# LEAD SCORING CASE STUDY

Building a Logistic Regression Model for X Education, an online education company using Python libraries.

# GROUP

By – Anupam Vaibhav, R V Gokul and Nisha Venkatesh

Cohort – DSC44

Batch ID- 3161

# Approach

- Datasets were imported and Python libraries – Pandas, Numpy, Matplotlib.plotly, Seaborn, Statsmodels and sklearn were used for analysis.

- Data was observed, cleaned and analysed using imputations, deletion, outlier treatments, data imbalance, Univariate and Multivariate analysis.

- Model was built using MultiClass Logistic Regression

# Data Understanding

1. There were 9240 rows and 37 columns initially.

2.  30 of the columns contained string/object type data, rest were numerical.

3. 'Converted' column as Target column

4. Initial conversion rate at 38%.

# Data Cleaning

- 5 columns, eg. Magazine with no variance were deleted.
- Columns with more than 30% values either missing or having 'Select' as value were deleted. Eg. How did you hear about X Education
- Columns having 2% to 30% missing values were imputed with 'Other' or 'Unknown' for categorical columns.
- Rows of variables such as TotalVisits with less than 2% missing values were deleted.

# Data Cleaning

- ▫ Categorical variables such as 'Country' having large number of unique value but having less than 8% contribution were merged as Other.
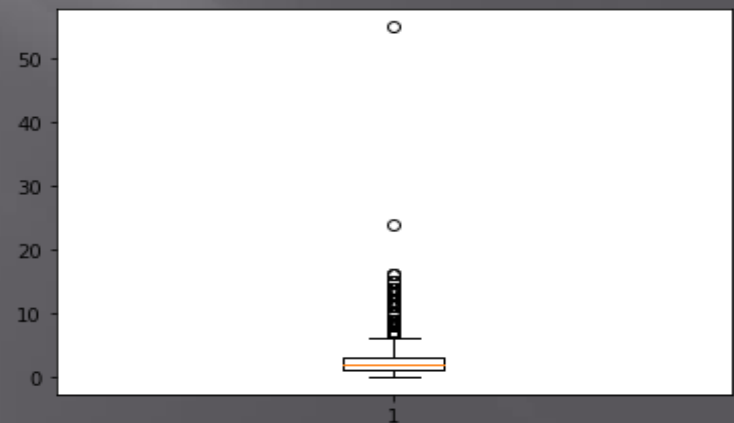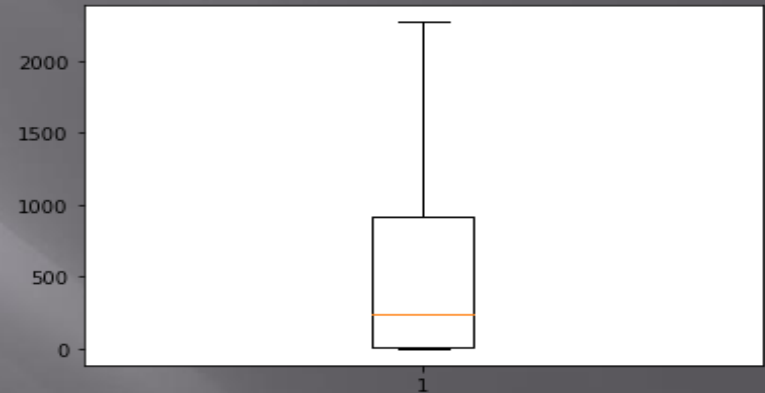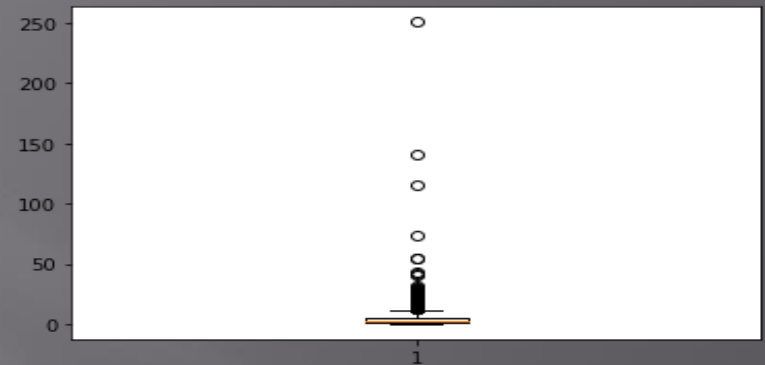
## Outliers

For numerical columns , Boxplot were made and outliers were analysed.

For, TotalVisits column, 25 was set as cutoff.

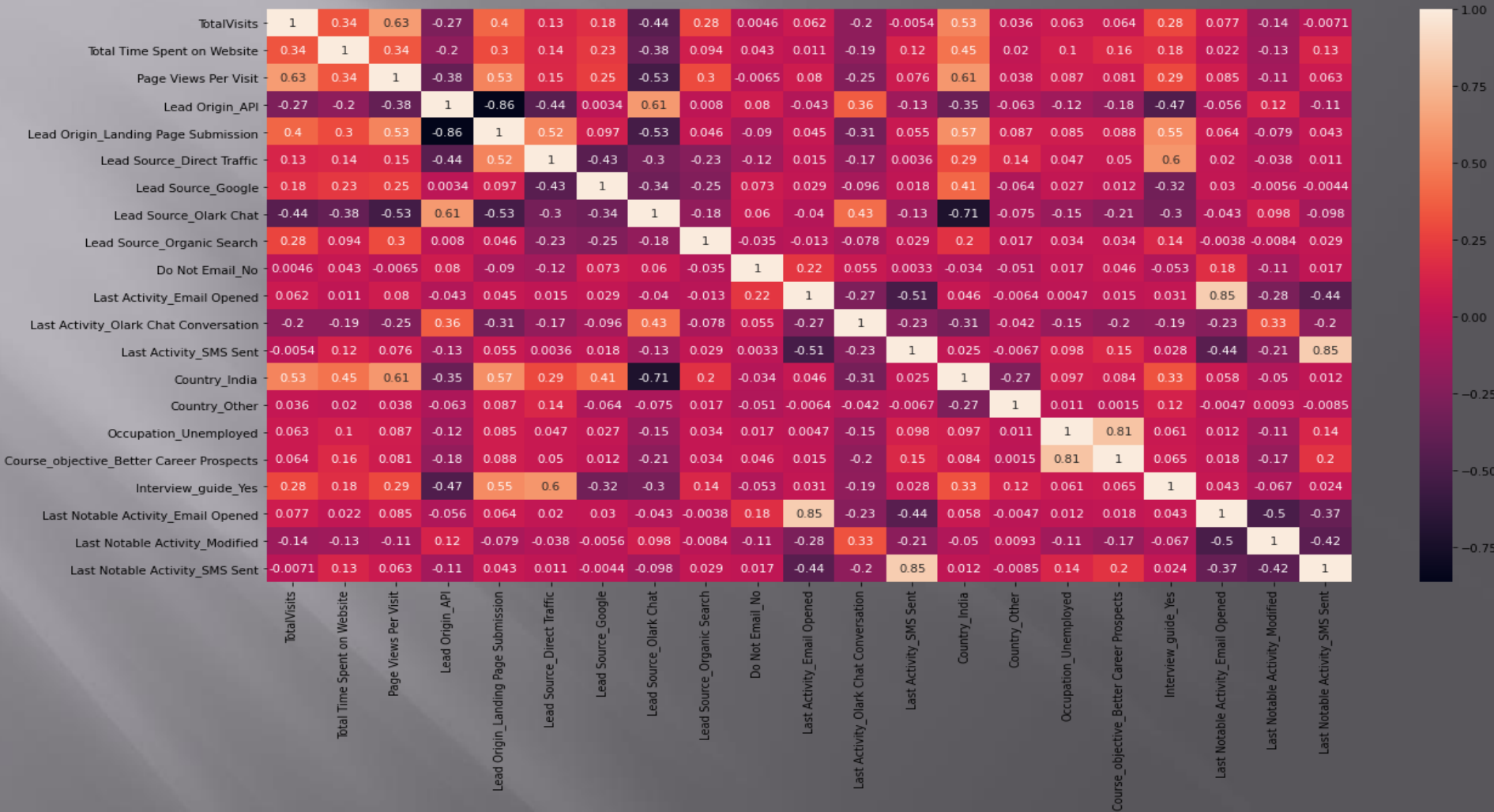For, Total Time Spent on Website , no outlier treatment was required.

For, Page Views Per Visit col, 10 was set as cutoff.

# Splitting, Dummy Creation and Scaling

- Test-Train split was done as 70-30%.

- Train set numerical variables were fit and scaled using Standardisation.

- Test set numerical variables were only scaled using Standardisation.

- Dummy variables were created for rest of the categorical columns and one was dropped from each.

# Correlations



Columns with correlation > |0.7| were deleted

# Statistical limits

- We were left with 17 columns.
- RFE was used to keep 12 columns and rest were dropped.
- P-value < 0.02 and VIF < 5 was used as cutoff to keep remaining variables thus, 3 more were deleted after a few iterations.
- Total 9 variables were left for model to be trained upon

# Model Building

- Logistic Regression was used to teach the model.
- Confusion matrix was made on the predicted and actual conversions at cutoff probability of 0.5 initially.
- Recall, precision and accuracy were calculated for multiple cutoff probability.
- Final cutoff was set at 0.4 while optimising the model for Recall.
- Score column was made as 100 times probability.

# Model Performance

- On Train Set: Recall = 75.4%, Accuracy = 77%, Precision = 68% and F1 Score = 71.5%.

- On Test Set: Recall = 75%, Accuracy = 76%, Precision = 68% and F1 Score = 71.8%.

# CONCLUSION

**Most important Variables for Model:**
1. Do Not Email_No
2. Occupation_Unemployed
3. Lead Source_Google
4. Lead Source_Direct Traffic
5. Last Notable Activity_Modified
6. Last Activity_SMS Sent
7. Lead Source_Organic Search
8. Last Activity_Olark Chat Conversation
9. Total Time Spent on Website

*Final F1 score was 0.71 with recall of 0.75, precision of 0.68 and accuracy of 0.76.*