

# Lead Scoring Case Study

By:-

Arjun Mukerji

Harinandan Praveen

Abhishek Sahebrao Bolke

# 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

- X Education wants to know most promising leads.
- For that they want to build a Model which identifies the 'Hot Leads'.
- Deployment of the Model for the future use.

# Solution Methodology

- EDA

- 1.Check and handle duplicate data

- 2.Check and handle NAN values and missing values

- 3.Drop columns, if it contains large amount of missing values and are not useful for the analysis

- 4.Imputation of the values, if necessary

- 5.Check and handle Outliers in data

- Feature Scaling and Dummy Variables

- Classification Techniques: Logistic Regression used for the model making and prediction

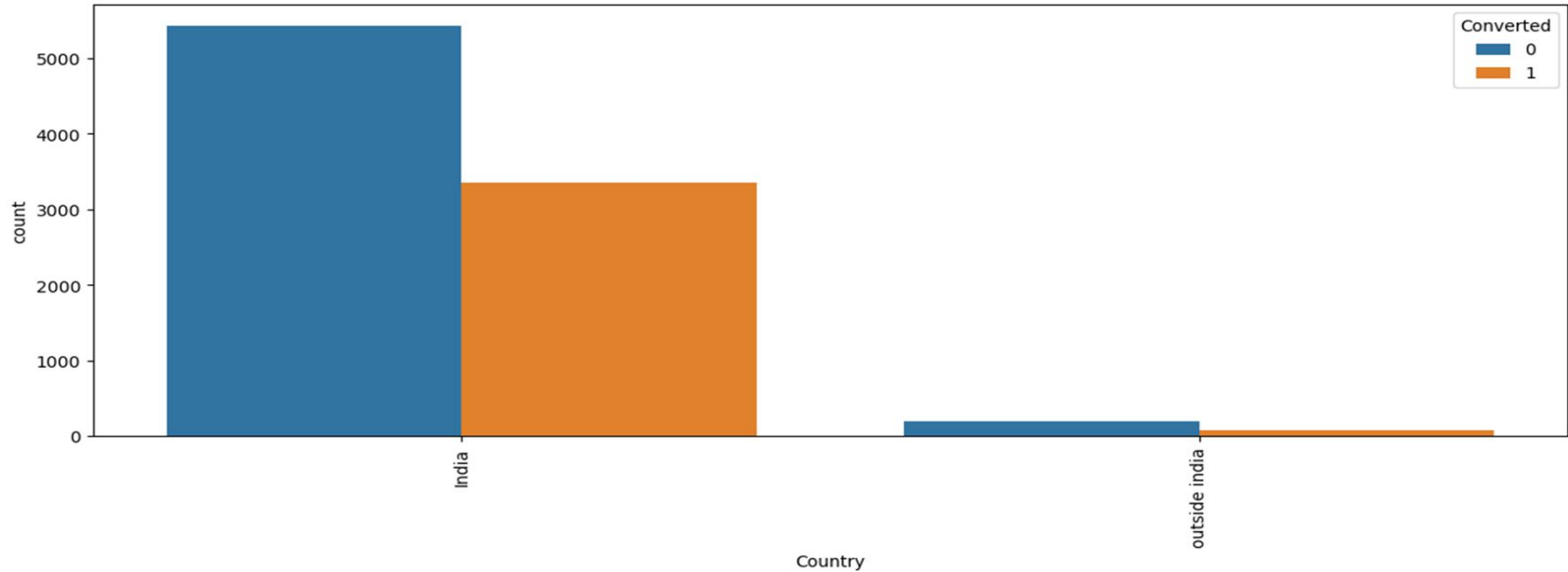
# Solution Methodology (Continued)

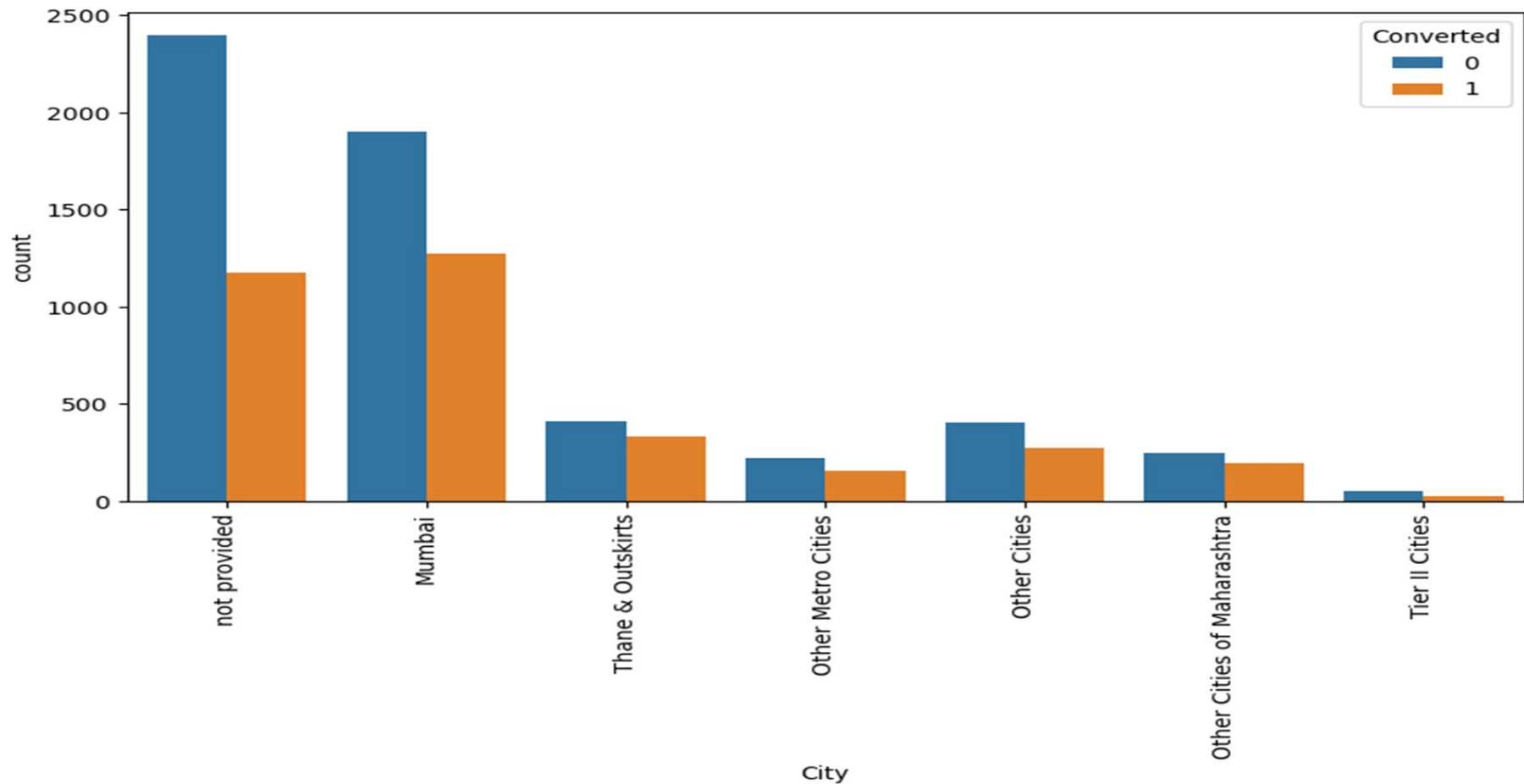
- Validation of the Model
- Model presentation
- Conclusions and Recommendations

# Data Cleaning and Data Manipulation

- Total number of Rows = 9240, Total number of Columns=37.
- Removing 'Lead Number' and 'Prospect ID' which is not necessary for analysis since they all have Unique values.
- Converting 'Select' values to 'NAN'.
- Dropping unique valued Columns.
- Dropping off the Columns with more than 45% missing values.
- Imputing 'Not Provided' values with 'India' since 'India' is the most common occurrence among the non-missing values.

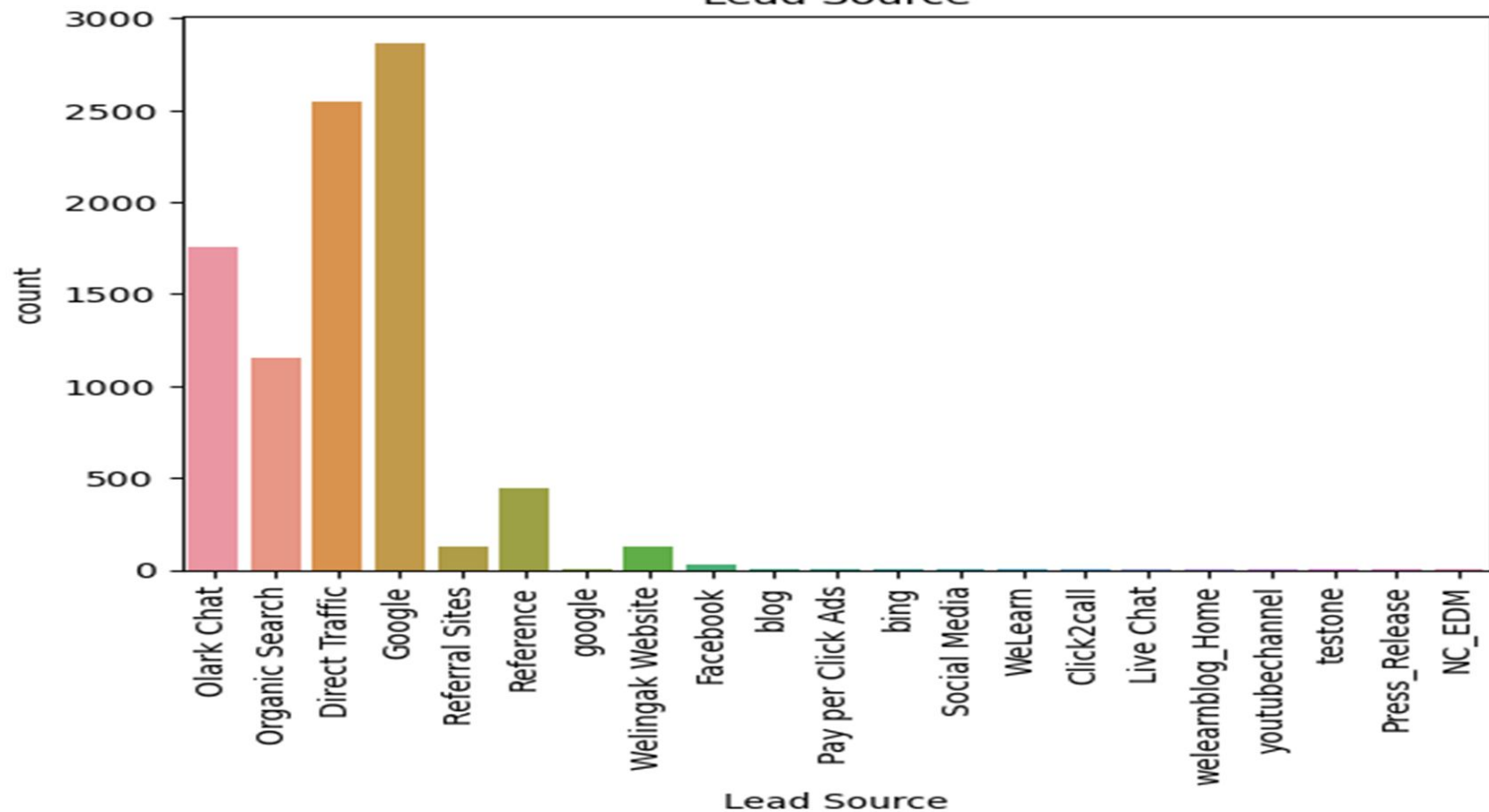
# EDA (Categorical Attributes Analysis)



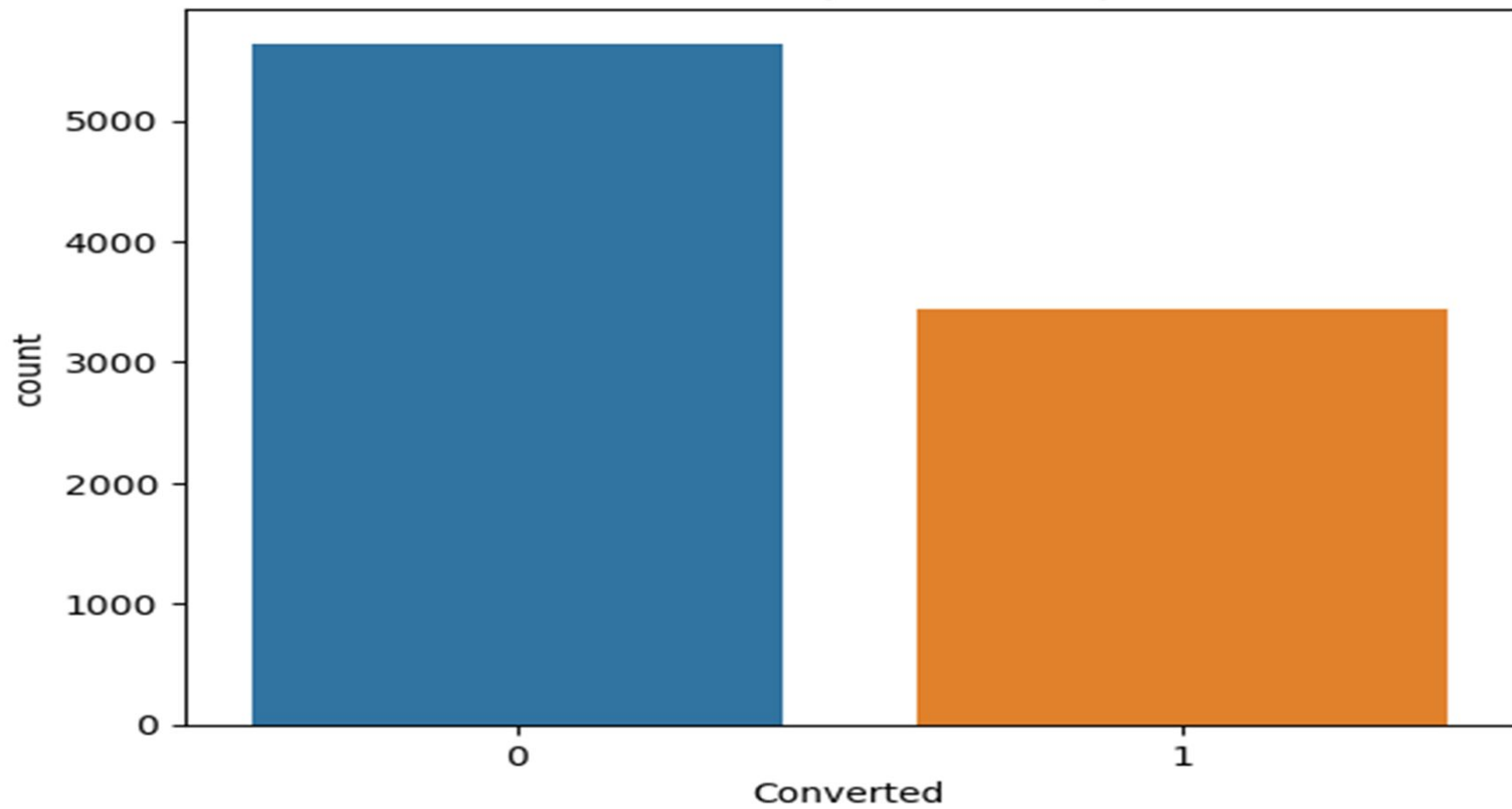




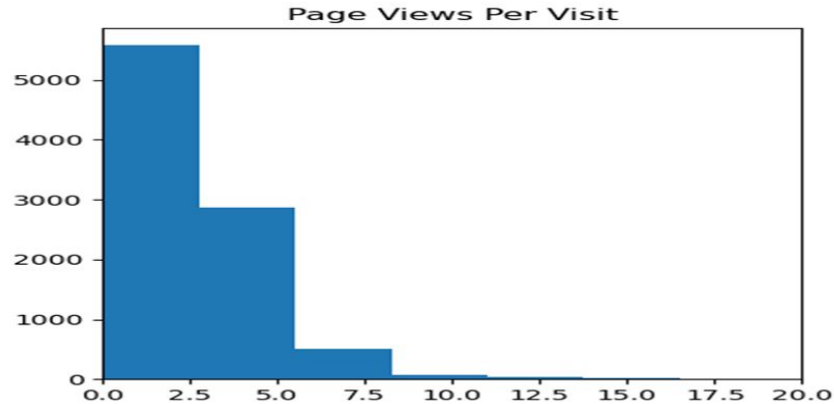
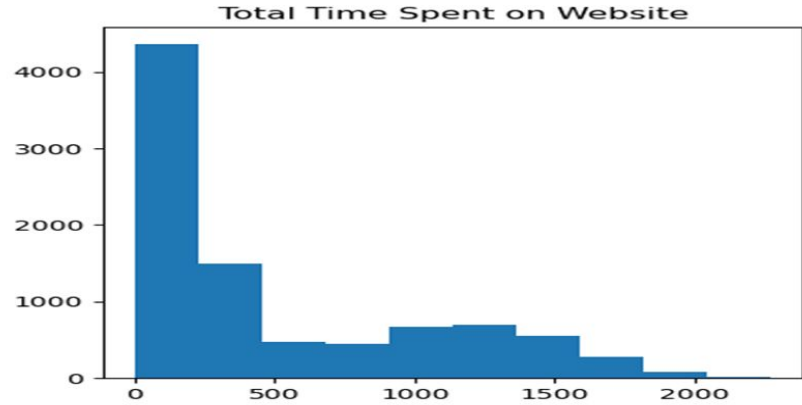
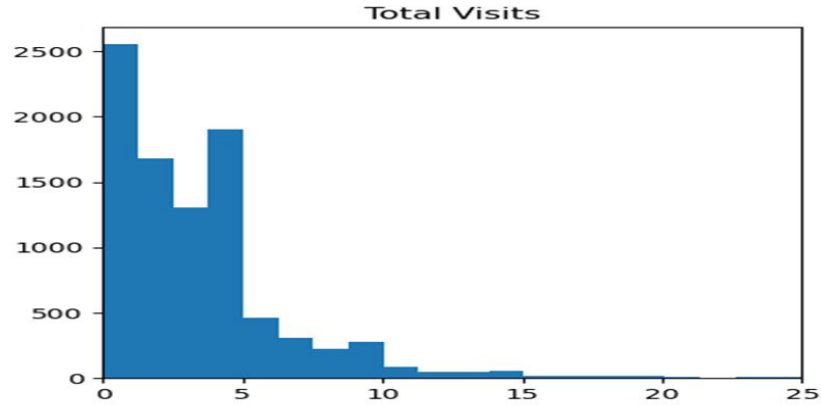
Lead Source



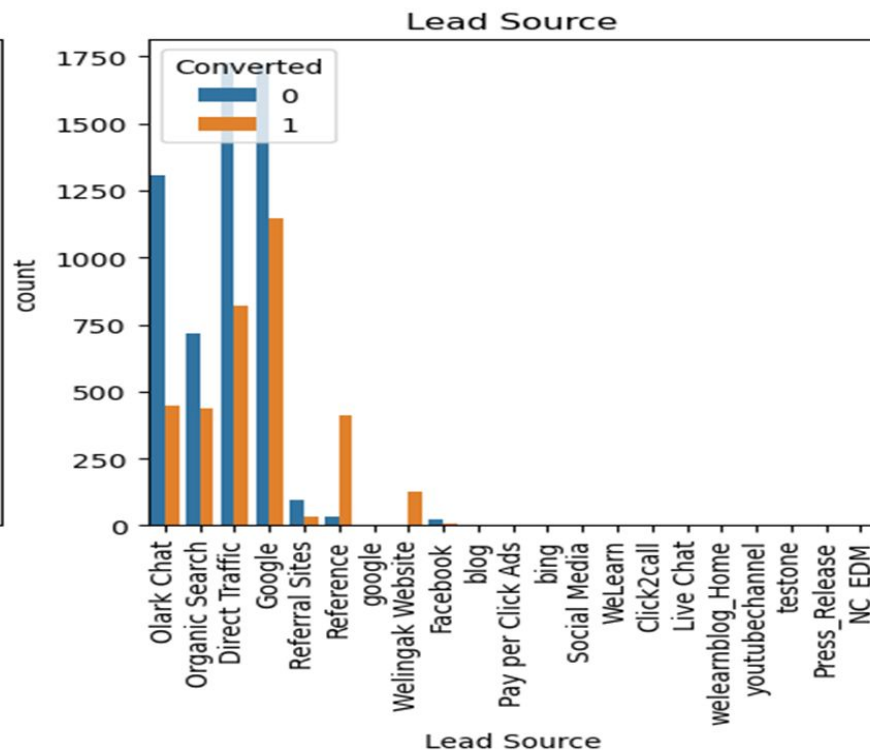
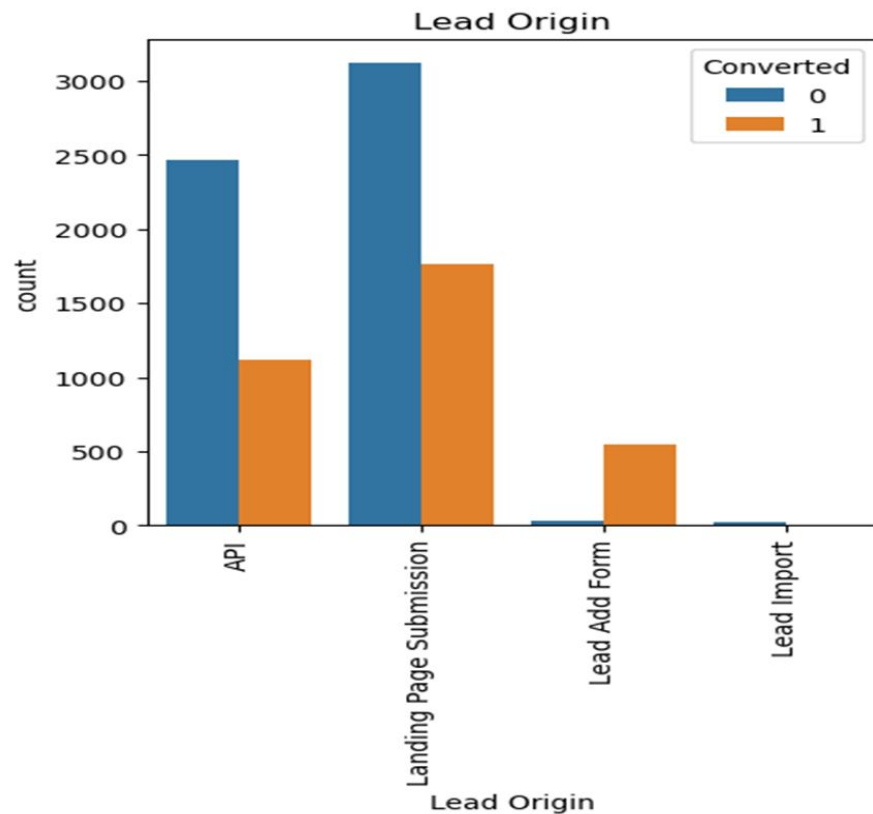
Converted("Y variable")



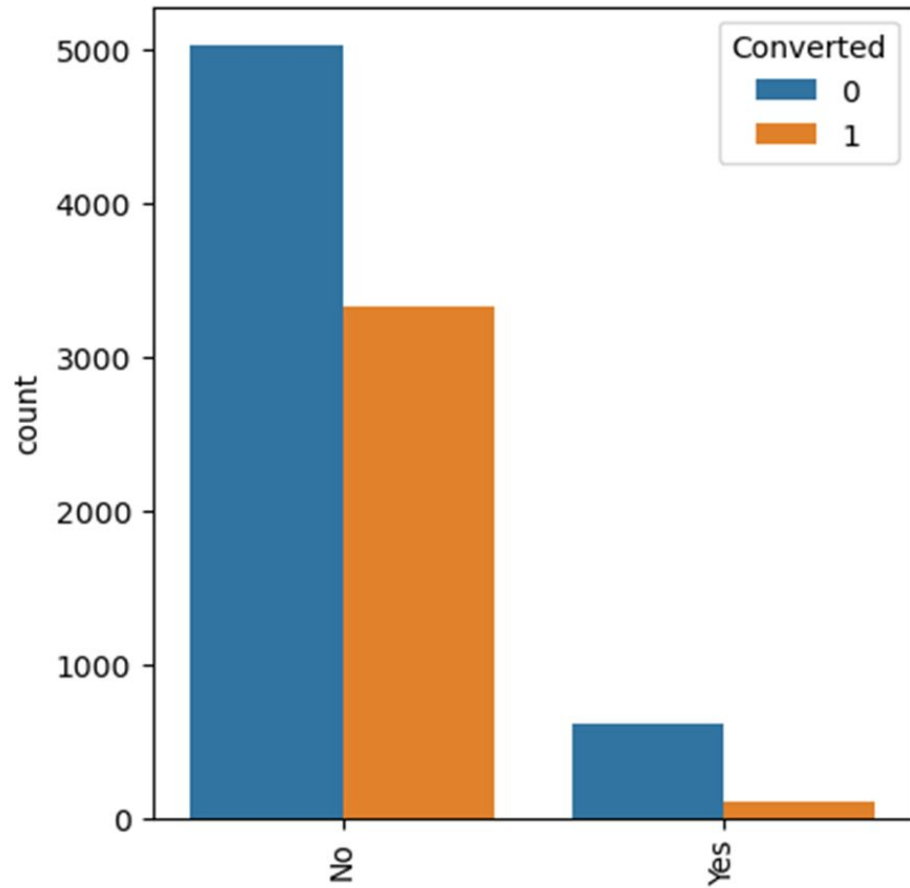
# EDA (Numerical Variables)



# EDA (Categorical Variable Relation)

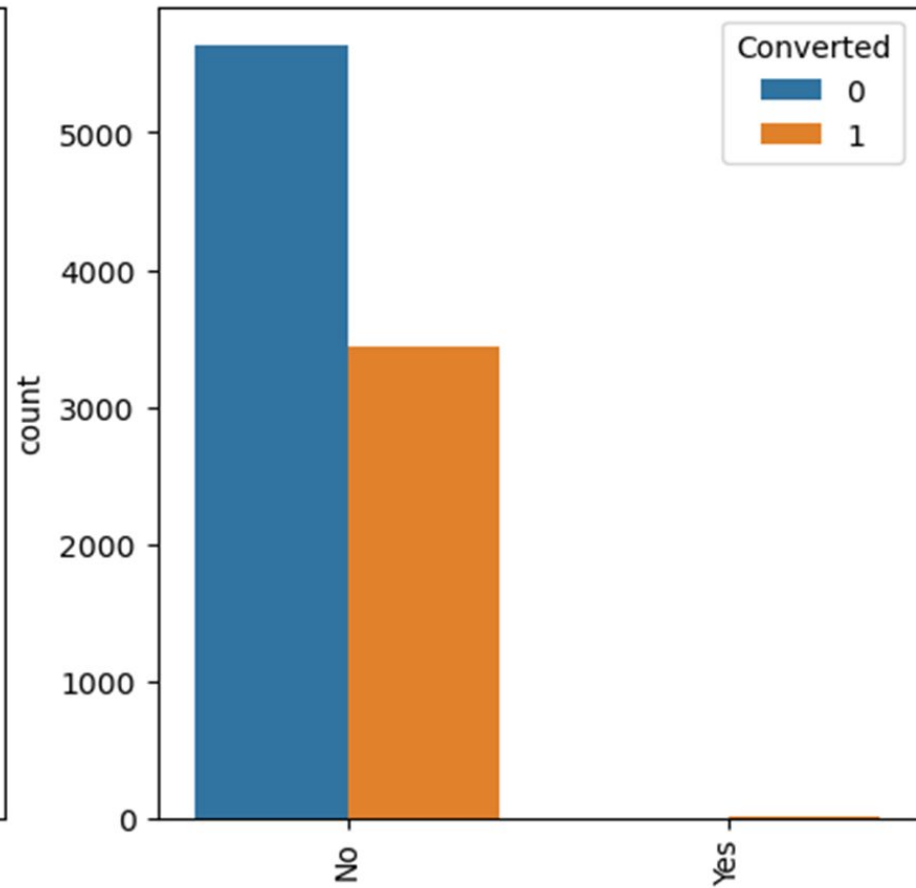


Do Not Email

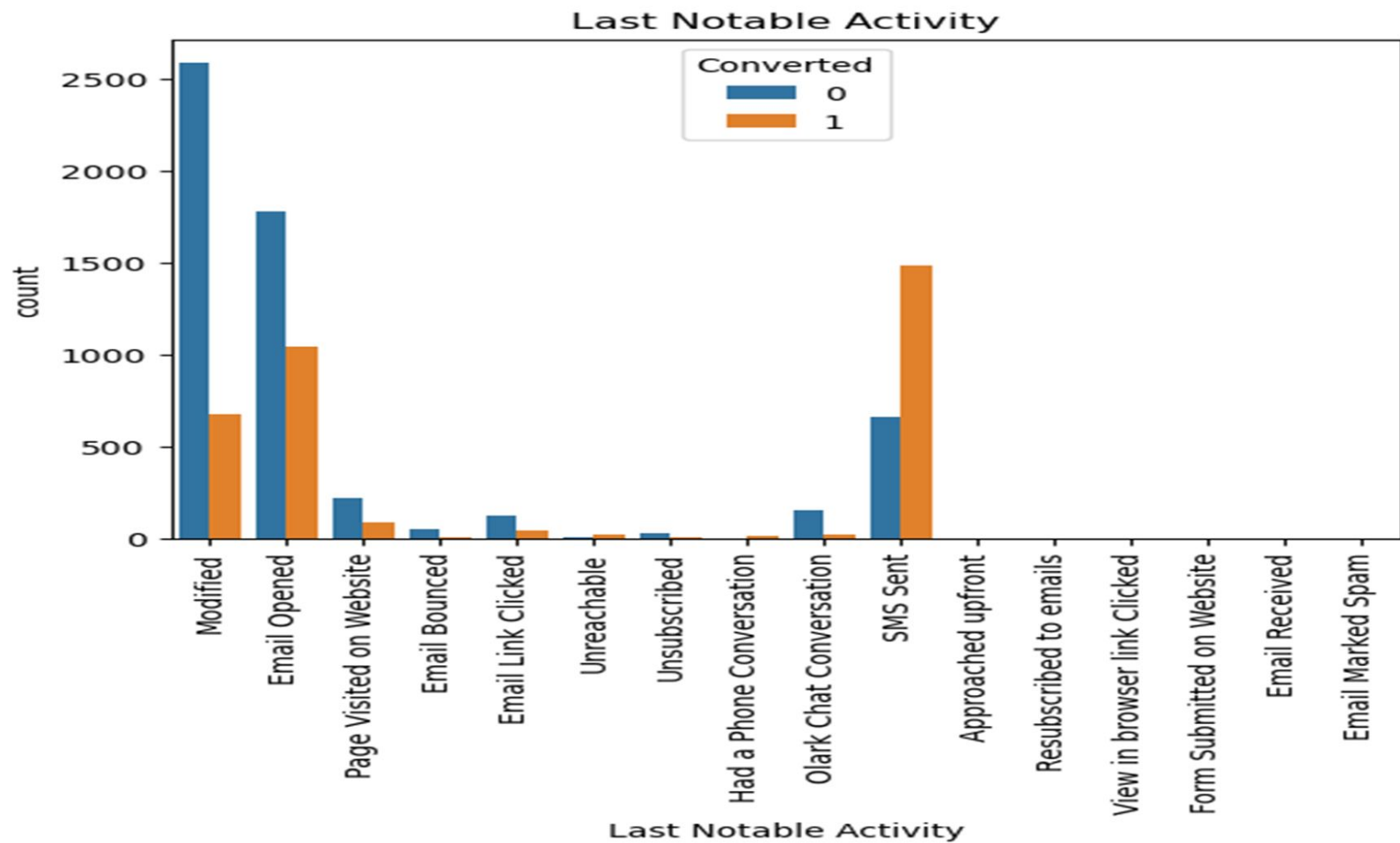


Do Not Email

Do Not Call



Do Not Call



# Data Conversion

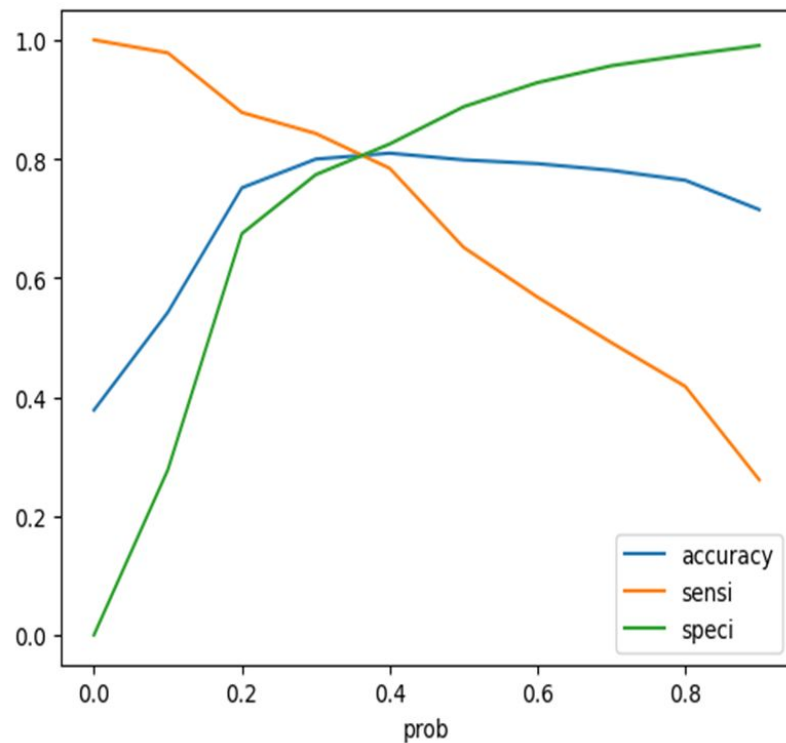
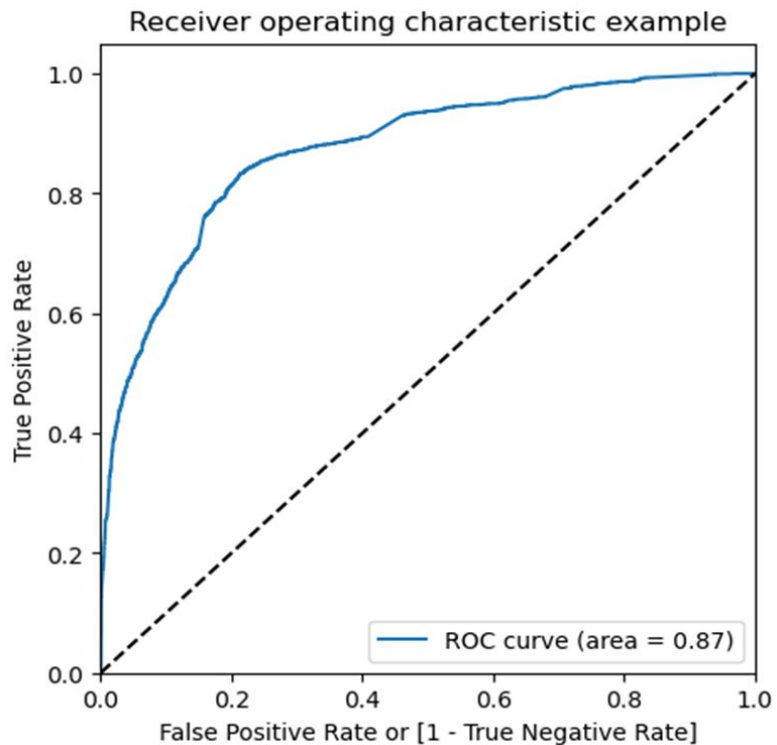
- Numerical Variables are Normalised.
- Dummy Variables are created for Object Type Variables.
- Total Rows for Analysis: 8991
- Total Columns for Analysis: 82

# Model Results

- Splitting the Data into Training and Testing Sets.
- The first basic step for Regression is performing a train-test split, we have chosen 70:30 ratio.
- Use RFE for Feature Selection.
- Running RFE with 15 variables as output.
- Building Model by removing the variable whose p-value is greater than 0.05 and VIF value is greater than 5.
- Predictions on Test Data Set.
- Overall accuracy is 78.9%.



# ROC Curve



## ROC Curve (Continued)

- Finding Optimal Cut Off Point.
- Optimal Cut Off Probability is that Probability where we get balanced Sensitivity and Specificity.
- From the second graph, it's visible that the optimal cut-off is approximately at 0.37.

## Conclusion and Recommendations

• Looking at the below variables, we should be able to help X-Education company with finding out hot leads which will be very high potential customers:-

1.TotalVisits

2.The total time spend on the Website

3.Lead Origin\_Lead Add Form

4.Lead Source\_Direct Traffic

5.Lead Source\_Google

6.Lead Source\_Welingak Website

## Conclusion and Recommendations (Continued)

7. Lead Source\_Organic Search

8. Lead Source\_Referral Sites

9. Lead Source\_Welingak Website

10. Do Not Email\_Yes

11. Last Activity\_Olark Chat Conversation

12. Last Activity\_Email Bounced

- Taking the above details into consideration and using this model, the X-Education company can maximise the sales of their course and also generate an excellent revenue.

**THANK YOU!!!**