

Summary

This analysis is performed for X Education and to find ways to get more industry professionals to join their courses. The dataset provided gave us a lot of information about how the potentials customers visit the site, the time they spend over there, then how they reached the site and the conversion rate

Below steps are performed:

1. Data cleaning :

1. First step to clean the dataset we choose to remove the redundant variables/features.
2. The data set was partially clean except for a few null values and the option 'Select' has to replace with a null value since it did not give us much information.
3. Dropped the high percentage of Null values more than 40%.
4. Checked for number of unique Categories for all Categorical columns.
5. From that Identified the Highly skewed columns and dropped them.
6. Treated the missing values by imputing the favourable aggregate function like (Mean, Median, and Mode).
7. Detected the Outliers

2. Exploratory Data Analysis:

1. A quick EDA was done to check the condition of our data. It was found that a lot of elements in the categorical variables were irrelevant. The numeric values seems good but found the outliers
2. Performed Univariate Analysis for both Continuous and Categorical variables. Performed Bivariate Analysis with respect to Target variable.

3. Dummy Variables: The dummy variables are created for all the categorical columns.

4. Scaling: Used *Standard scalar* to scale the data for Continuous variables.

6. Model Building:

1. The train test Split was done at 70% and 30% for train and test the data respectively.
2. By using RFE with provided 20 variables. It gives top 20 relevant variables. Later the irrelevant features was removed manually depending on the VIF values and p-value (The variables with VIF < 5 and p-value 0.05 were kept).

4. Model Evaluation:

A confusion matrix was made. Later on the optimum cut-off value by using ROC curve was used to find the accuracy, sensitivity and specificity which came to be around 80%.

8. Prediction:

Prediction was done on the test data frame an optimum cut-off as 0.34 with accuracy, sensitivity and Specificity of 80%.

9. Precision-Recall:

The method was also used to recheck and a cut-off of 0.41. The precision-recall method was used to re check and cut off of 0.41 was found with precision around 79% and recall around 70% on test data frame.

10. Conclusion :

Thus we have achieved our goal of getting a ballpark of the target lead conversion rate to be around 80% . The Model seems to predict the Conversion Rate very well and we should be able to give the CEO confidence in making good calls based on this model to get a higher lead conversion rate of 80%.

Recommendations:

- a. The company should make calls to the leads coming from the lead sources "Welingak Websites" and "Reference" as these are more likely to get converted.
- b. The company should make calls to the leads who are the "working professionals" as they are more likely to get converted.
- c. The company should make calls to the leads who spent "more time on the websites" as these are more likely to get converted.
- d. The company should make calls to the leads coming from the lead sources "Olark Chat" as these are more likely to get converted.
- e. The company should make calls to the leads whose last activity was SMS Sent as they are more likely to get converted.
- f. The company should not make calls to the leads whose last activity was "Olark Chat Conversation" as they are not likely to get converted.
- g. The company should not make calls to the leads whose lead origin is "Landing Page Submission" as they are not likely to get converted.
- h. The company should not make calls to the leads whose Specialization was "Others" as they are not likely to get converted.
- i. The company should not make calls to the leads who chose the option of "Do not Email" as "yes" as they are not likely to get converted.