

“

Lead Score Case Study

”

Submitted By:

Adnan Mirza

Viren Pawar

Shraddha Nerlekar



Problem Statement



- X Education sells online courses to industry professionals.
- 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 communicating with the potential leads rather than making calls to everyone. We call these 'Hot Leads'.
- If they successfully identify this set of leads, the lead conversion rates should go up as the sales team will now be focusing



Business Objective

- X education wants to know most promising leads.
- For that they want to build a Model which identifies the hot leads.
- Deployment of the model for the future use.



Solution Methodology



Data cleaning and data manipulation

1. Check and handle duplicate data.
2. Check and handle NA values and missing values.
3. Drop columns, if it contains large amount of missing values and not useful for the analysis.
4. Imputation of the values, if necessary.
5. Check and handle outliers in data.



➤ EDA

1.Univariate data analysis: value count, distribution of variable etc.

2.Bivariate data analysis: correlation coefficients and pattern between the variables etc.

➤ Feature Scaling & Dummy Variables and encoding of the data.

➤ Classification technique: logistic regression used for the modelmaking and prediction.

➤ Validation of the model.

➤ Model presentation.

➤ Conclusions and recommendations.







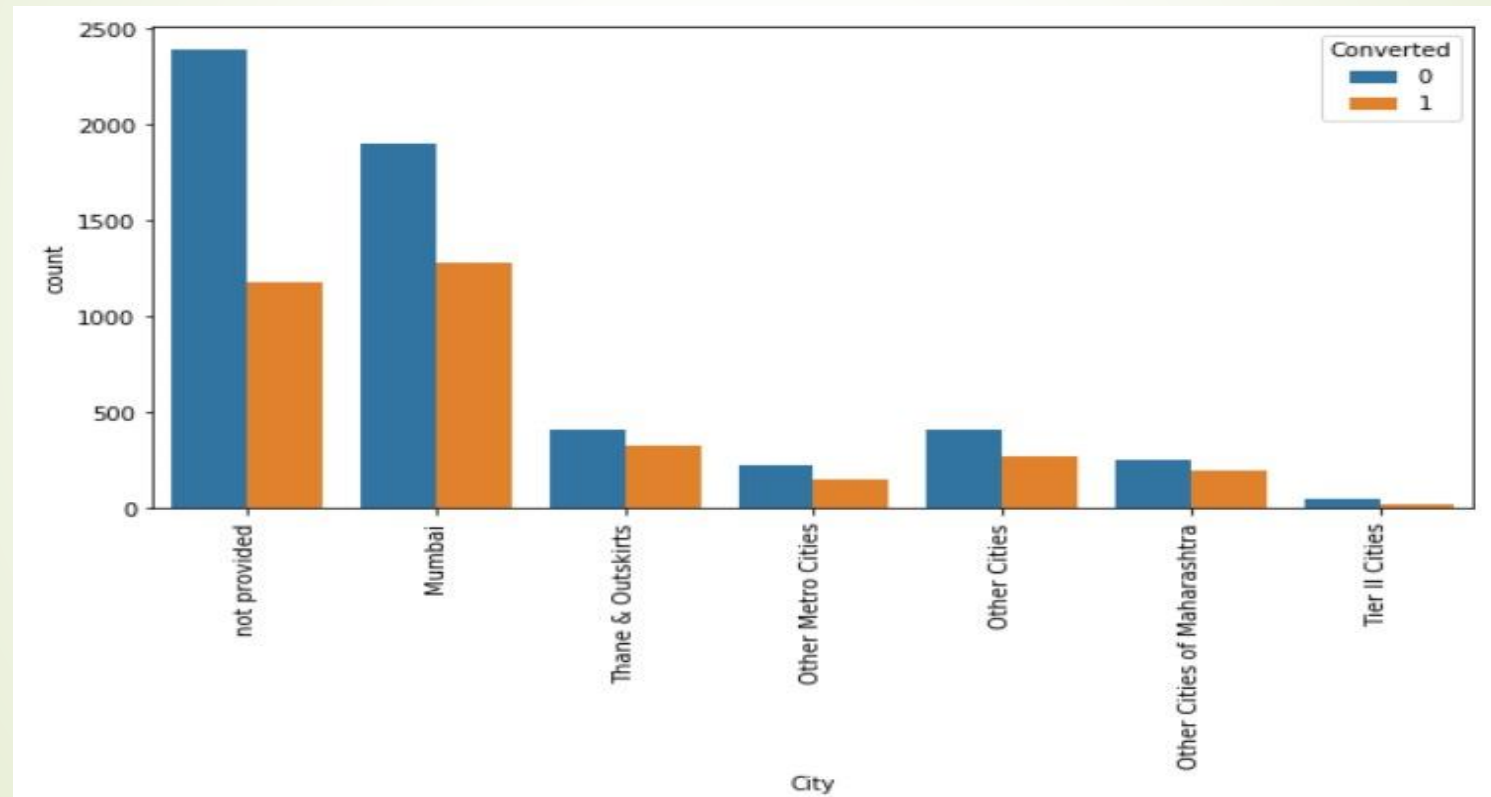
Data Manipulation

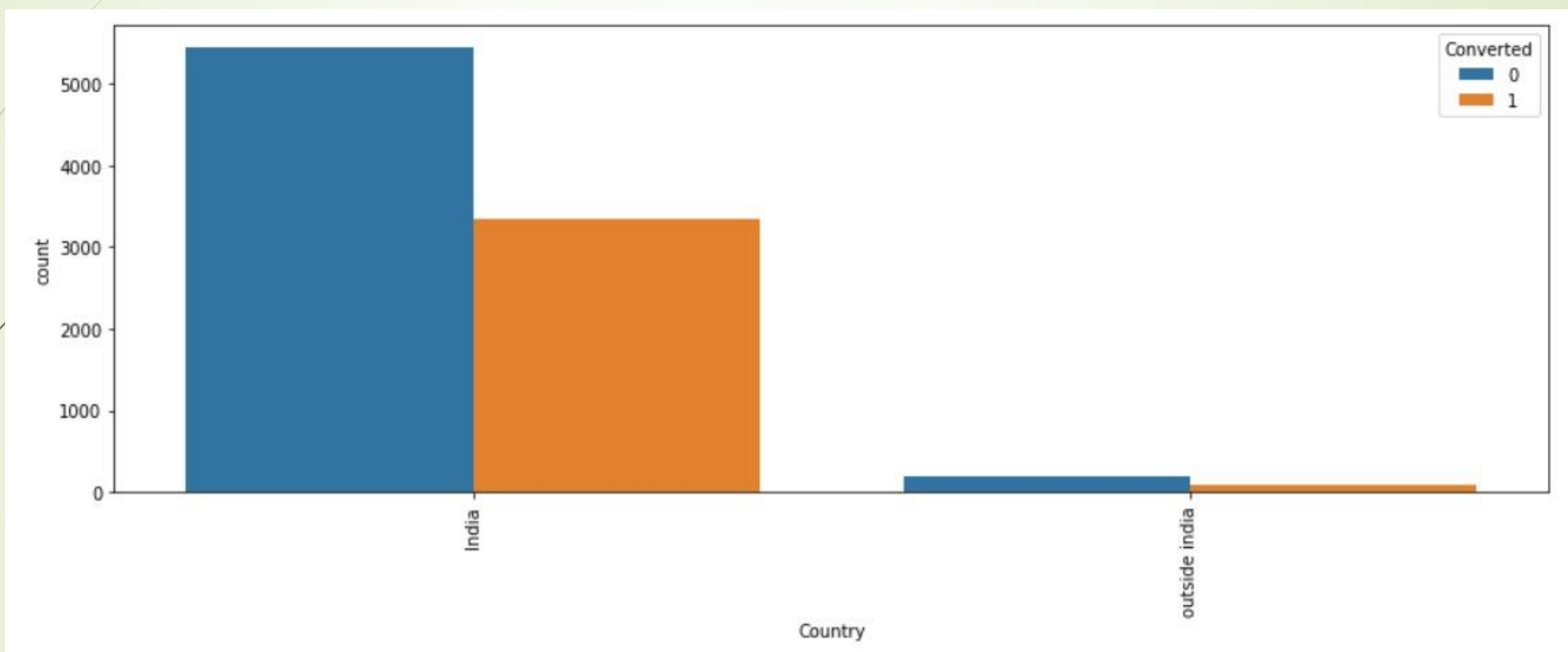


- Total Number of Rows =37, Total Number of Columns =9240.
- Single value features like “Magazine”, “Receive More UpdatesAbout Our Courses”, “Update me on Supply”.
- “Chain Content”, “Get updates on DM Content”, “I agree to pay the amount through cheque etc. have been dropped.
- Removing the “Prospect ID” and “Lead Number” which is not necessary for the analysis.

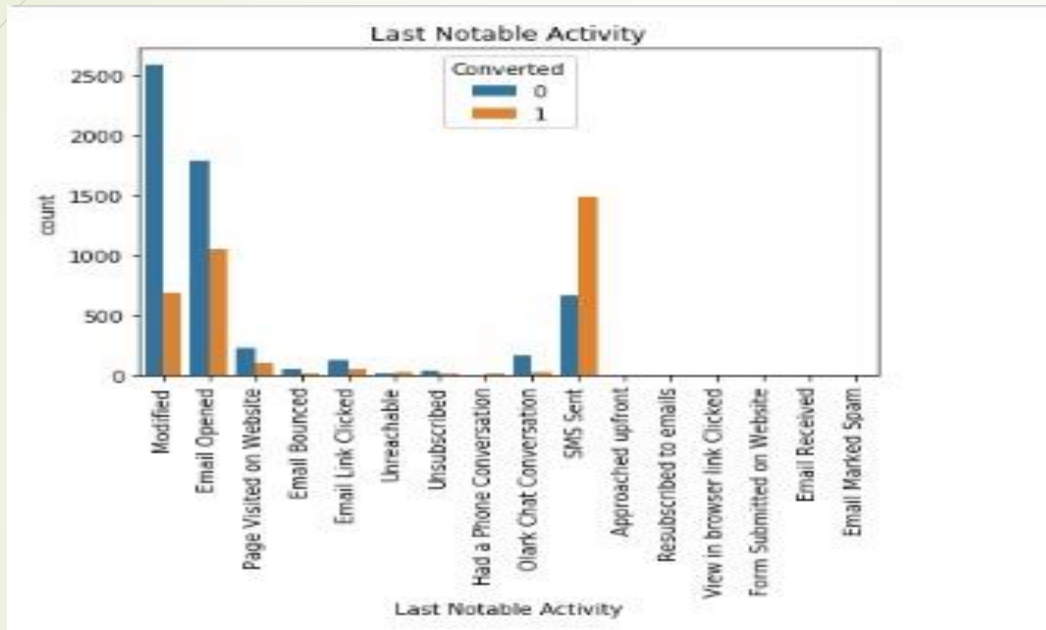
- 
- 
- After checking for the value counts for some of the object type variables, we find some of the features which has no enough
 - variance, which we have dropped, the features are: “Do Not Call”, “What matters most to you in choosing course”, “Search”, “Newspaper Article”, “X Education Forums”, “Newspaper”, “DigitalAdvertisement” etc.
 - Dropping the columns having more than 35% as missing value such as ‘How did you hear about X Education’ and ‘Lead Profile’.

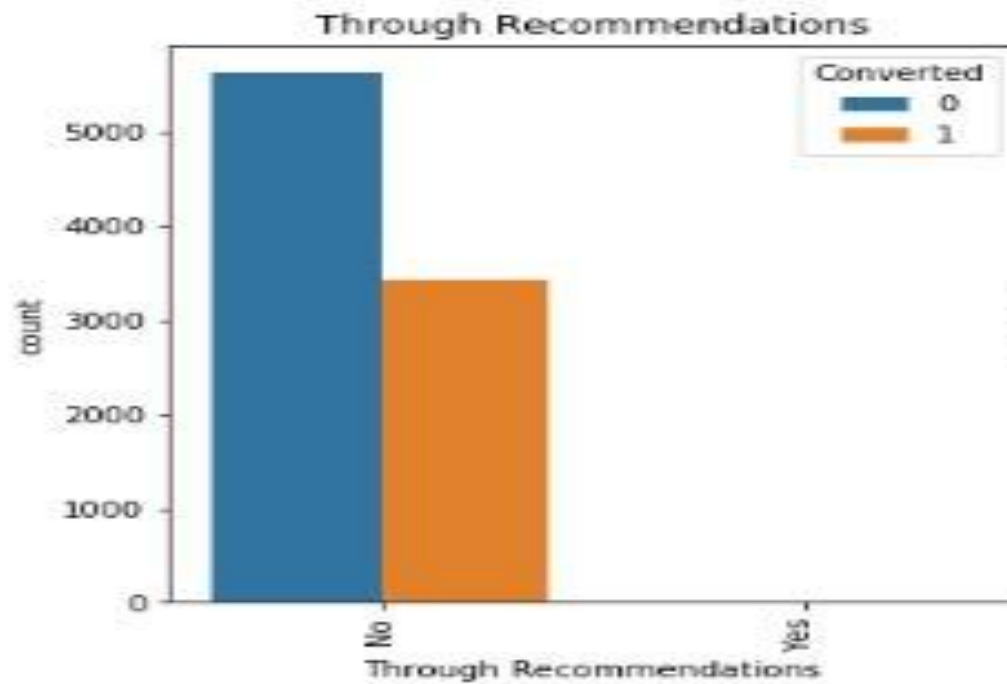
Exploratory Data Analysis

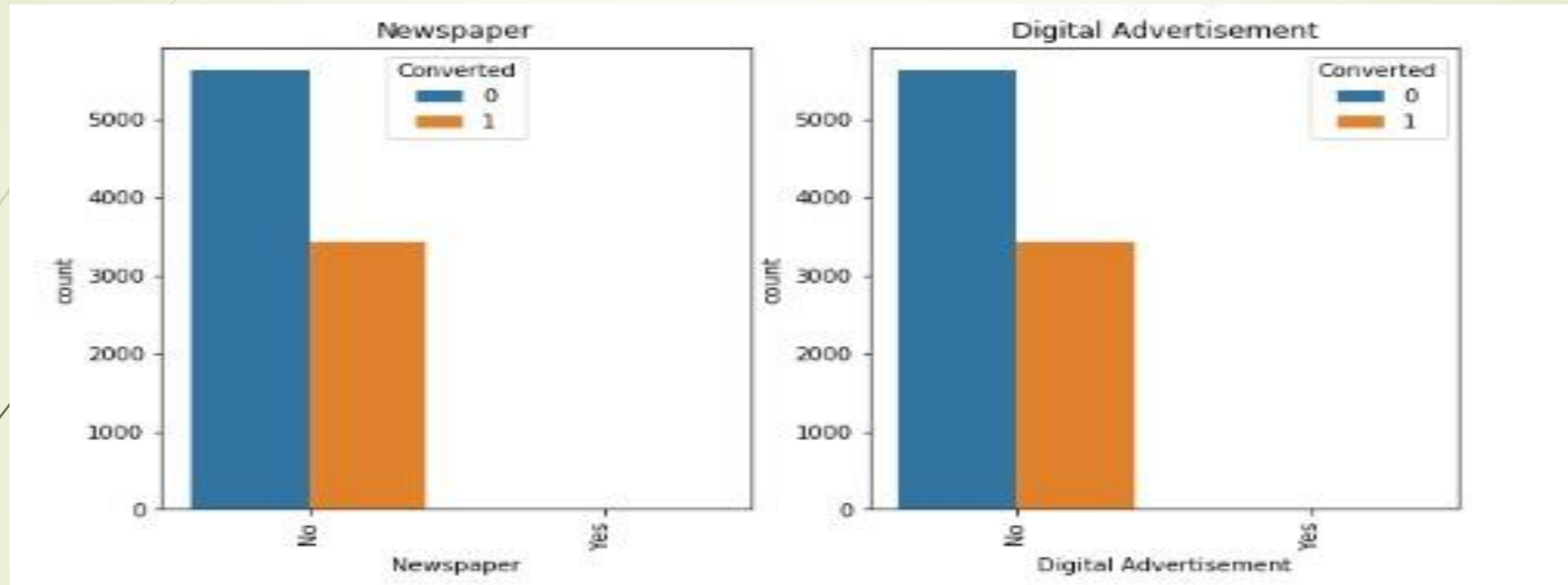


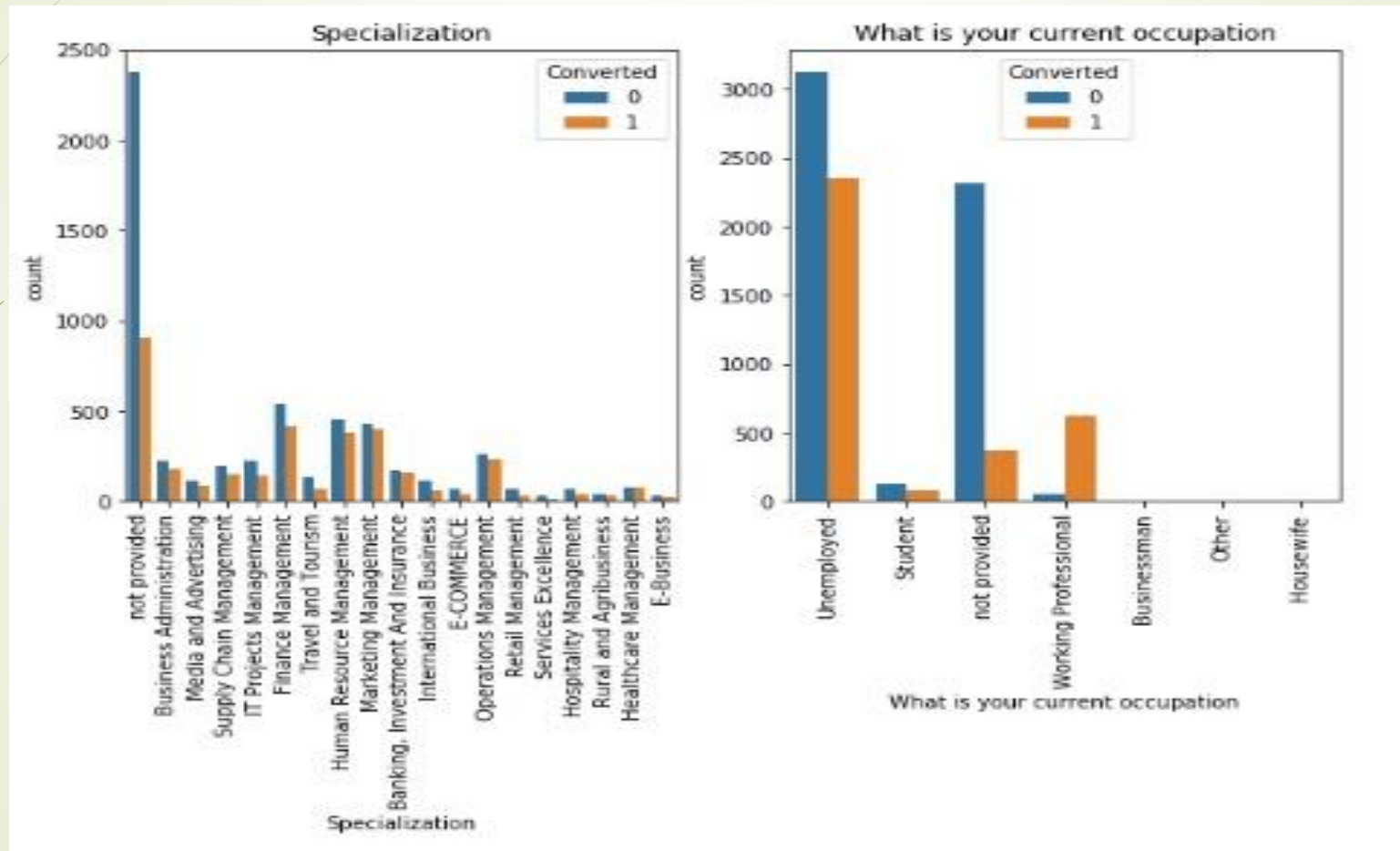


Categorical Variable Relation

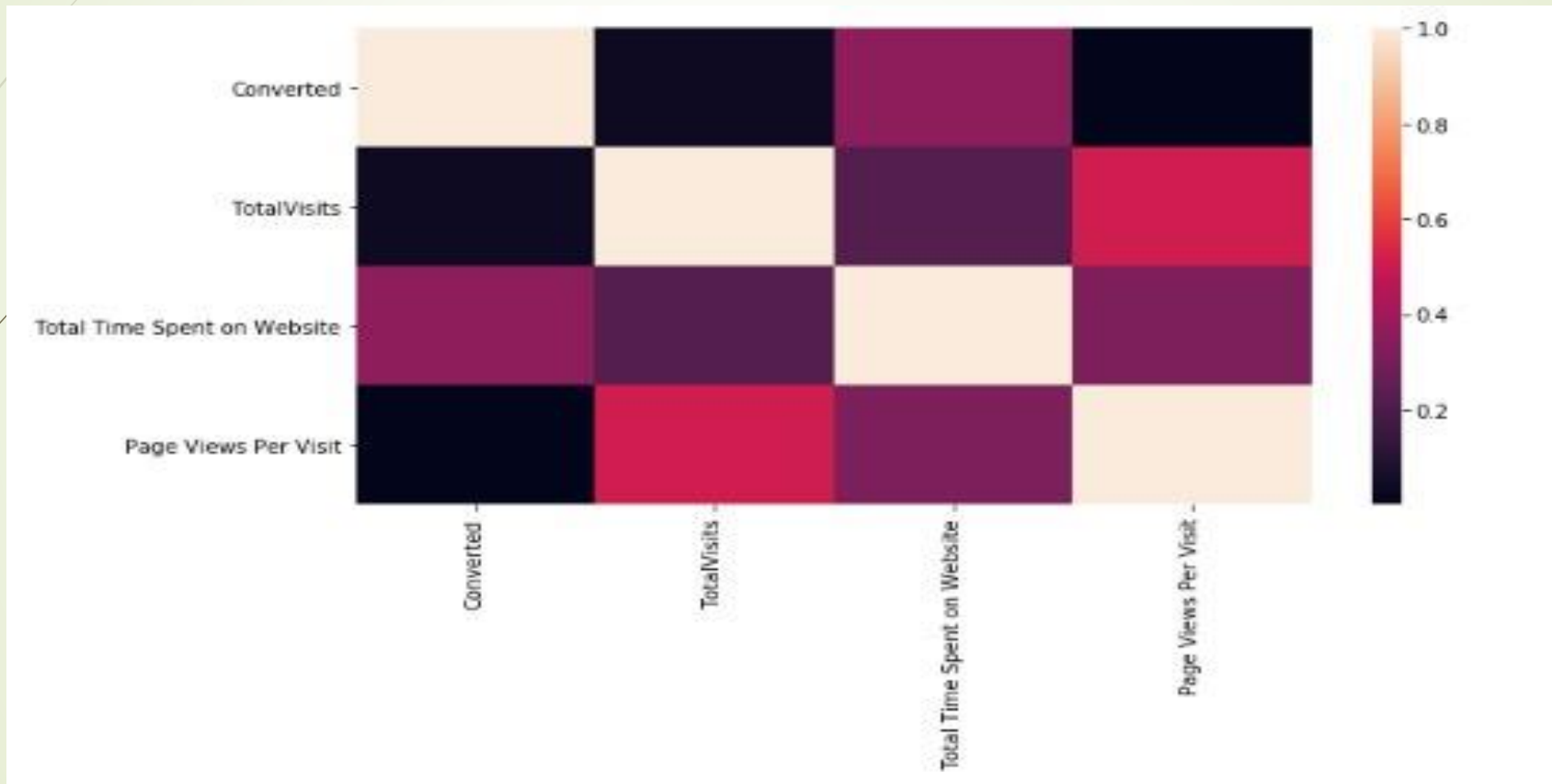




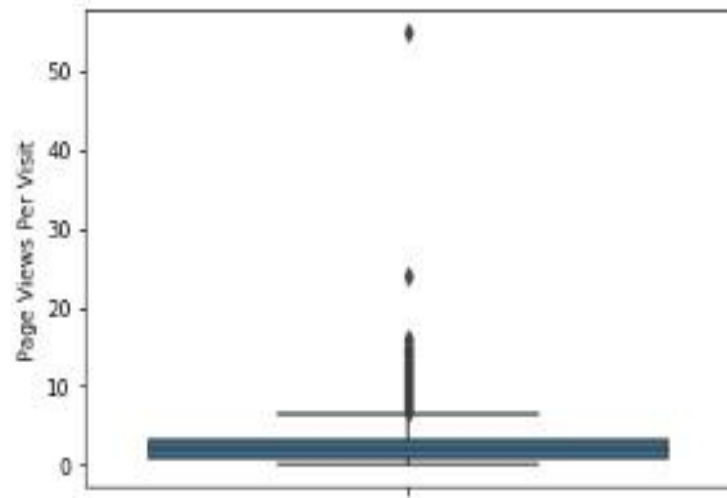




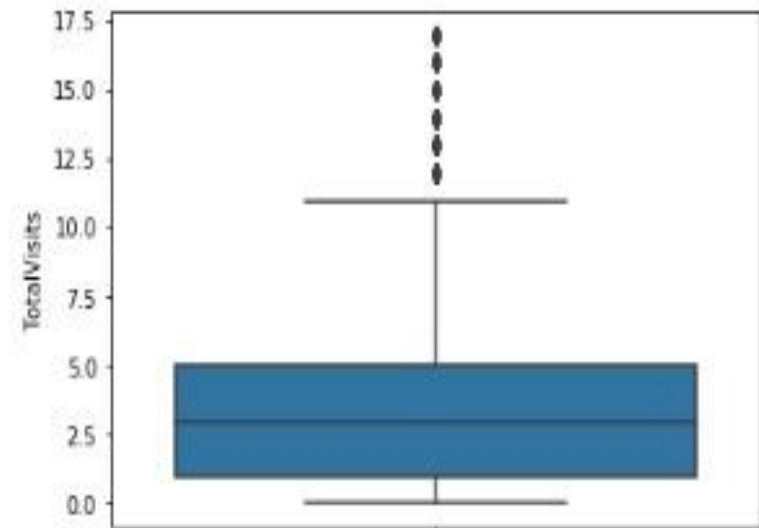
Corelation Among Variables

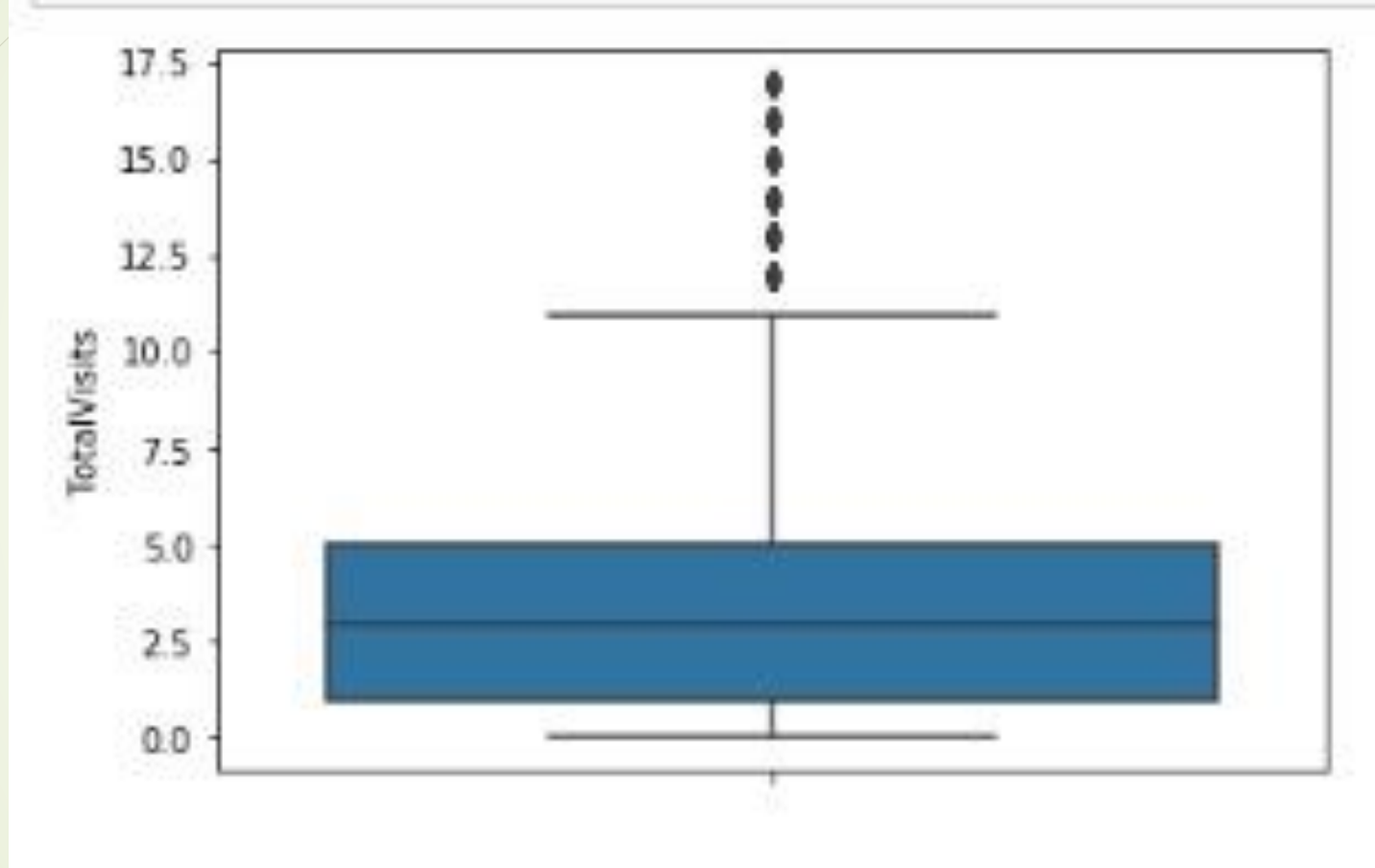


Outliers Detected

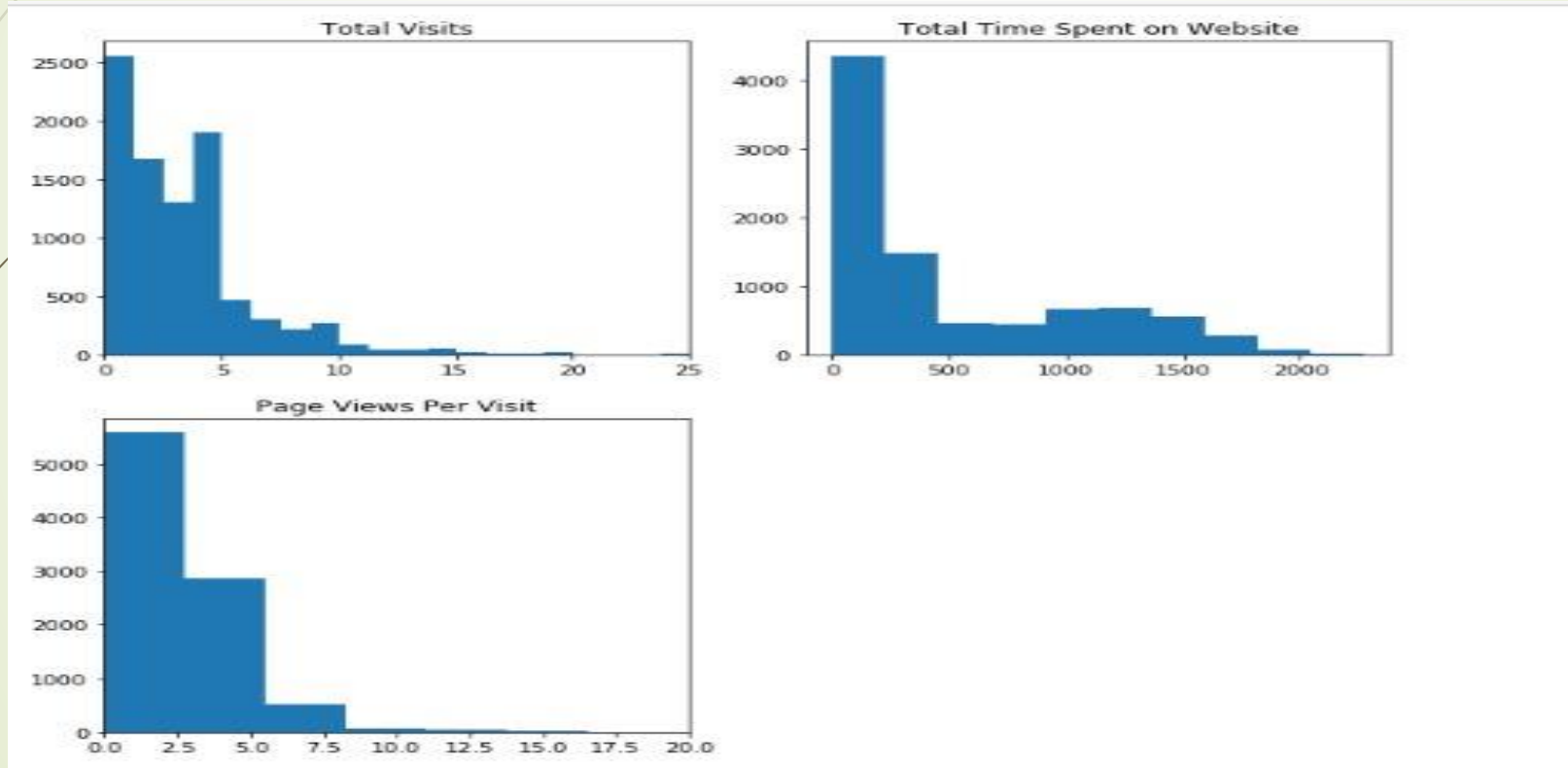


We can see presence of outliers in TotalVisits





Numerical Variable Relation





Data Conversion

- Numerical Variables are Normalised.
- Dummy Variables are created for object type variables.
- Total Rows for Analysis: 9240
- Total Columns for Analysis: 37



Model Building

- Splitting the Data into Training and Testing Sets.
- The first basic step for regression is performing a train-test split, we have chosen 70:30 ratio.
- Use RFE for Feature Selection.
- Running RFE with 15 variables as output.
- Building Model by removing the variable whose p-value is greater than 0.05 and VIF value is greater than 5.
- Predictions on test data set.
- Overall accuracy 81%.

Conclusion

- It was found out that the variables that mattered the most in the potential buyers are :
- The Total Time Spent on Website.
- Total number of Visits.
- When The Lead source was :a. Olarkchat b. WellingakWebsite
- When the Last Activity was :a. SMS b. OlarkChat Conversation
- When the lead origin is Lead add Form.
- When the current occupation was :
 - a. Working Professionals
 - b. Student
 - c. Unemployed
 - d. Other
- Keeping The above mentioned points in mind the X education can increase all thepotential buyers to change their mind and buy their courses.