

# **Lead Scoring Case Study**

# 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.
- The CEO, in particular, has given a ballpark of the target lead conversion rate to be around 80%.

# Steps Taken

## **Data cleaning, preparation and 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.

# Steps Taken

## Data Manipulation

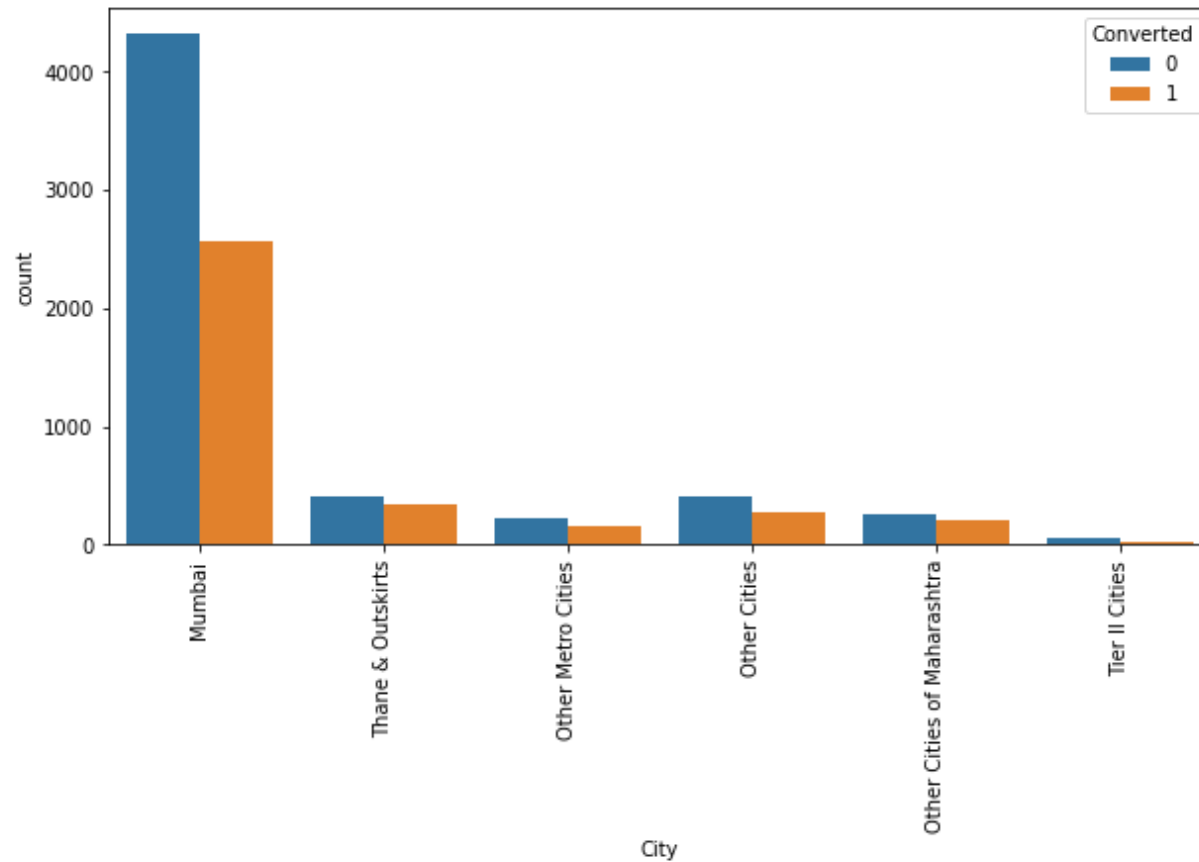
- Total Number of Rows =37, Total Number of Columns =9240.
- Removing the “Prospect ID” and “Lead Number” which is not necessary for the analysis.
- After checking for the value counts for some of the object type variables, we find some of the features which has no enough variance, which we have dropped, the features are: “Do Not Call”, “Last Activity”, “What matters most to you in choosing course”, “Search”, “Newspaper Article”, “X Education Forums”, “Newspaper”, “Digital Advertisement” etc.
- Dropping the columns having more than 40% as missing value such as ‘How did you hear about X Education’ and ‘Lead Profile’.

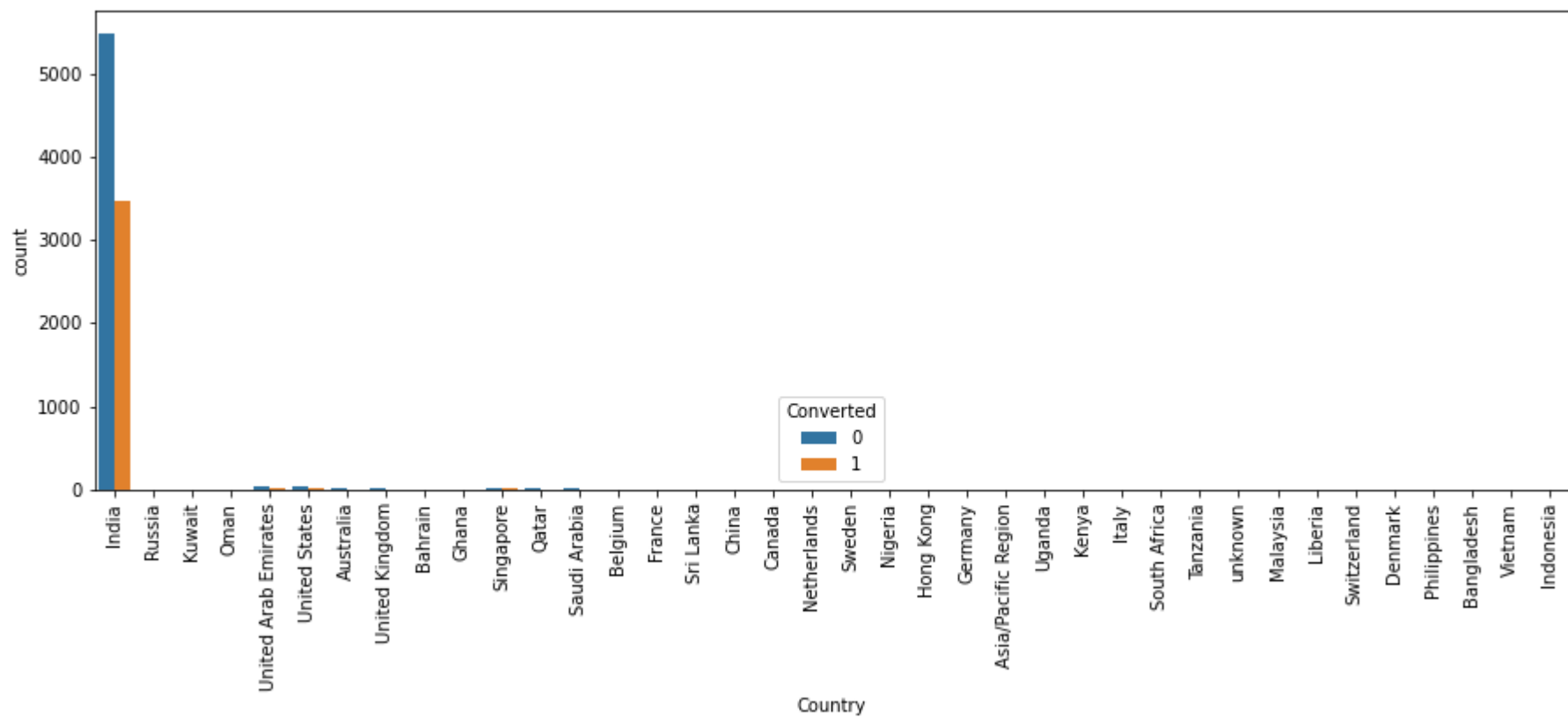
# Steps Taken

## EDA

- 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.
- Model presentation.
- Conclusions and recommendations.

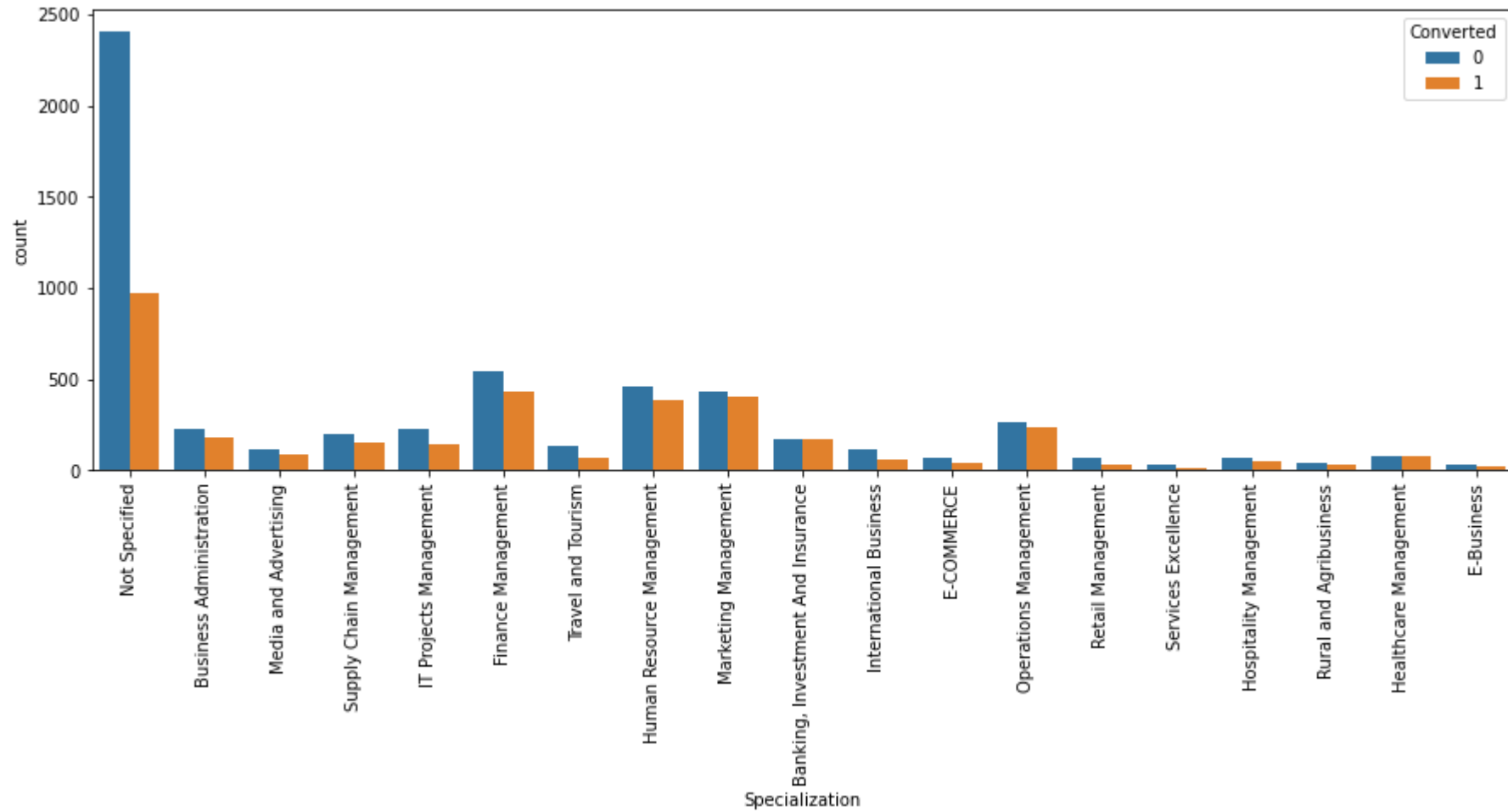
# EDA Visualization





Country

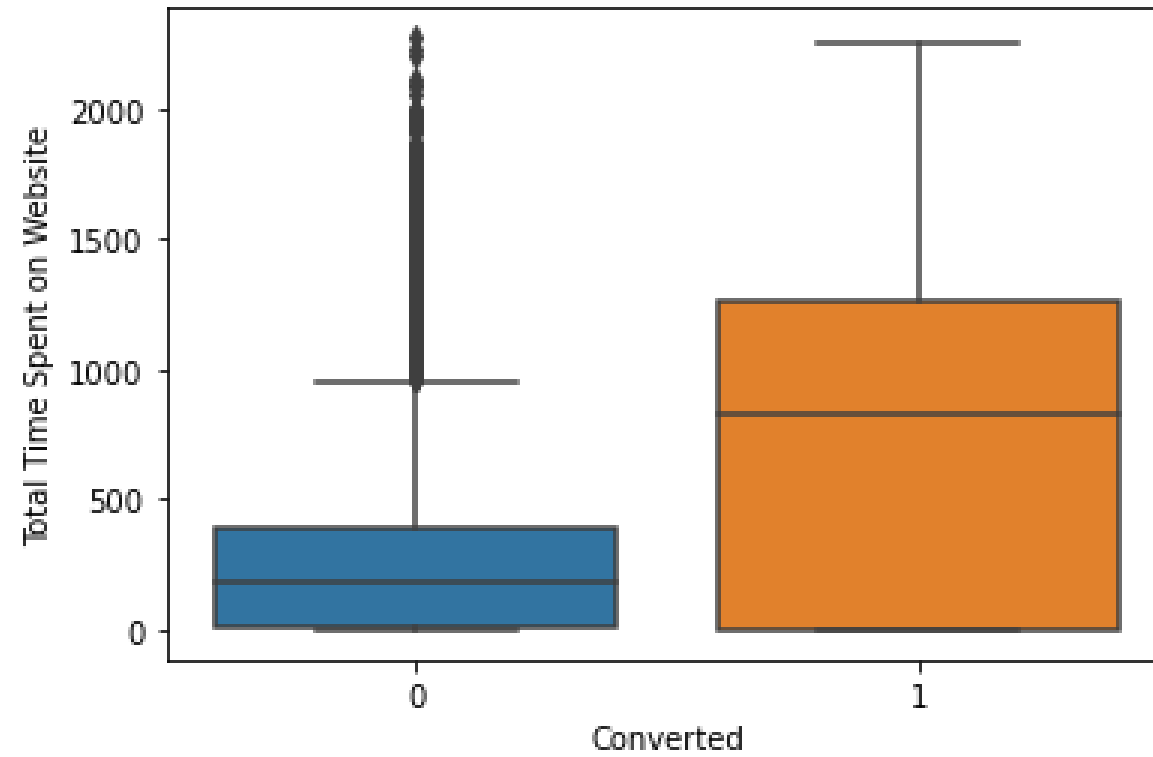




**Specialization**

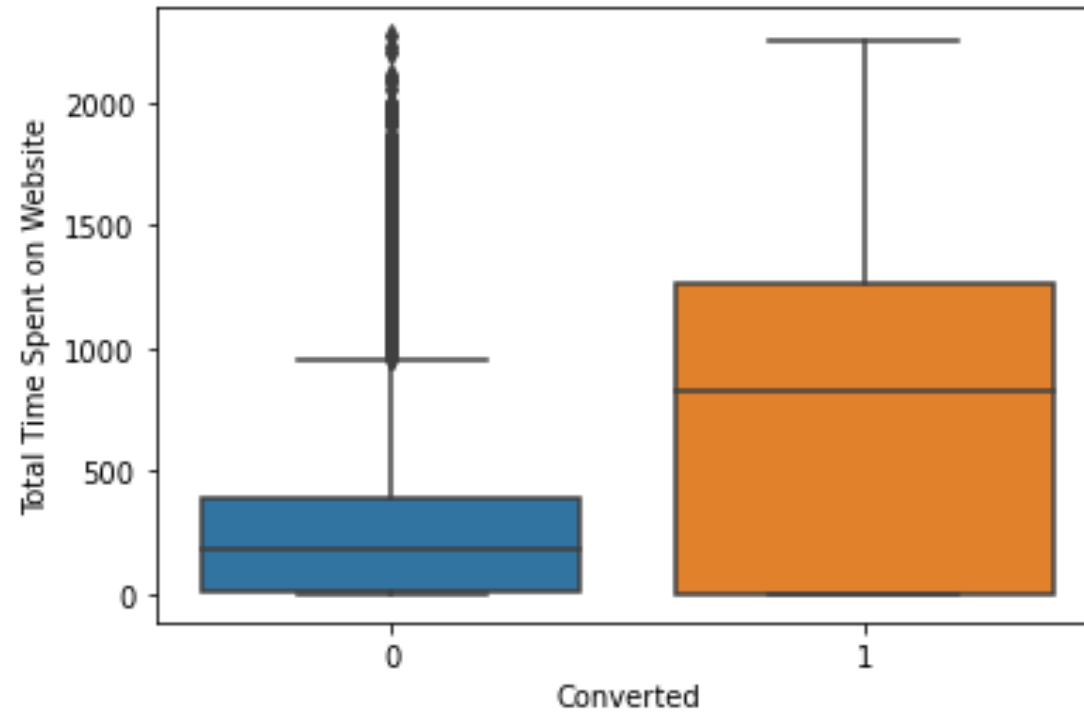
## Time Spent on Website vs Converted variable

Exploratory /



## Page Views Per Visit vs Converted

Page Views Per Visit



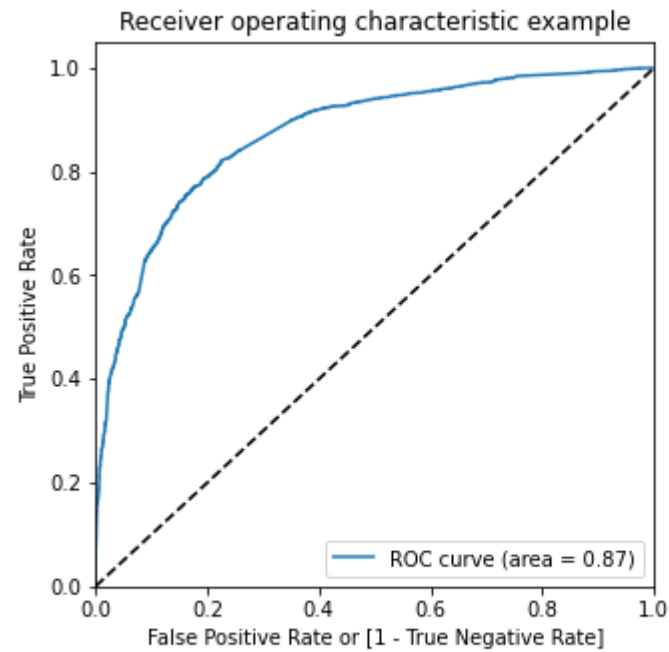
# Data Conversion

- Numerical Variables are Normalised
- Dummy Variables are created for object type variables
- Total Rows for Analysis: 8792
- Total Columns for Analysis: 43

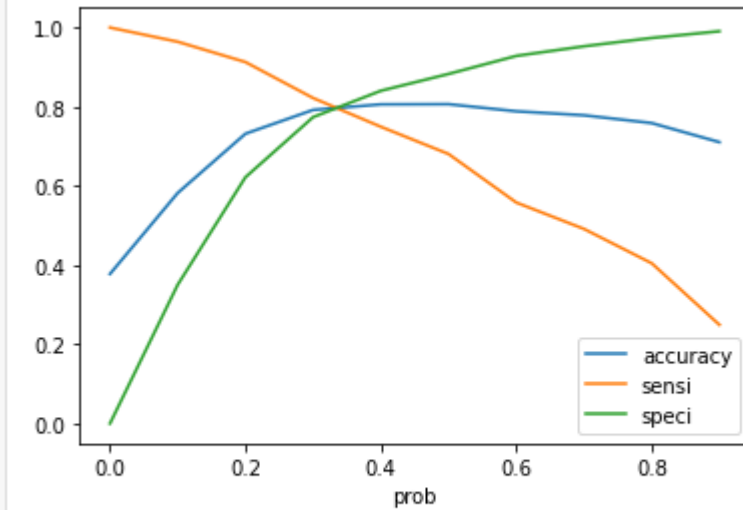
# Model Building

- 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 81%

# ROC Curve



The area under ROC curve is 0.87 which is a very good value.



From the graph it is visible that the optimal cut off is at 0.35

# Conclusion

It was found that the variables that mattered the most in the potential buyers

1. The total time spend on the Website.
2. Total number of visits.
3. When their current occupation is as a working professional.
4. Page Views Per Visit
5. When the lead source was:
  - a) Google
  - b) Direct traffic
  - c) Organic search
  - d) Welingak website
6. When the last activity was:
  - a) SMS
  - b) Olark chat conversation