

# Lead Score Case Study

---

Atanu Bhowmick

Ashish Ahuja

Abheer Brahme

# Contents

---

- Problem statement
- Problem approach
- EDA
- Model building
- Model evaluation
- Observation
- Conclusion



# Problem Statement

---

- 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. Moreover, the company also gets leads through past referrals. 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%.
- Now, although 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. The company looking to build a model where 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, in particular, has given a ballpark of the target lead conversion rate to be around 80%.

# Problem approach

---

Below steps are involved to approach the solution.

- EDA
- Model building
- Model evaluation
- Observation
- Conclusion

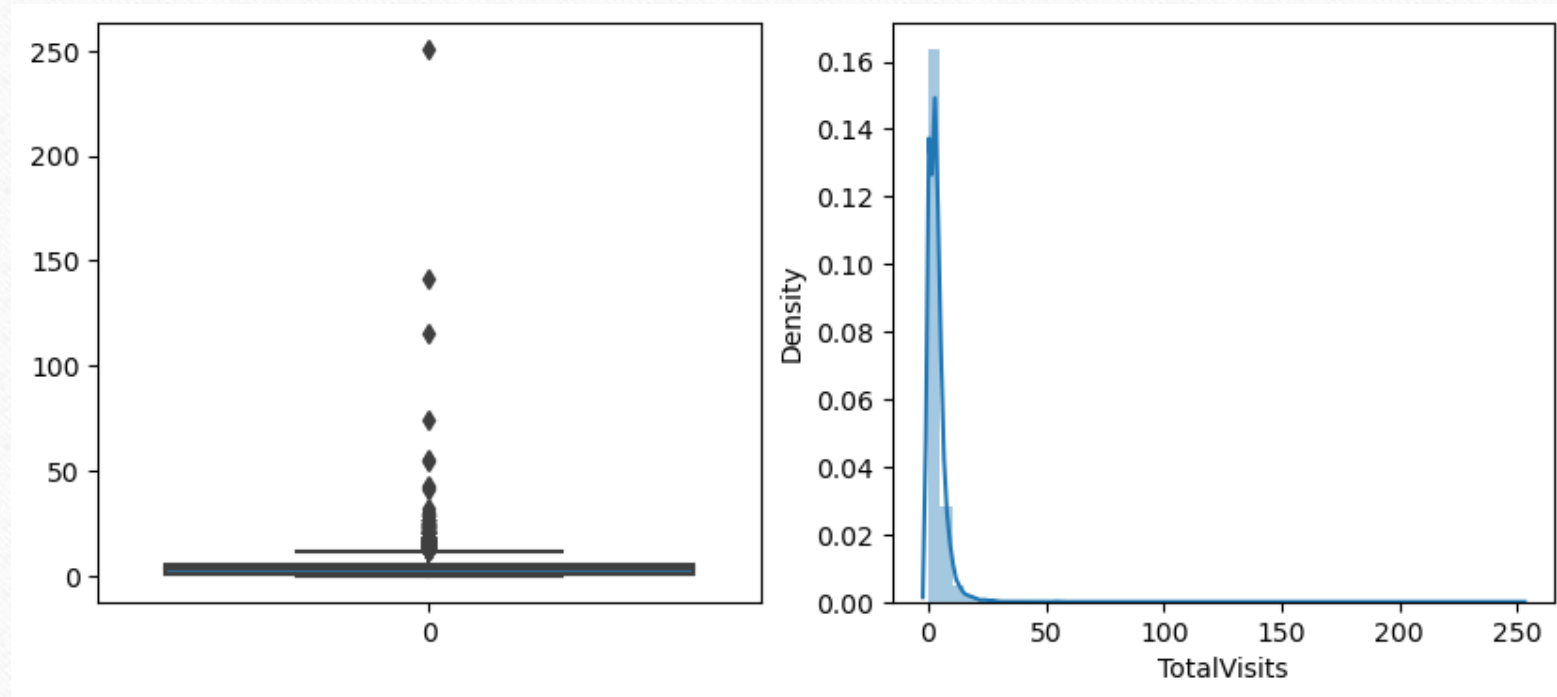
# Reading data & Cleanup

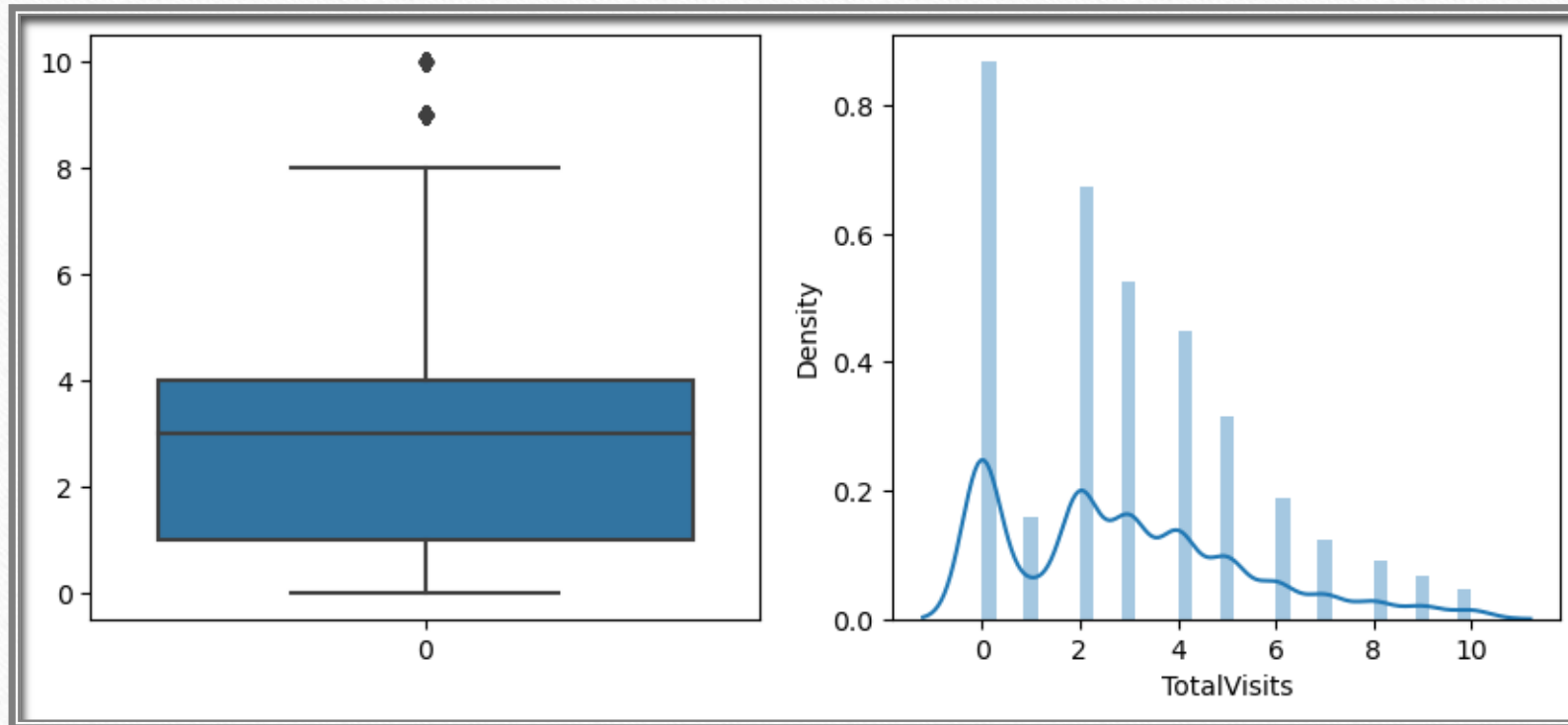
---

- There are total 37 rows and 9240 columns
- Columns which contain unique value (Prospect ID, Lead Number) are dropped
- Replaced 'Select' value with Null
- Columns for which missing data percentage is more than 33% are dropped
- The columns which represent the ad channels are dropped. Columns are (Search, Magazine, Newspaper Article, X Education Forums, Newspaper, Digital Advertisement, Through Recommendations)
- Country column is dropped as most of the user's country is India or missing

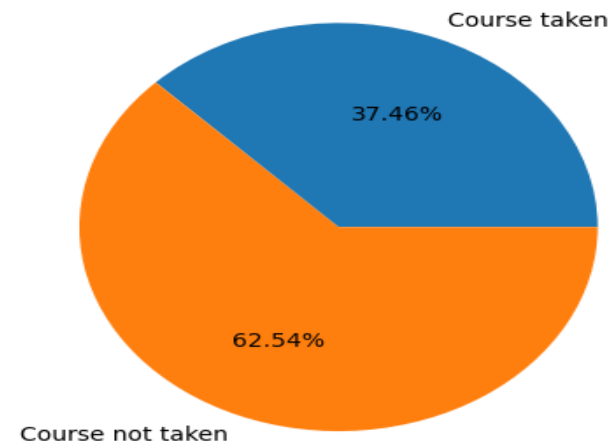
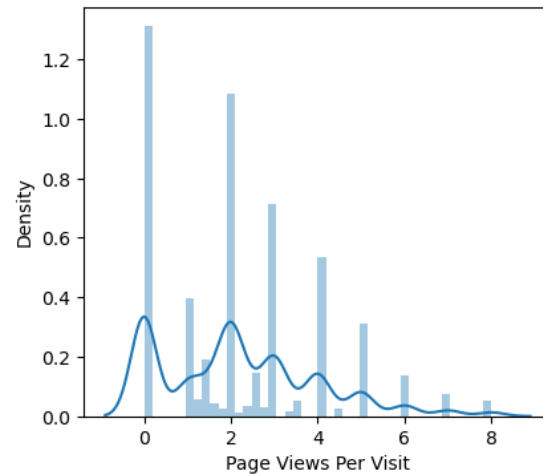
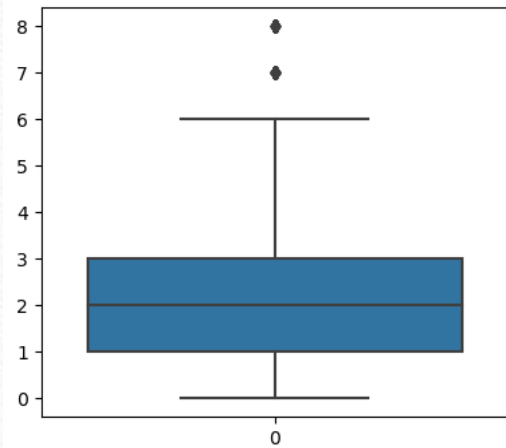
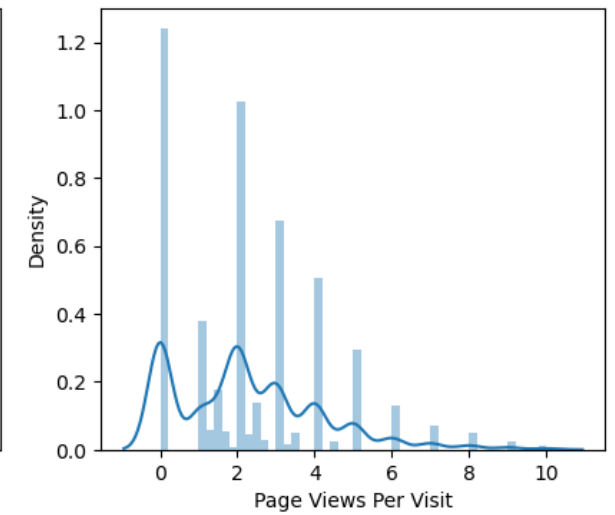
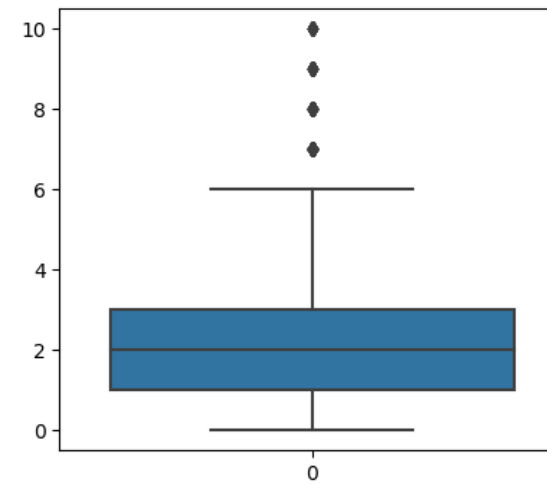
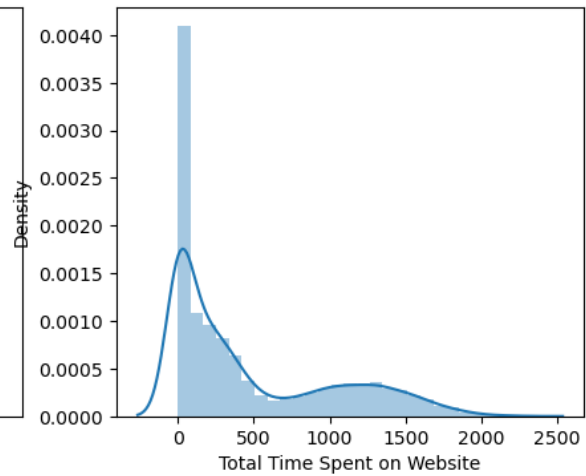
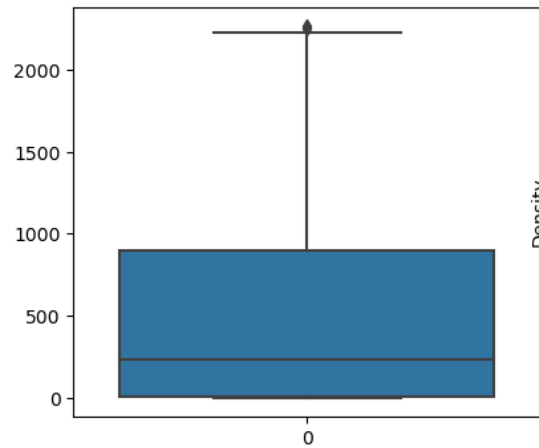


# EDA

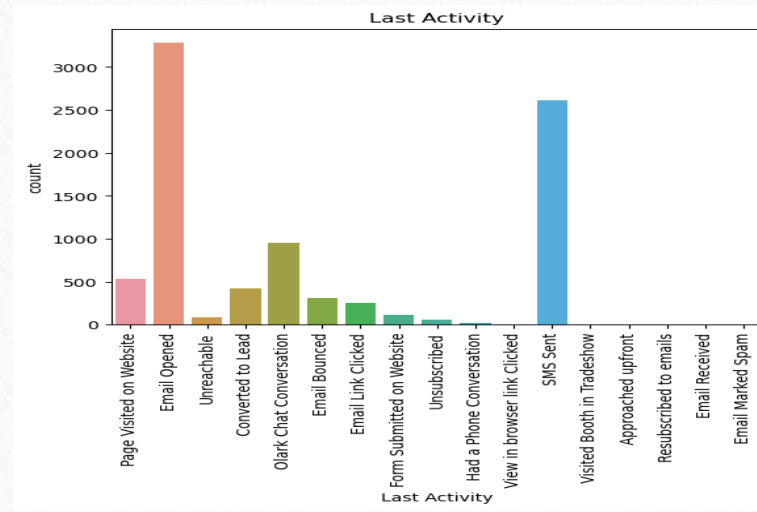
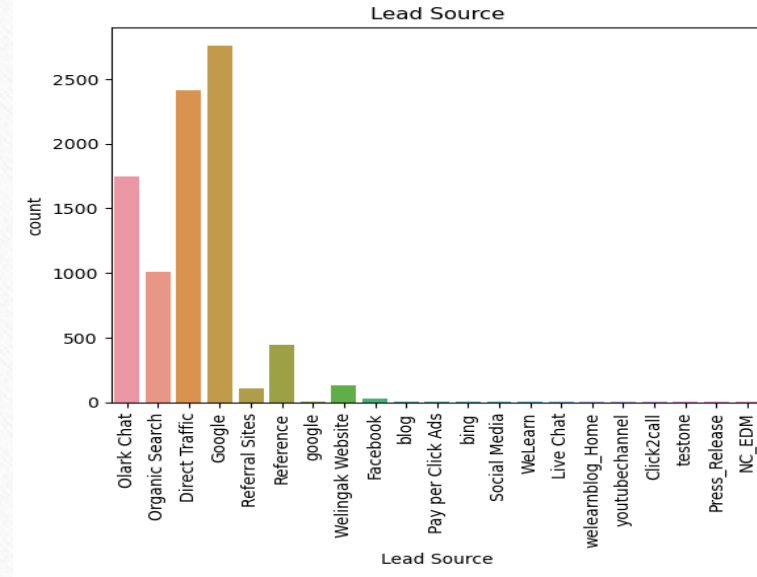
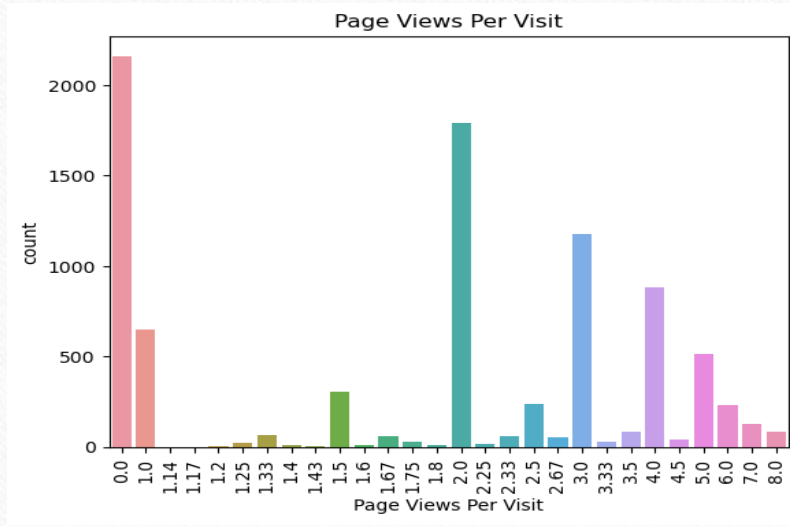
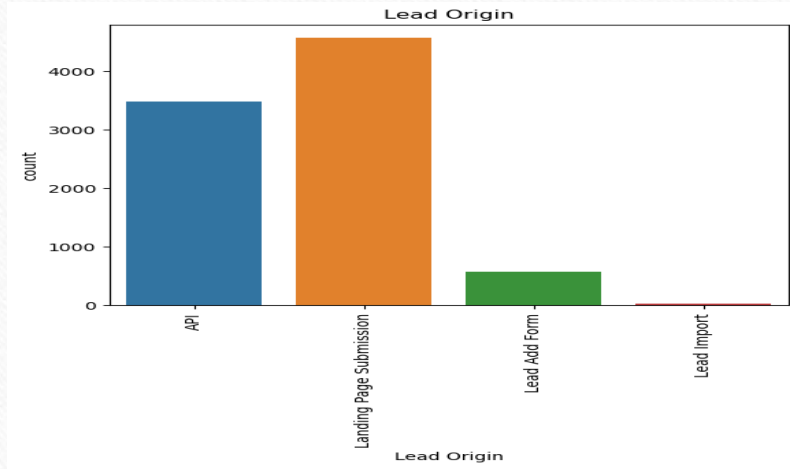




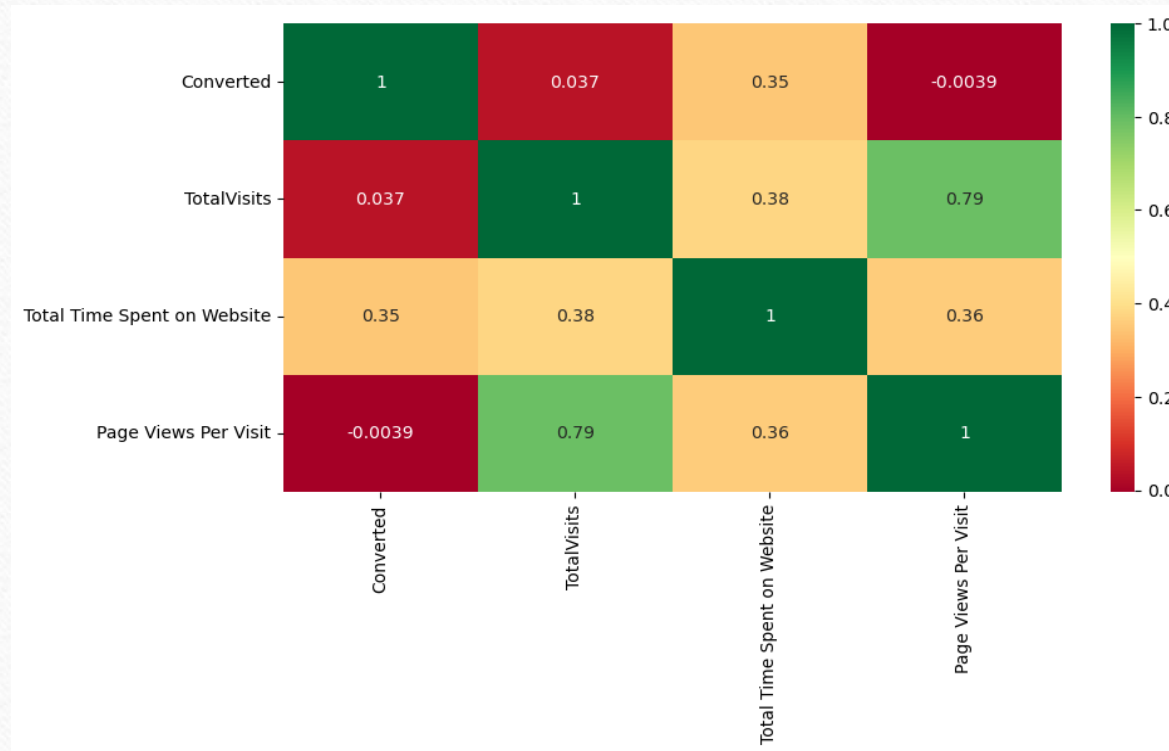
After removing outliers







# Correlation



# Data Conversion

---

- Numerical variables are normalized
- Dummy variables are created for object type variables
- Total rows for analysis: 8653
- Total columns for analysis: 11

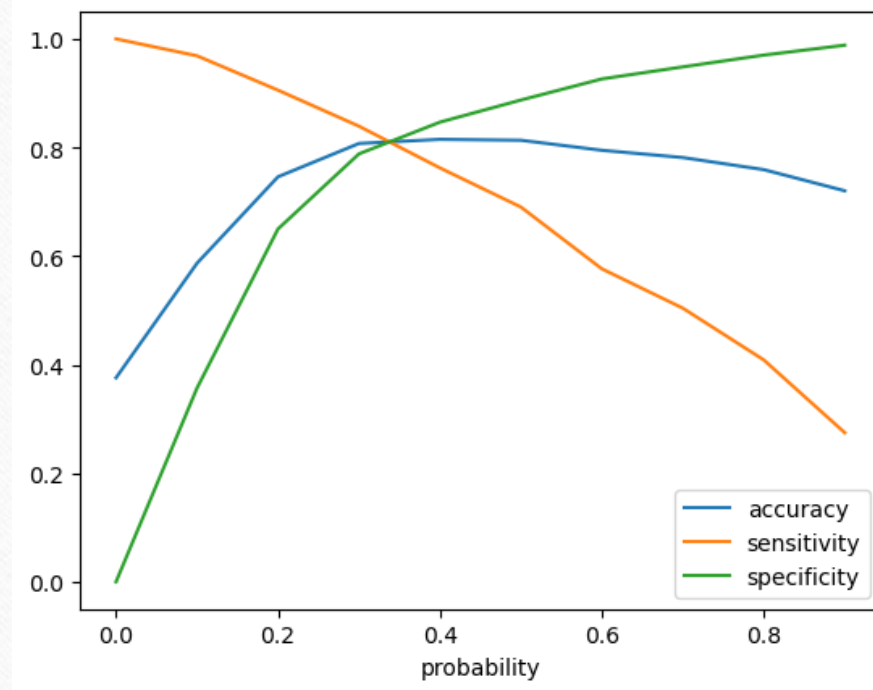


# Model Building

---

- Split the data into training and test set
- Train-Test split with 70:30 ratio
- The RFE for feature selection is used in Logistic Regression model
- Choosing 10 features out of all the features
- Model selection by dropping features which has p-value greater than 0.00 and VIF value greater than 2.5
- Predicting the test data
- Achieved the model accuracy is 80.20
- Model predicted with a optimum cutoff of 0.36

# ROC Curve



- 
- From the model we can see that the variables that mattered the most in the conversion are
    - Last notable activity
    - Lead source
    - Total time spent
    - Current Occupation



Thank You