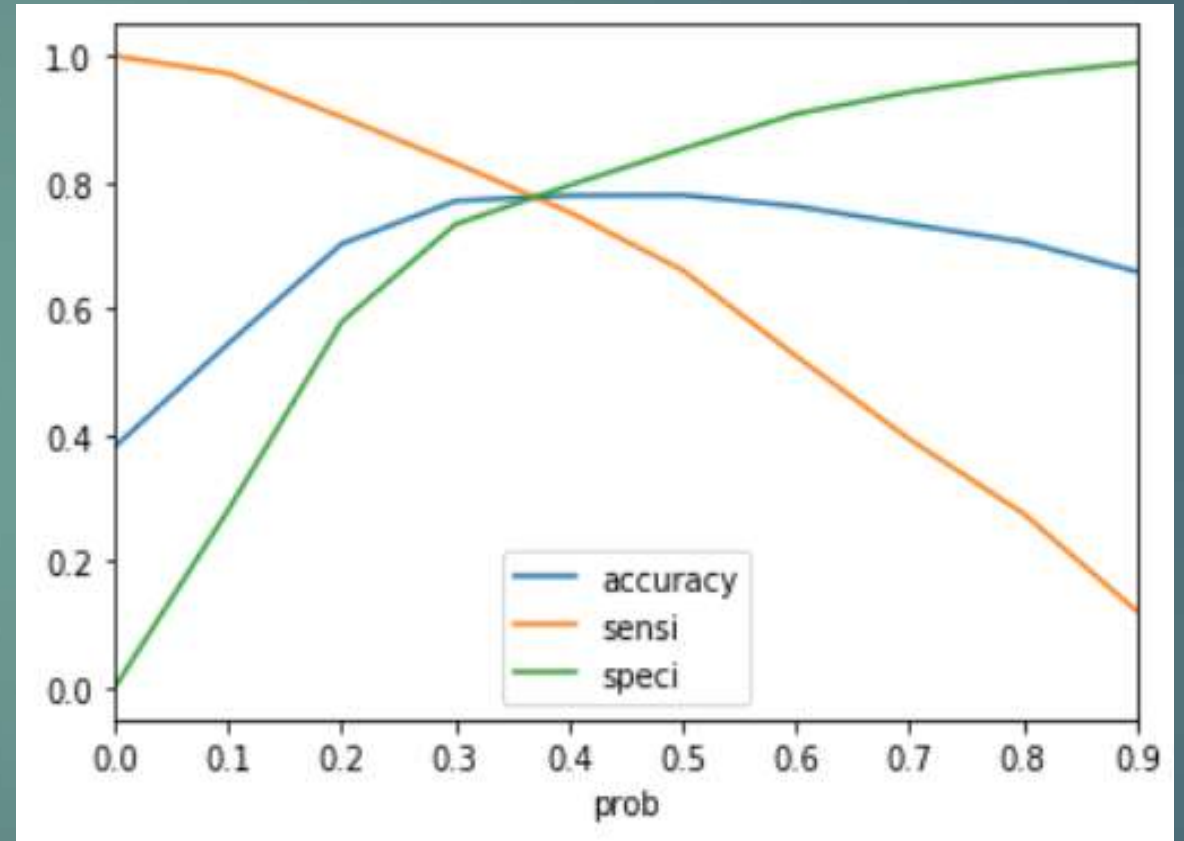
# ROC Curve to get the best cut-off

▶ From the ROC curve we can see that significant amount of area is under the curve which suggests that the built model is a good model.
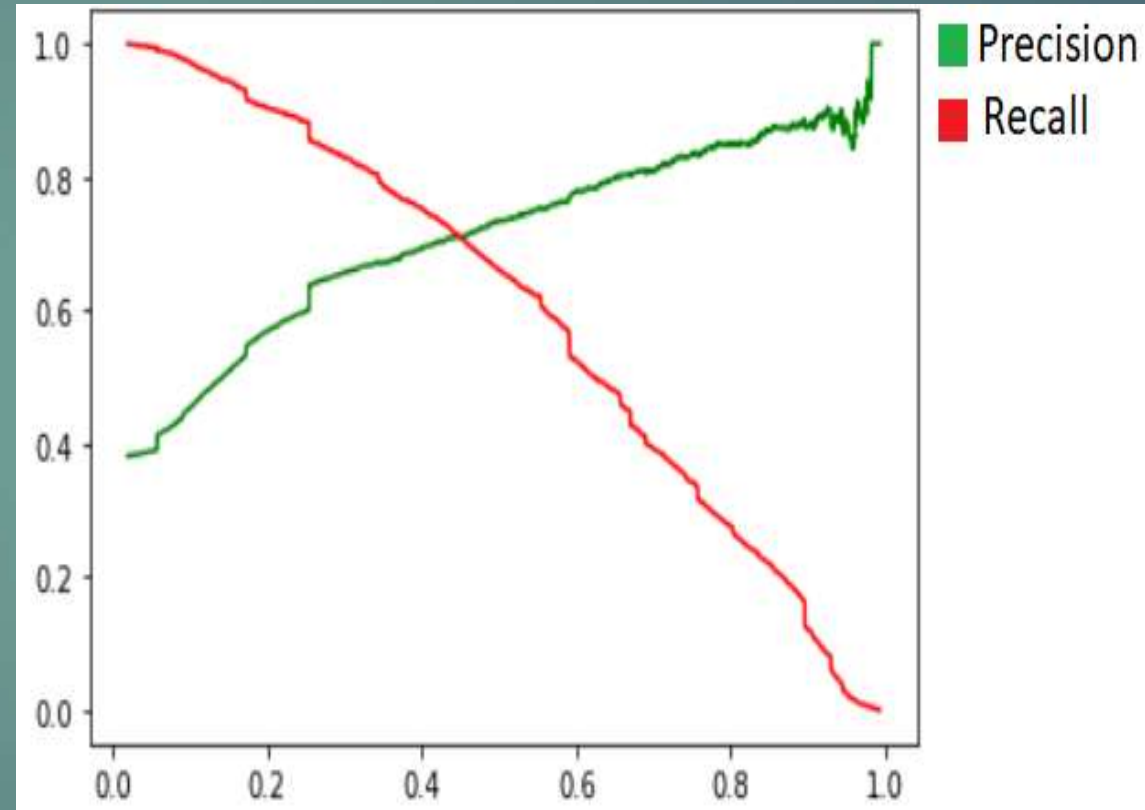


Receiver operating characteristic example

# Optimum Cut-off ?

- From the plot between accuracy, sensitivity and specificity it is clear that the cut-off lies between 0.3 and 0.4

- Chose the cut-off of 0.3 and updated the final predicted values

- For this model:

  - Accuracy ➔ 77%

  - Sensitivity ➔ 73%

  - Specificity ➔ 83%

# Precision – Recall trade off

- From the plot between precision and recall we can see that they intersect at around 0.4

# Making Predictions

▶ Used the learnt model from training data set and applied on the test dataset

▶ Following are the quality metrics after making predictions on test dataset:

 ▶ Accuracy ➜ 78%

 ▶ Sensitivity ➜ 82%

# Conclusion

- The top 3 variables affecting the lead conversion are:
  - Lead Source
  - Last Notable Action
  - Last Action
- The top 3 dummy variables affecting the lead conversion are:
  - LS_Direct Traffic
  - LS_Organic Search
  - LS_Google