

A decorative graphic on the left side of the slide, consisting of a network of white lines and circles on a blue gradient background, resembling a circuit board or a neural network.

# LEAD SCORING CASE STUDY

BY:

Sayantan Singha

Sagar Rishi

Shreeragh C V

# 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 US TO HELP THEM 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 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. THE CEO, IN PARTICULAR, HAS GIVEN A BALLPARK OF THE TARGET LEAD CONVERSION RATE TO BE AROUND 80%.



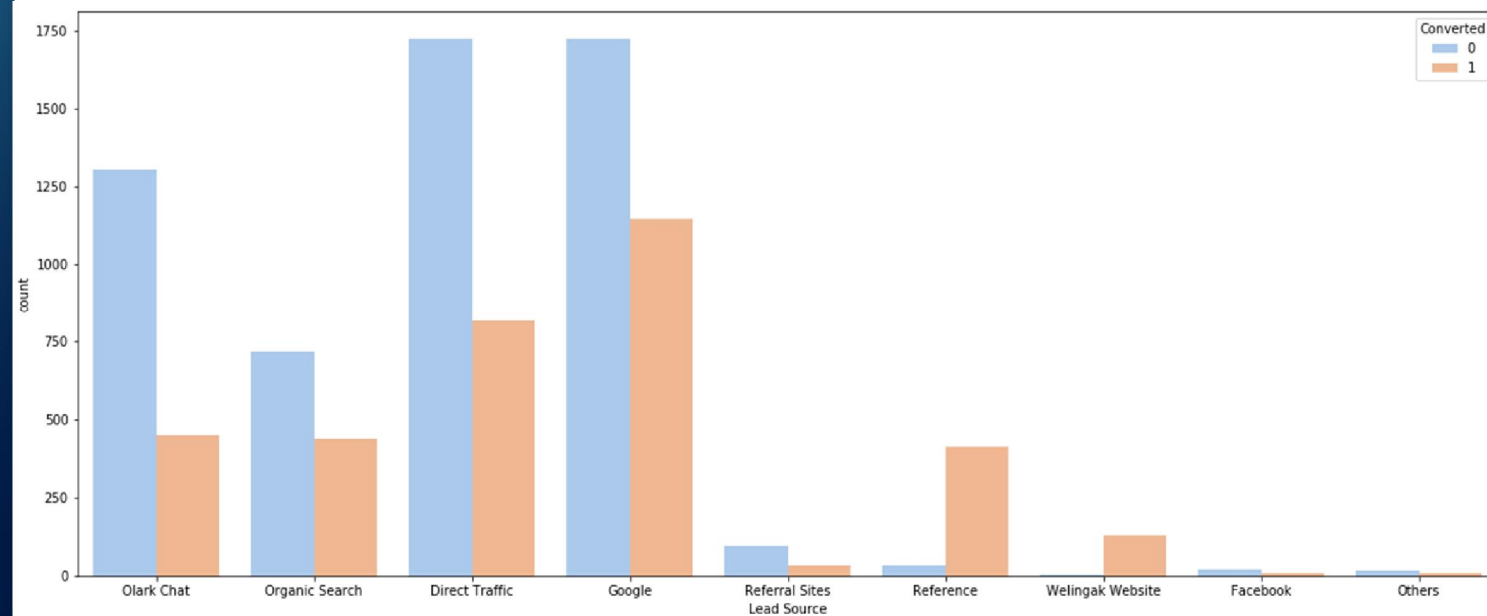
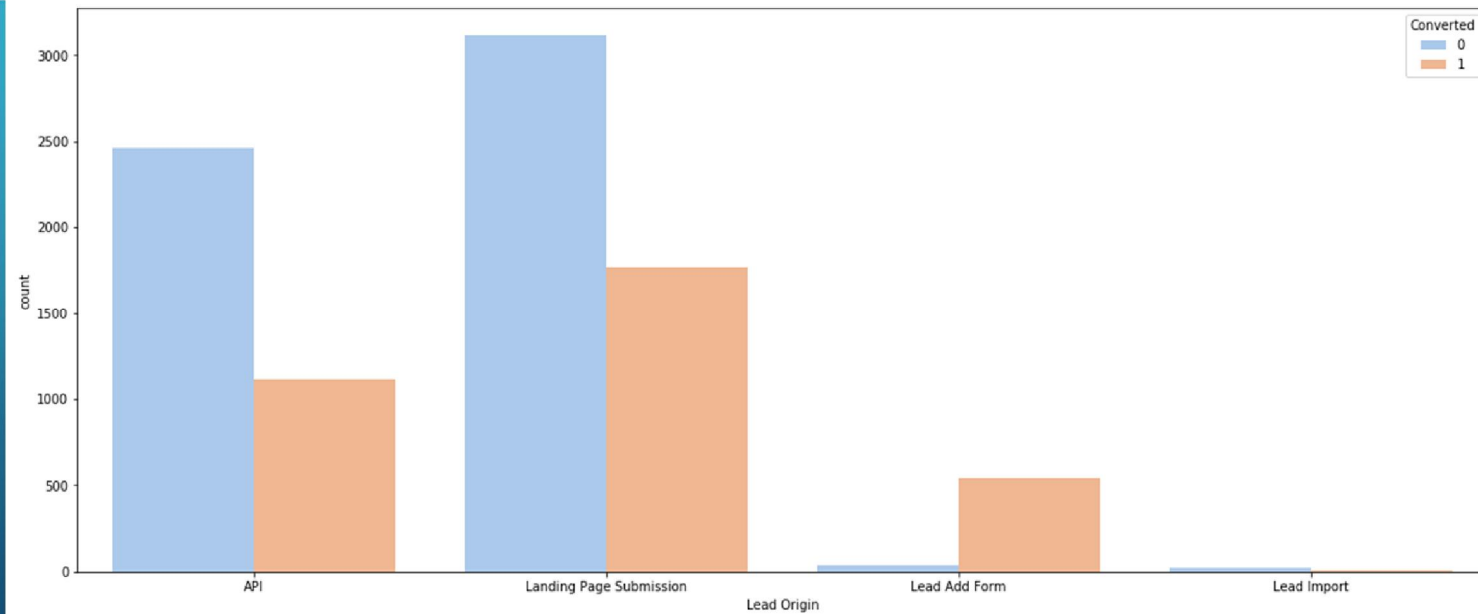
# ANALYSIS APPROACH

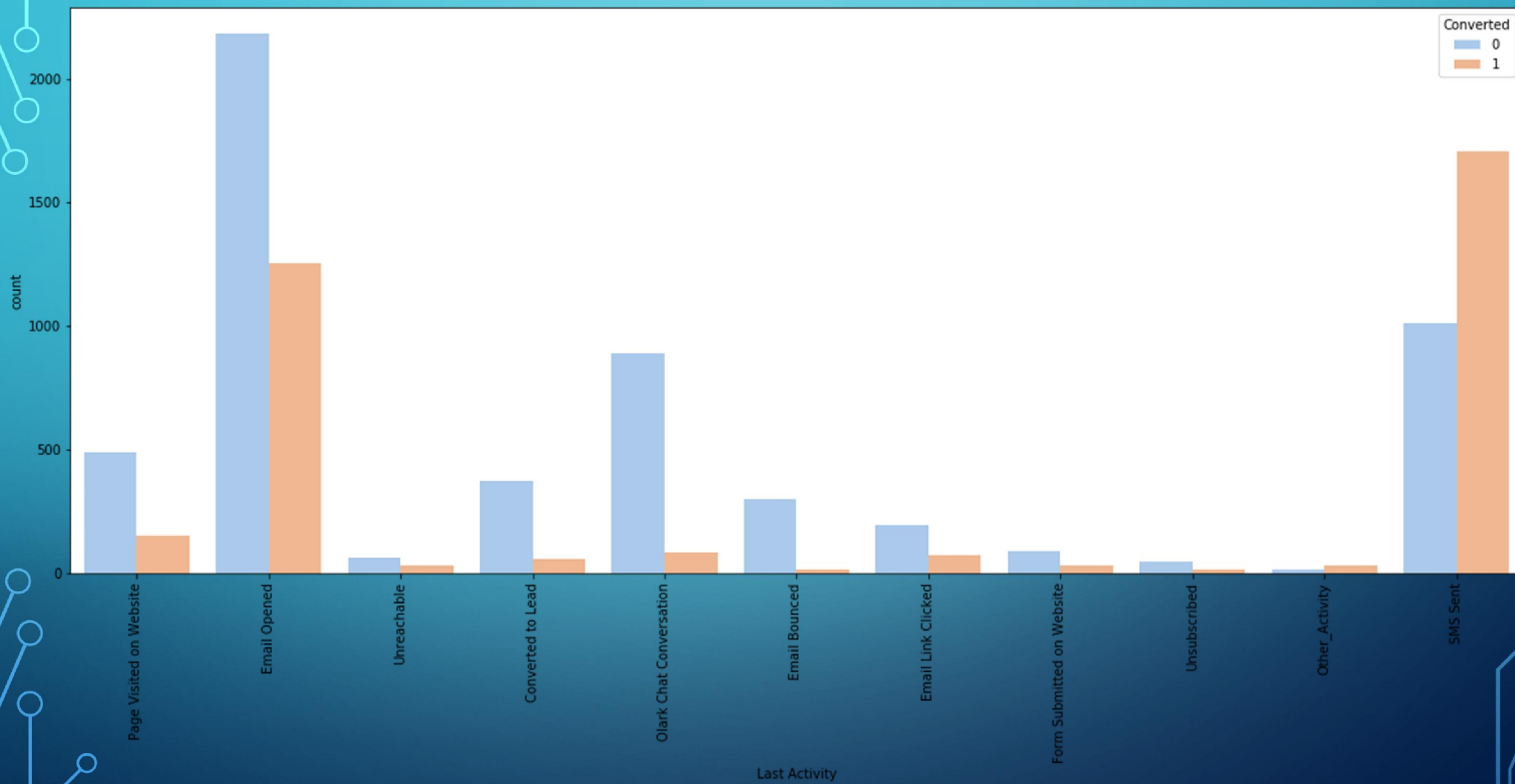
## STEP 1: MISSING VALUE TREATMENT

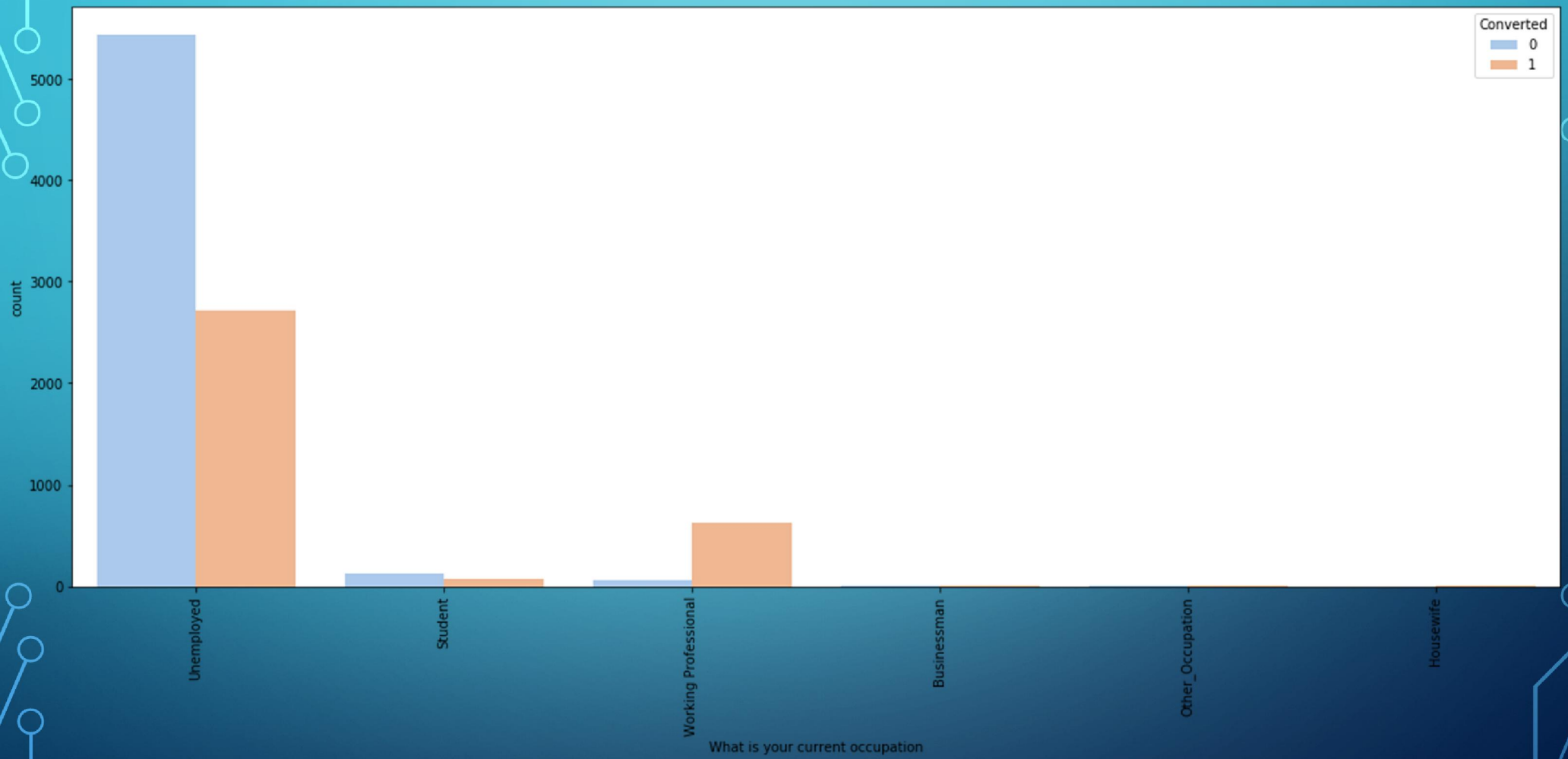
- Missing Values in each columns were treated.
- Columns with missing value  $> 70\%$  were dropped straight away
- Column 'Lead Quality' had nearly 50% of missing values. As this column had category 'Not Sure', all missing values were imputed with this category
- Other columns with missing values  $> 40\%$  were also dropped
- For rest of the columns, missing values were imputed with most frequently occurring value.

## STEP 2: EDA

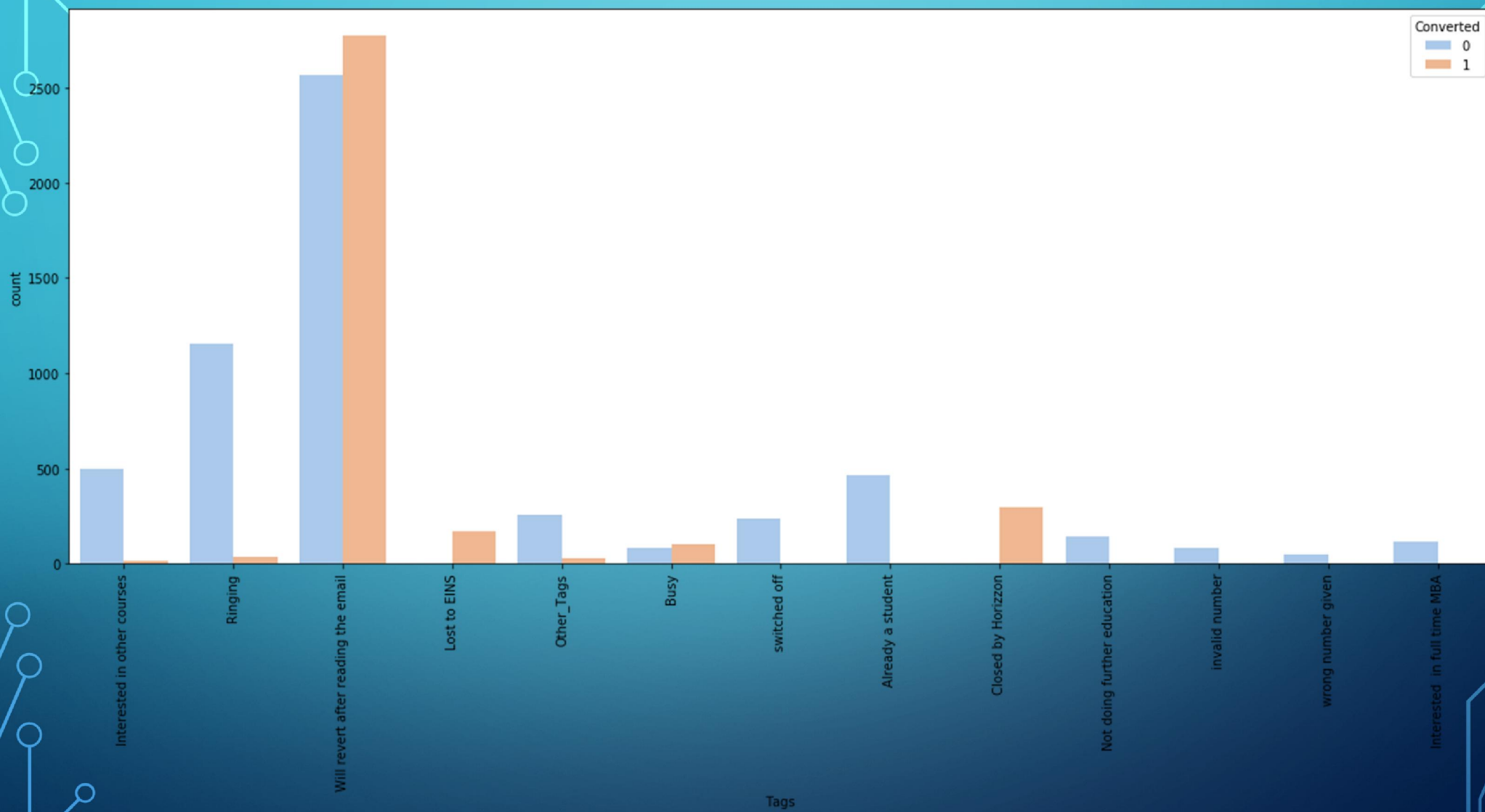
- Each columns was analyzed with respect to the target column, i.e., 'Converted' column
- Importance of the column for building the logistic model was determined by analyzing each graph
- Those column which did not have any significant affect on the target column were dropped at the end of EDA before moving on to model building



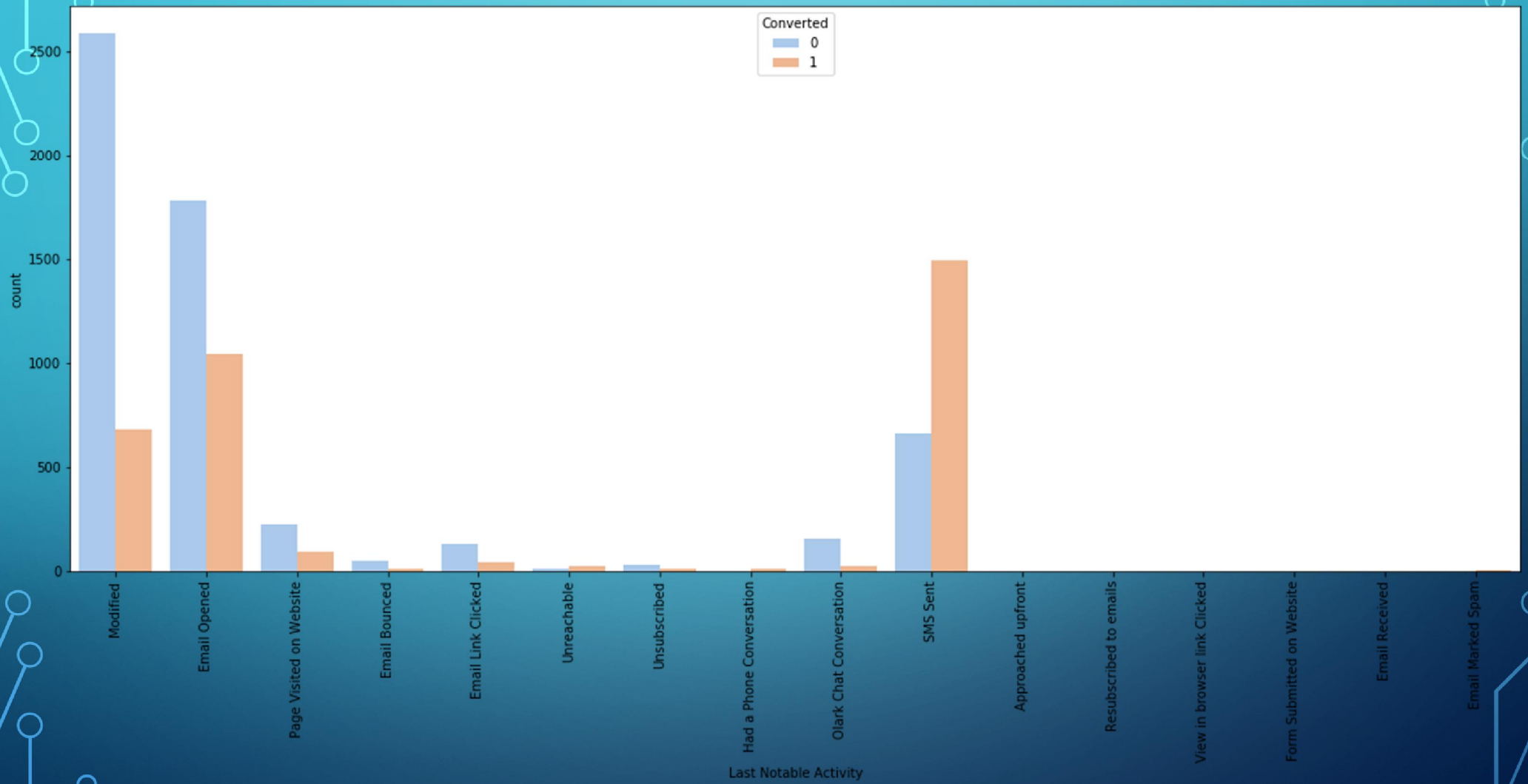












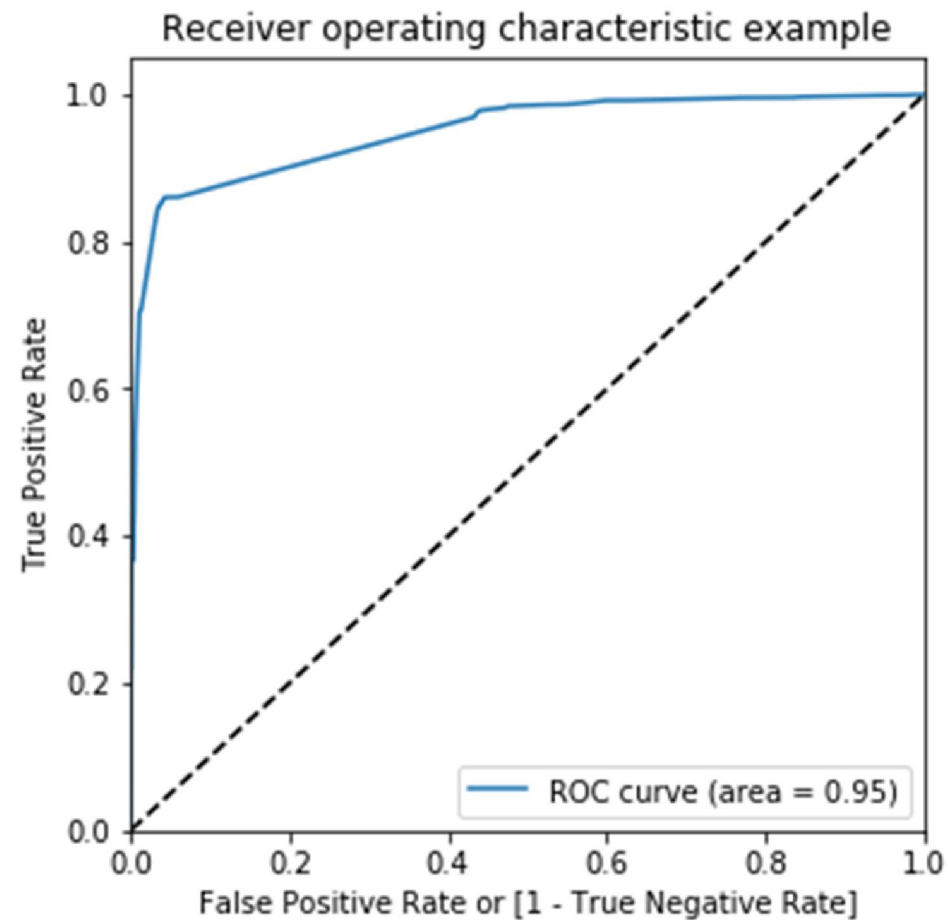


## STEP 3: DATA PREPARATION & MODEL BUILDING

- For all categorical column, dummy columns were created, and original column was dropped.
- Dataset was split into Train and Test.
- For train dataset, numerical columns were standardized.
- Model was build using RFE, considering Top-15 columns.
- After that, each column's p-value was checked and those column with higher p-values were dropped manually.
- Final model was arrived with 13 feature variable.
- VIF of the feature variables were also checked and found ok.

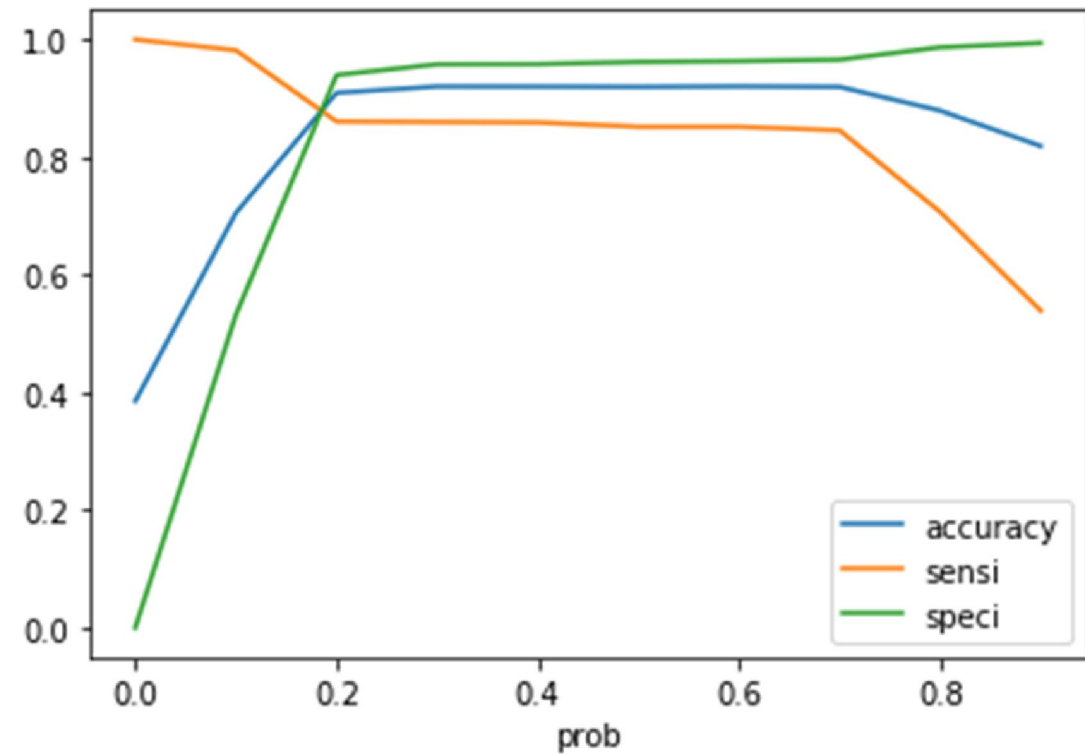
## STEP 4: ROC CURVE PLOT

- 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



## STEP 5: FINDING OPTIMAL CUTOFF POINT

- Optimal cutoff probability is that prob where we get balanced sensitivity and specificity



## STEP 6: MAKING PREDICTION ON TEST SET

- Numerical columns of Test dataset were standardized.
- Final model developed was used to predict target variable of the test dataset

### Comparison of Accuracy, Sensitivity & Specificity of Train and Test dataset

|       | Sensitivity | Specificity |
|-------|-------------|-------------|
| Train | 86.05       | 93.95       |
| Test  | 84.42       | 93.88       |

# CONCLUSION

Following are the variables / dummy variables which should be focused on

- 'Do Not Email'
- 'Lead Origin - Lead Add Form'
- 'Lead Source - Welingak Website'
- 'What is your current occupation - Unemployed'
- 'Tags - Busy'
- 'Tags - Closed by Horizzon'
- 'Tags - Lost to EINS'
- 'Tags - Ringing'
- 'Tags - Will revert after reading the email'
- 'Tags - switched off'
- 'Lead Quality - Not Sure'
- 'Lead Quality - Worst'
- 'Last Notable Activity - SMS Sent'



The background is a blue gradient. In the corners, there are white line art designs resembling circuit boards or neural networks, with lines and small circles.

THANK YOU