# LEAD SCORING CASE STUDY

# Problem statement

An education company named X Education sells online courses to industry professionals. On any given day, many professionals who are interested in the courses land on their website and browse for courses.

The company markets its courses on several websites and search engines like Google. Once these people land on the website, they might browse the courses or fill up a form for the course or watch some videos. When these people fill up a form providing their email address or phone number, they are classified to be a lead. Moreover, the company also gets leads through past referrals. Once these leads are acquired, employees from the sales team start making calls, writing emails, etc. Through this process, some of the leads get converted while most do not. The typical lead conversion rate at X education is around 30%.

Now, although X Education 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 rate should go up as the sales team will now be focusing more on communicating with the potential leads rather than making calls to everyone.

# Goal

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.

There are some more problems presented by the company which your model should be able to adjust to if the company's requirement changes in the future so you will need to handle these as well. These problems are provided in a separate doc file. Please fill it based on the logistic regression model you got in the first step. Also, make sure you include this in your final PPT where you'll make recommendations.
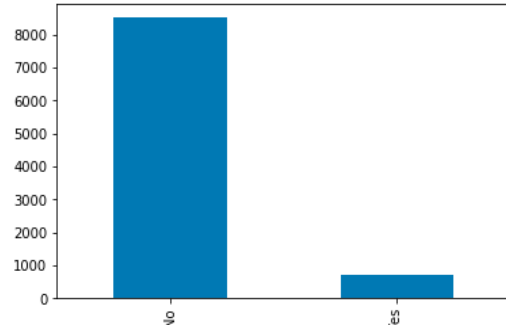
# Steps Performed

- Data Import
  - Reading and Visualizing the Data
- Data Preparation and Cleaning
  - Treating Missing Values (Dropping columns with high missing percentage)
  - Removing Highly Skewed Columns
  - Treating Outliers
- Data Modelling
  - Train Test Split
  - Scaling of data
  - Iterative Modelling using coarse tuning (RFE + Manual Tuning)
  - Modelling using Logistic Regression's GLM
  - Checking key parameters like Accuracy, Sensitivity, Specificity, ROC Curve, Optimal Cutoff, Precision and Recall
- Prediction using the derived model
  - Running the model on Test Data Set
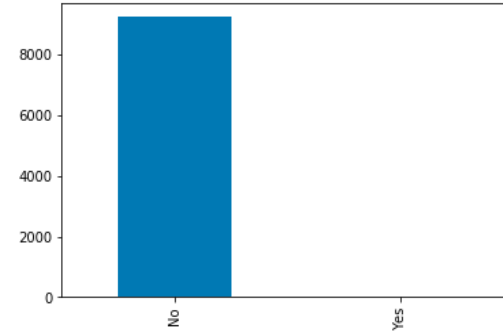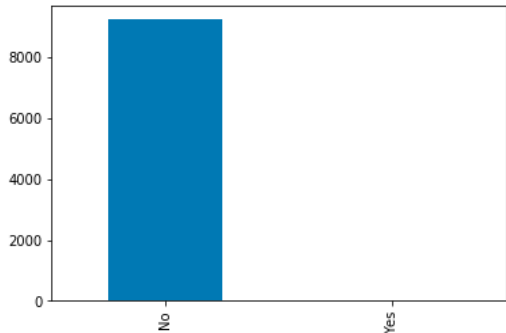  - Checking the key parameters whether model performs good on unseen data

Name: Do Not Email, dtype: int64
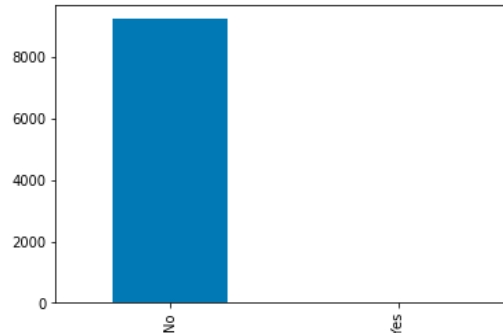
Name: Newspaper Article, dtype: int64

Name: X Education Forums, dtype: int64

Name: Do Not Call, dtype: int64

# Skewness Check

▶ As seen in the adjacent images, some of the columns shown here are highly skewed.

▶ Hence it is best to drop them
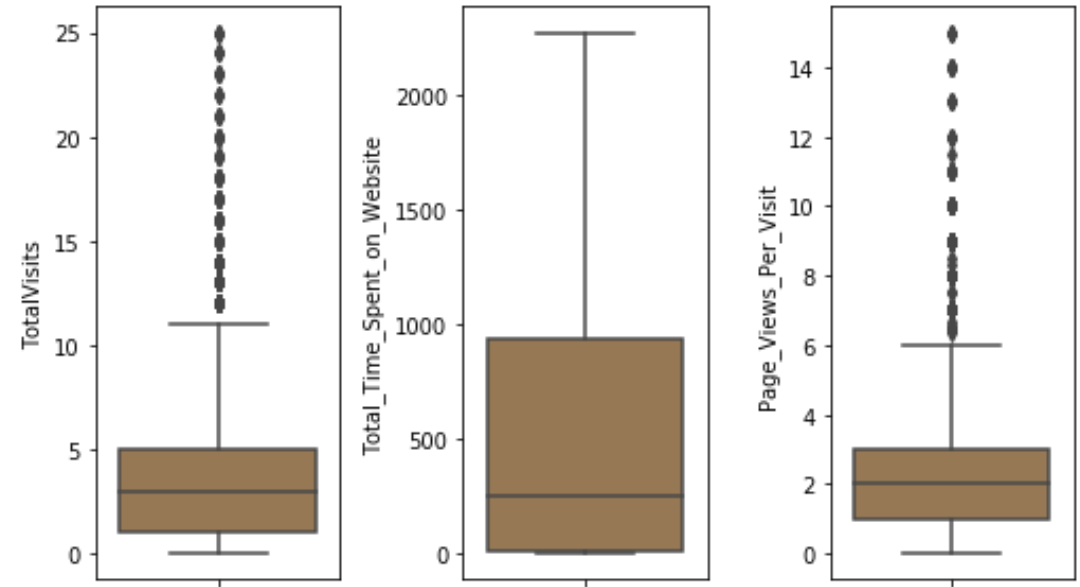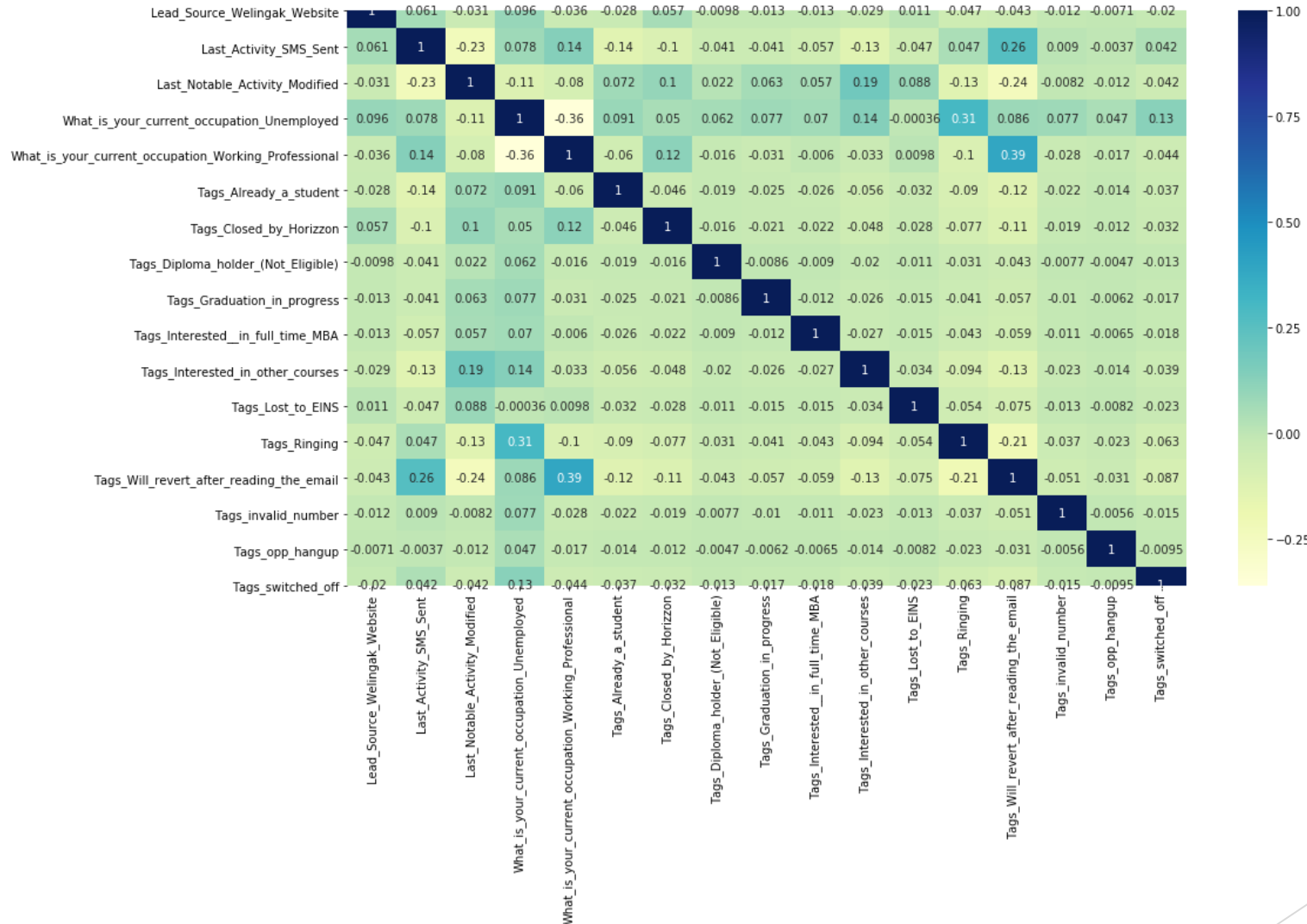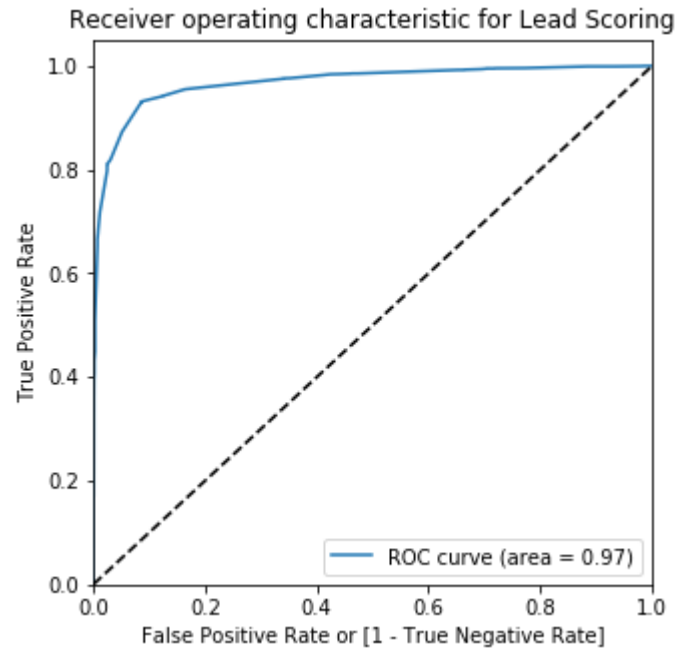
# Outlier Treatment

▶ As seen from adjacent box plots, Total Visits and Page Views Per Visit have outliers.

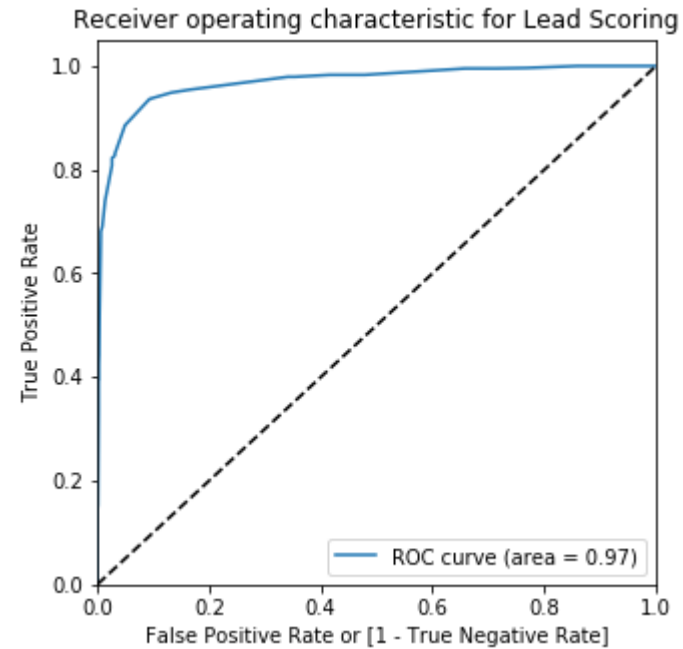▶ It is imperative to remove them.

# Correlations on Chosen Features
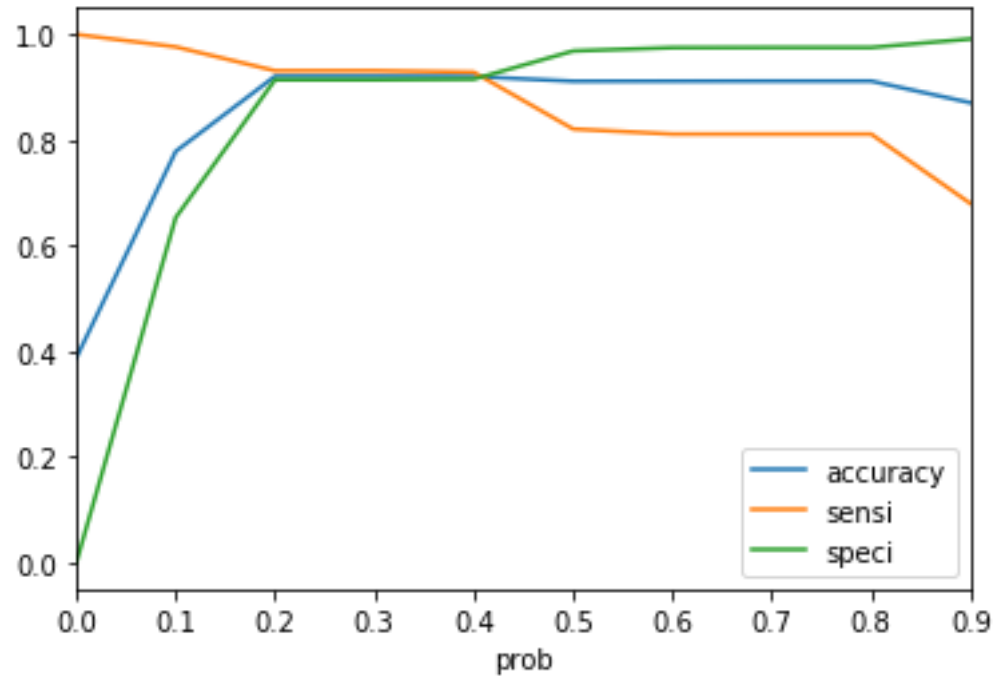
# ROC Curve



Train Data Set

Test Data Set

ROC Curve is pretty good for both train and test data

# Optimal Cutoff



*As the sensitivity, specificity and accuracy intersects from x: 0.2 to x:0.4, we are choosing 0.3 as optimal cutoff point*

# Key Parameters Compared

| Parameters | Train Data | Test Data |
|---|---|---|
| Accuracy | 92.07% | 91.78% |
| Sensitivity | 93.12% | 93.64% |
| Specificity | 91.42% | 90.64% |
| Precision | 94.30% | 86.03% |

F1 Score = 0.8967

# Feature Importance Decoded

# Recommendations

- As seen from the bar plot, Tags_Closed_by_Horizzon feature has highest conversion rate

- It is followed by Tags_Lost_to_EINS

- Lead_Source_Welingak_Website contributes the next most to the successful conversion rate

- On the other hand, whoever provided an invalid number, they tend to never convert into a student.

- Most of the respondents whose contact number is invalid or the phone is switched off or never take a ring, may not be converted.

- So the X Education marketing representatives should not call these respondents more than once, which results in better focus on others

- Instead they should focus more on people, who report 'Closed by Horizzon' or 'Lost to EINS' for better conversion rate.

- They can also promote the Welingak Website on various platforms including Social media, which will significantly increase their conversion rate