

# Lead Score Case Study

Group members

1-Saba gul

2-Gayatri patil

3- Brijesh Dwivedi

# Problem Statement

- ▶ XEducation sells online courses to industry professionals.
- ▶ XEducation 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 rates 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.

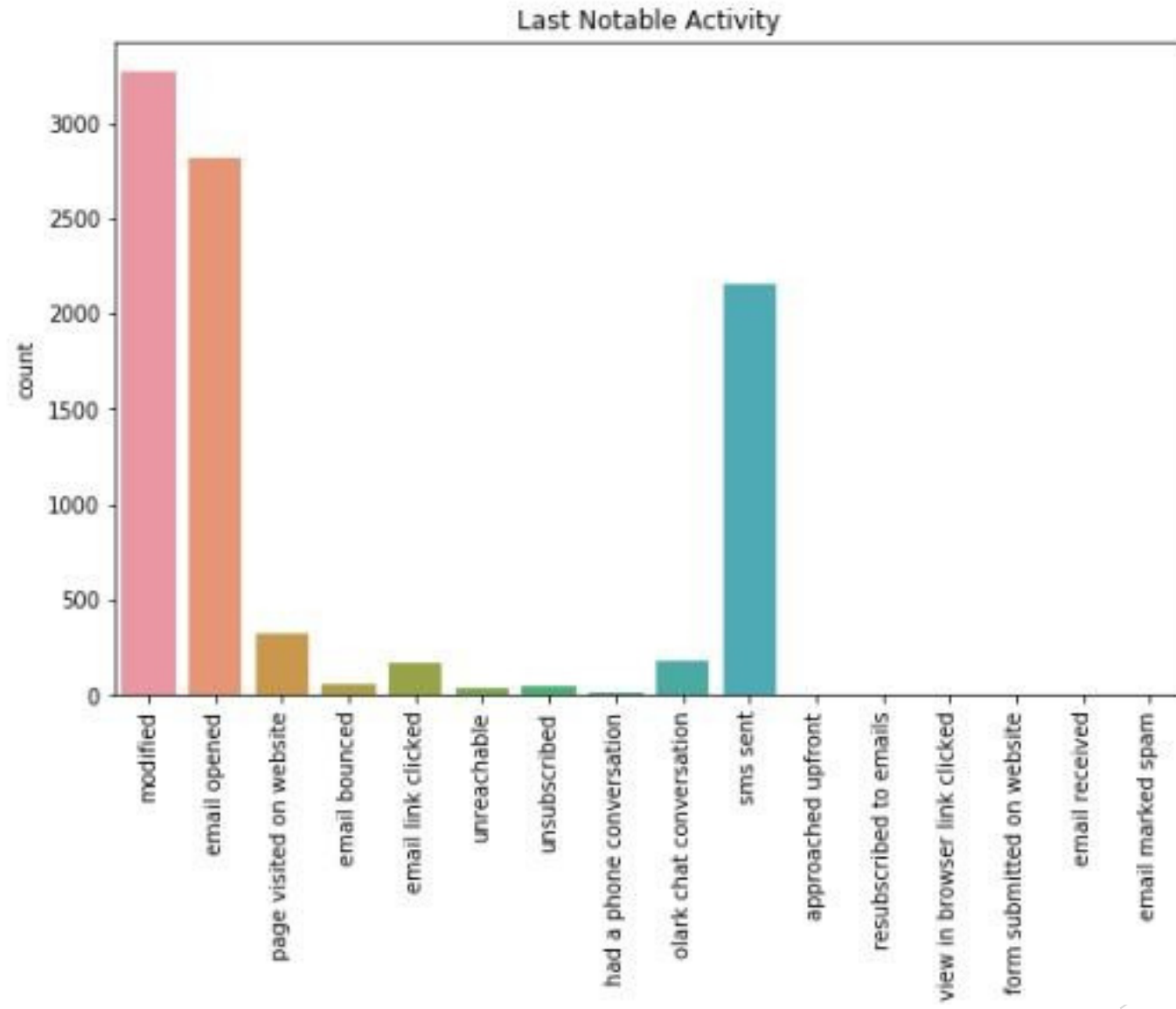
# Solution Methodology

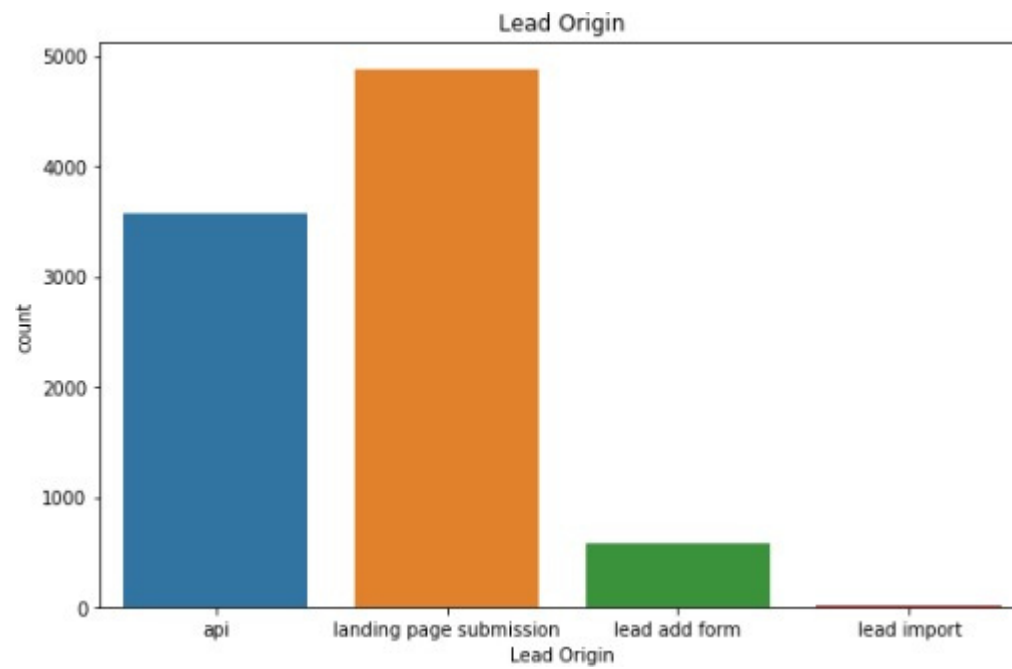
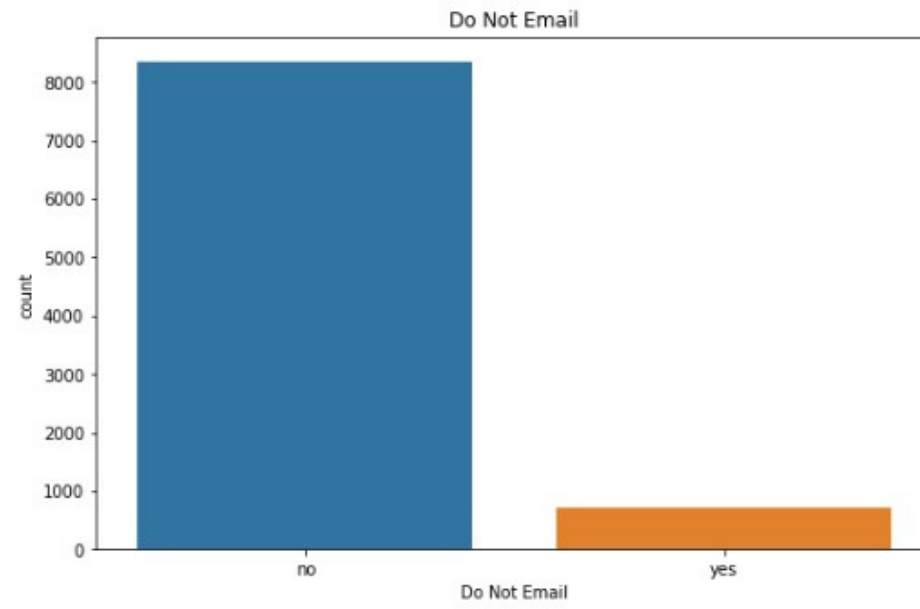
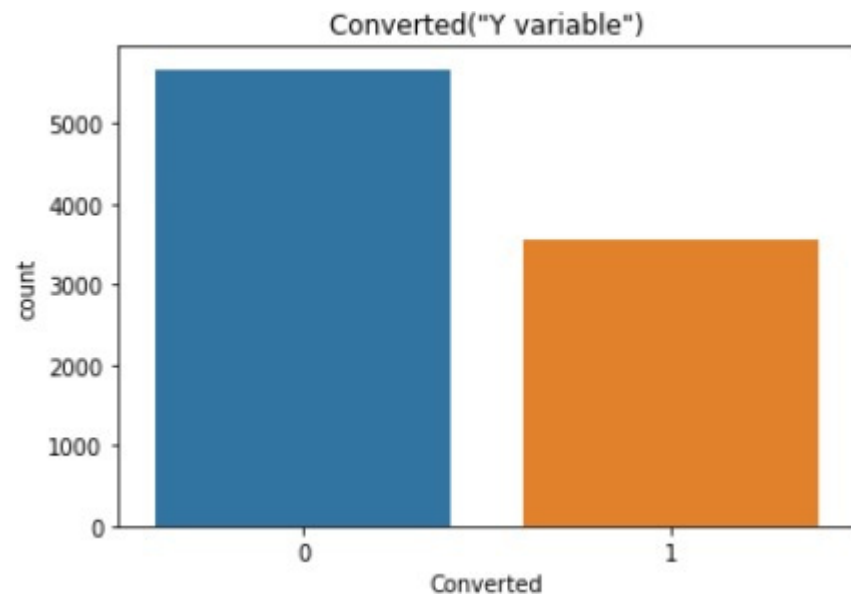
- ▶ Data cleaning and data manipulation.
  1. Check and handle duplicated data.
  2. Check and handle NA values and missing values.
  3. Drop columns, if it contains large amount of missing values and not useful for the analysis.
  4. Imputation of the values, if necessary.
  5. Check and handle outliers in data.
- ▶ EDA
  1. Univariate data analysis: value count, distribution of variable etc.
  2. 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.

# Data Manipulation

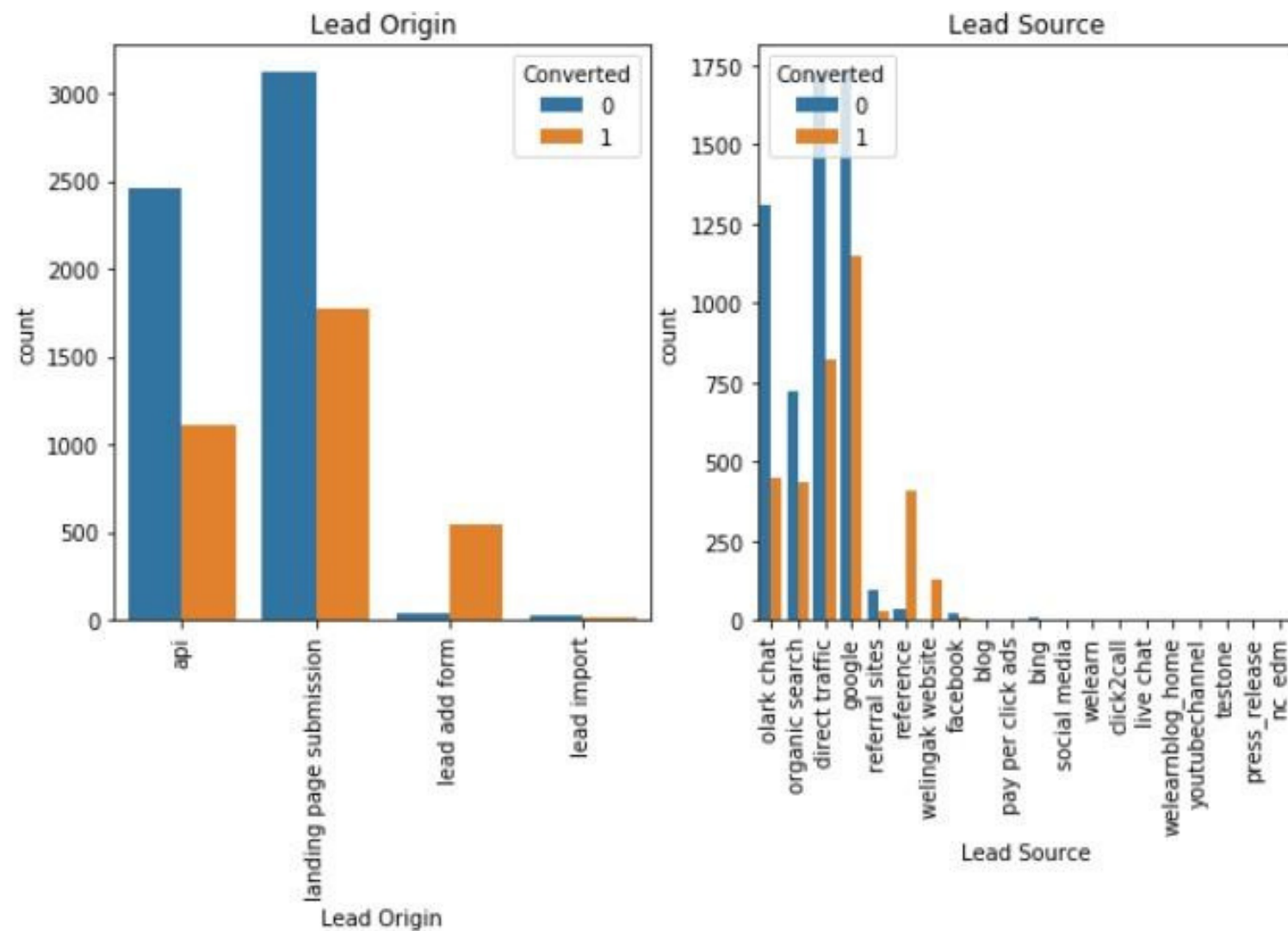
- ▶ TotalNumberOfRows=37, TotalNumberOfColumns=9240.
- ▶ Single value features like “Magazine”, “Receive More Updates About Our Courses”, “Update me on Supply”
- ▶ ChainContent”, “Get updates on DM Content”, “I agree to pay the amount through cheque” etc. have been dropped.
- ▶ Removing the “ProspectID” and “LeadNumber” 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 have no enough variance, which we have dropped, the features are: “DoNotCall”, “What matters most to you in choosing course”, “Search”, “Newspaper Article”, “XEducation Forums”, “Newspaper”, “Digital Advertisement” etc.
- ▶ Dropping the columns having more than 35% as missing values such as ‘How did you hear about XEducation’ and ‘Lead Profile’.

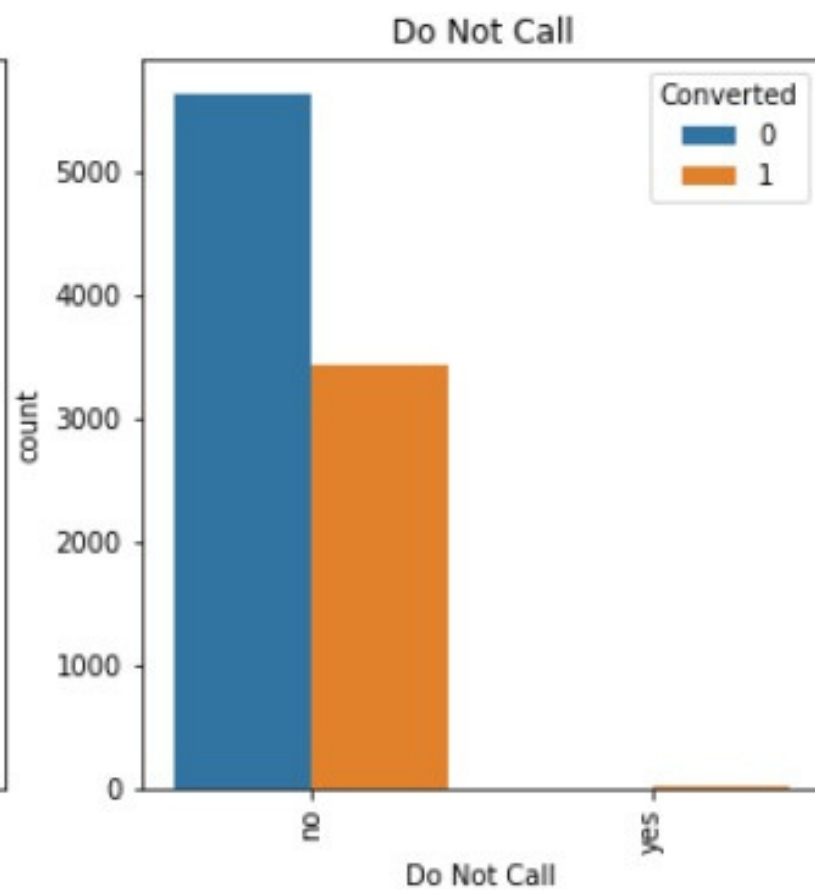
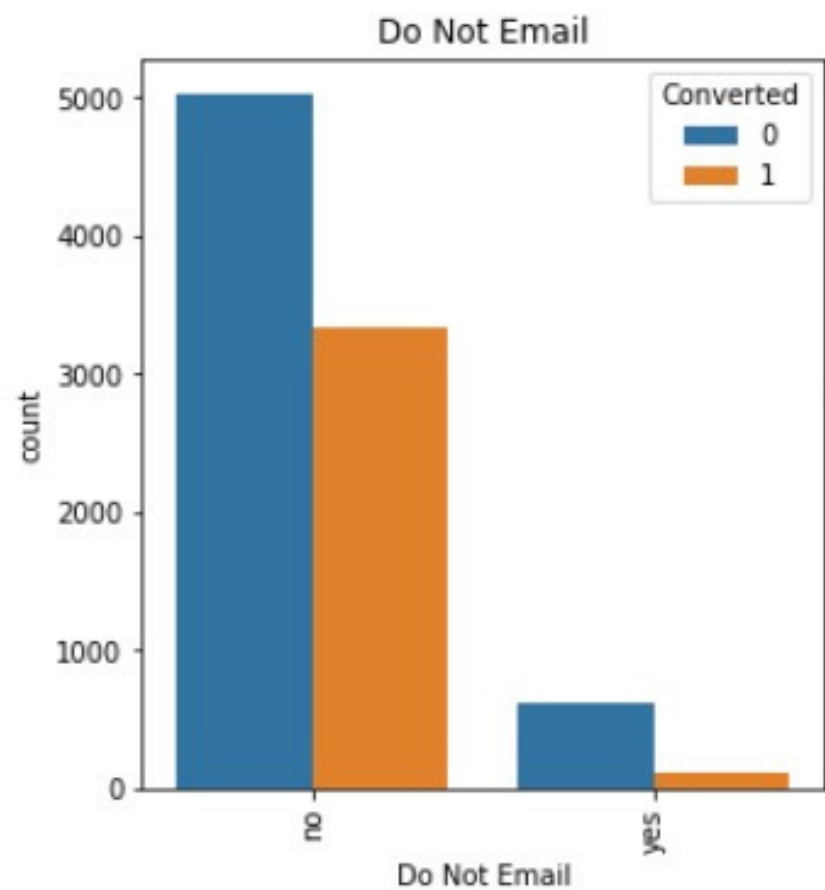
# EDA



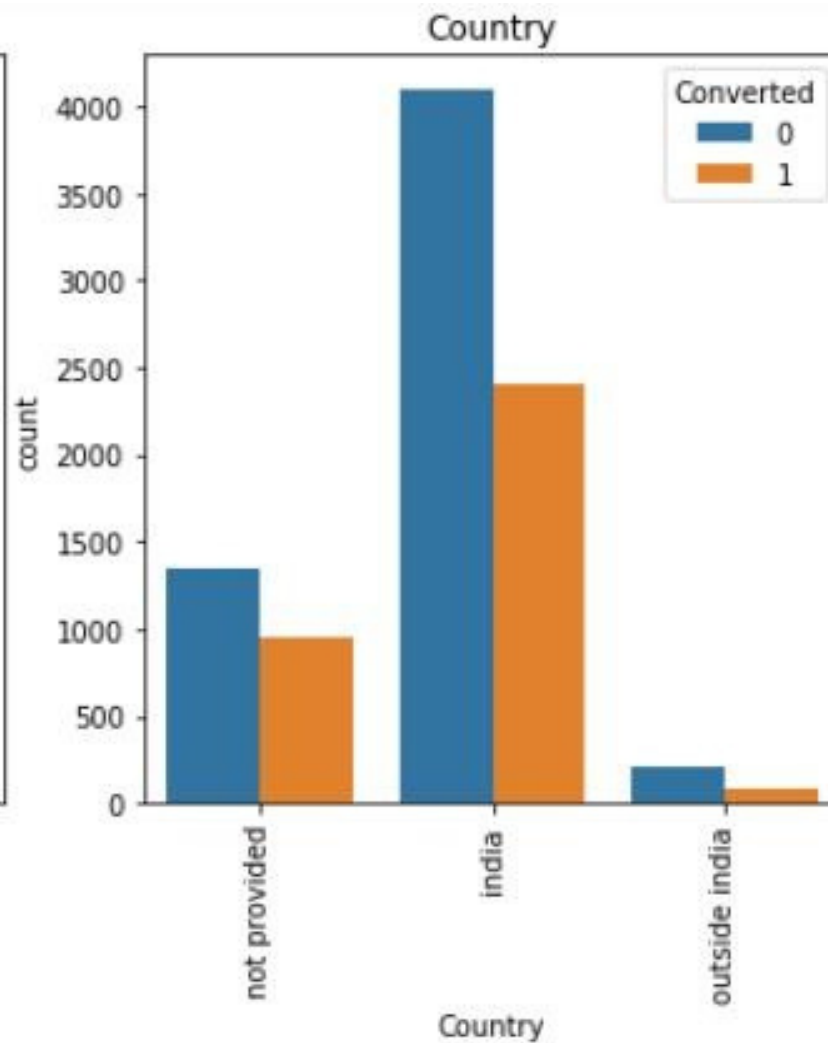
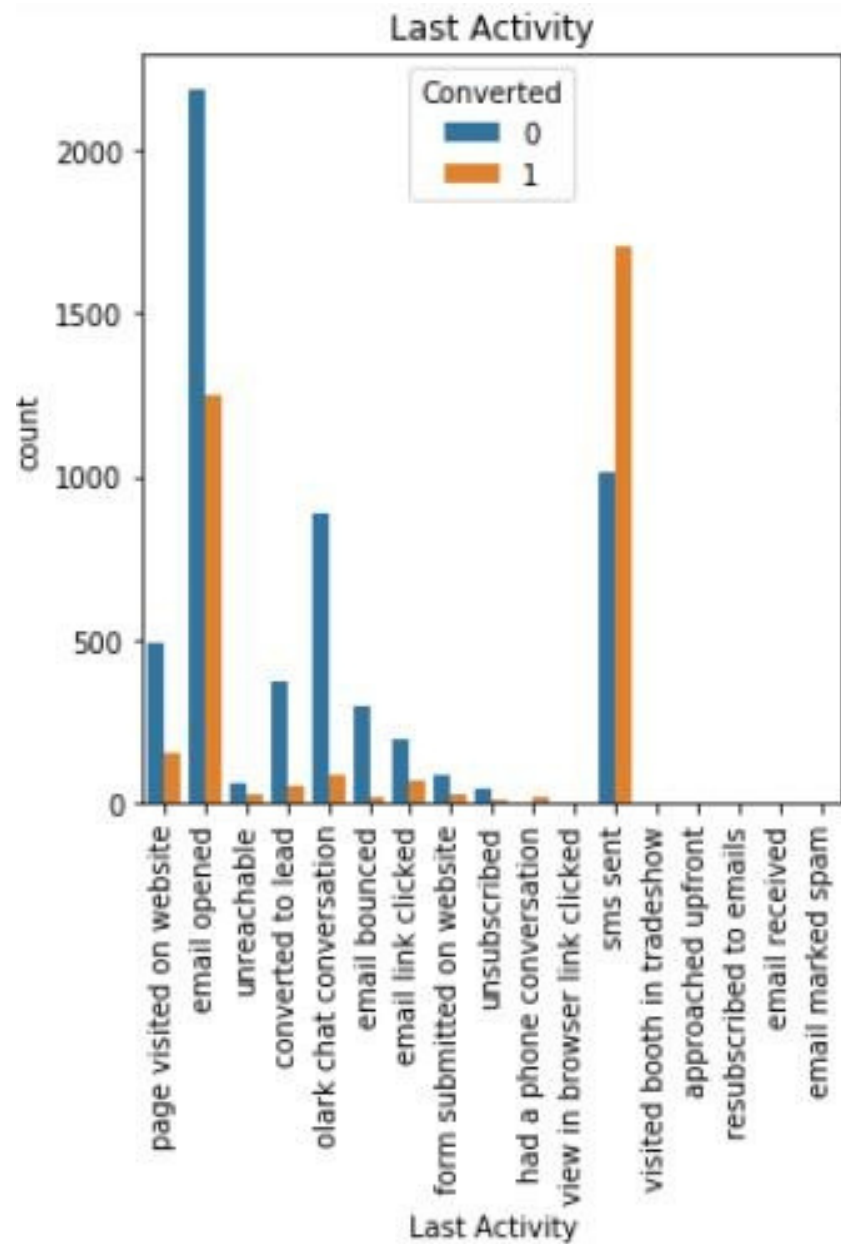


# CategoricalVariableRelation









# DataConversion

- ▶ NumericalVariablesareNormalised
- ▶ DummyVariablesarecreatedforobjecttypevariables
- ▶ TotalRows forAnalysis:8792
- ▶ TotalColumnsforAnalysis:43

# ModelBuilding

- ▶ 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

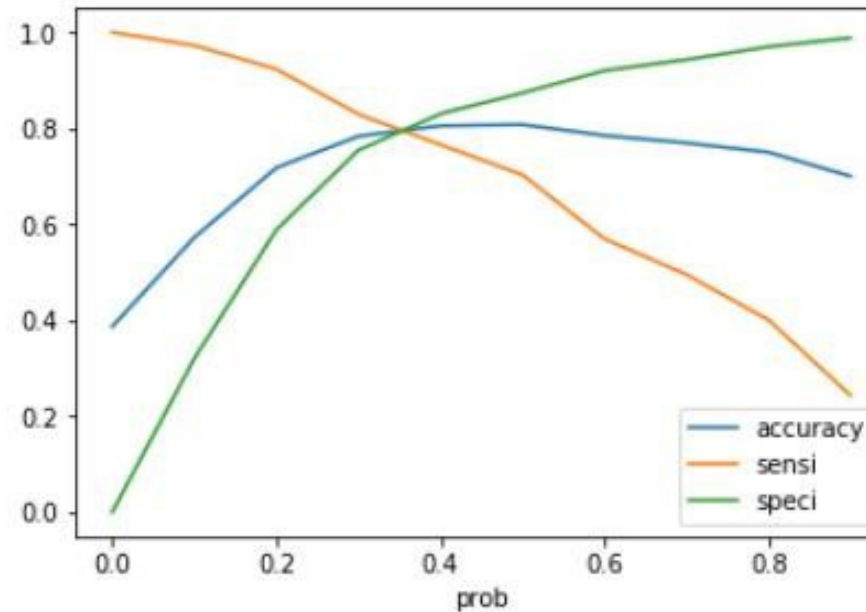
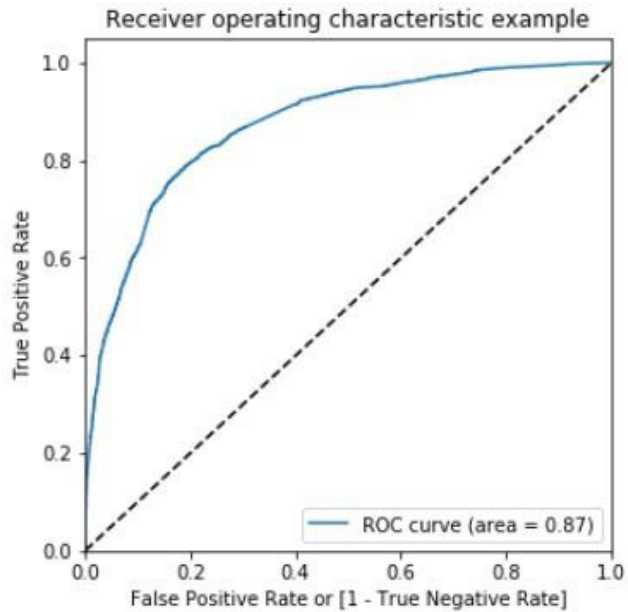
- ▶ Running RFE with 15 variables as output

- ▶ Building Model by removing the variable whose  $p$ -value is greater than 0.05 and if value is greater than 5

- ▶ Predictions on test dataset

- ▶ Overall accuracy 77%

# ROCCurve



- ▶ **Finding Optimal Cut off Point**
- ▶ Optimal cutoff probability is that
- ▶ probability where we get balanced sensitivity and specificity.
- ▶ From the second graph it is visible that the optimal cutoff is at 0.35.

# Conclusion

It was found that the variables that mattered the most in descending order):

- ▶ The total times spend on the Website.
  - ▶ Total number of visits.
  - ▶ When the lead source was:
    - a. Google
    - b. Direct traffic
    - c. Organic search
    - d. We link a website
  - ▶ When the last activity was:
    - a. SMS
    - b. Olark chat conversation
  - ▶ When the lead origin is Lead add format.
  - ▶ When their current occupation is as a working professional.
- Keeping these in mind the X Education can flourish as they have a very high chance to get almost all the potential buyers to change their mind and buy their courses.

the potential buyers are (In