

LEAD SCORE CASE STUDY SUMMARY

Summary contains the analysis about the X Education company which is an online course seller. Company needs to know the factors which can identify the potential buyers, so that sales team can target them. All the information related to the past buyers and their conversion is provided.

DATA PROVIDED :-

1. Leads.csv
2. Leads Data Dictionary.csv
3. The Target Variable which is 'Converted' which tells whether a past lead converted or not. 0 indicates that it was not converted whereas 1 indicates that it was converted .

The following are the steps taken to build an efficient model which are highlighted below :

Data Cleaning: There were many SELECT values which are considered as missing ones . If the missing values counts is more than 40% in a column, then drop the column. Dropping columns with unique value like Prospect ID and lead ID. Dropping columns with only one value (ex-Magazine). Imputing missing values with mode for categorical data and median for continuous data.

EDA (Exploratory Data Analysis): Created the list for categorical and continuous data type separately for easy graph plotting and analysis. There were many

columns having 99% of single value, dropped them (ex- Do not call). Outliers were not found.

METHODS TO PERFORM EDA :-

1. UNIVARIATE ANALYSIS

2. BIVARIATE ANALYSIS

3. MULTIVARIATE ANALYSIS

Dummy variables: For logistic regression model all categorical data are split into dummy variables to change the values to numeric form (1,0).

Train –Test data: Data is split into train and test dataset by 70:30 ratio. Train and test data are normalised using standard scaler.

Feature selection and Model building: Using hybrid approach to select feature. RFE is used to select 15 features and then manually features are dropped by analysing p-value (less than 0.05) and VIF (less than 5). Finally, 11 features are obtained in Model5.

Final Features are as follow :

1. Lead Origin_Lead Add Form
2. Do Not Email_Yes
3. Last Activity_Olark Chat Conversation
4. What is your current occupation_Working Professional
5. Tags_Busy

6. Tags_Closed by Horizzon
7. Tags_Lost to EINS
8. Tags_Ringing
9. Tags_Will revert after reading the email
- 10.Tags_in touch with EINS
- 11.Last Notable Activity_SMS Sent

Model evaluation: Predicting from train data and evaluating Accuracy, Recall and Precision. Plotting ROC curve to check the model. Checked the metrics at cutoff value 0.35 and 0.38. We got Optimum cutoff at 0.35. As stated, the TPR is maintained above 80%

Prediction on test set :- Normalising data in test set. Predicting from test data and evaluating Accuracy, Recall and Precision. On comparing the metrics from test and train data set were almost same depicting the model efficiency.

Key learnings to identify the hot leads are:

- When the customer is a working professional, it has high chance of conversion.
- When the Lead origin is Lead add form.
- Last activity is identified as Olark Chat Conversation and SMS sent.
- When customer has permitted for email.
- When customer is tagged as 'lost to EINS', 'closed by horizon','Will revert after reading mail', 'In touch with EINS'.

- Through EDA we can see the 'time spent in websites', 'total visits', 'lead source as Google' seems to give fruitful result.

X Education company can refer to above factors and go for potential buyers.

By – Saurabh Taneja

Sargam Agarwal

Shail Shikhar.....