



# X Education - Lead Scoring Case Study

Detection of Hot Leads to concentrate more of marketing efforts on them, improving conversion rates for X Education

Team Members: Ganapathi Chitturi, Anwesha Kundu & Tabassum Sayeed

# Table of Contents

- Background of X Education Company
- Problem Statement & Objective of the Study
- Suggested Ideas for Lead Conversion
- Analysis Approach
- Data Cleaning
- EDA
- Data Preparation
- Model Building (Regression - Binominal Class ,VIF & Manual fine tuning)
- Model Evaluation
- Recommendations

# Background of X Education Company

- 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.
- 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%.

# Problem Statement & Objective of the Study

## Problem Statement:

- X Education gets a lot of leads, its lead conversion rate is very poor at around 30%
- X Education wants to make lead conversion process more efficient by identifying the most potential leads, also known as Hot Leads
- Their sales team want to know these potential set of leads, which they will be focusing more on communicating rather than making calls to everyone.

## Objective of the Study:

- To help X Education 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 a higher lead score have a higher conversion chance and the customers with a lower lead score have a lower conversion chance.
- The CEO has given a ballpark of the target lead conversion rate to be around 80%.

# Suggested Ideas for Lead Conversion



## Leads Grouping

- Leads are grouped based on their propensity or likelihood to convert.
- This results in a focused group of hot leads.



## Better Communication

- We could have a smaller pool of leads to communicate with, which would allow us to have a greater impact.



## Boost Conversion

- We would have a greater conversion rate and be able to hit the 80% objective since we concentrated on hot leads that were more likely to convert.



Since we have a target of 80% conversion rate, we would want to obtain a high **sensitivity** in obtaining hot leads.

# Analysis Approach



## Data Cleaning:

Loading Data Set,  
understanding &  
cleaning data



## EDA:

Check imbalance,  
Univariate &  
Bivariate analysis



## Data Preparation

Dummy variables,  
test-train split,  
feature scaling



## Model Building:

VIF & Manual  
Feature Reduction  
& finalizing model



## Model Evaluation:

Confusion matrix,  
Cutoff Selection,  
assigning Lead  
Score



## Predictions on Test Data:

Compare train vs  
test metrics, Assign  
Lead Score and get  
top features



## Recommendation:

Suggest top 3  
features to focus for  
higher conversion &  
areas for  
improvement

# Data Cleaning

- **"Select"** level represents null values for some categorical variables, as customers did not choose any option from the list so convert to most appearing value.
- Columns with over 40% null values were dropped.
- Missing values in categorical columns were handled based on value counts and certain considerations.
- Drop columns that don't add any insight or value to the study objective (tags, country)
- Imputation was used for some categorical variables.
- Additional categories were created for some variables.
- Columns with no use for modeling (Prospect ID, Lead Number) or only one category of response were dropped.
- Category Variables are minimized to smaller number of values
- Numerical data was imputed with mode after checking distribution.

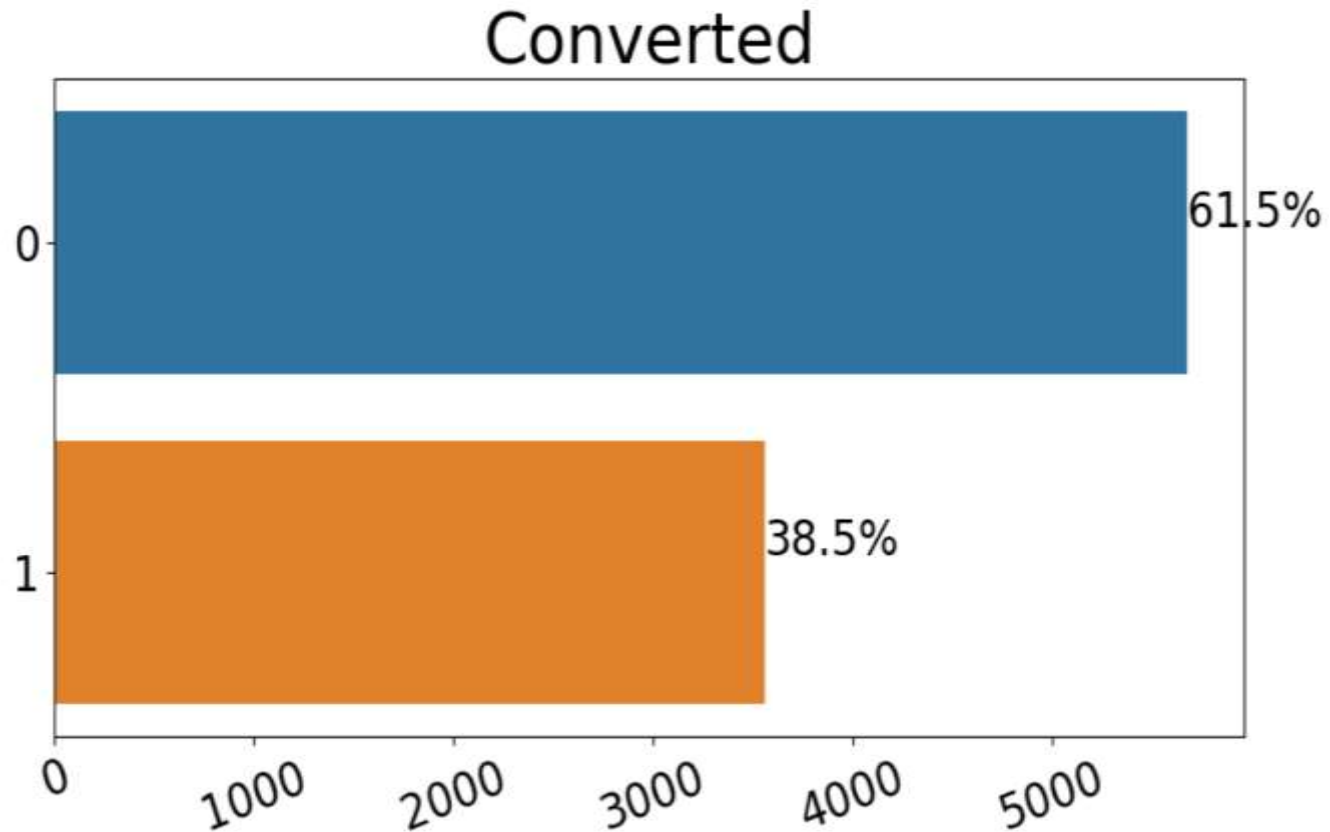
# Data Cleaning

- Skewed category columns were checked and dropped to avoid bias in logistic regression models.
- Outliers in **TotalVisits** and **Page Views Per Visit** were treated and capped.
- Invalid values were fixed and data was standardized in some columns, such as lead source.
- Low frequency values were grouped together to “Others” in all category variables.
- Binary categorical variables were mapped.
- Other cleaning activities were performed to ensure data quality and accuracy.
  - Fixed Invalid values & Standardizing Data in columns by checking casing styles, etc. (lead source has Google, google)



# EDA

- Data is imbalanced while analyzing target variable.

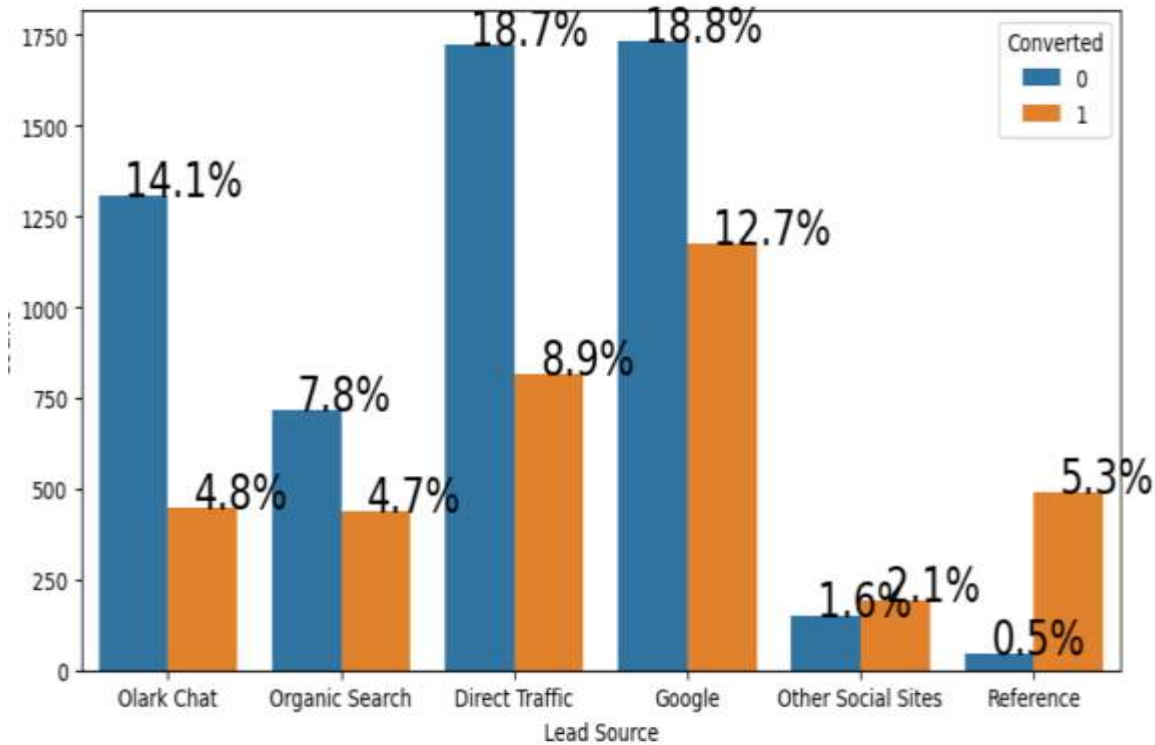


Conversion rate is of 38.5%, meaning only 38.5% of the people have converted to leads. (Minority)

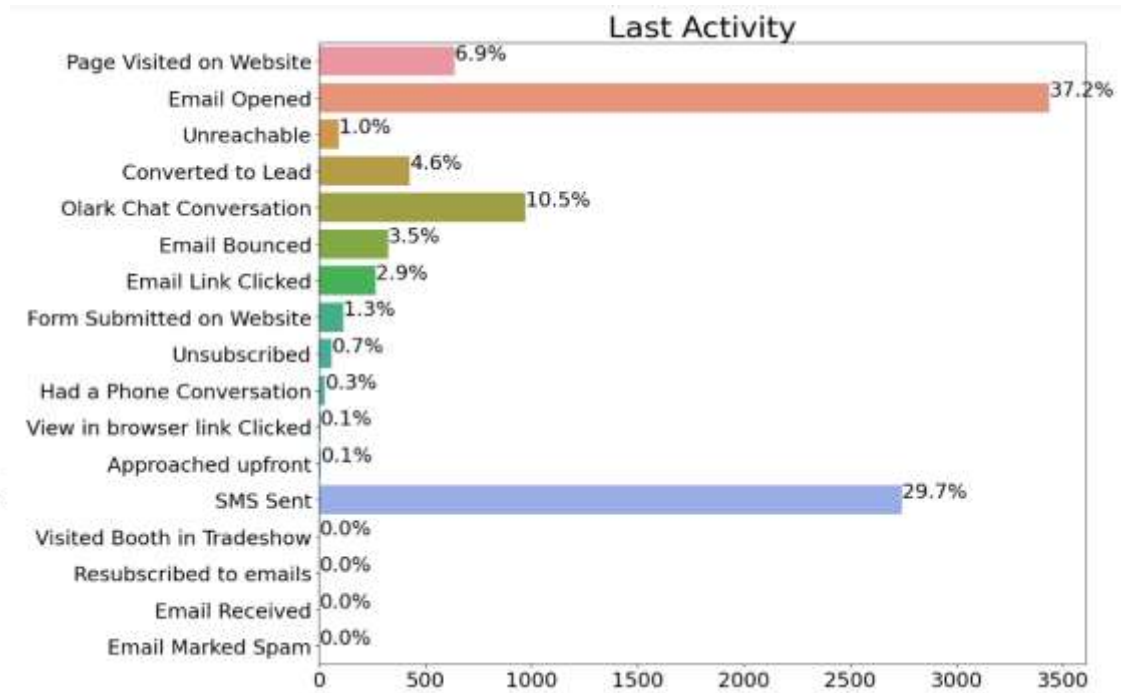
While 61.5% of the people didn't convert to leads. (Majority)

# EDA

## Univariate Analysis - Categorical Variables



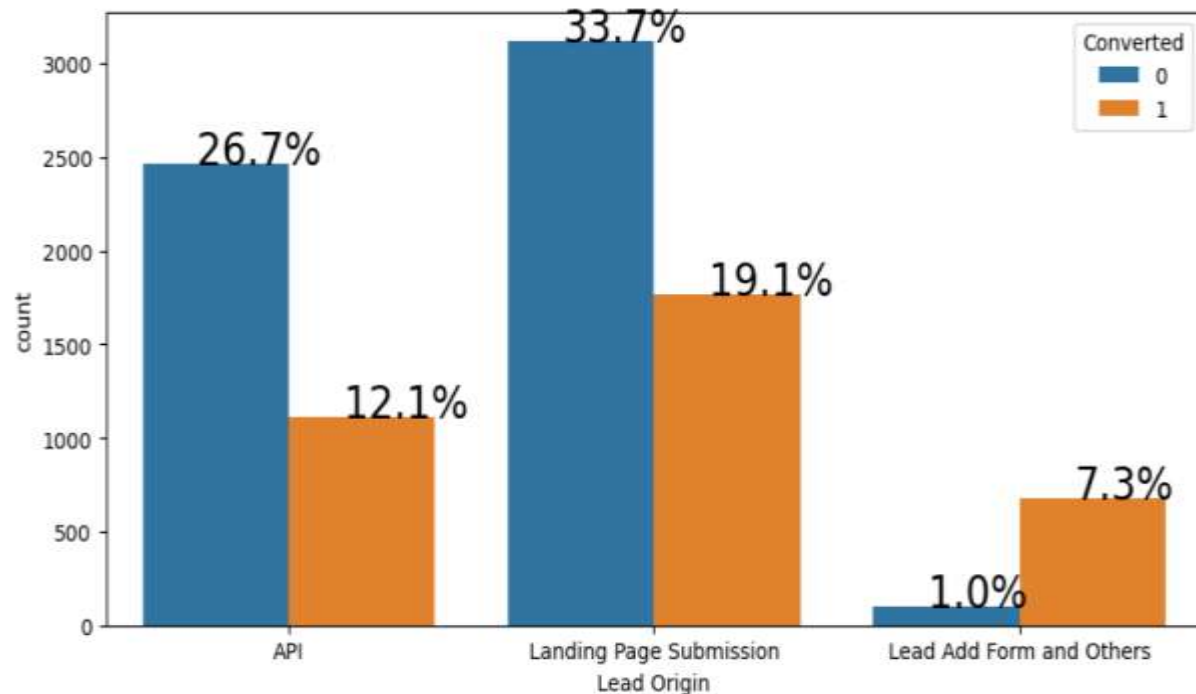
- **Lead Source:** 59.1% Lead source is from Google & Direct Traffic combined.



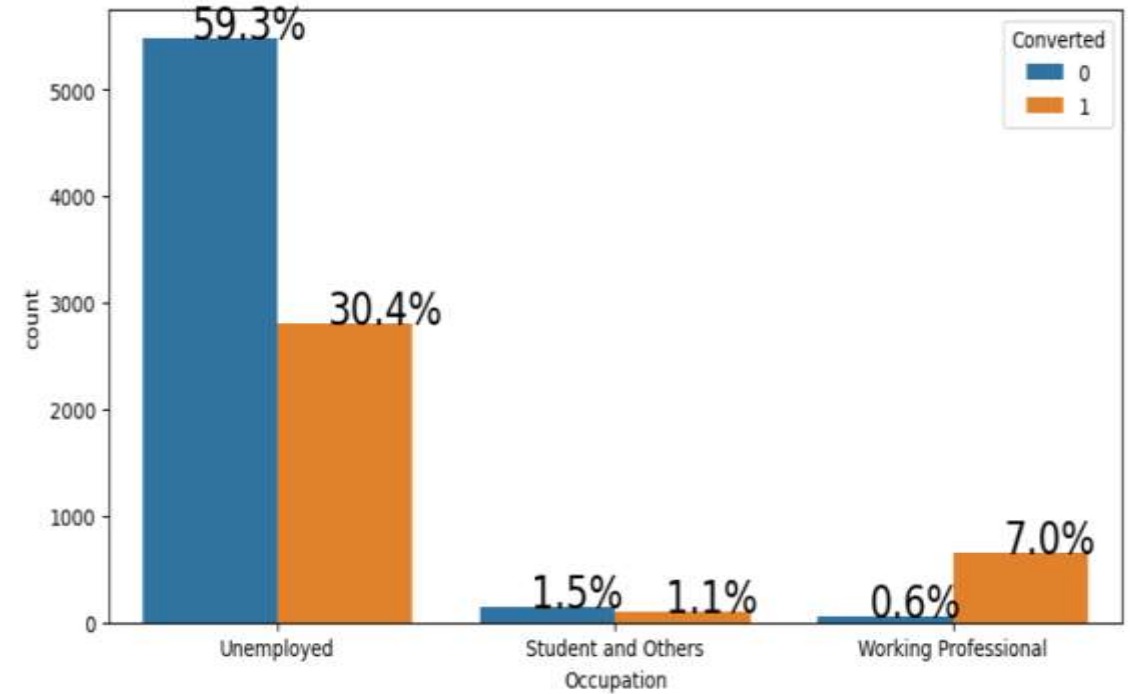
- **Last Activity:** 66.9% of customers contribution in SMS Sent & Email Opened activities.

# EDA

## ● Univariate Analysis – Categorical Variables

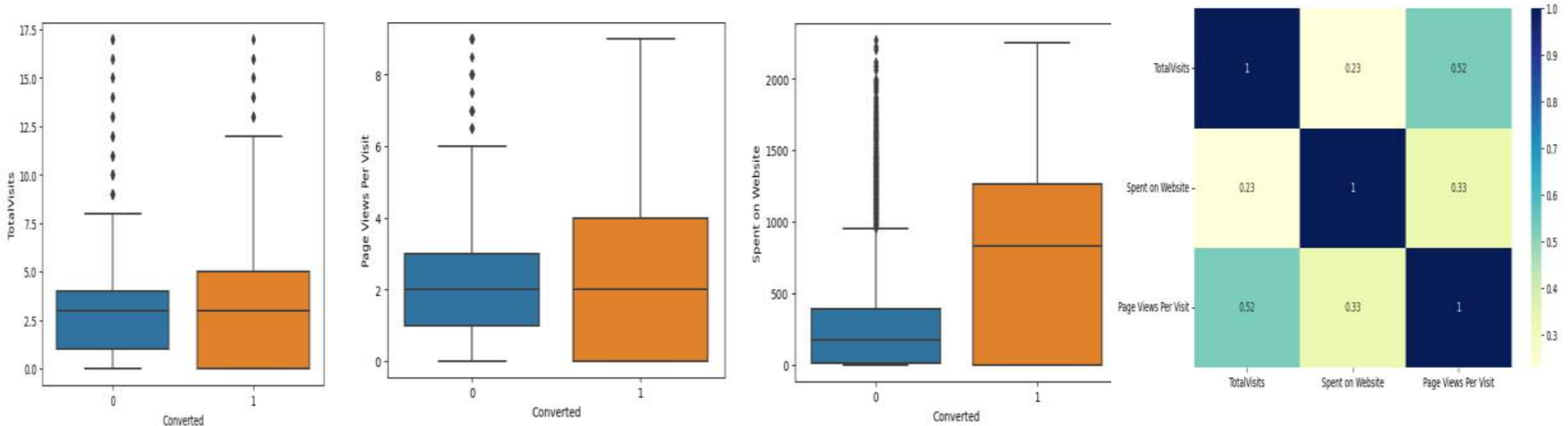


- **Lead Origin:** "Landing Page Submission" identified 52.8% of customers, "API" identified 38.8%.



- **Current\_occupation:** It has 89.7% of the customers as Unemployed.

# EDA - Bivariate Analysis for Numerical Variables



- Past Leads who **spends more time on the Website** have a higher chance of getting successfully converted than those who spends less time as seen in the **box-plot**

# Data Preparation before Model building

- Binary level categorical columns were already mapped to 1 / 0 in previous steps
- Created dummy features (one-hot encoded) for categorical variables - Lead Origin, Lead Source, Last Activity, Specialization, Current\_occupation
- Splitting Train & Test Sets
  - 70:30 % ratio was chosen for the split
- Feature scaling
  - Standardization method was used to scale the features
- Checking the correlations
  - Predictor variables which were highly correlated with each other were dropped (Lead Origin\_Lead Import and Lead Origin\_Lead Add Form).

# Model Building

## Feature Selection

- The data set has lots of dimension and large number of features.
- This will reduce model performance and might take high computation time.
- Hence it is important to perform **variance\_inflation\_factor** (VIF) & RFE and to select only the important columns.
- Then we can manually fine tune the model.
- RFE outcome
  - Pre RFE - 34 columns & Post RFE - 17 columns

# Model Building

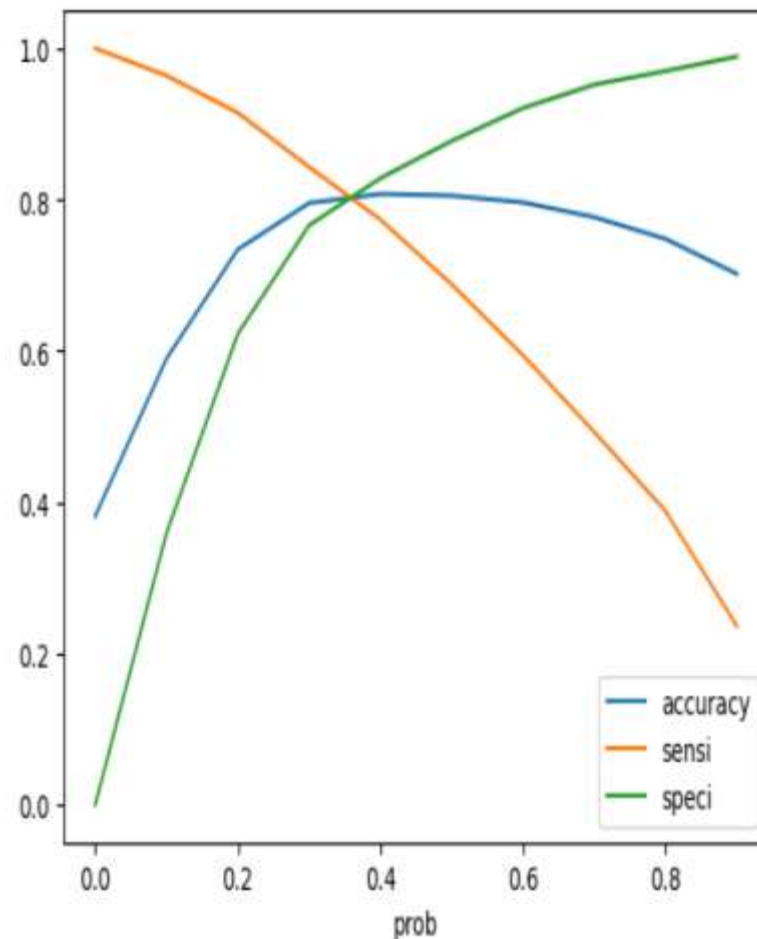
- Manual Feature Reduction process was used to build models by dropping variables with p - value greater than 0.05.
- Model 4 looks stable after four iteration with:
  - significant p-values within the threshold (p-values < 0.05) and
  - No sign of multicollinearity with VIFs less than 5
- Hence, **reg5** will be our final model, and we will use it for Model Evaluation which further will be used to make predictions.

# Model Evaluation

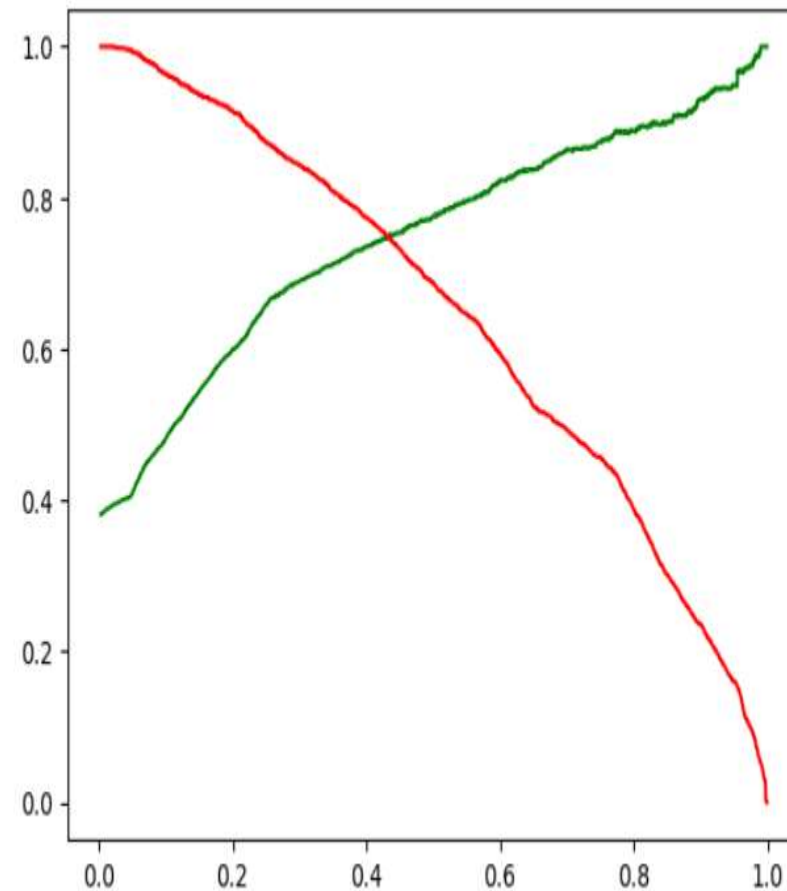
Confusion Matrix & Evaluation Metrics  
with 0.355 as cutoff

## Train Data Set

It was decided to go ahead with 0.35  
as cutoff after checking evaluation  
metrics coming from both plots



Confusion Matrix & Evaluation Metrics  
with 0.41 as cutoff

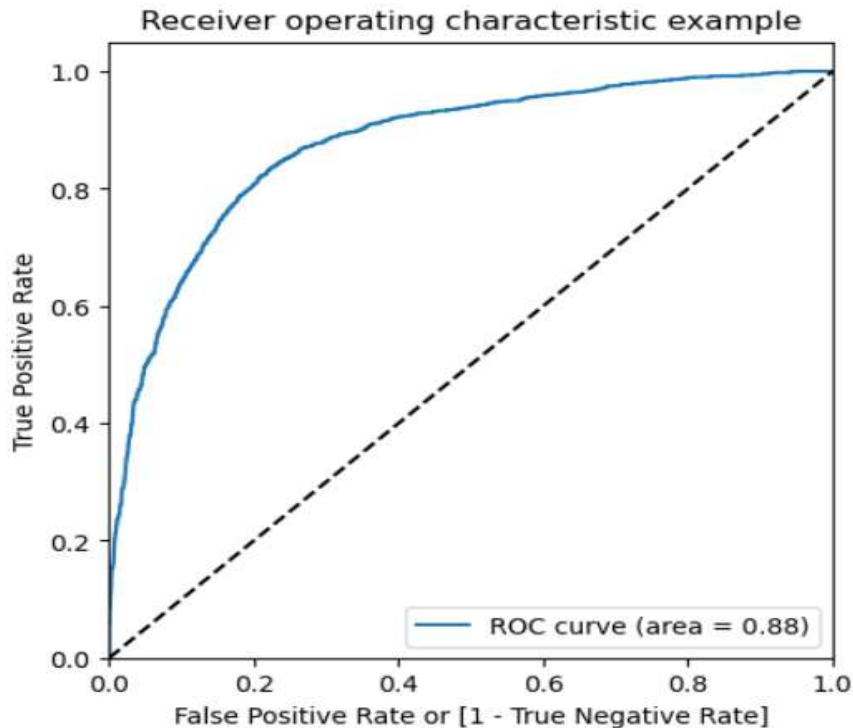




# Model Evaluation

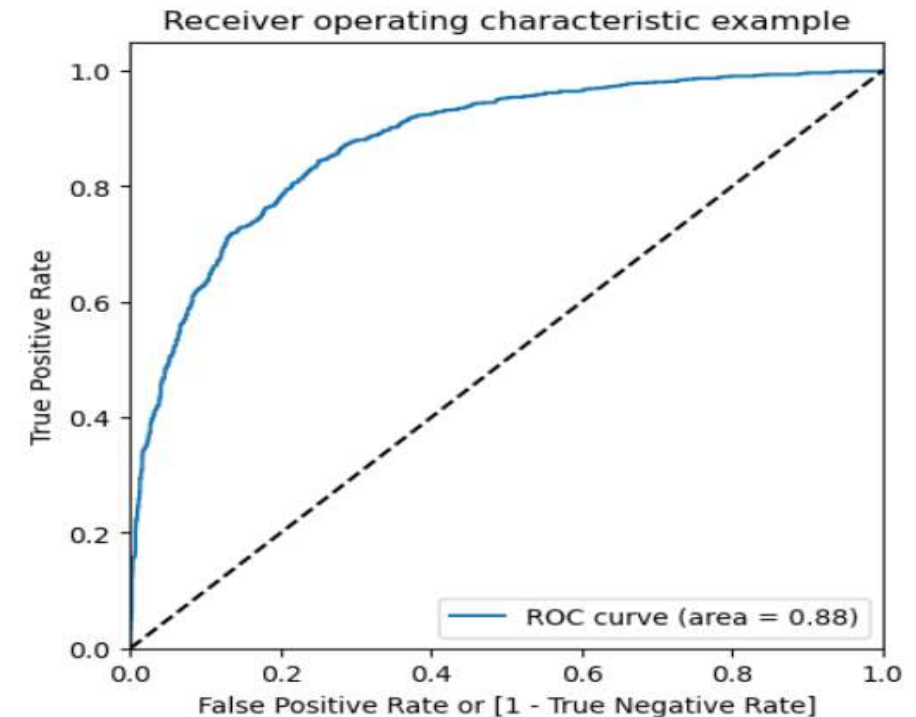
## ROC Curve - Train Data Set

- Area under ROC curve is 0.88 out of 1 which indicates a good predictive model.
- The curve is as close to the top left corner of the plot, which represents a model that has a high true positive rate and a low false positive rate at all threshold values.



## ROC Curve - Test Data Set

- Area under ROC curve is 0.87 out of 1 which indicates a good predictive model.
- The curve is as close to the top left corner of the plot, which represents a model that has a high true positive rate and a low false positive rate at all threshold values.



# Recommendation based on Final Model

- As per the problem statement, increasing lead conversion is crucial for the growth and success of X Education. To achieve this, we have developed a regression model that can help us identify the most significant factors that impact lead conversion.
- We have determined the following features that have the highest positive coefficients, and these features should be given priority in our marketing and sales efforts to increase lead conversion.
  - Lead Source\_Reference: 3.4907
  - Current\_occupation\_Working Professional: 2.47
  - Last Activity\_SMS Sent: 0.4457
  - Total Time Spent on Website: 1.0919
  - Last Activity\_Email Opened: 0.94
  - Lead Source\_Olark Chat: 0.6880

# Recommendation based on Final Model

## To increase our Lead Conversion Rates

- Focus on features with positive coefficients for targeted marketing strategies.
- Develop strategies to attract high-quality leads from top-performing lead sources.
- Optimize communication channels based on lead engagement impact.
- Engage **working professionals** with tailored messaging.
- More budget/spend can be done on **Google Website** in terms of advertising, etc.
- Incentives/discounts for providing reference that convert to lead, encourage providing more references.
- Working professionals to be aggressively targeted as they have high conversion rate and will have better financial situation to pay higher fees too.

## To identify areas of improvement

- Analyze negative coefficients in specialization offerings.
- Review landing page submission process for areas of improvement.



*Thank  
You!*