

Lead Scoring Case Study

Group Members

Jinal Gohil

Jayant Kashid

Jayaram K R

Problem Statement

- X – Education sells online courses to industry professionals
- X – Education gets several leads, unfortunately its rate of conversion of leads is pretty low example:- out of 100 leads only 30 of them gets converted
- Company wants to make its process more efficient by identifying the most potential leads (Hot Leads)
- If they are successful in identifying this set of leads company expects to grow the conversion rate to 80%

Business objectives

- X – Education wants to know about the most promising (Hot Leads)
- For this they are looking at building a logistic regression model to identify those hot leads
- Development of the model is for the future use

Solution Methodology

Data cleaning and data manipulation.

1. Check and handle duplicate 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.

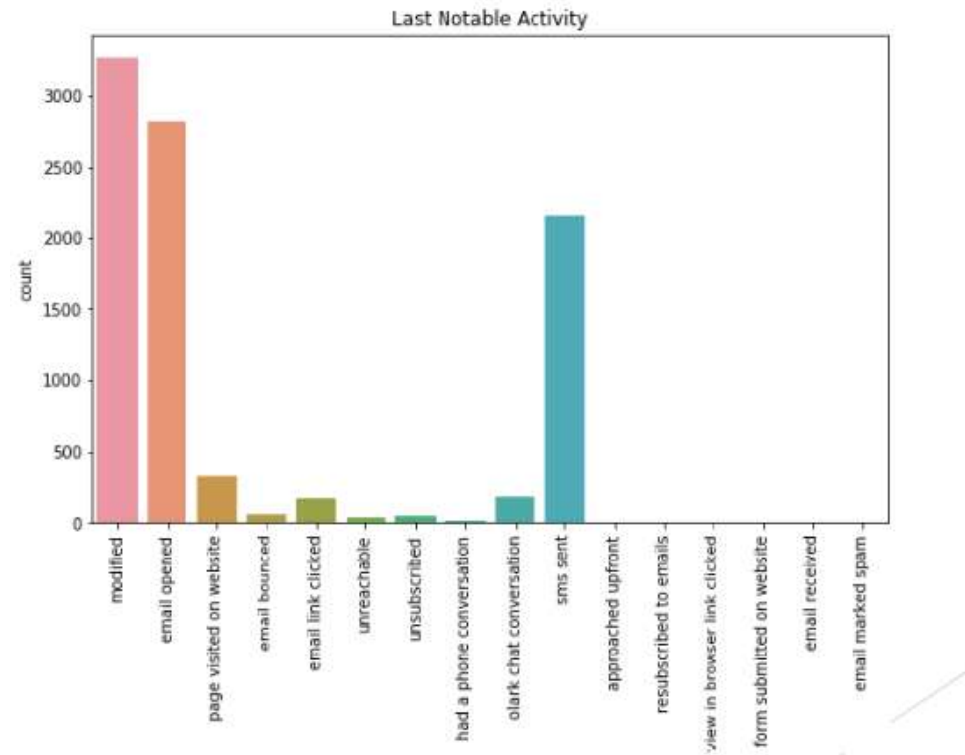
EDA

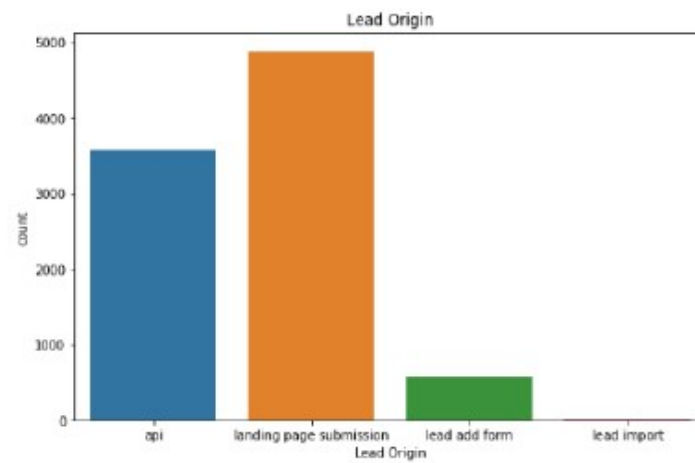
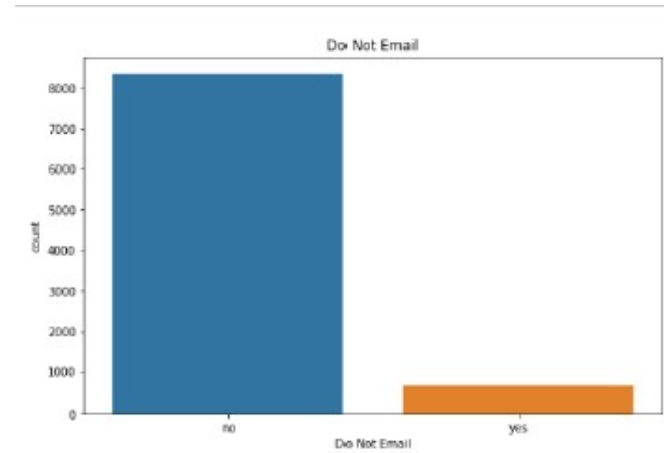
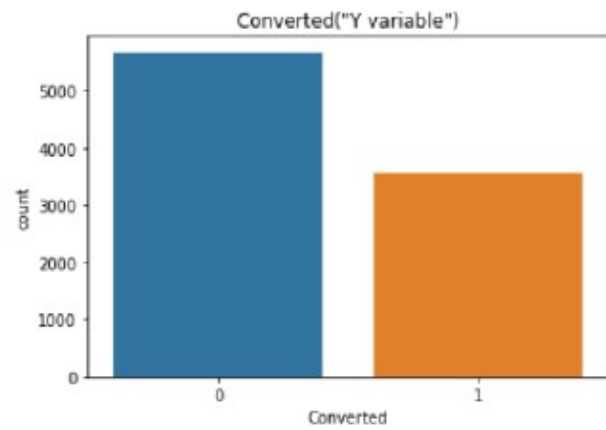
1. Univariate data analysis: value count, distribution of variable etc.
2. Bivariate data analysis: correlation coefficients and pattern between the variables etc.
- 3 Feature Scaling & Dummy Variables and encoding of the data.
- 4 Classification technique: logistic regression used for the model making and prediction.
- 5 Validation of the model.
- 6 Model presentation.
- 7 Conclusions and recommendations.

Data Manipulation

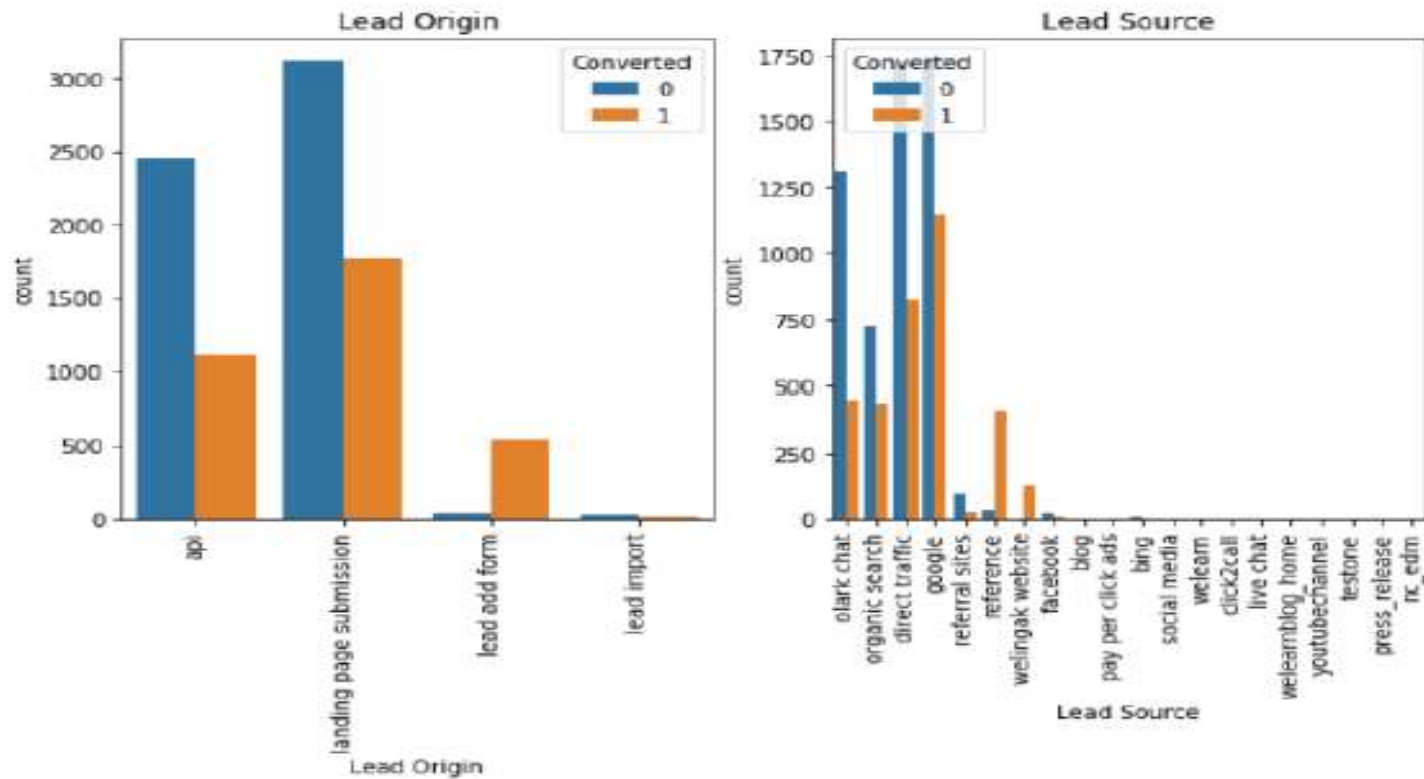
- Total Number of Rows =37, Total Number of Columns =9240.
- Single value features like “Magazine”, “Receive More Updates About Our Courses”, “Update me on Supply”
- Chain Content”, “Get updates on DM Content”, “I agree to pay the amount through cheque” etc. have been dropped.
- Removing the “Prospect 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”, “What matters most to you in choosing course”, “Search”, “Newspaper Article”, “X Education Forums”, “Newspaper”, “Digital Advertisement” etc.
- Dropping the columns having more than 35% as missing value such as ‘How did you hear about X Education’ and ‘Lead Profile’.

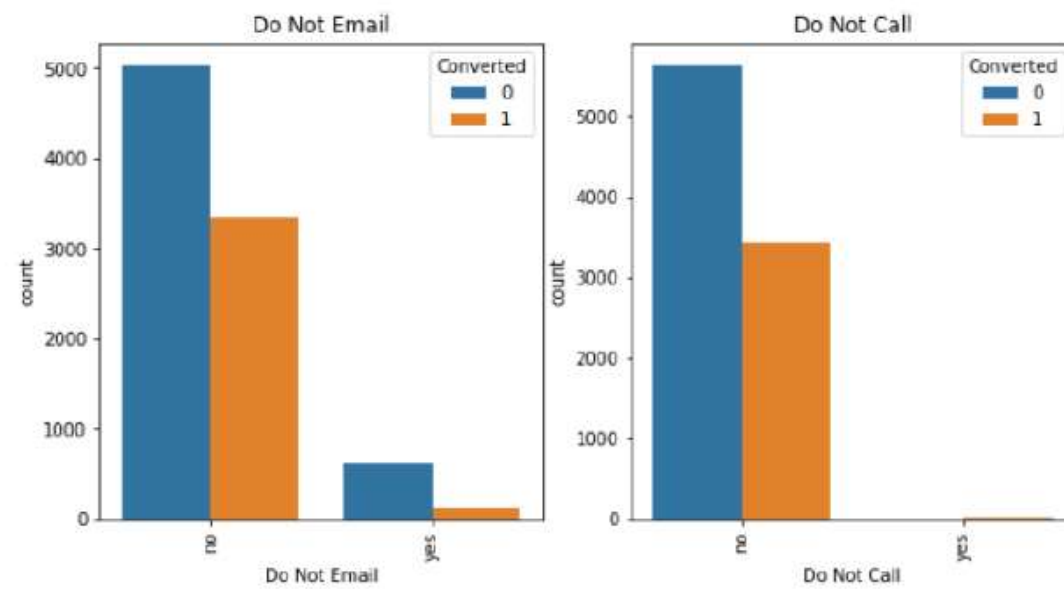
EDA

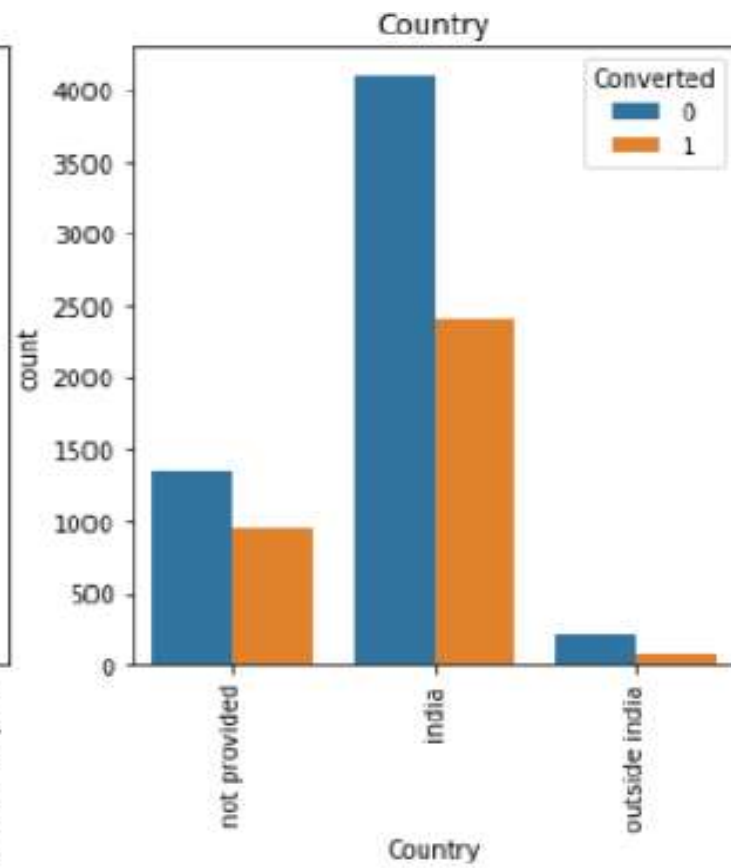
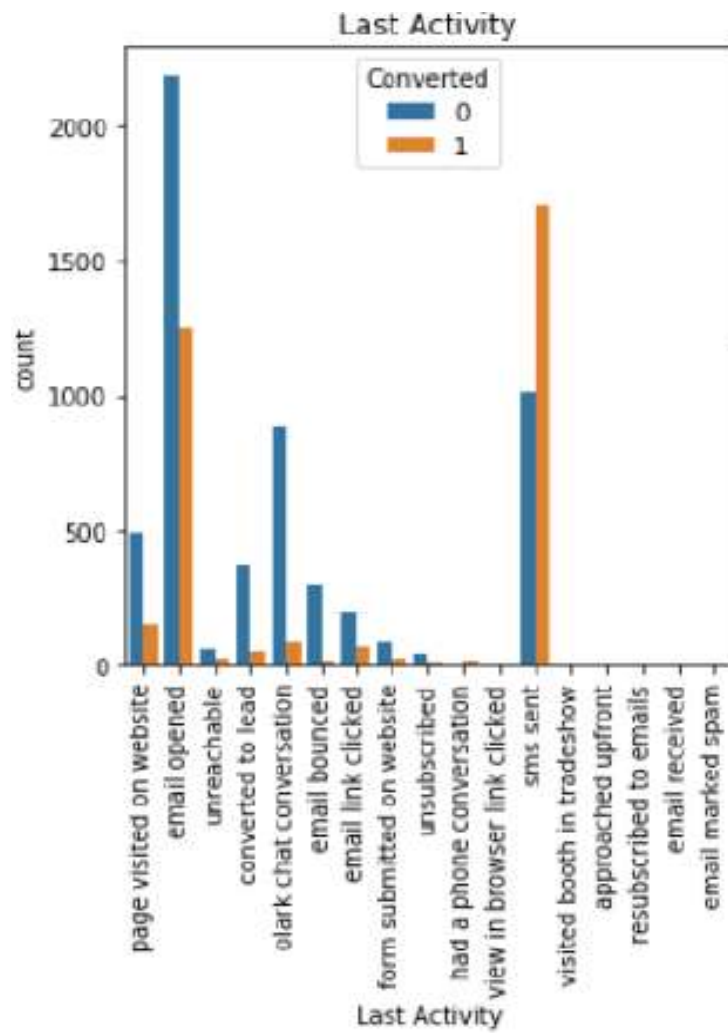




Categorical Variable relation







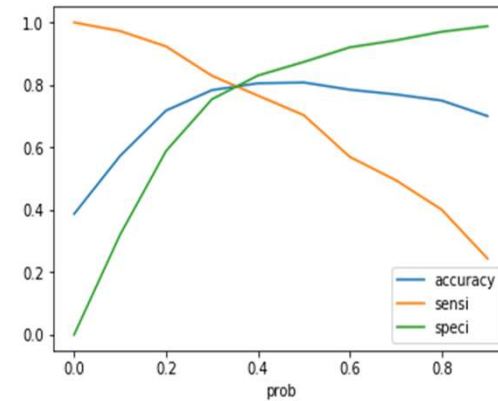
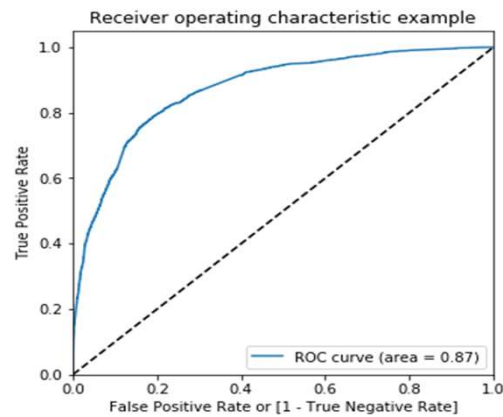
Data Conversion

- Numerical Variables are normalized
- 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 70:30ratio.
- 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



Finding Optimal Cut off Point

Optimal cut off probability is that probability where we get balanced sensitivity and specificity.

From the second 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 are (In descending order) :
- The total time spend on the Website.
- Total number of visits.
- When the lead source was:
 - a. Google
 - b. Direct traffic
 - c. Organic search
 - d. Welingak 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