# LEAD SCORING CASE STUDY

Kumar Mayank

Mayank Jain

Mangesh Deepak Bagul

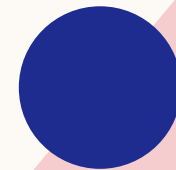# AGENDA

# 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.

❑ Building a suitable model which identifies the hot leads.

# METHODOLOGY

❑ **Data cleaning and data manipulation.**

- ▪ Check and handle duplicate data.
- ▪ Check and handle NA values and missing values.
- ▪ Drop columns, if it contains large amount of missing values and not useful for the analysis.
- ▪ Imputation of the values, if necessary.
- ▪ Check and handle outliers in data.

❑ **Exploratory data Analysis**

- ▪ Univariate data analysis: value count, distribution of variable etc.
- ▪ Bivariate data analysis: correlation coefficients and pattern between the variables etc.

❑ **Feature Scaling & Dummy Variables and encoding of the data**.

❑ **Classification technique:** Logistic regression used for the  model making and prediction.
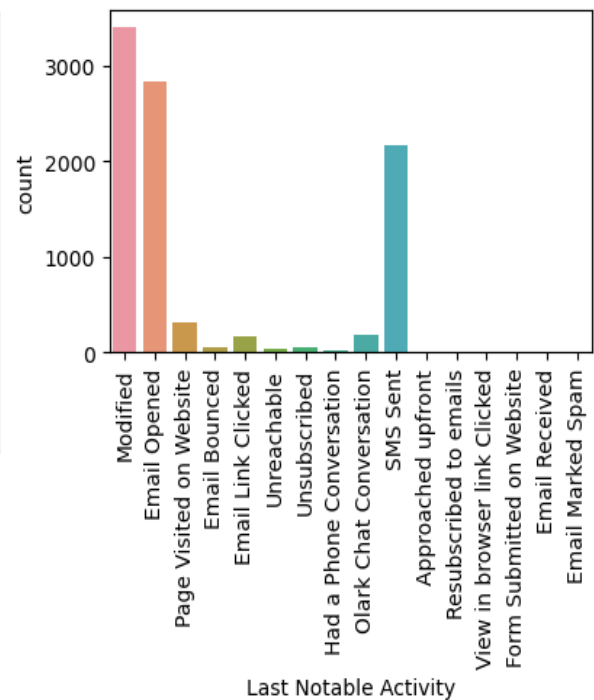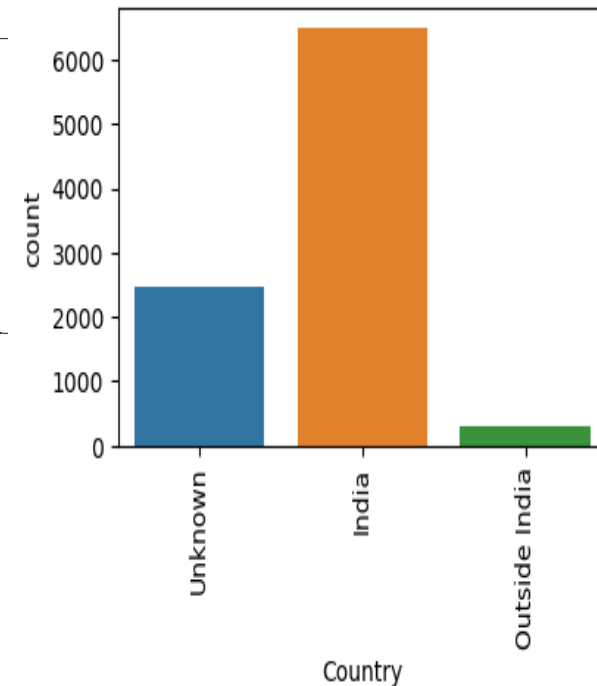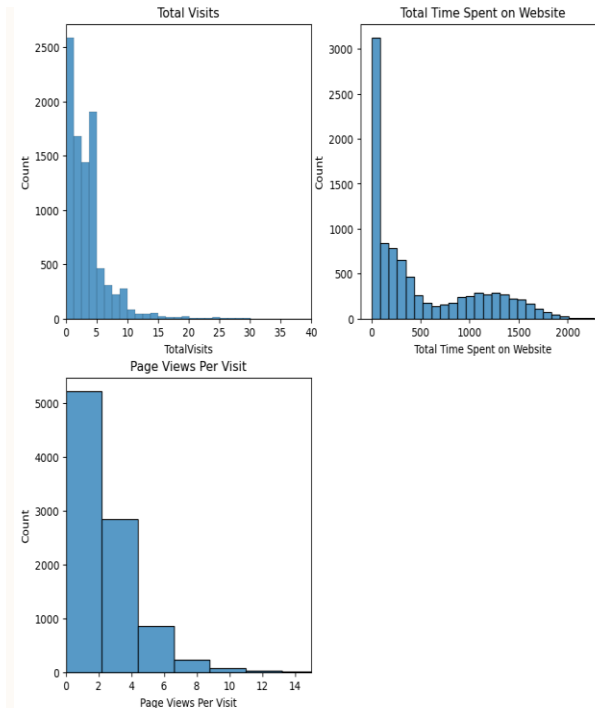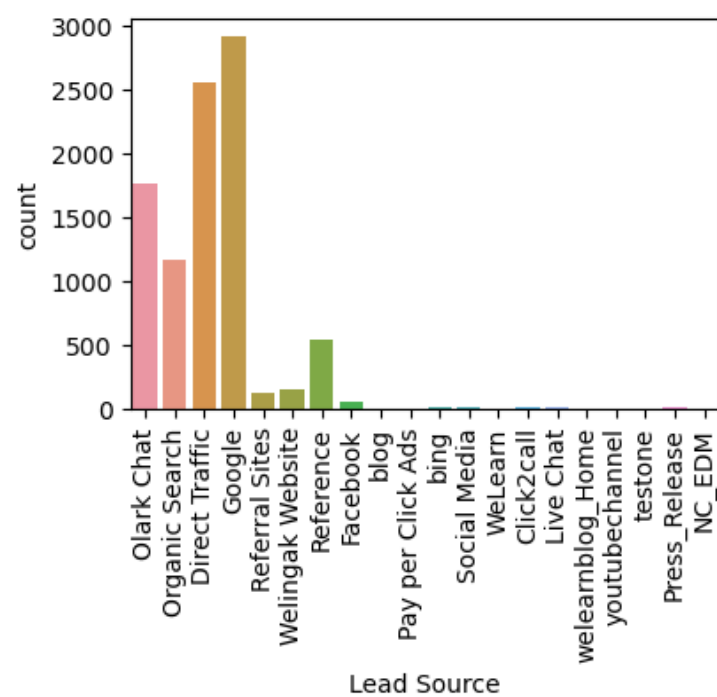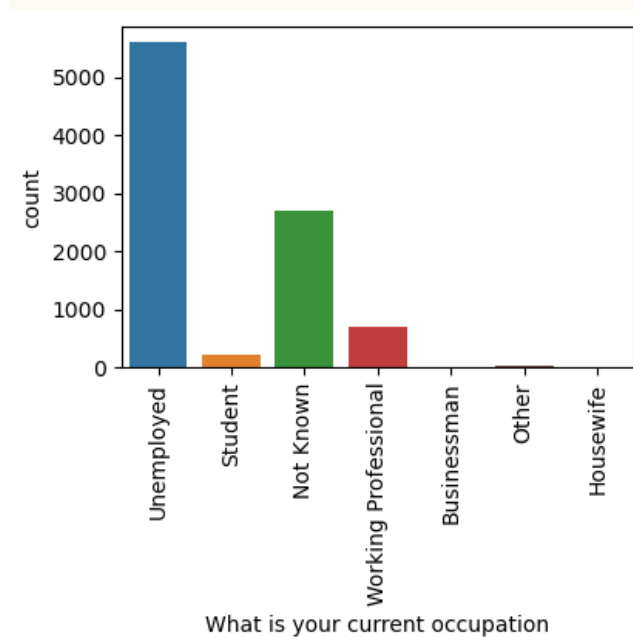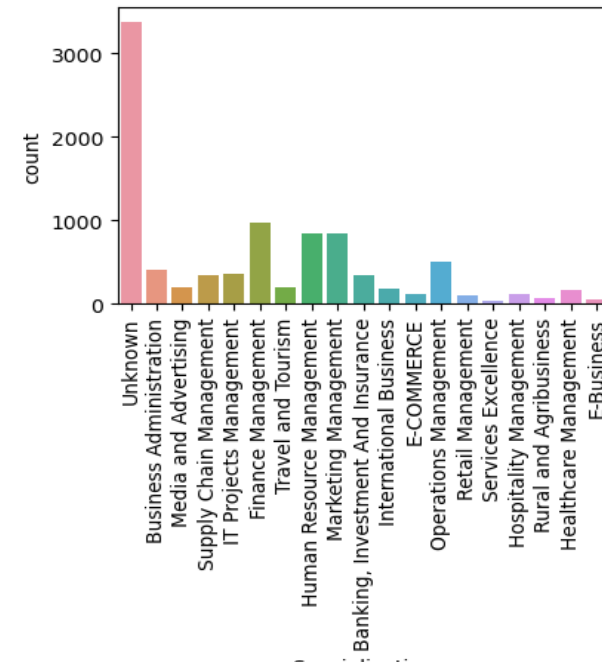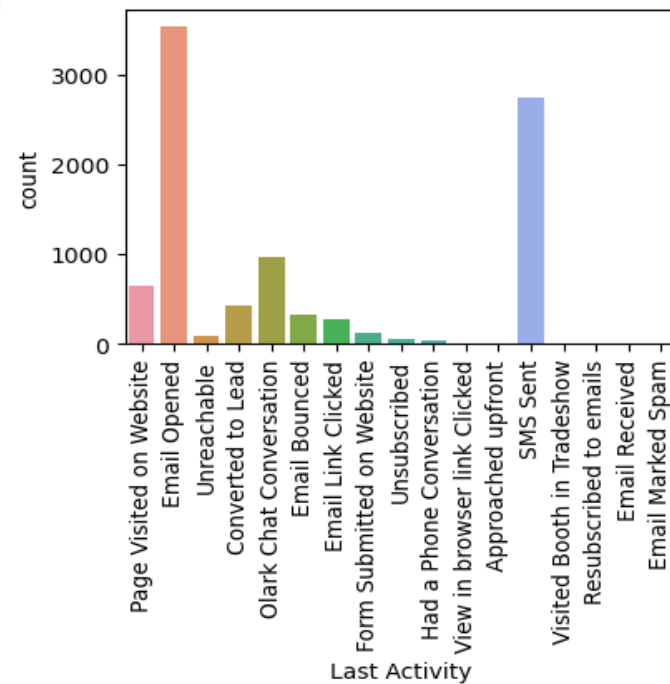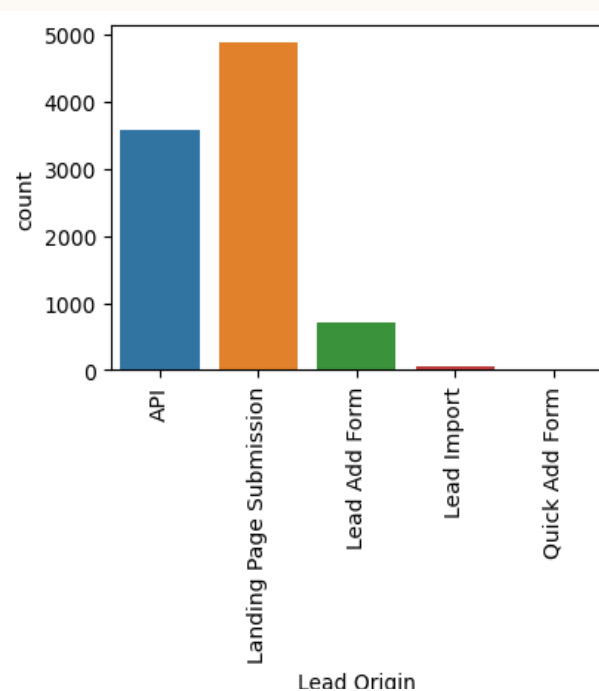
❑ **Validation of the model.**

# DATA CLEANING AND MANIPULATION

- 'Select' values were present which we have replaced with Null Values.

- Columns having minimum 40% Null values were dropped off

- Columns which had only one unique value were dropped off

- Imputing values with mean and mode in which less than 5% null values are present.

- Checking for Duplicate data points and rectifying the same.
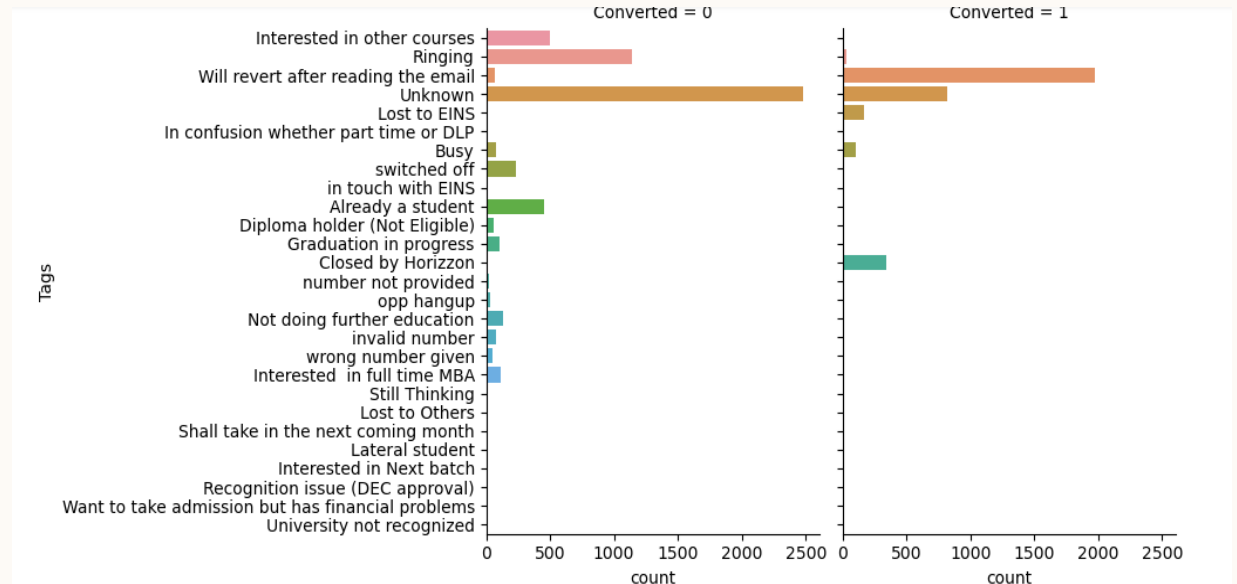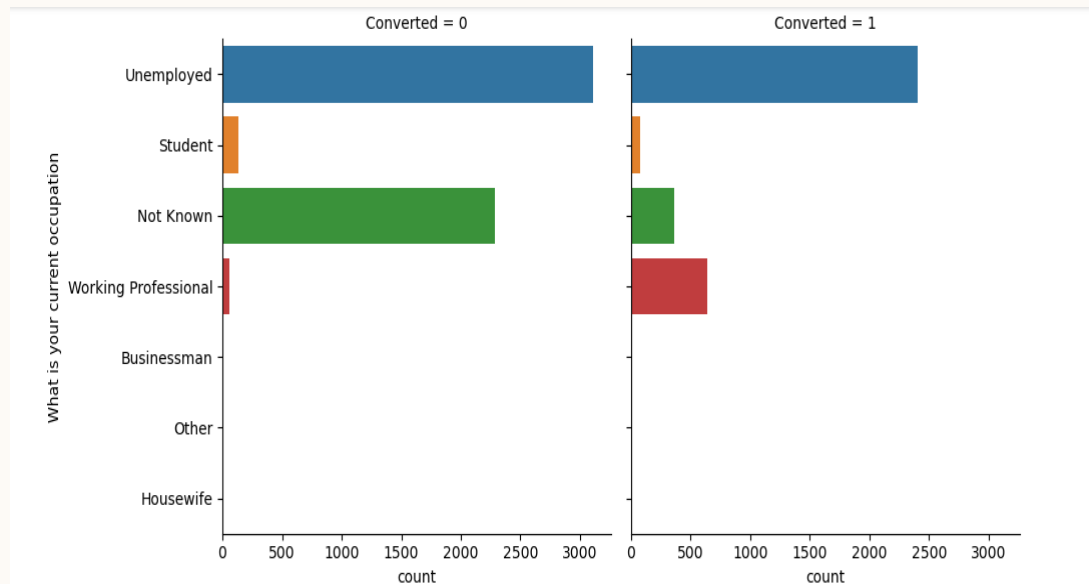
# PERFORMING EDA

**Univariate Analysis**

In Univariate analysis, mainly count plots and histograms were plotted to ge the understanding of the counts and frequencies of  the different variables present in the data
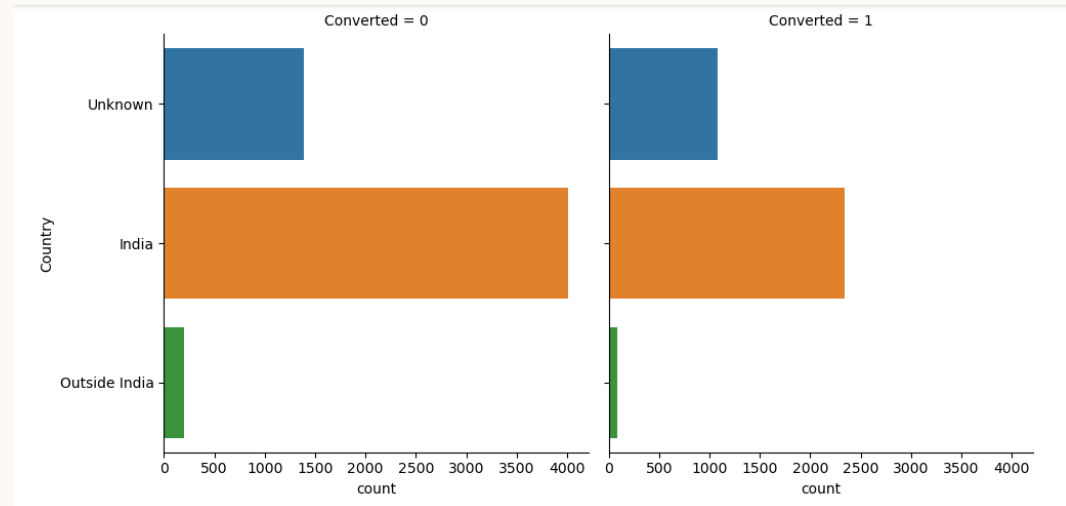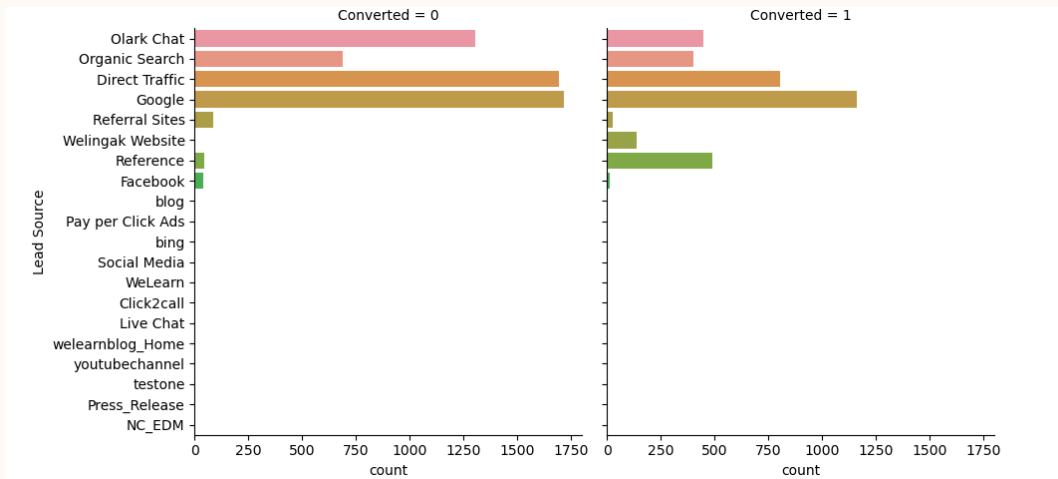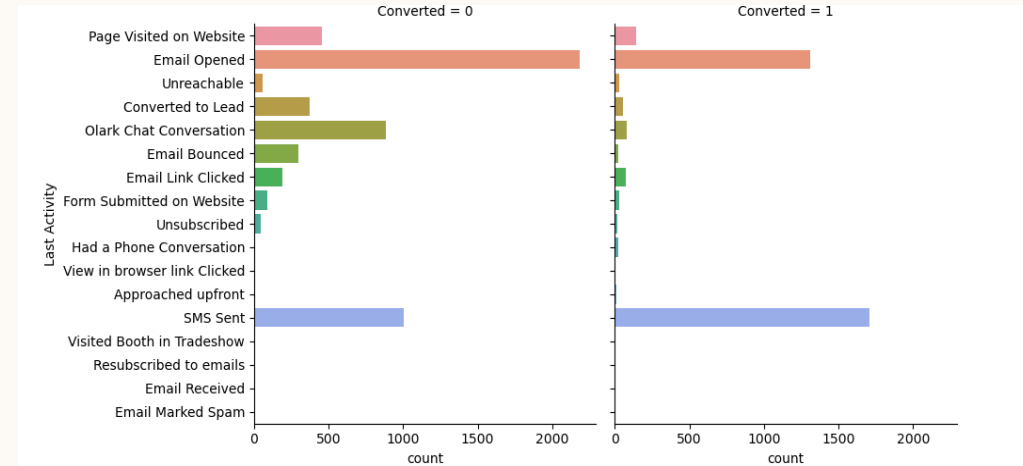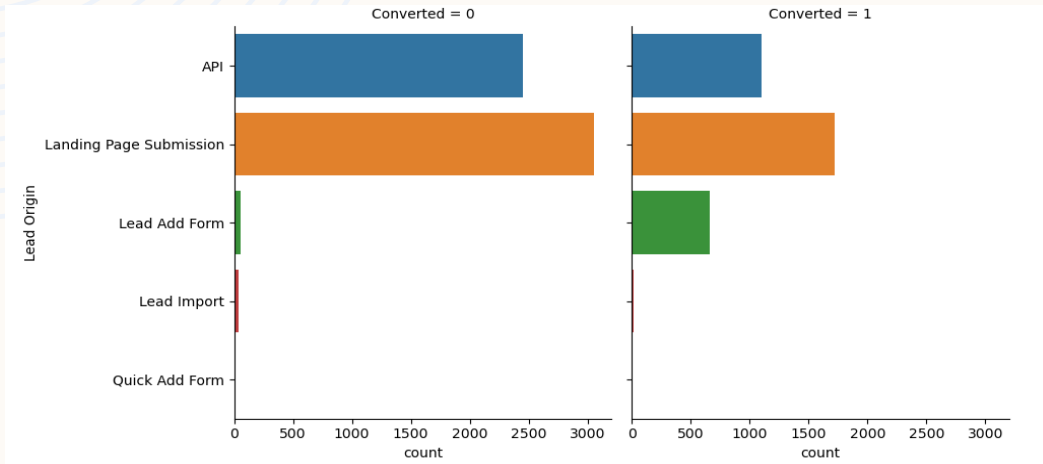
# PERFORMING EDA

**Bivariate Analysis**

In Bivariate analysis, count plots were plotted keeping 'Converted' as the target variable and other independent variables
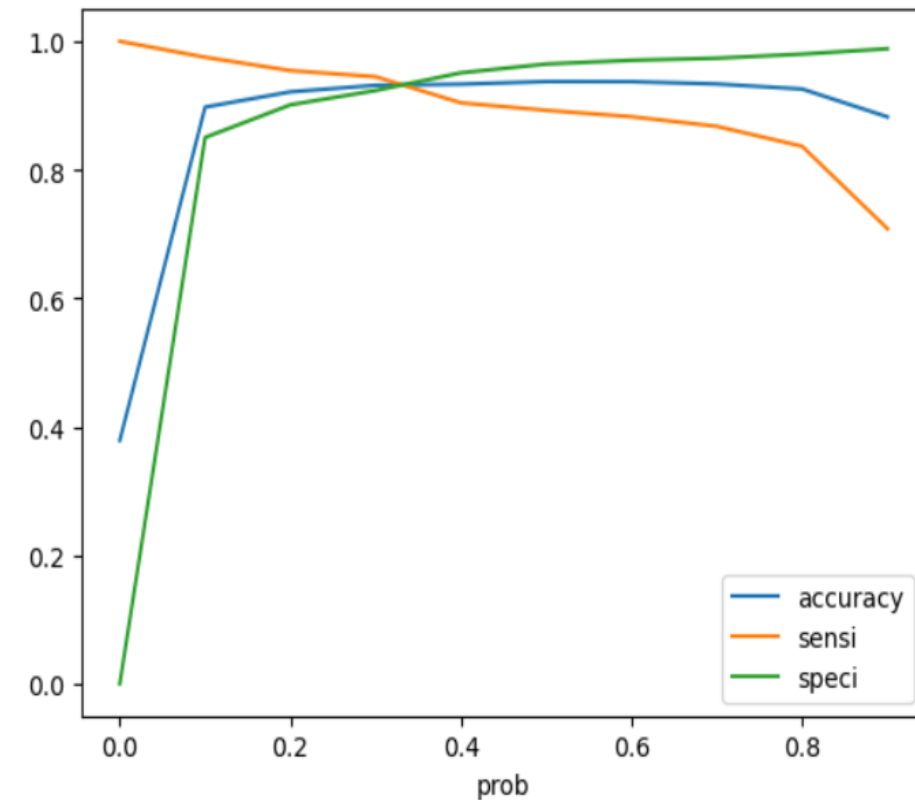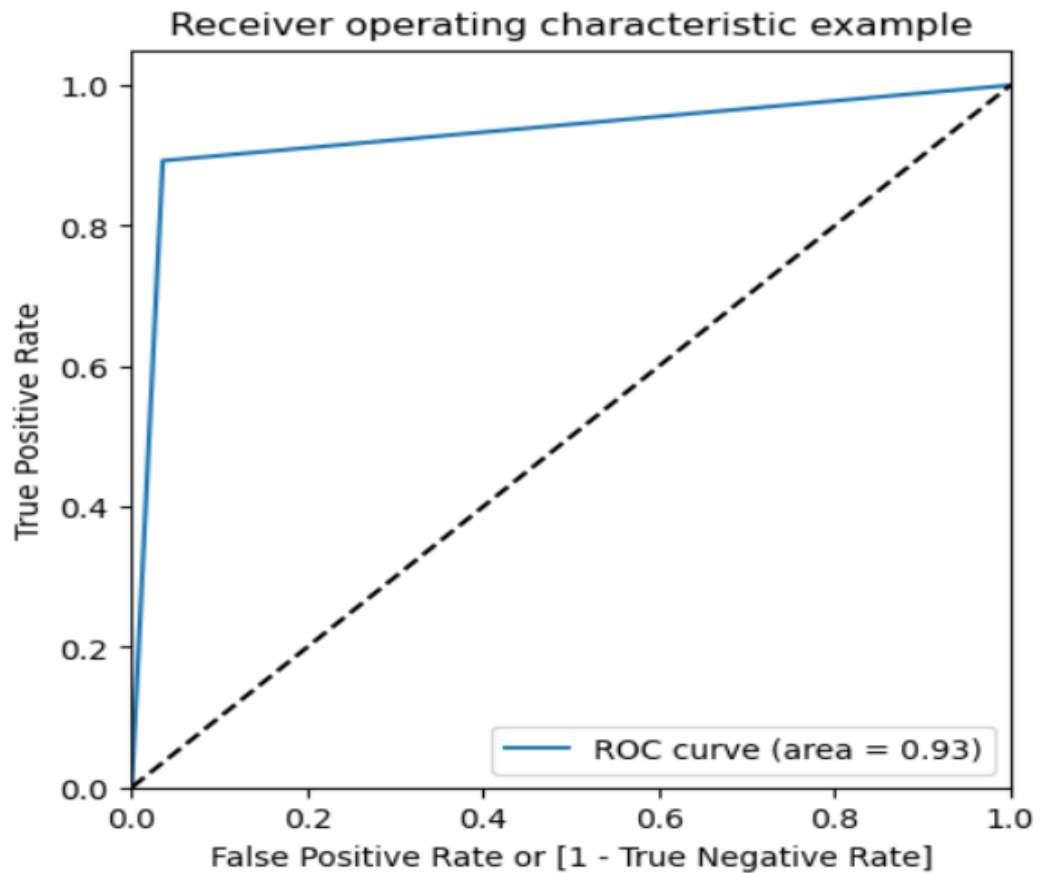
# DATA CONVERSION

- **Before initiating the model building, the dummy variables were created using 'get dummies'**

- **The data type of Boolean columns was also converted into integer type**

# MODEL BUILDING (1/3)

- 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.

- Feature scaling

- Checking correlation of variables using heat maps

- 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

- Predicting the probabilities on train set and evaluating the model using accuracy, sensitivity specificity.

- Predictions on test data

# MODEL BUILDING (2/3)

# MODEL BUILDING (3/3)

# MODEL BUILDING (4/4)

```
## Calculating Specificity and Specificity
print("Sensitivity with threshold as 0.35 is :", round(TP/(TP+FN)*100,2),"%")
print("Specificity with threshold as 0.35 is :", round(TN/(TN+FP)*100,2),"%")
```

```
Sensitivity with threshold as 0.35 is : 91.52 %
Specificity with threshold as 0.35 is : 92.2 %
```

## Precision-Recall

```
print("Precision with threshold as 0.35 is :", round(precision_score(y_train_pred_f.Converted,y_train_pred_f.Final_Predicted)*100,2),"%")
print("Recall with threshold as 0.35 is :", round(recall_score(y_train_pred_f.Converted,y_train_pred_f.Final_Predicted)*100,2),"%")
```

```
Precision with threshold as 0.35 is : 88.93 %
Recall with threshold as 0.35 is : 93.36 %
```

```
Sensitivity of test set is : 91.52 %
Specificity of test set is : 92.2 %
Precision of test set is : 88.58 %
Recall of test set is : 91.52 %
```

# KEY INSIGHTS

The variables identified as most influential in predicting potential buyers, ranked in descending order of importance, are:

- Tags

- Lead Source

- Total Time Spent on Website

- Last Notable Activity

- Last Activity.

- Current occupation as a working professional.

In conclusion, focusing efforts on these key variables presents an opportunity for X Education to significantly enhance its conversion rates by effectively engaging and persuading potential buyers.

# THANK YOU