# Lead Score
# Logistic Regression

- Mohan T
- Unnikrishnan
- Vanshika

# Business Problem & Approach

## Business Problem:

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

X Education offers online courses for industry professionals. Despite receiving numerous leads, their conversion rate is quite low. (30%).

To make this process more efficient, the company wishes to identify the most potential leads, also known as 'Hot Leads'

## Approach:

**Lead Data:** files has been provided with leads information containing approximately 9,000 data points

Perform Data Cleaning and EDA

Develop a model to identify leads with a conversion probability greater than 80%

# Detailed Approach

Import data & Perform data checks

EDA

Dummy Variable creation

Test- Train Split

Feature scaling

Model building

- Features selection – RFE
- Checking correlation between variables i.e., VIF and p-values
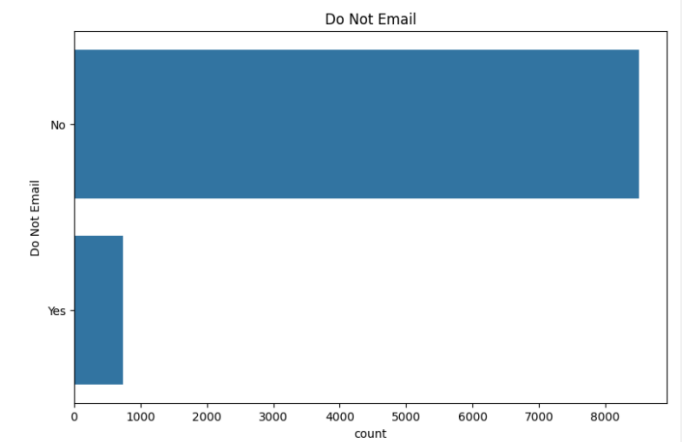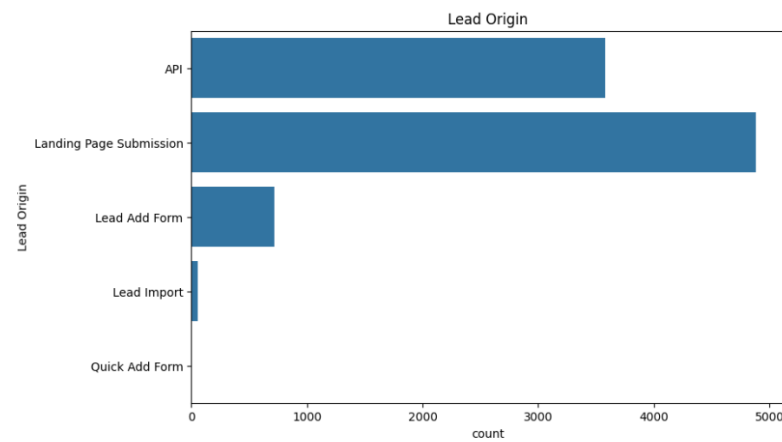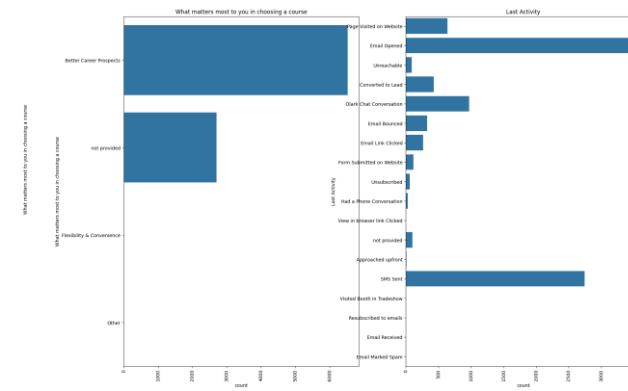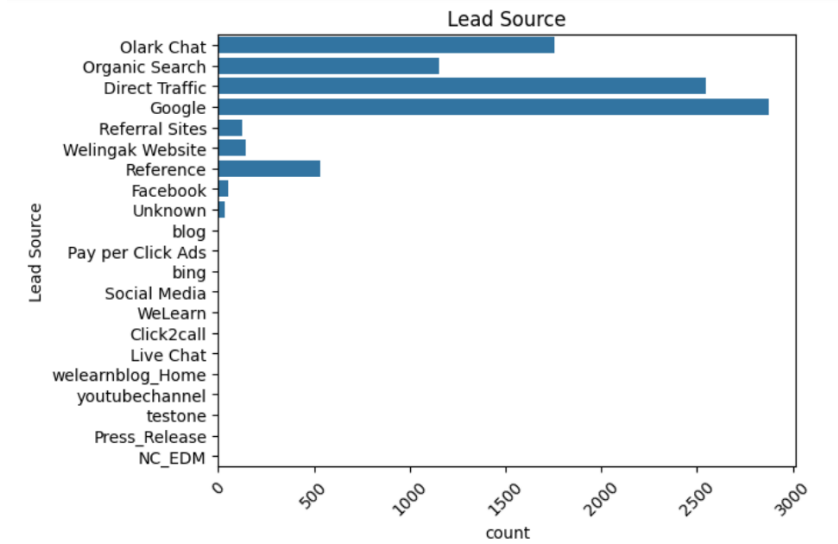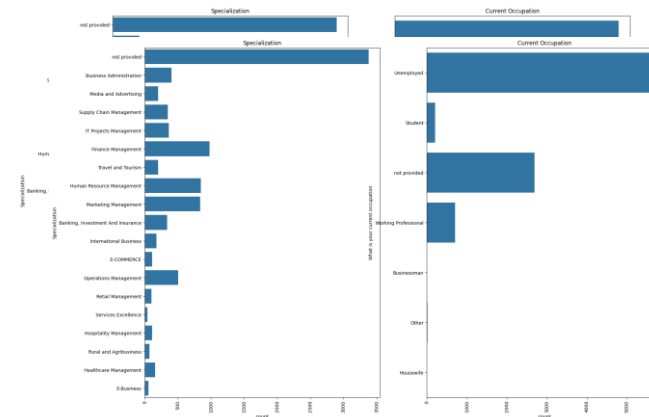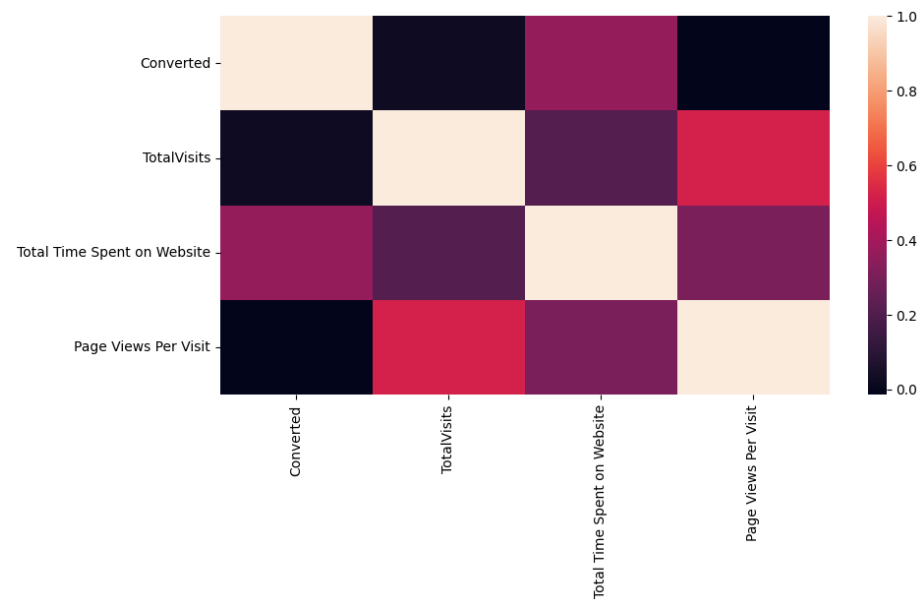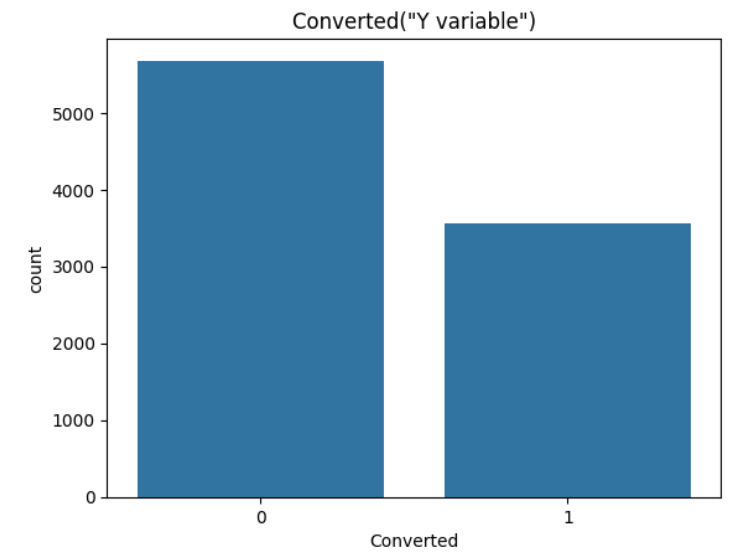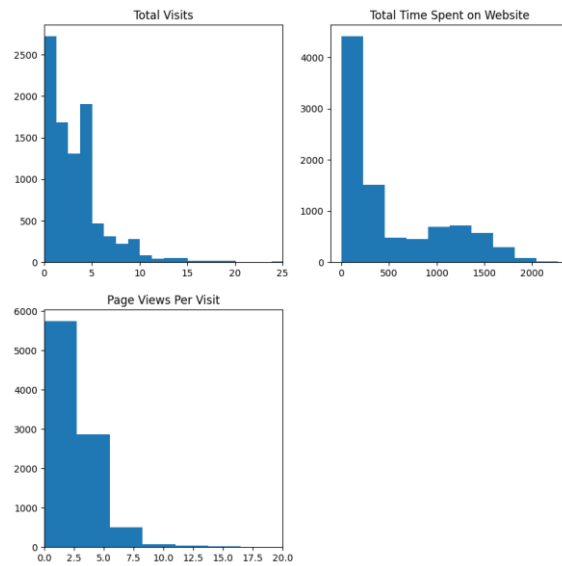
Model Evaluation

Making predictions

# Clean and Transform Data

- Started with loading data and analyzing data for null/missing values. Columns that have more than 35% of missing values, excluded from analysis and model building.

- Columns with unique values can be removed, as they do not impact the overall solution.

- Data with "select" has been replaced with "Not provided".

- Numerical variable missing values are replaced with 0. I.e., Total Time Spent on Website, Page Views Per Visit and total visits.

- Since the majority of leads are from India, with only a few from other countries, the data has been grouped into two categories: India and Others
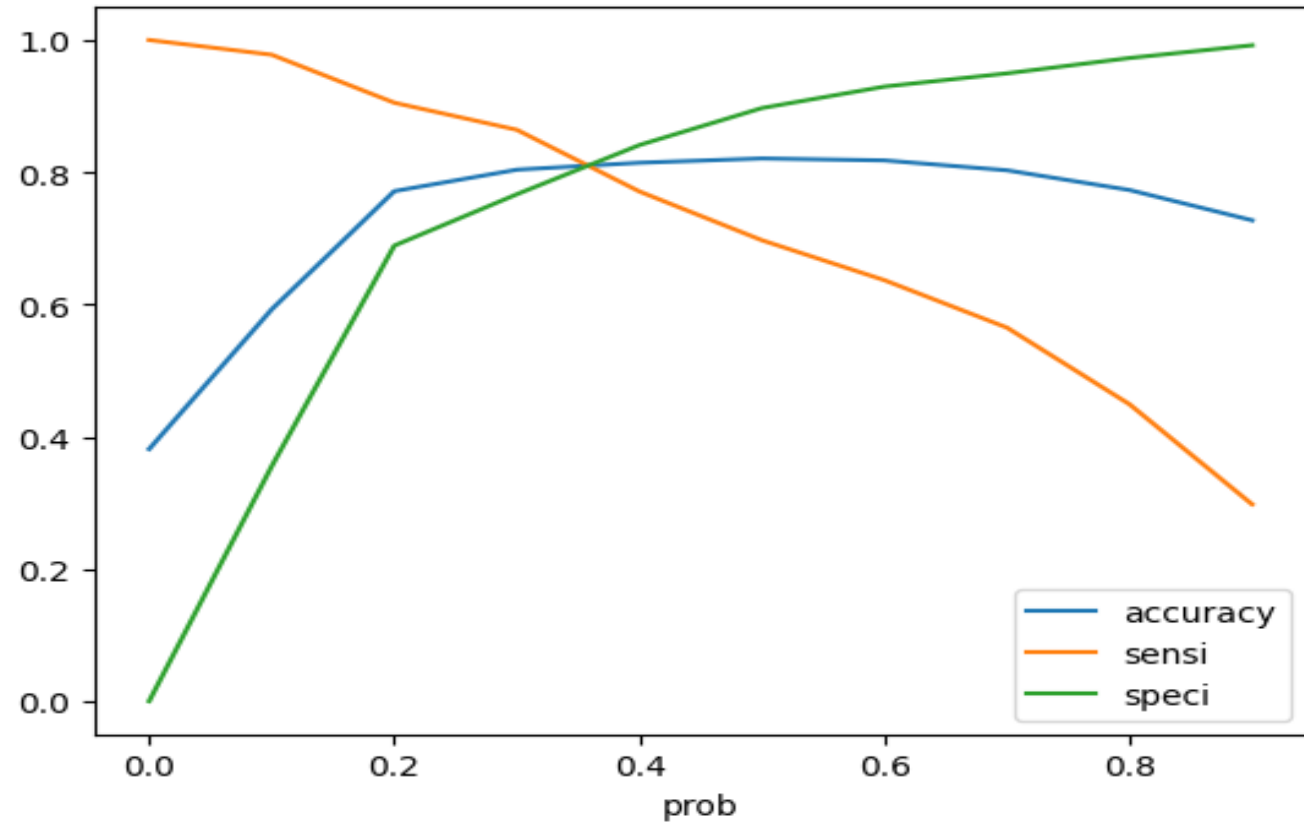
# EDA

# EDA

# ROC Curve



**Finding Optimal Cut off Point:**
- Probability where we get balanced sensitivity and specificity
- Optima cut off point is .35

# Observations



**Train test Data**

Accuracy: 82%

Sensitivity: 69%

Specificity: 89%



**Test test Data**

Accuracy: 82%

Sensitivity: 75%

Specificity: 86%

- Feature list
  - Do Not Email
  - Total Time Spent on Website
  - Lead Origin_Lead Add Form
  - Lead Source_Welingak Website
  - Last Activity
    - Had a Phone Conversation
    - SMS Sent