# LEAD SCORE CASE STUDY

Group members :-
1. Pramod Poojary
2. Shivani Gupta

# PROBLEM STATEMENT

X Education sells online courses to industry professionals

X Education gets a lot of leads, its lead conversion rate is very poor. For example, if, say, they acquire 100 leads in a day, only about 30 of them are converted.

To make this process more efficient, the company wishes to identify the most potential leads, also known as 'Hot Leads'. If they successfully identify this set of leads, the lead conversion rate should go up as the sales team will now be focusing more on communicating with the potential leads rather than making calls to everyone.

# BUSINESS OBJECTIVE

- Select the most promising leads, i.e. the leads that are most likely to convert into paying customers.

- The company requires to build a model wherein we 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.

- Target lead conversion rate to be around 80%.

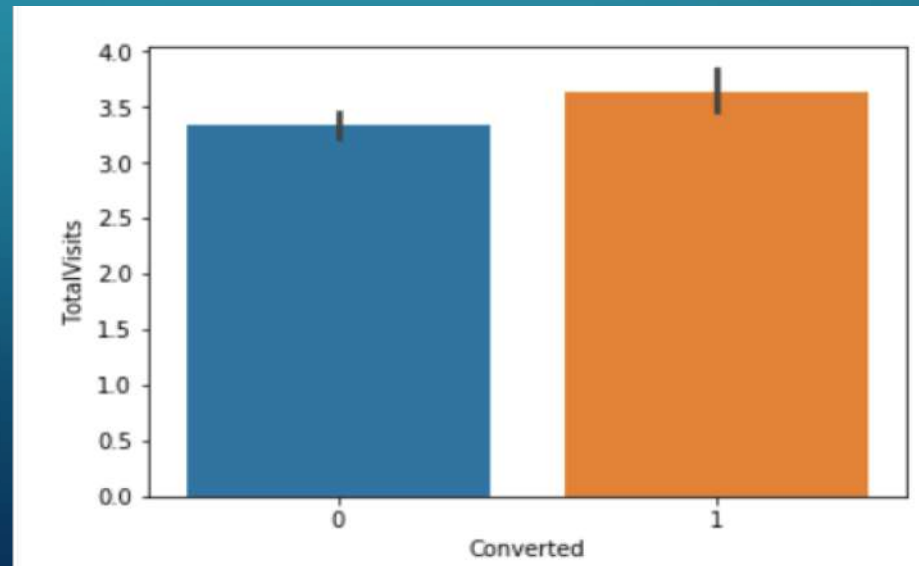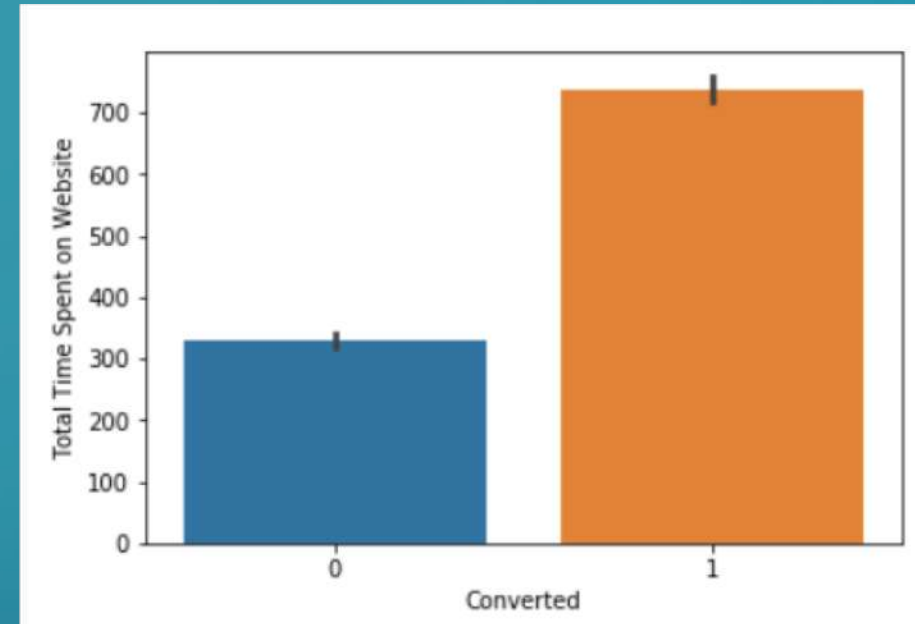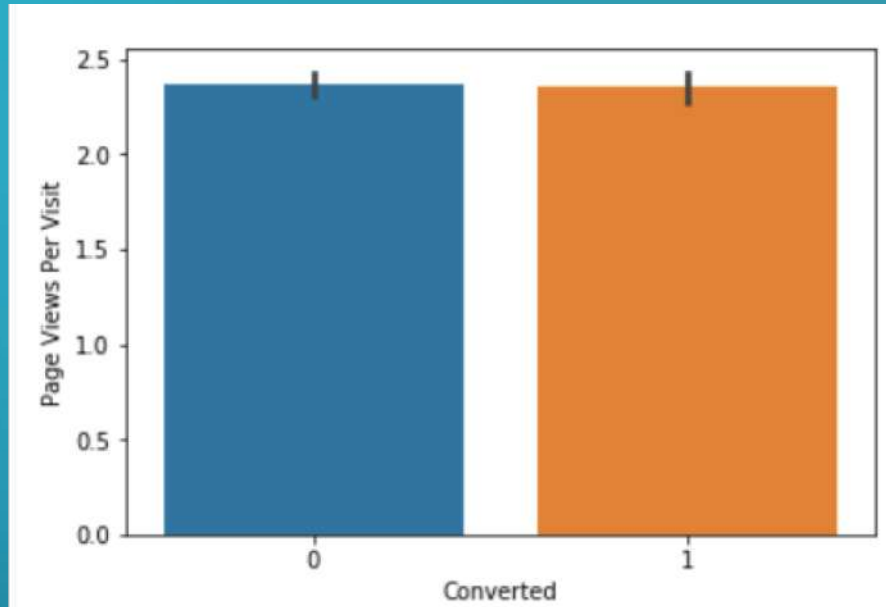- Deployment of the model for the future use.

# SOLUTION METHODOLOGY

1. Data cleaning and Data manipulation (Check and handle duplicate data, null values, missing values and outliers in data).

2. Mapping Yes/No values to 0/1, dummy variable creation and splitting data into Test-Train.

3. Feature scaling and finding correlations between variables.

4. Building Logistic Regression model

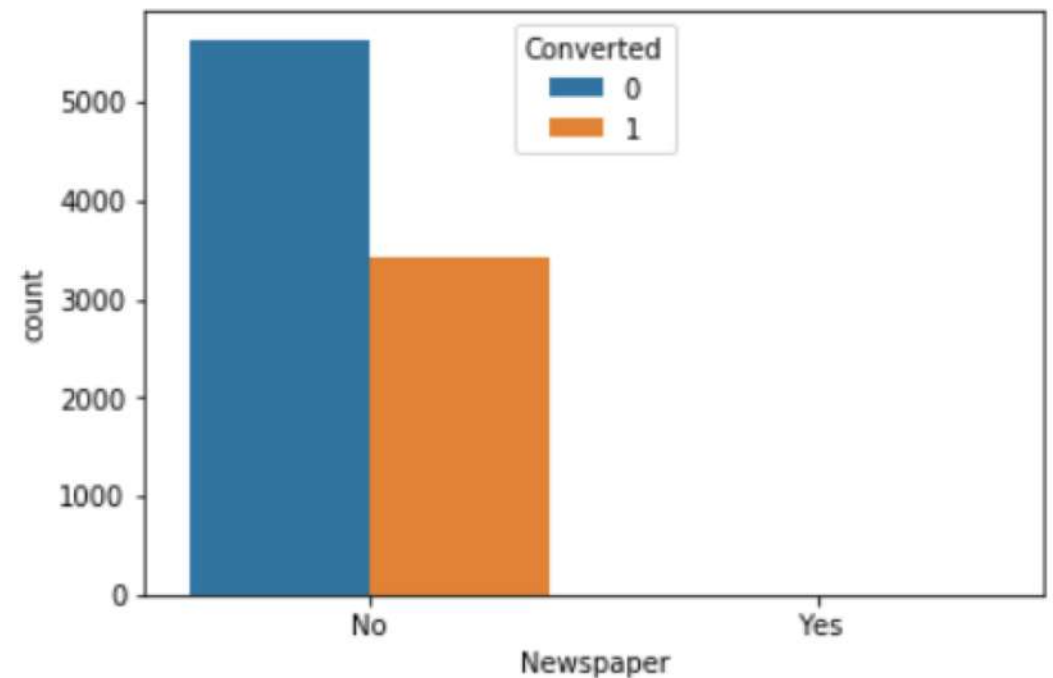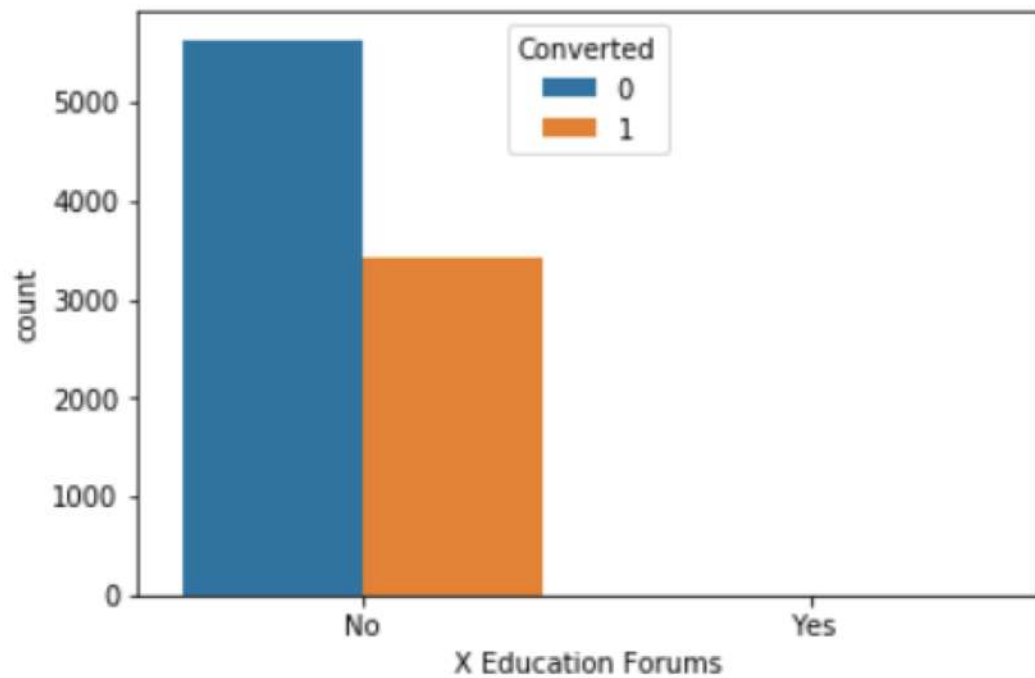5. Validation of the model

6. Conclusion

# DATA CLEANING AND MANIPULATION

- 'Select' present in many columns replaced with NaN value as it is useless.

- Insignificant columns like "City", "Country" and "Prospect ID" are dropped.

- Columns containing more than 3000 columns has been dropped.

- Retained the rows having <=5 NaN's.

- After removing above all and outliers , we left with 8997 Rows and 24 Columns.
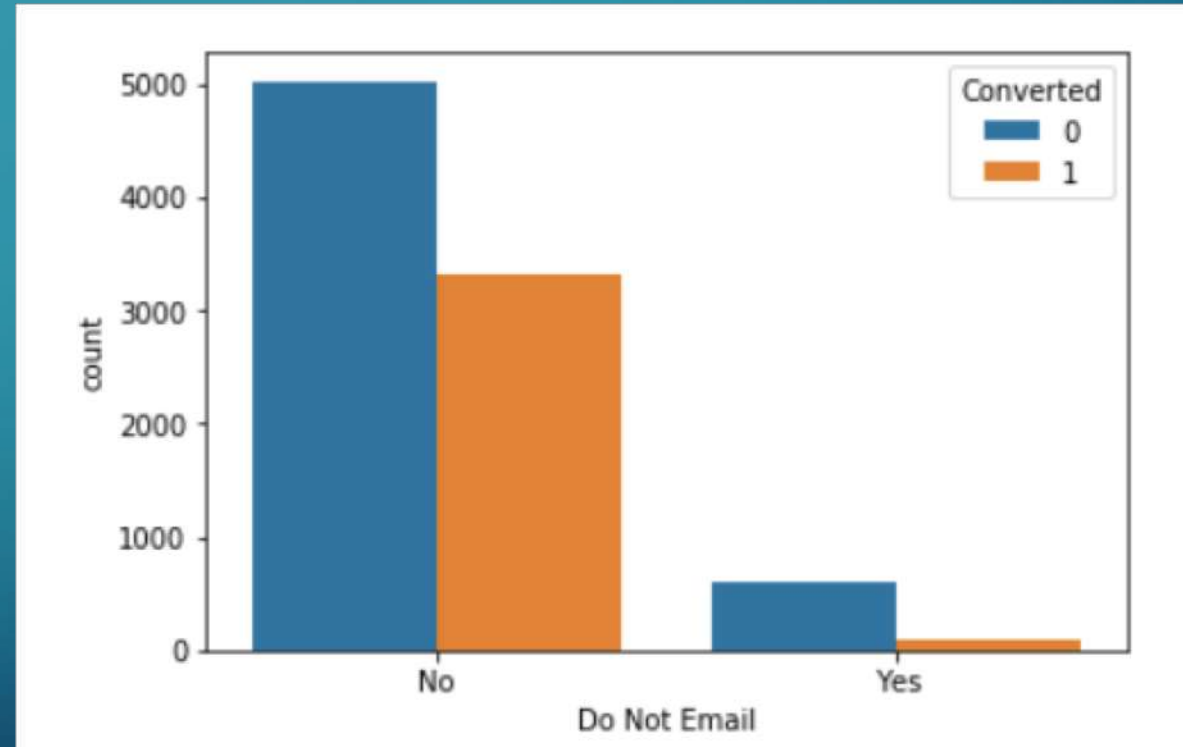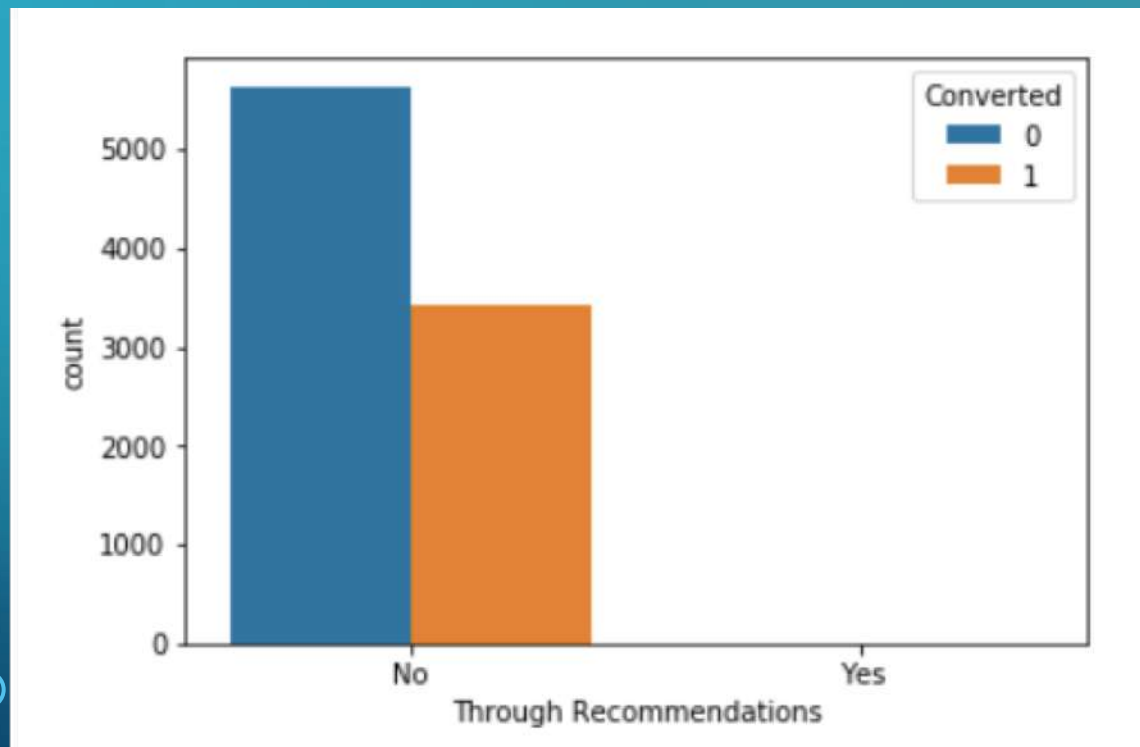
# DATA VISUALIZATION OF CONTINUOUS VARIABLES
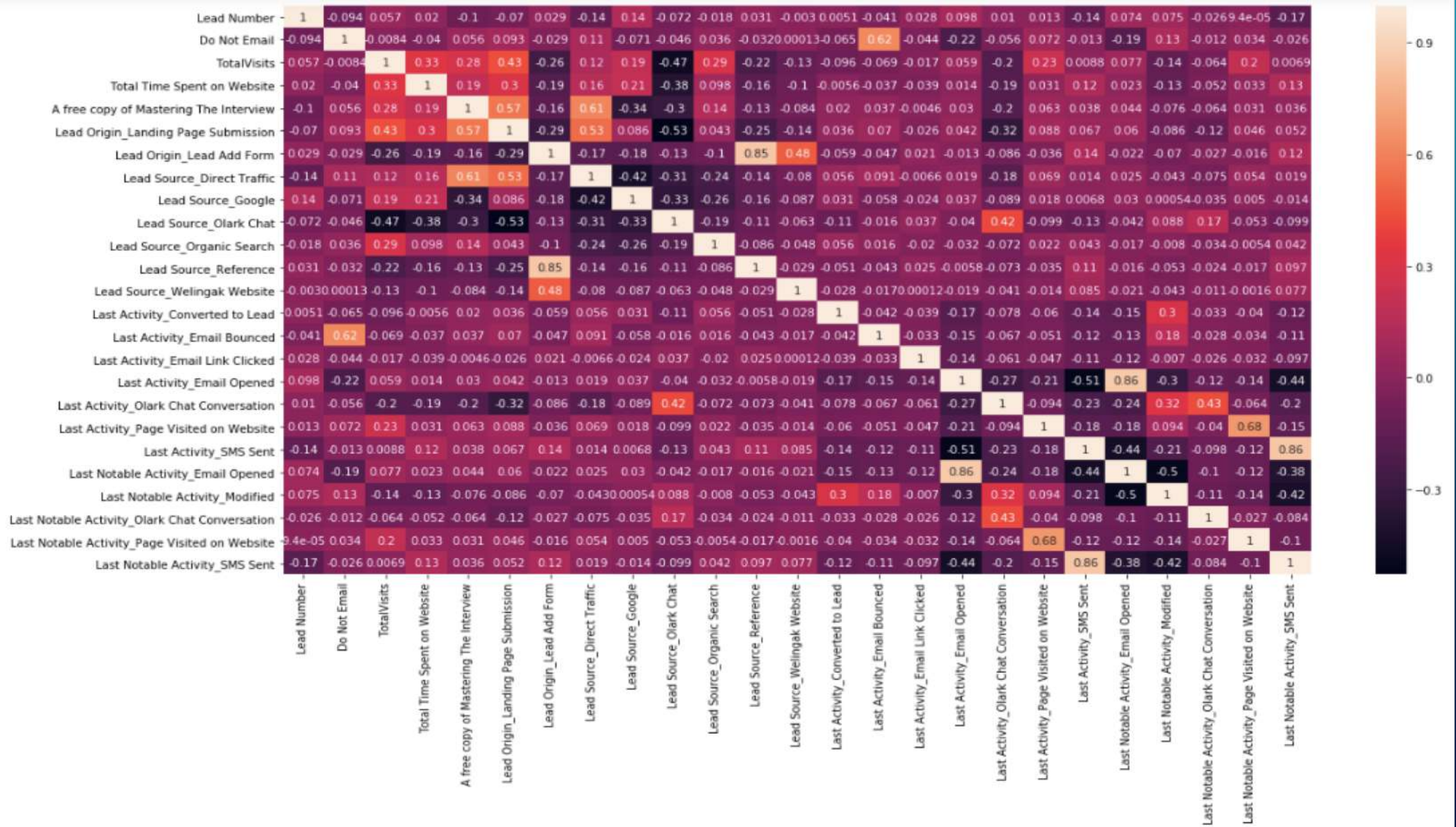
# VISUALIZATION OF VARIABLES

# VISUALIZATION CONTINUED..

# MODEL BUILDING

- Basic step to splitting the data into training and test sets and we choosen 70:30 ratio.

- Used RFE for feature selection.

- Running RFE with 24 variables as output.

- Removing the variable whose P-value is greater than 0.05 and VIF is greater than 0.02.

- Prediction on test data set.
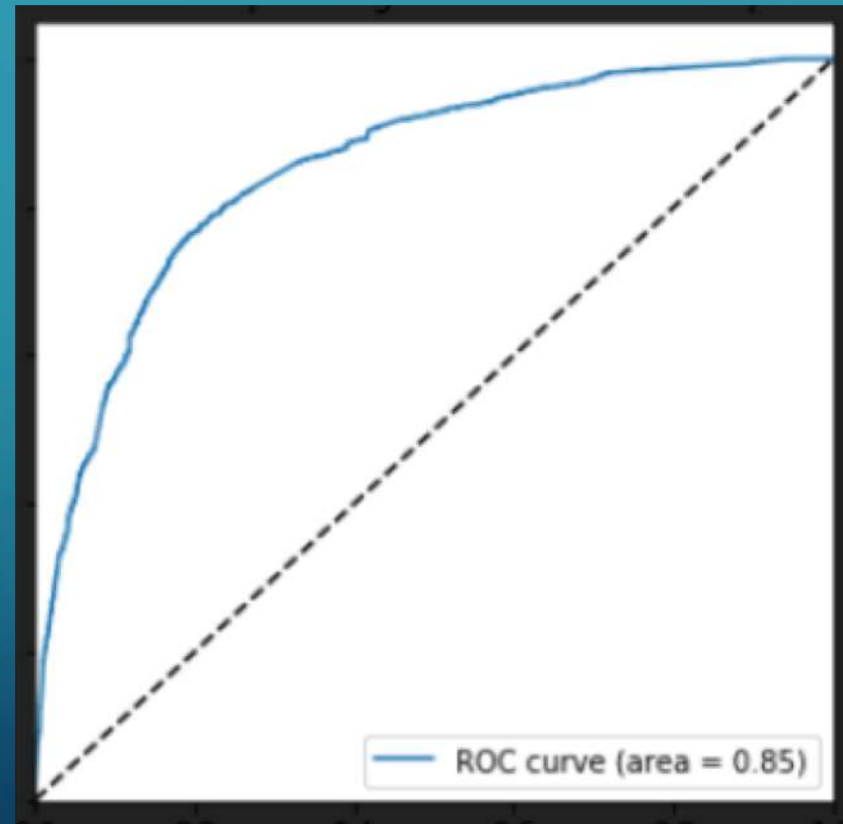
# CORRELATION BETWEEN VARIABLES

# RESULTS

After building the stable model ,we made the prediction and checked the statistics of the model -

- Sensitivity – 66%

- Specificity – 86%

- False positive rate – 13%

- Positive predicted value – 75%

- Negative predicted value – 80.5%

# ROC CURVE

ROC curve is Area under the curve is around 0.85 which shows a good model.



ROC curve (area = 0.85)

# CONCLUSION

- Main variables that contribute to analysis are "**Lead Source_Reference**", "**Total Time Spent on Website**" and "**Last Notable Activity_SMS Sent**".

- For more conversion we need to consider lower lead score, If we want to cover 100% of the leads, then we should consider score from 0. And If we want to cover around 97% of the leads then we should consider lead scores above 10.

- And we observed from the model if we want to precise conversion we need to consider higher lead score, If we want to be 98% accurate , then we should consider score from 90 or above. And to be around 100% accurate leads then we should consider lead scores above 99.