# LEAD SCORING CASE STUDY
## LOGISTIC REGRESSION

PRESENTED BY: DISHARI DEBNATH

# PROBLEM STATEMENT

- X Education sells online courses to industry professionals. The company markets its courses on several websites and search engines like Google. When people fill up a form providing their email address or phone number, they are classified to be a lead. Although X Education gets a lot of leads, its lead conversion rate is very poor.

# BUSINESS GOAL

- The company wishes to identify the most potential leads, also known as 'Hot Leads'. The goal is to build a model which assigns 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 target is to achieve a lead conversion rate of 80%
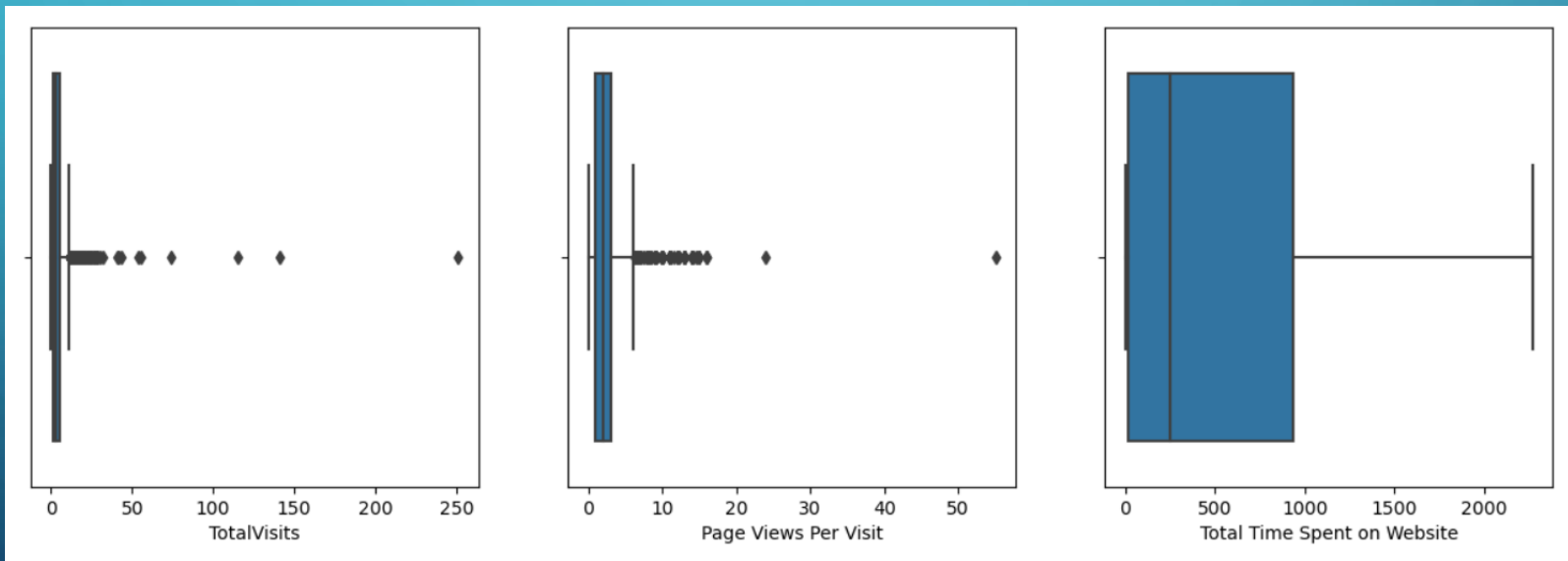
# APPROACH

- Data importing/loading
- Data preparation
  - encoding categorical variable
  - handling null values
- EDA
  - Univariate Analysis
  - Outlier detection
  - Checking data imbalance
- Dummy Variable Creation
- Test Train Split
- Feature Scaling
- Correlations
- Model Building
  - Feature Selection Using RFE
  - Model Improvisation
- Model Finalisation
- Model Evaluation
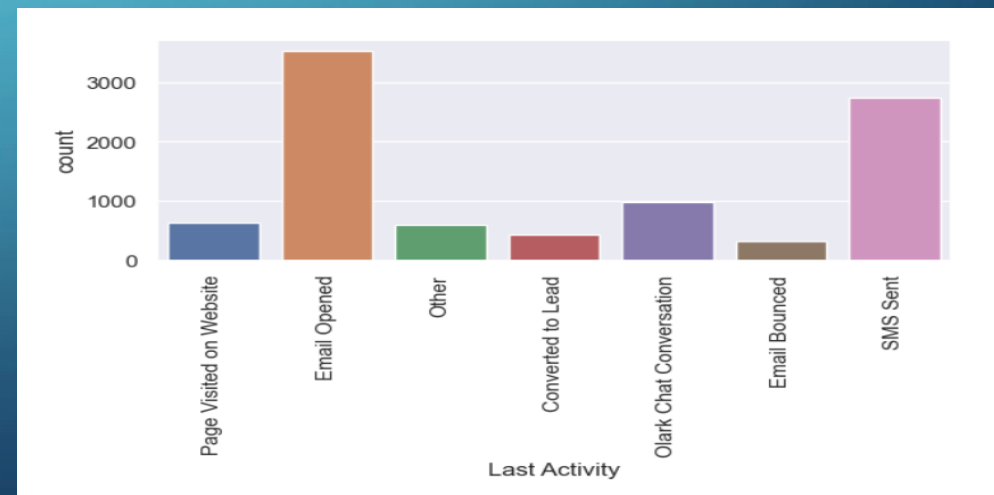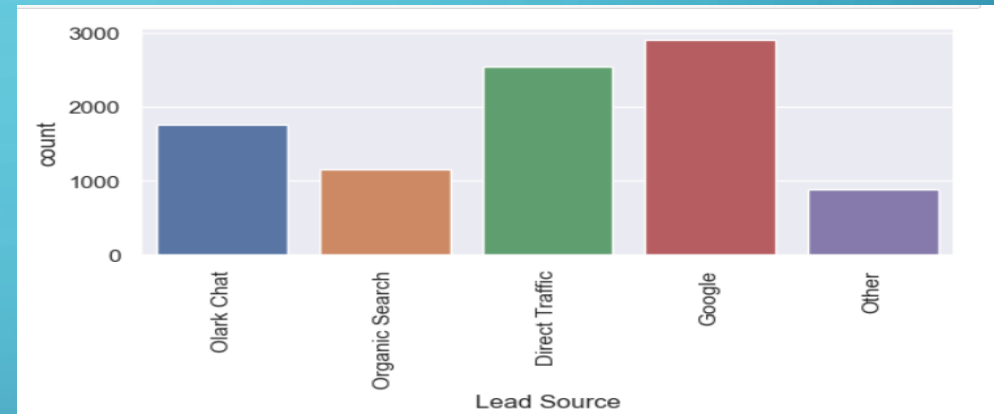- Recommendation

# OUTLIER DETECTION

To detect the outliers we performed univariate analysis on numerical dtype columns.

Outliers are Present in both the Variables 'TotalVisits' and 'Page Views Per Visit' it should be treated and the value are spreaded above median highly in 'Total Time Spent on Website'.
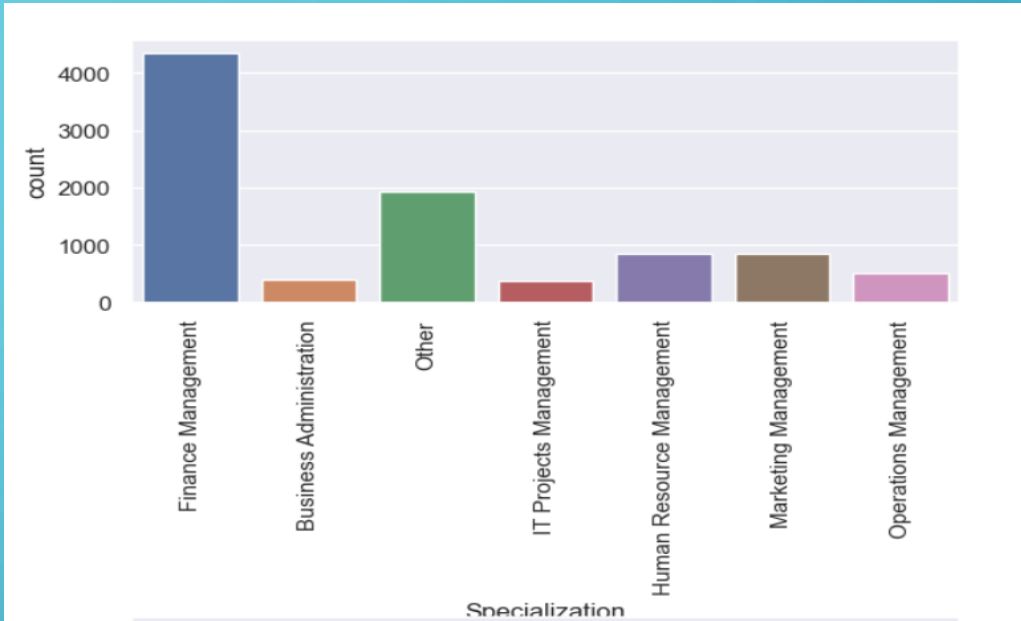
# UNIVARIATE ANALYSIS

- Lead Source Vs Converted
  - In Lead Source Direct Traffic and Google are the two main source for Leads



- Last Activity Vs Converted
  - The Number of values is High in Email Opened and SMS Sent in Last Activity
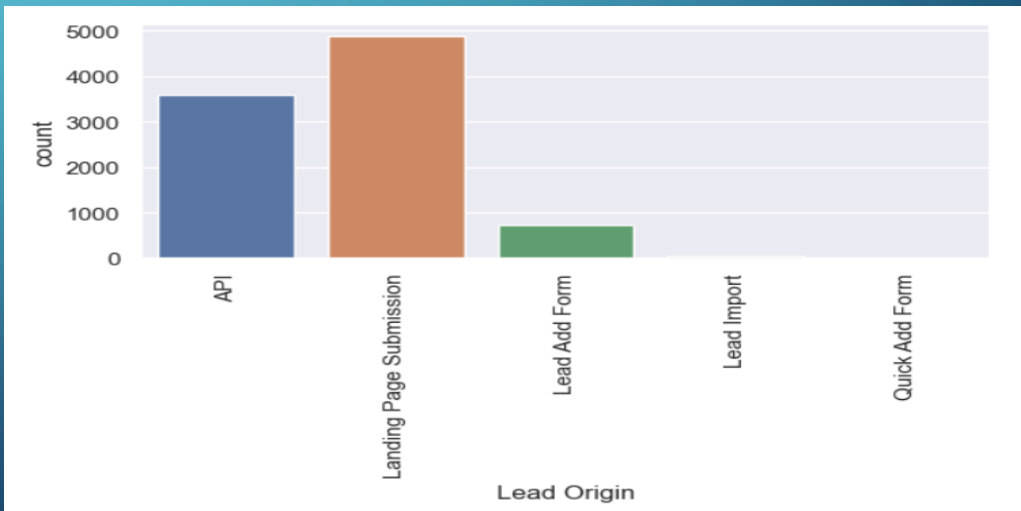
- Specialization Vs Coverted
  - Most of the people chooses Finance Management Specialization rather than other Specialization
  - The IT Project management have very lees so that most of the People not prefered this Specialization
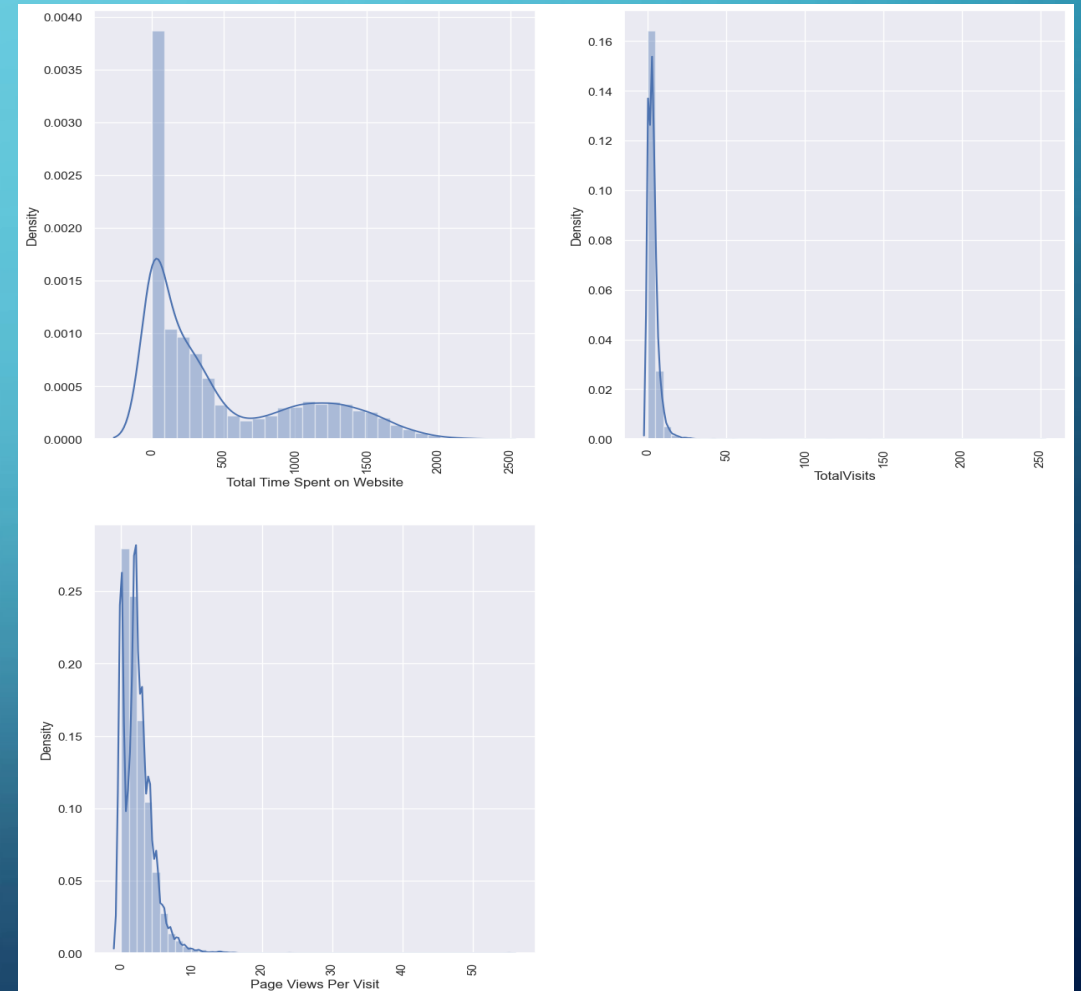


- Lead Source Vs Converted
  - Landing Page Submission and API to be a promising Lead Source for high conversion.
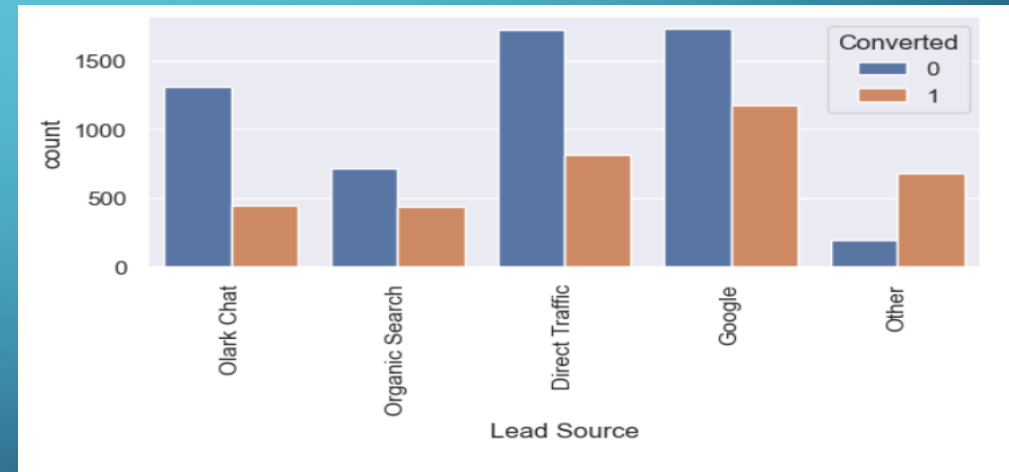
# CONTINUOUS UNIVARIATE ANALYSIS

- None of the Continueous Variables are in Normal distribution

- Presence of Outliers in Total Visits and Page Views Per Visit

- In total visits more values is between 0-50 and page views per visits 0-20
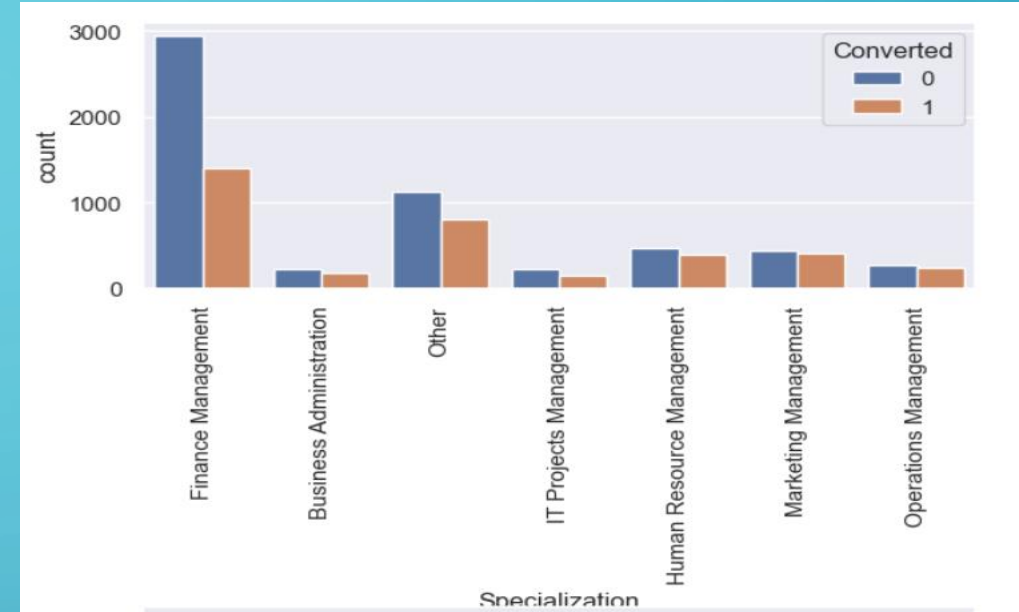
# BIVARIATE ANALYSIS
# ANALYSIS WITH RESPECT TO TARGET COLUMN CONVERTED

- Lead SourceVs Converted
  - In Lead Source The number of Hot leads is higher in Direct Traffic and Google less in Other Category
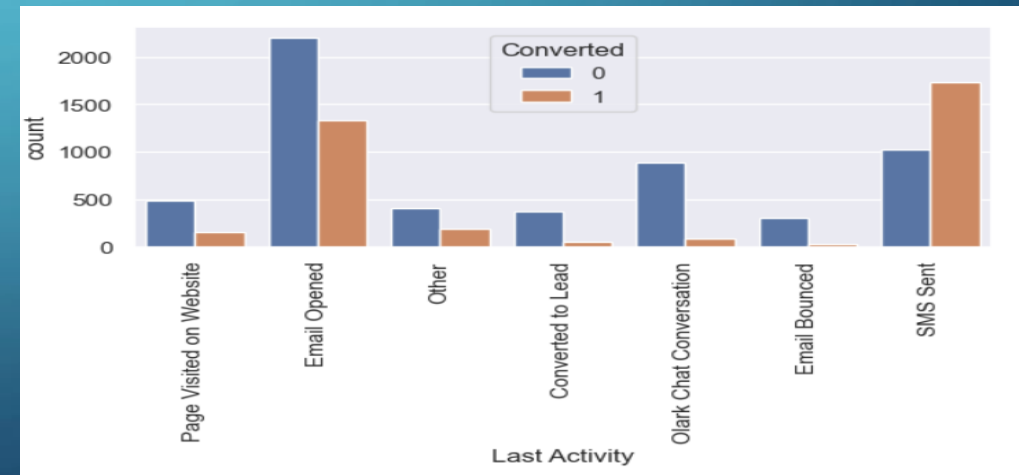
- Specialization Vs Converted
  - In Specialization the most of the leads are comes from Finance management but here Hot leads are lesseer than Cold leads.
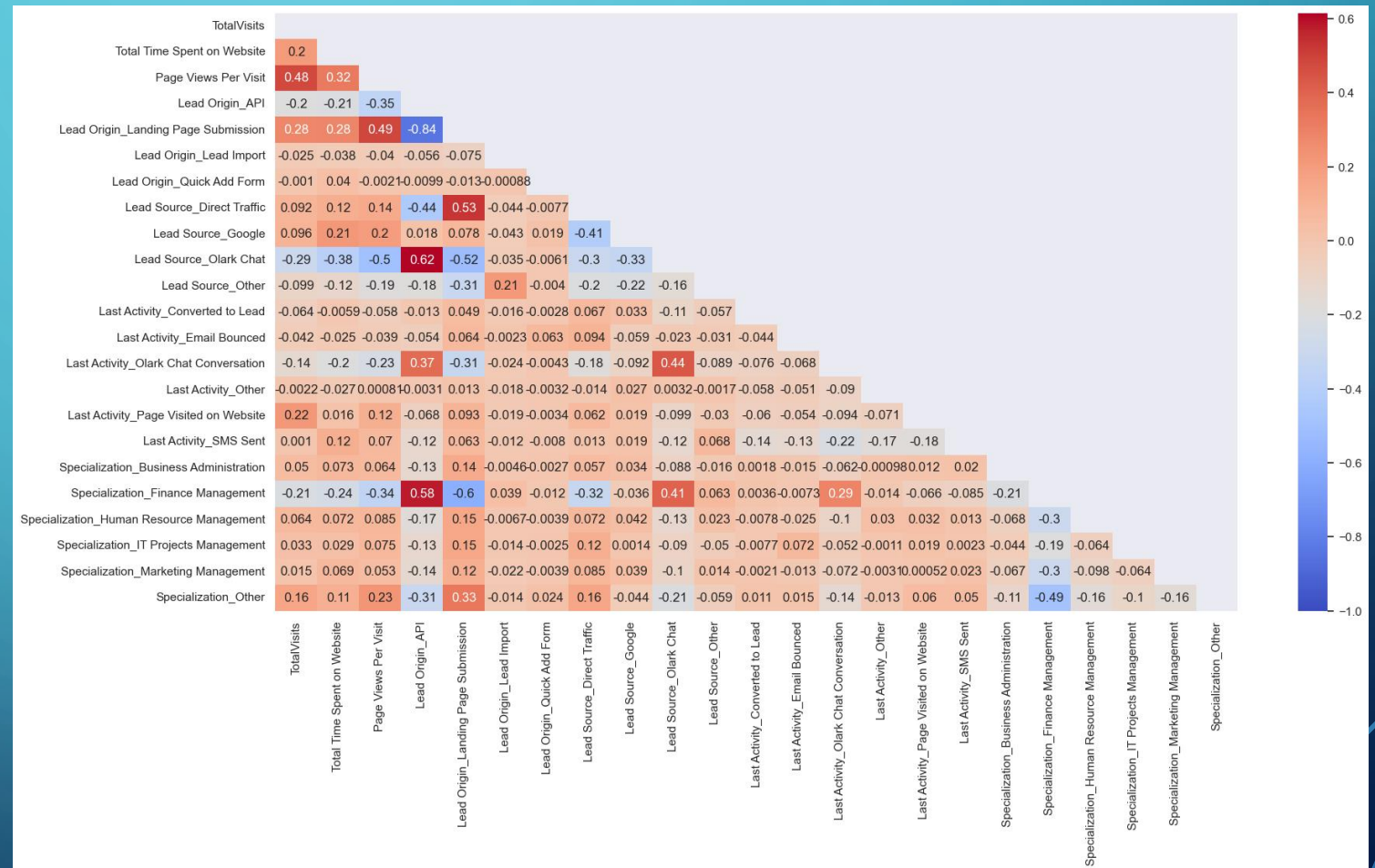
- Last Activity Vs. Converted
  - In Last Activity the number of Hot leads is higher in SMS and in EMAIL cold leads is higher than hot leads.
  - In Last Notable Activity it's mostly same as Last Activity.

# FINDING THE CORRELATION USING HEATMAP

- We can see some columns are highly correlated

- We will let RFE to decide to drop columns or not

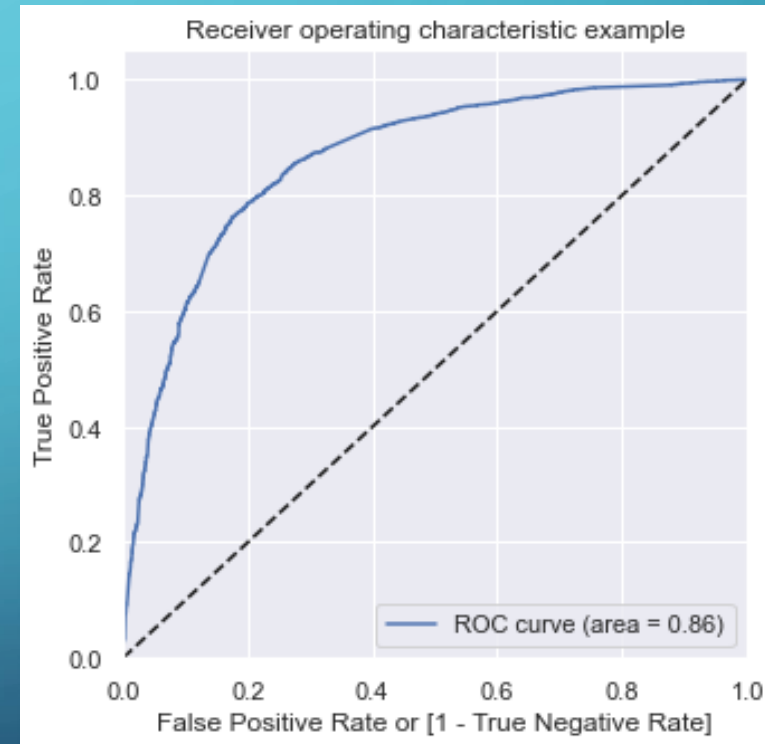# MODEL BUILDING

- Splitting into Train-Test Set

- Scale variables in train set

- Build the first model

- Use RFE to  eliminate less relevant variable

- Build next model

- Eliminate variables based on high p-value

- Check VIF value for all existing columnsPredict using train set

- Evaluate accuracy and other metric

- Predict using text set

- Precision and Recall analysis on test predictions

# ROC CURVE

- An ROC curve shows tradoff between sensitivity and specificity (increase in one will cause decrease in other). The closer the curve follows the y-axis and the top border of the ROC space then it means more area under the curve and the more accurate the test prediction will be. The Closer the curve comes to the 45degree diagonal of the ROC space the lesser the area under the curve and less the test prediction will be.
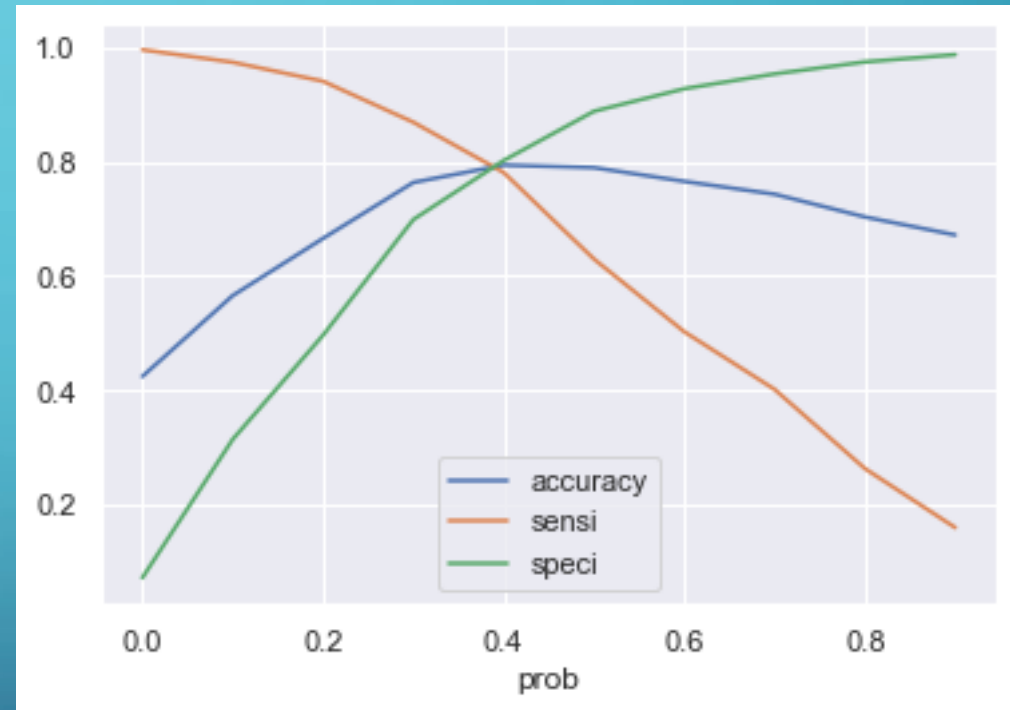
# MODEL EVALUATION

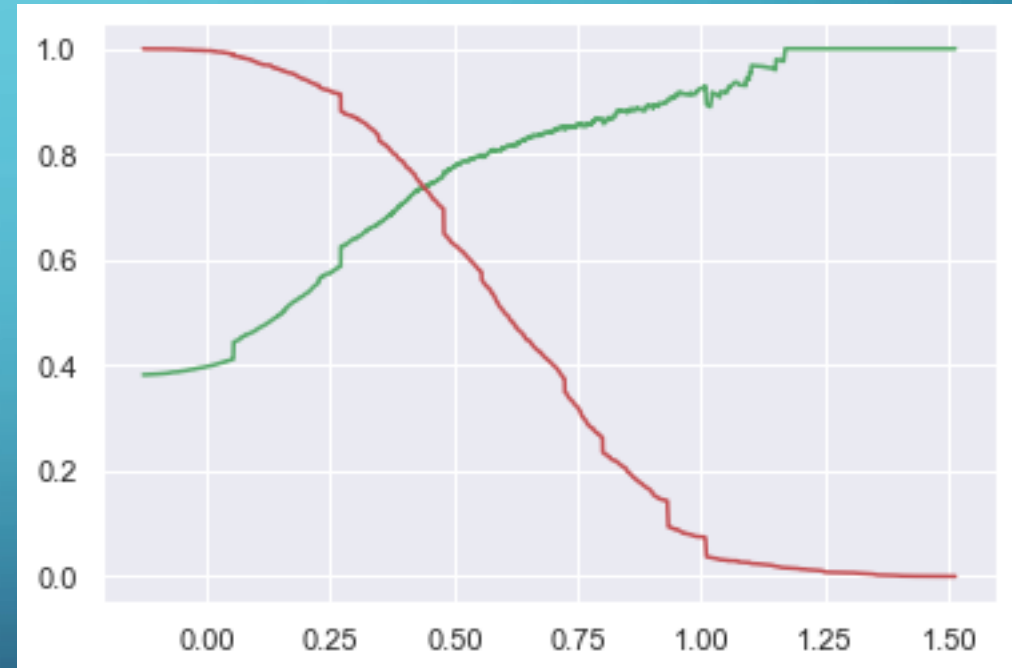- Accuracy, sensitivity and specificity for various probability cutoffs

|  | prob | accuracy | sensi | speci |
|---|---|---|---|---|
| 0.0 | 0.0 | 0.424088 | 0.996350 | 0.071464 |
| 0.1 | 0.1 | 0.565708 | 0.974453 | 0.313843 |
| 0.2 | 0.2 | 0.665894 | 0.941200 | 0.496252 |
| 0.3 | 0.3 | 0.763915 | 0.869424 | 0.698901 |
| 0.4 | 0.4 | 0.794372 | 0.781427 | 0.802349 |
| 0.5 | 0.5 | 0.789889 | 0.629765 | 0.888556 |
| 0.6 | 0.6 | 0.765770 | 0.502433 | 0.928036 |
| 0.7 | 0.7 | 0.743352 | 0.401460 | 0.954023 |
| 0.8 | 0.8 | 0.703463 | 0.262774 | 0.975012 |
| 0.9 | 0.9 | 0.671923 | 0.159367 | 0.987756 |



- From this we can have value 3.7 as the cutoff value

# PRECISION VS RECALL

- Green-Precision

- Red-Recall

- Here we got 0.37 as the Cut-off as Precesion-Recall Threshold

# ACCURACY, SENSITIVITY, SPECIFICITY

- **Sensitivity** (true positive rate) is the probability of a positive test result, conditioned on the individual truly being positive.

- **Specificity** (true negative rate) is the probability of a negative test result, conditioned on the individual truly being negative.

# CONFUSION MATRIX

For our model

| | |
|---|---|
| 3556 | 446 |
| 913 | 1553 |

Accuracy: 0.7857142857142857

Sensitivity 0.8102189781021898

Specifitiy 0.7706146926536732



| | Actual Values | |
|---|---|---|
| | Positive | Negative |
| Predicted Values Positive | TP | FP |
| Negative | FN | TN |

Figure 1

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

$$Sensitivity = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{TN + FP}$$

# CONCLUSION FROM LOGISTIC REGRESSION MODEL

- We can see that our model is doing well in test set also

- Sensitivity means how our model is telling that actually converted and model prdecited them as as converted.

- We can see that our model is giving about .80 sensitivity.

- it means that 80 percent time our model is able to predict (actually)converted as (prdicited)converted.