

Lead Scoring Case Study

By - Sanjay Yadav

Business Understanding

X Education is an education company that provides online courses to professionals in various industries. They attract a significant number of potential customers who visit their website to browse the available courses. When these visitors provide their email address or phone number on the website, they are considered as leads. Additionally, X Education also receives leads through referrals from their past customers. The sales team then proceeds to contact the leads through phone calls and emails, among other methods, in an attempt to convert them into paying customers. However, only a fraction of the leads get converted, with the typical conversion rate standing at 40%.

X Education has conducted an analysis to identify the characteristics of leads that are more likely to convert into paying customers. With this knowledge, they aim to increase their conversion rate to 80%.

Data Summary

We have been provided with a leads dataset from the past with around 9000 data points. This dataset consists of various attributes such as Lead Source, Total Time Spent on Website, Total Visits, Last Activity, etc. which may or may not be useful in ultimately deciding whether a lead will be converted or not. The target variable, in this case, is the column 'Converted' which tells whether a past lead was converted or not wherein 1 means it was converted and 0 means it wasn't converted.

Objective

Build a logistic regression model to assign a lead score between 0 and 100 to each of the leads which can be used by the company to target potential leads.

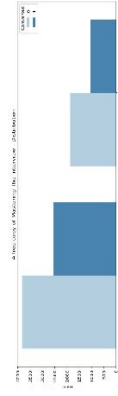
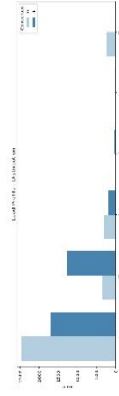
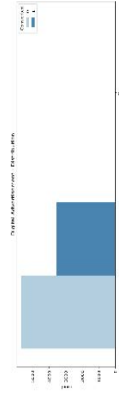
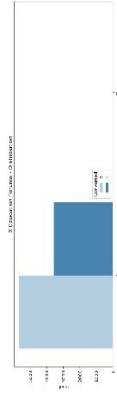
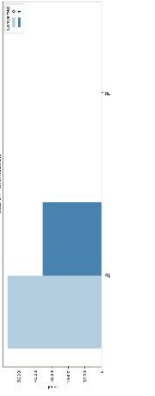
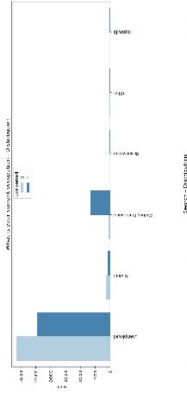
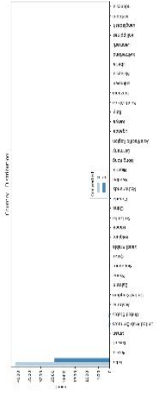
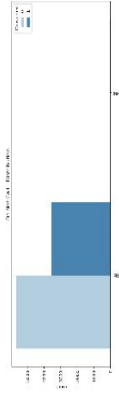
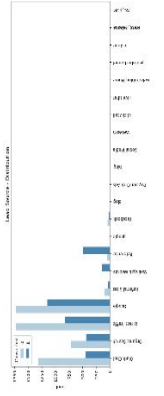
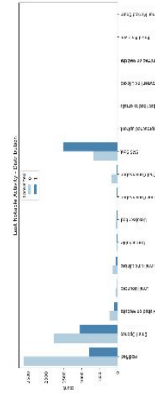
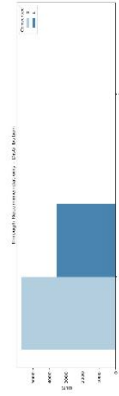
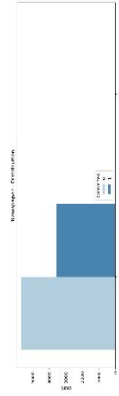
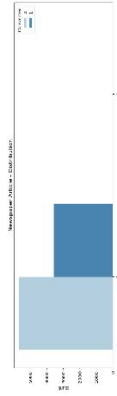
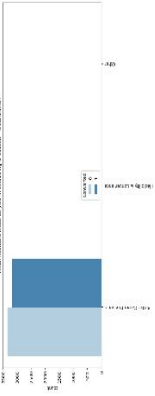
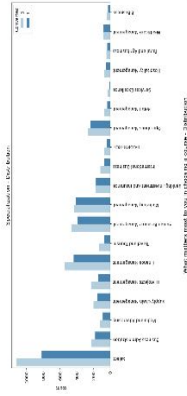
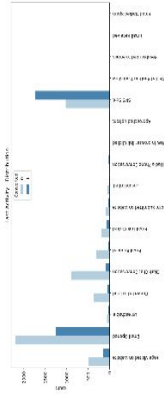
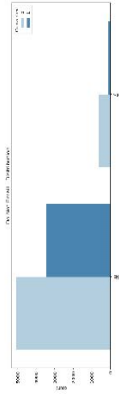
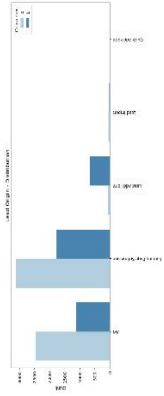
A higher score would mean that the lead is hot, i.e. is most likely to convert whereas a lower score would mean that the lead is cold and will mostly not get converted.

EDA

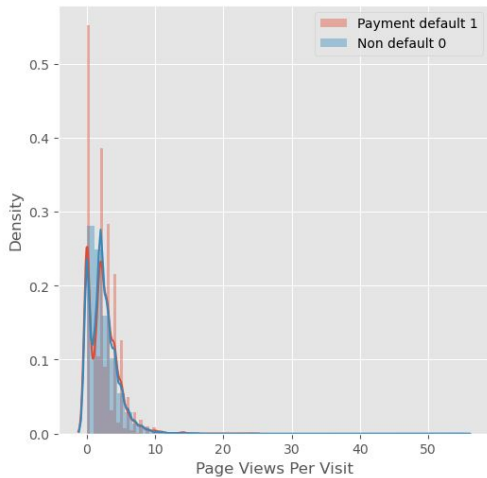
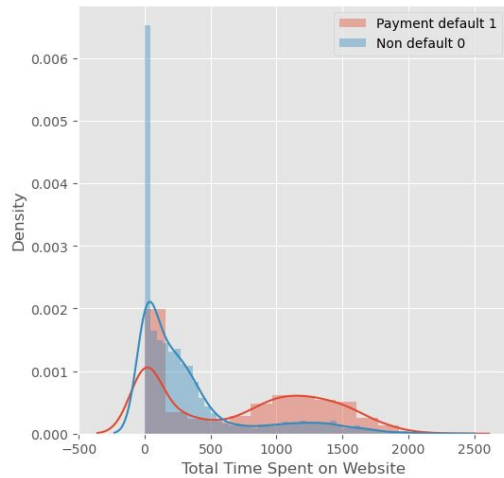
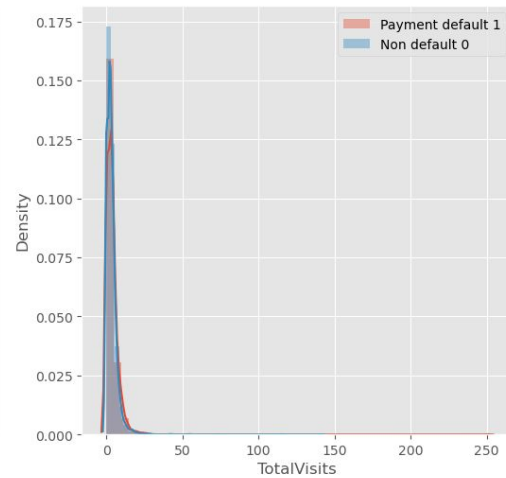
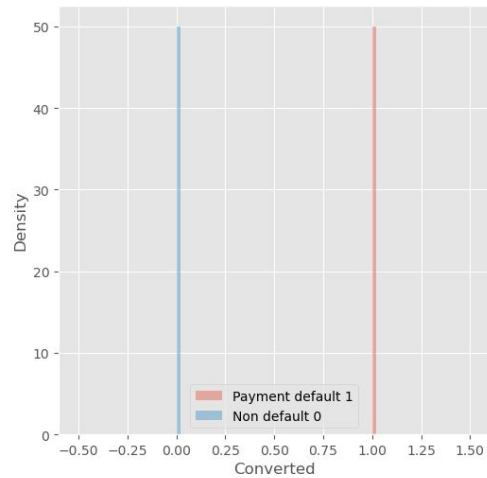
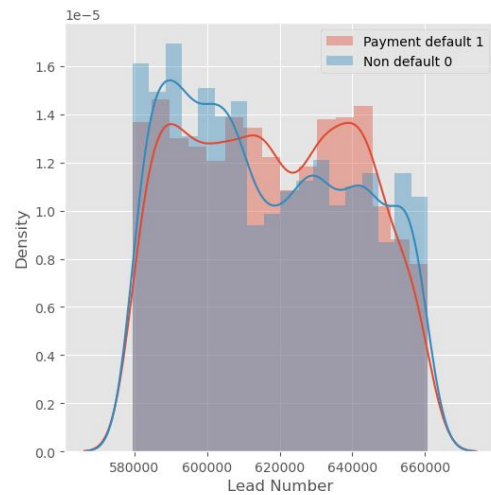
Steps Performed:

Data cleaning - removing columns with more than 30% missing values, replacing values with median for numerical columns and mode for categorical columns, checking for proper data types, making the columns positive with numerical attributes for proper analysis, replacing features with less than 100 value counts by 'others' category in particular feature.

Distribution of categorical features

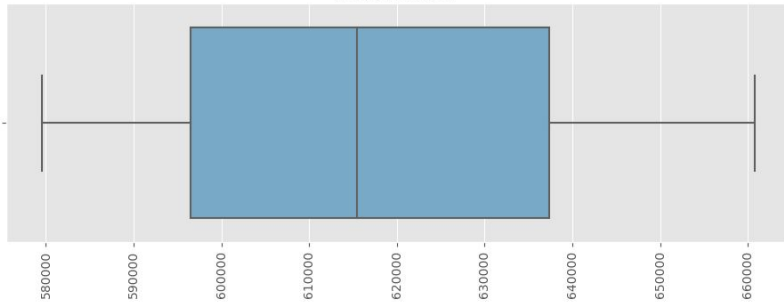


Distribution of Numerical Features

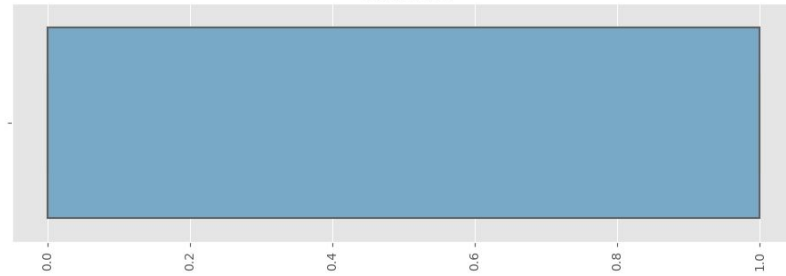


Boxplot of Numerical Features

Lead Number



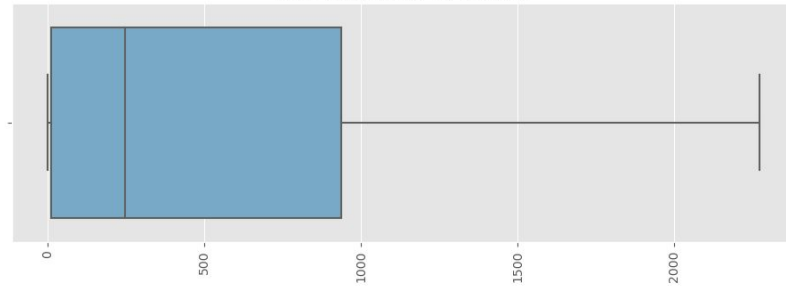
Converted



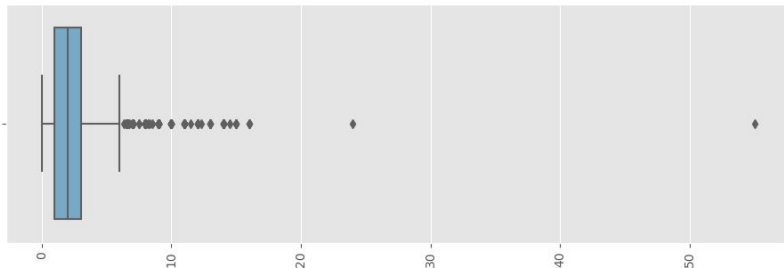
TotalVisits



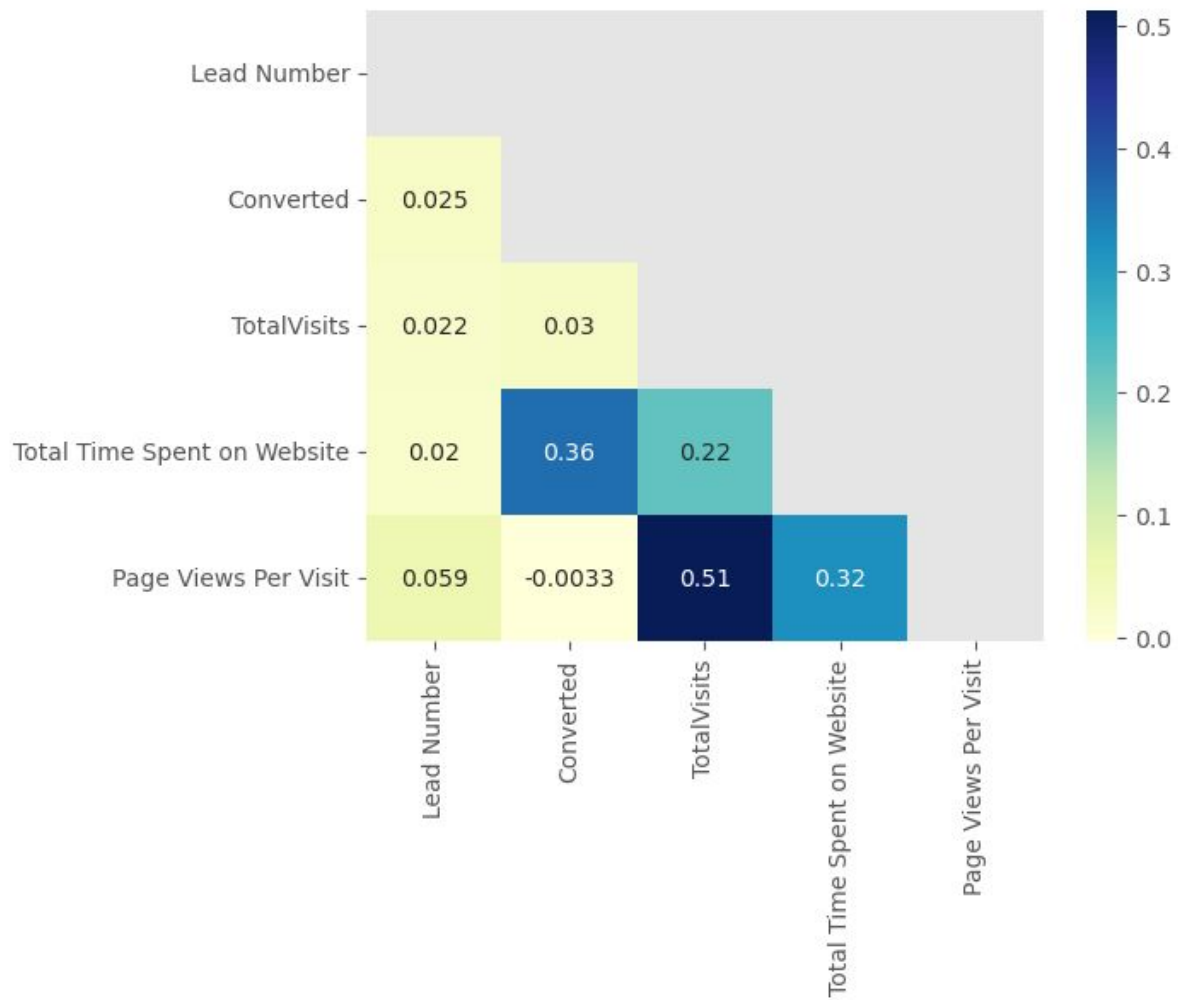
Total Time Spent on Website



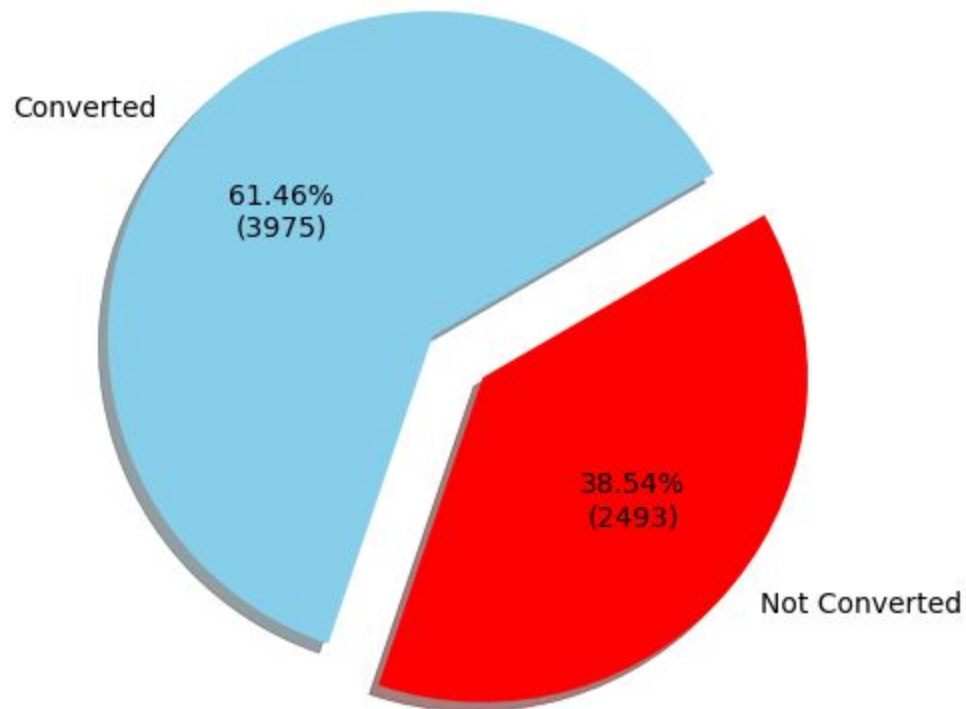
Page Views Per Visit



Correlation analysis for numerical features



Class Imbalance For Final Data



Model Evaluation

Accuracy on train data : 0.831

Accuracy on test data : 0.828

Precision on test data : 0.816

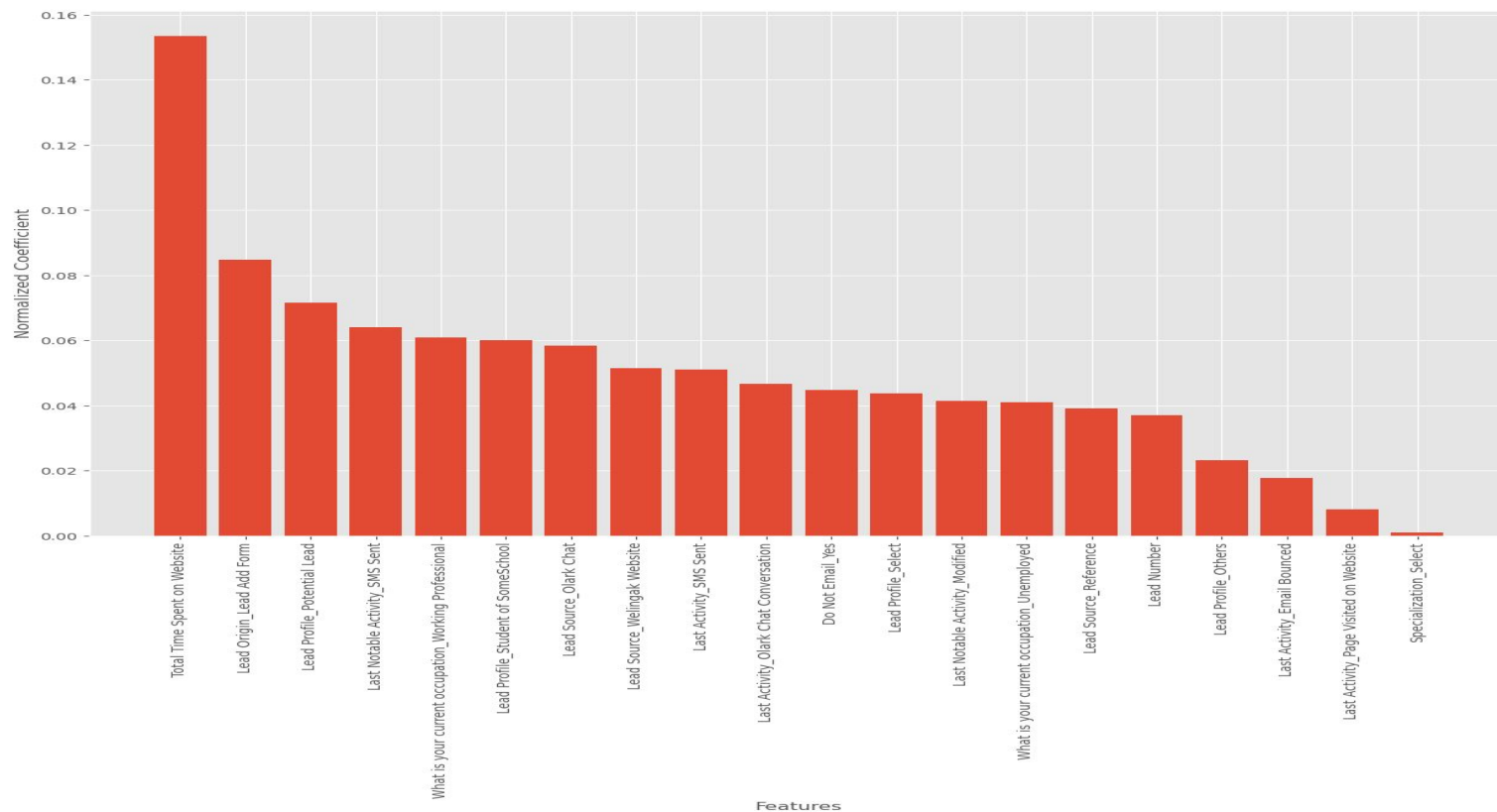
Recall on test data : 0.785

ROC score on test data : 0.818

Confusion Matrix



Feature Importance



Recommendations

- Individuals who are unemployed tend to prioritize courses that offer improved career opportunities.
- Targeting working professionals could be the next step in attracting individuals interested in courses that enhance career prospects.
- Utilizing Google search engine can significantly increase business for the company, as the lead conversion rate is higher compared to other sources, followed by referrals.
- Those who respond to email inquiries are more likely to enroll in courses.