

Lead Score Case Study

Group Members

Yasaswi Racha

Varun Duggal

Rajat Gupta

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 '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. Business Objective:
- X education wants to know most promising lead

Problem Statement (contd)

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

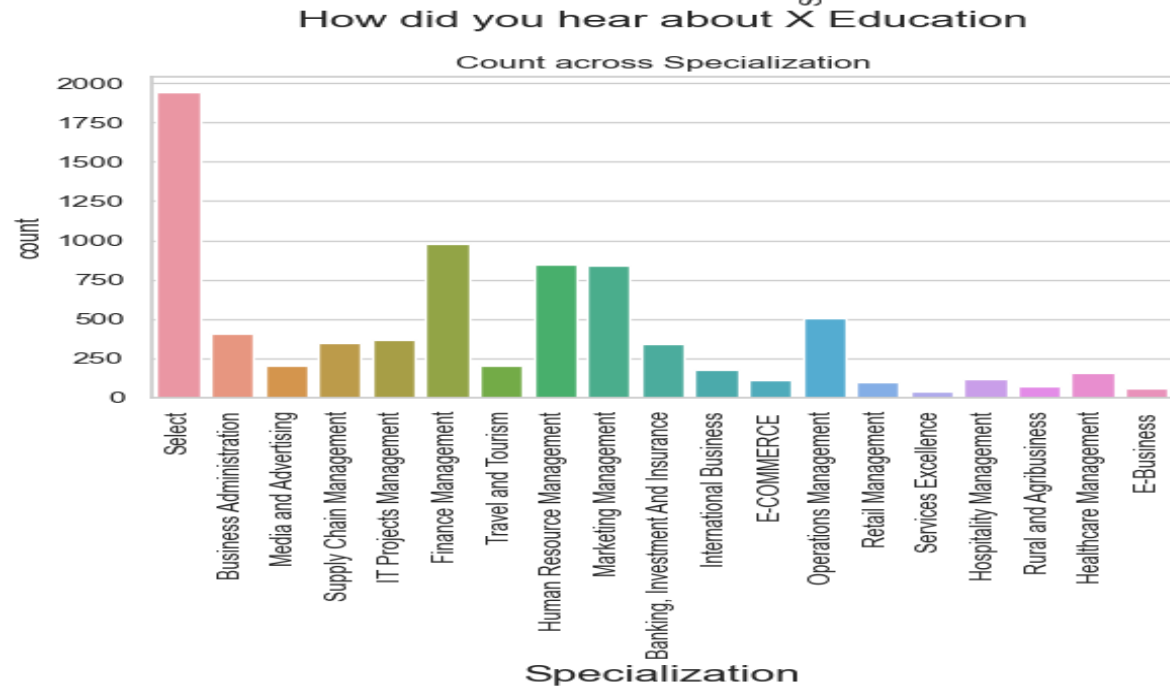
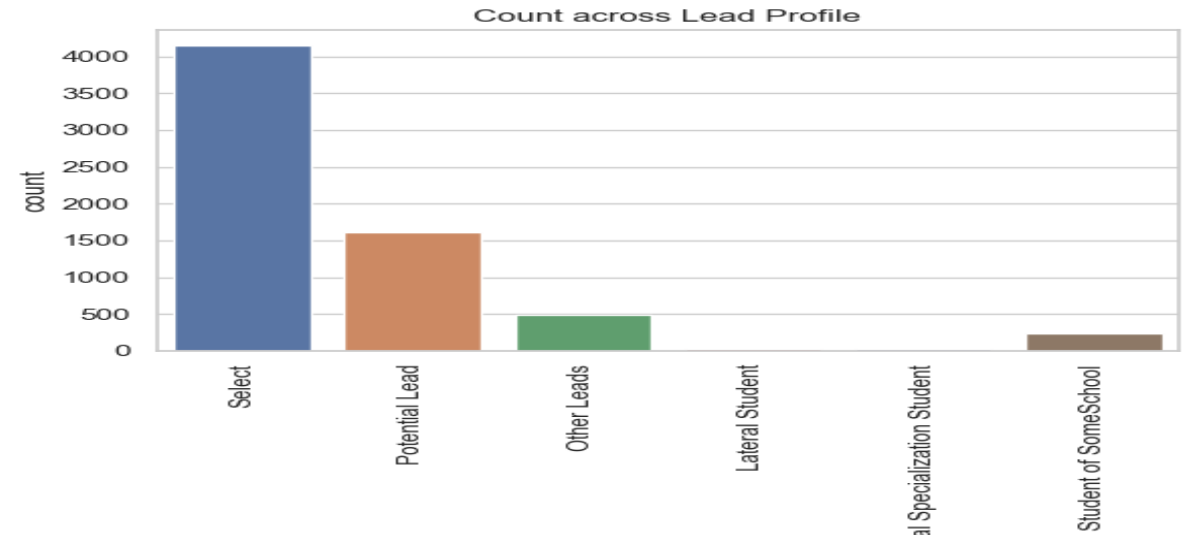
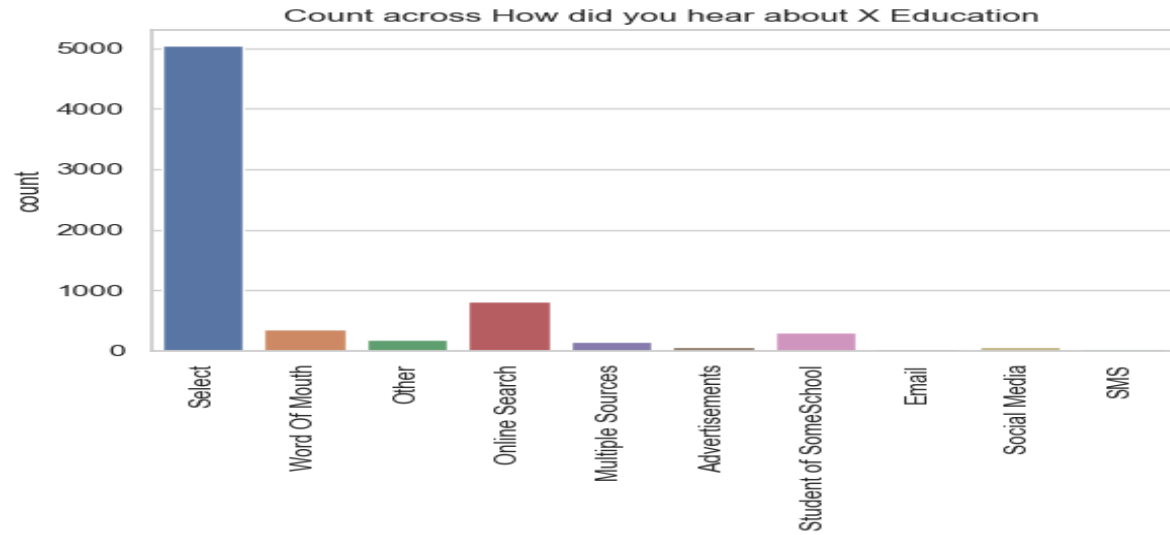
1. Univariate data analysis: value count, distribution of variable etc.
2. Bivariate data analysis: correlation coefficients and pattern between the variables etc.
 - Feature Scaling & Dummy Variables and encoding of the data.
 - Classification technique: logistic regression used for the model making and prediction.
 - Validation of the model. Model presentation.
 - Conclusions and recommendations.

Data Manipulation

➤ Data Manipulation

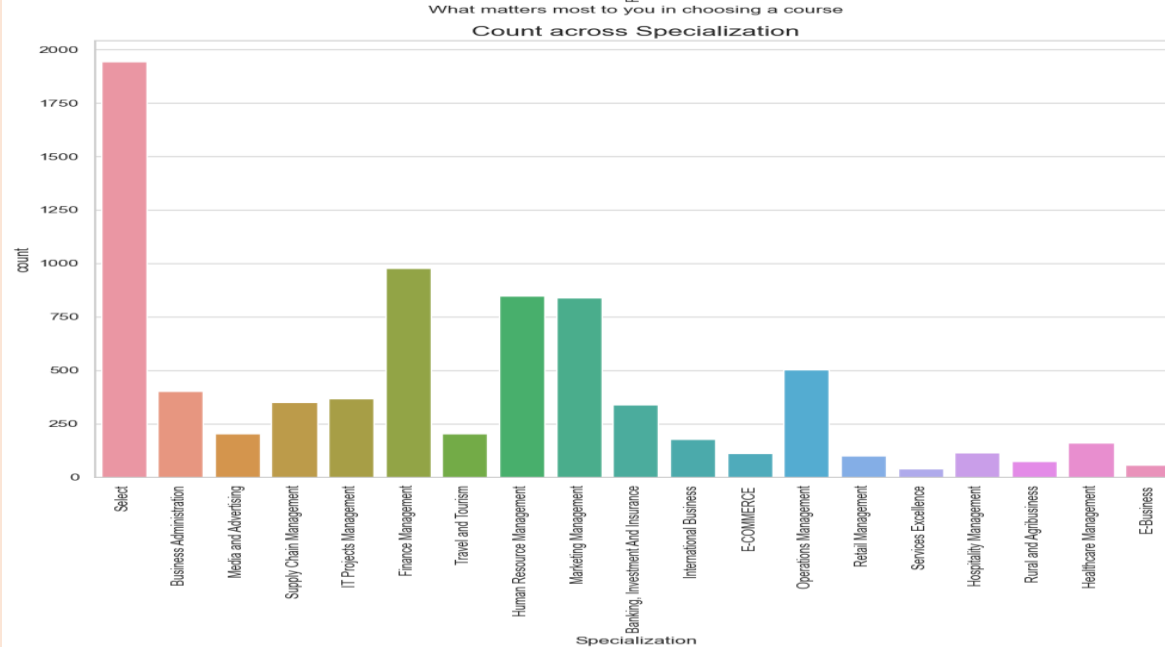
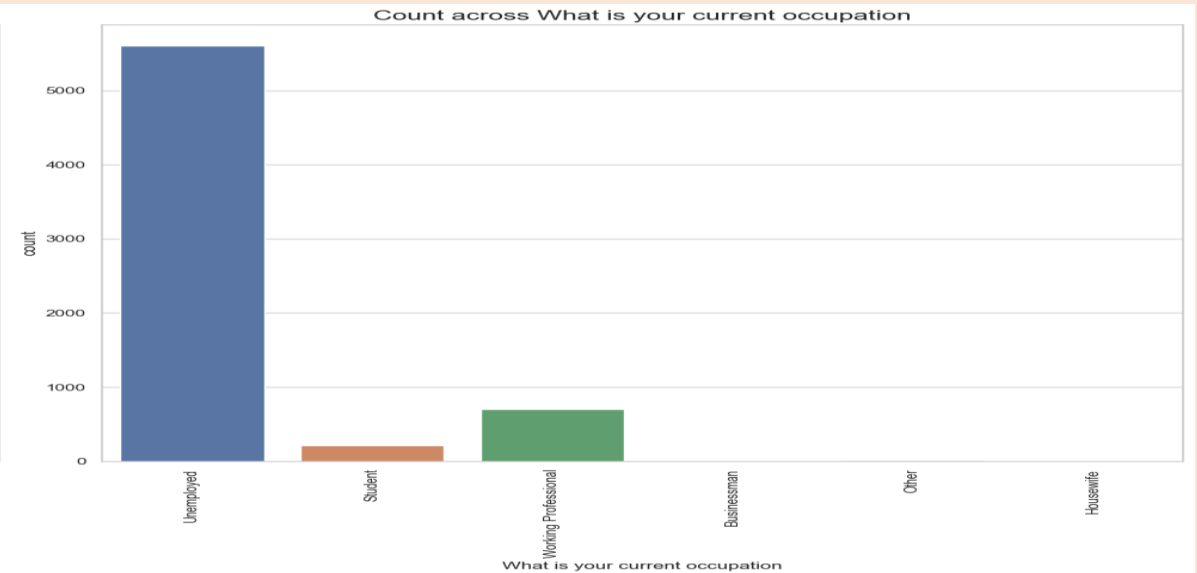
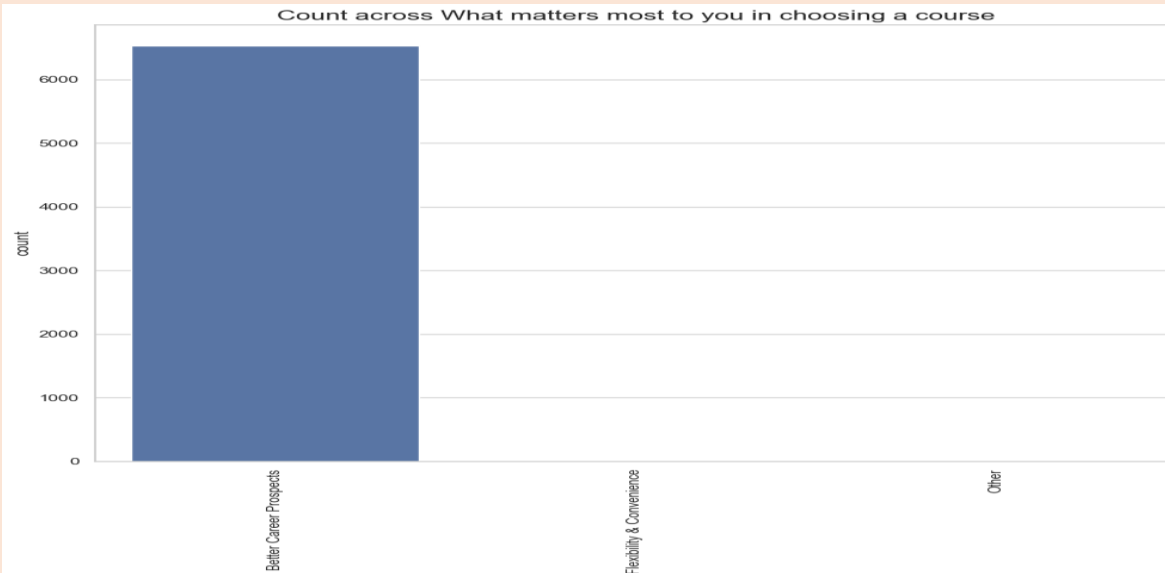
- Total Number of Rows =37
- Total Number of Columns =9240.
- Single value features like “Magazine”, “Lead Profile", 'How did you hear about X Education, City, Country, etc. have been dropped.
- Removing the “Prospect ID” and “Lead Number” which is not necessary for the analysis. After checking for the 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', 'Search', 'Magazine', 'Newspaper Article', 'X Education Forums', 'Newspaper', 'Digital Advertisement', 'Through Recommendations', 'Receive More Updates About Our Courses’.

EDA Graphs

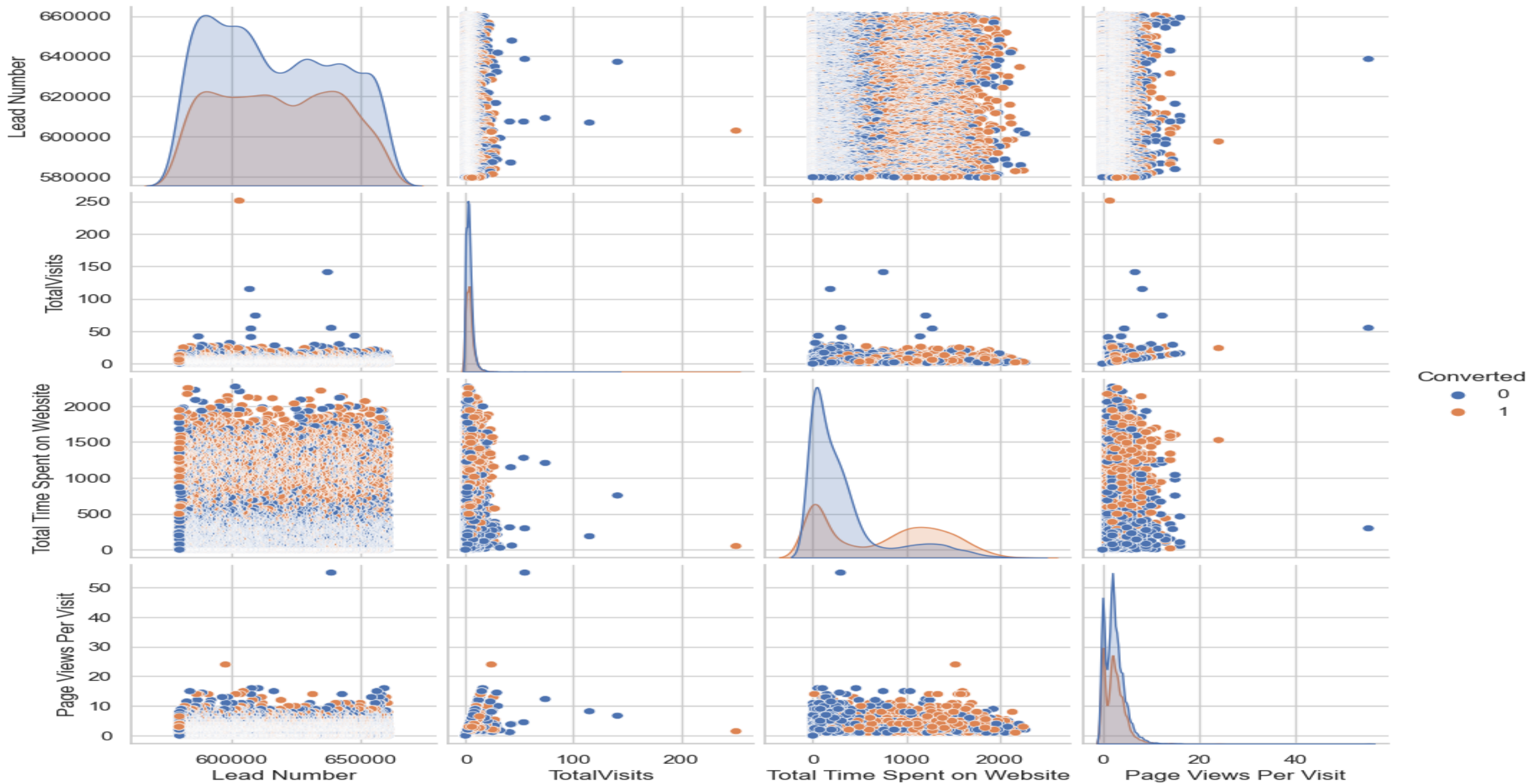


Lead Profile

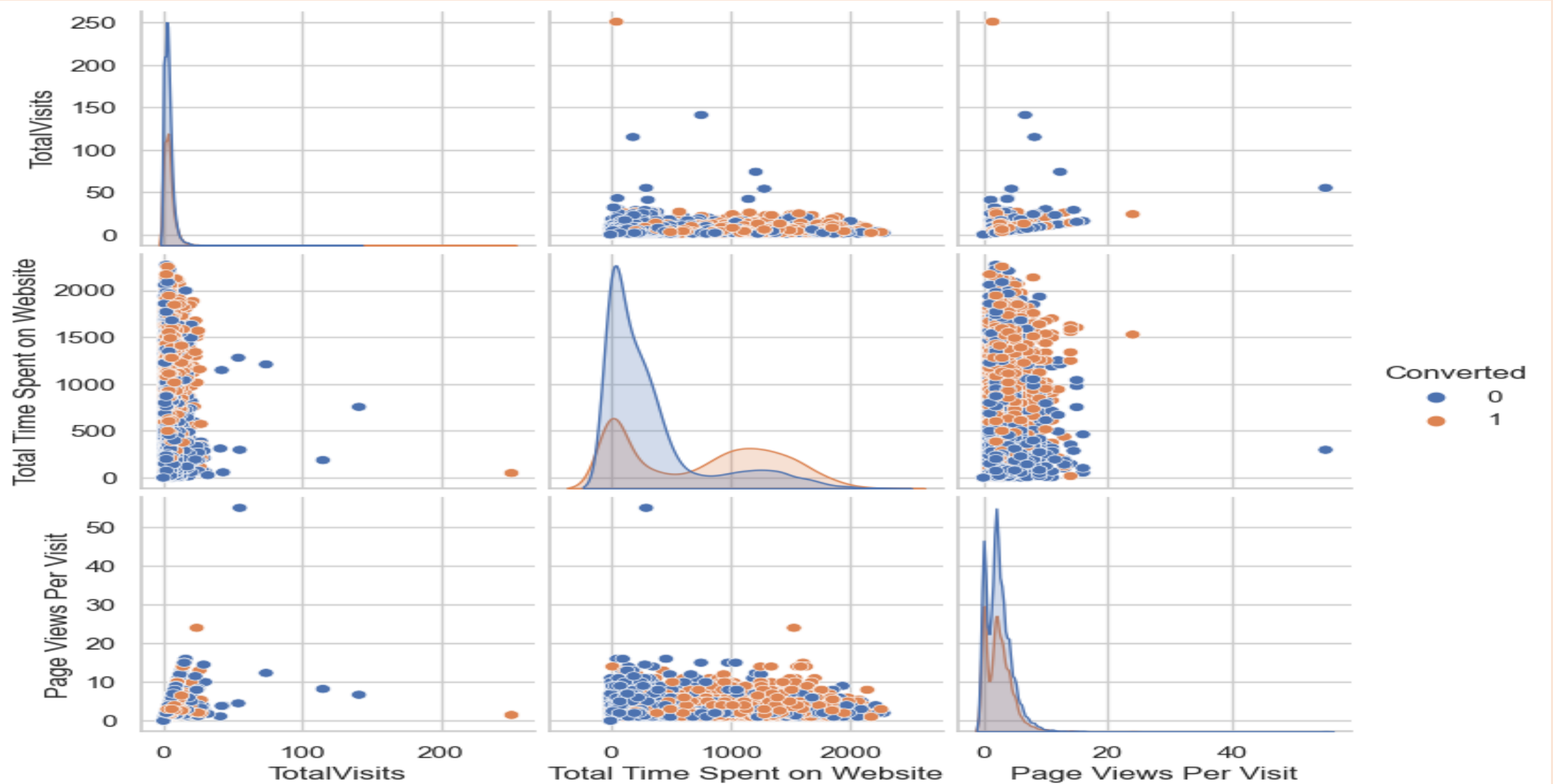
EDA Graphs



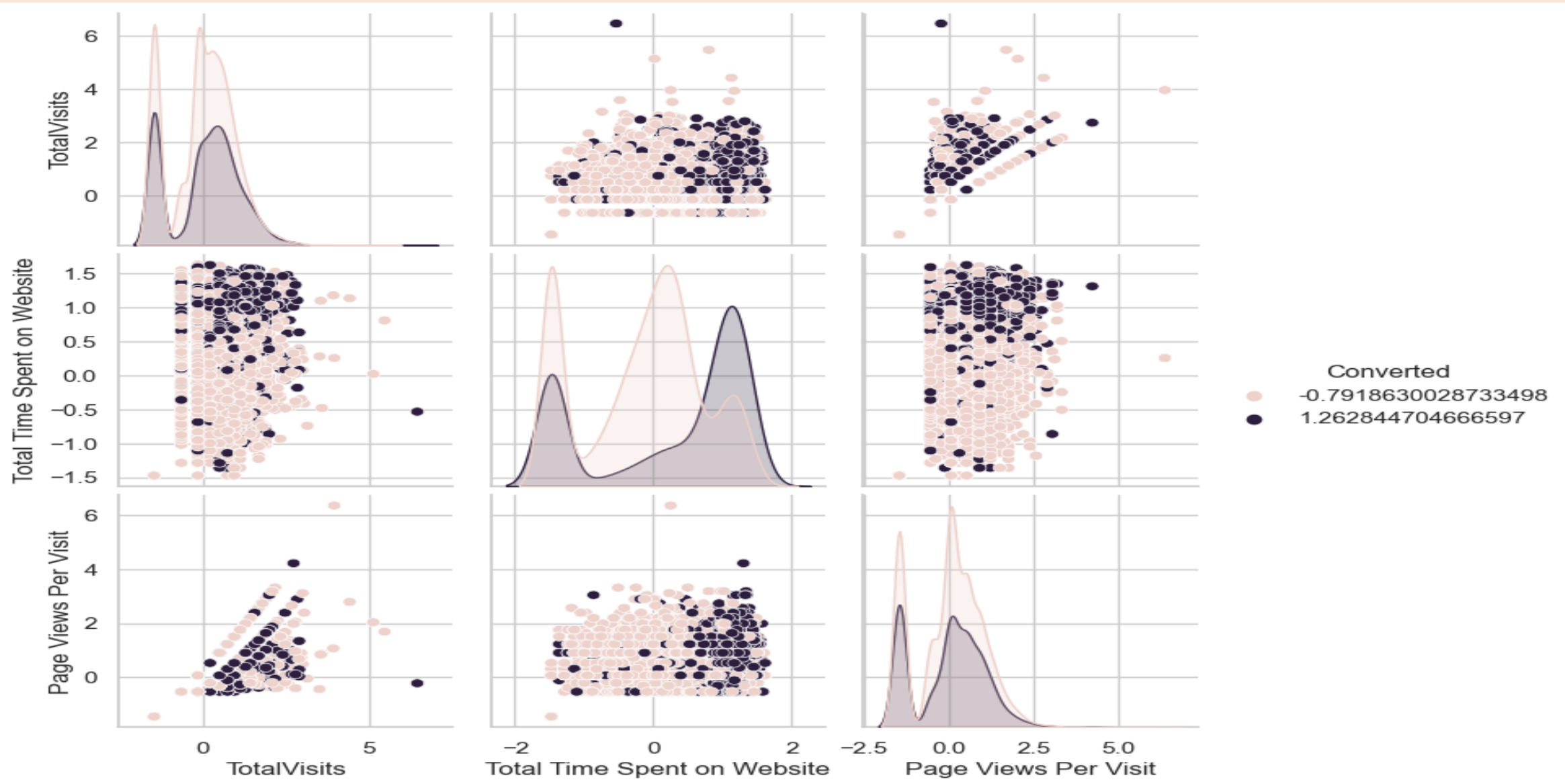
EDA Graphs



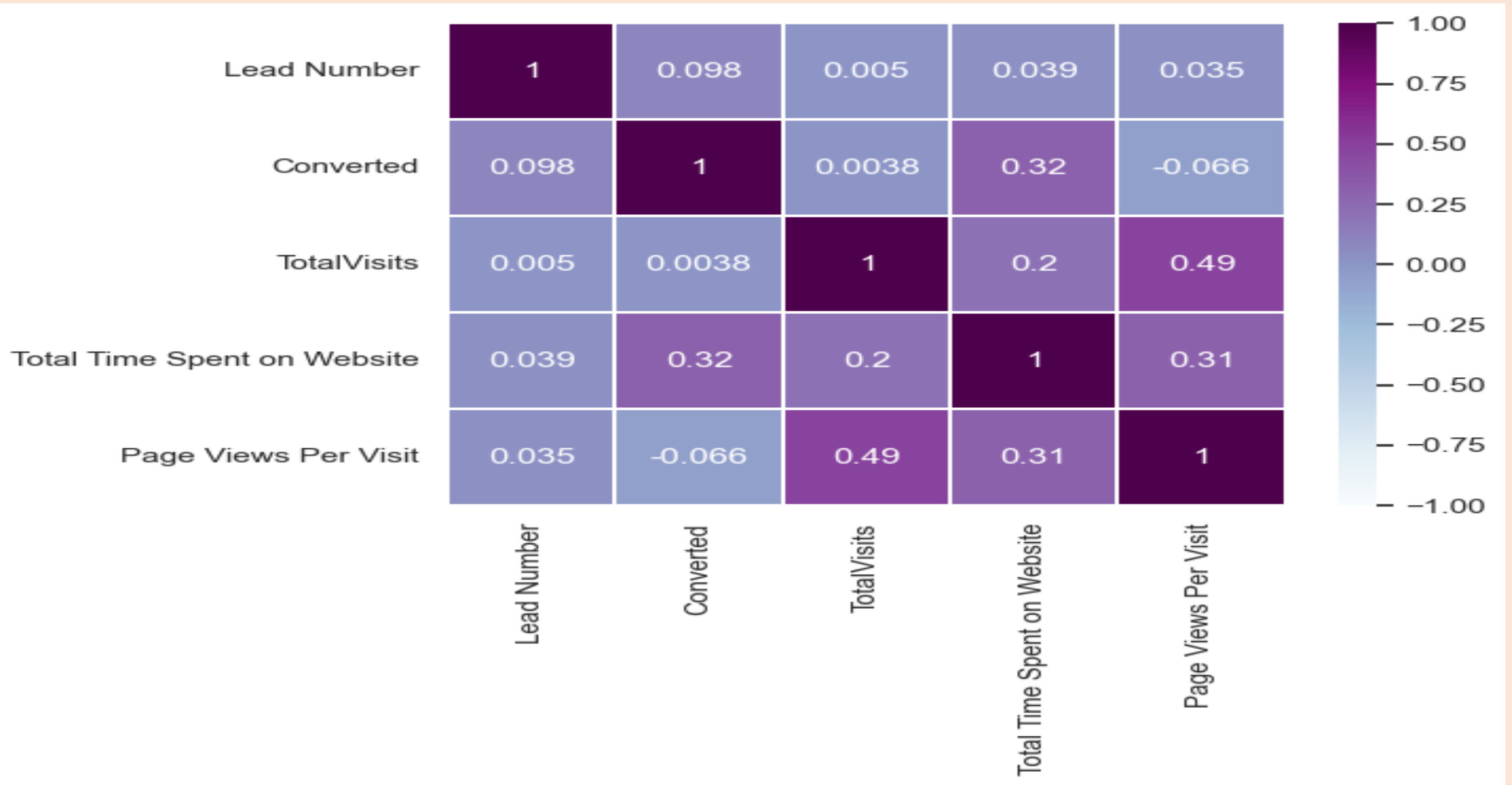
EDA Graphs



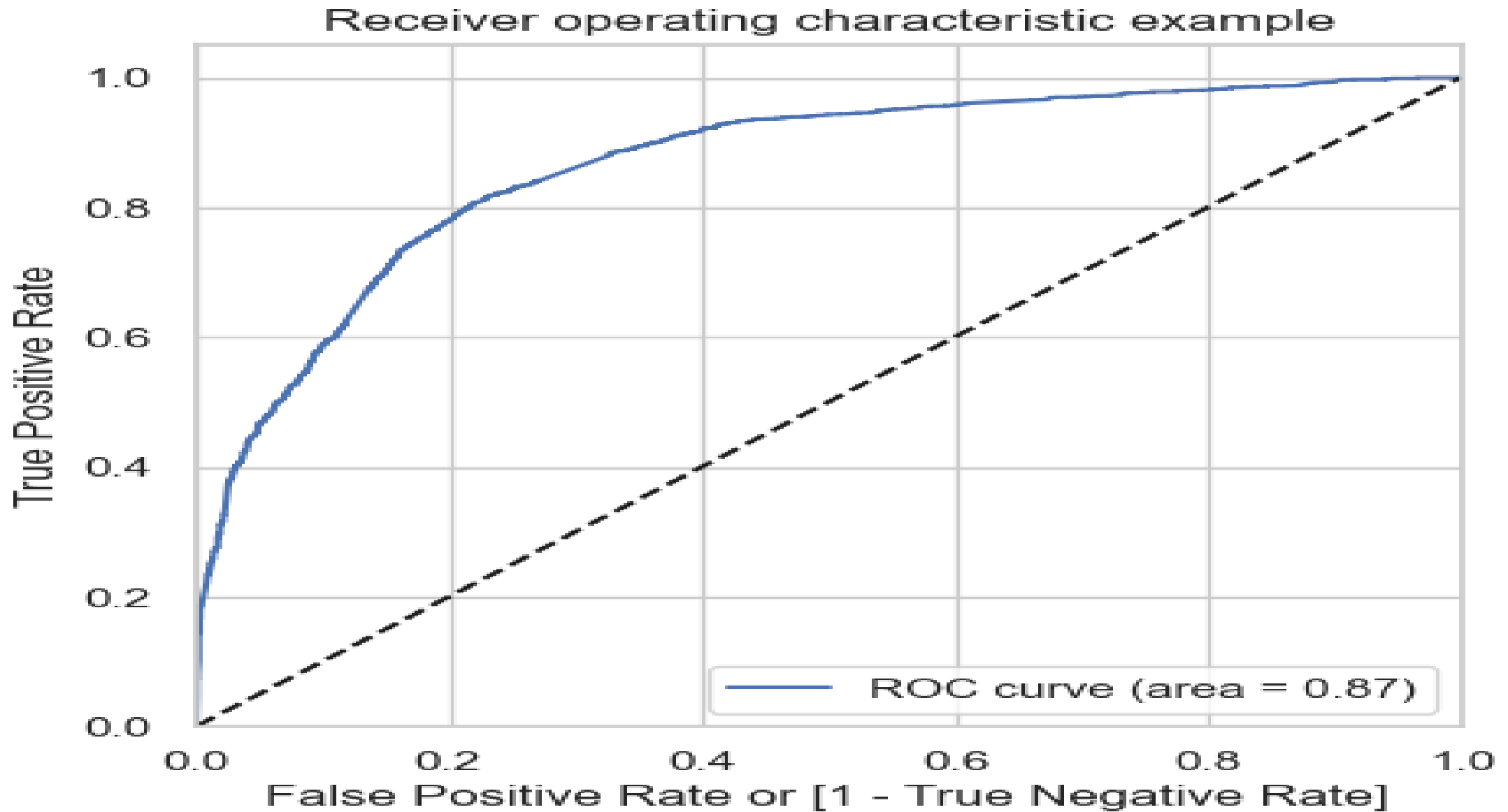
EDA Graphs



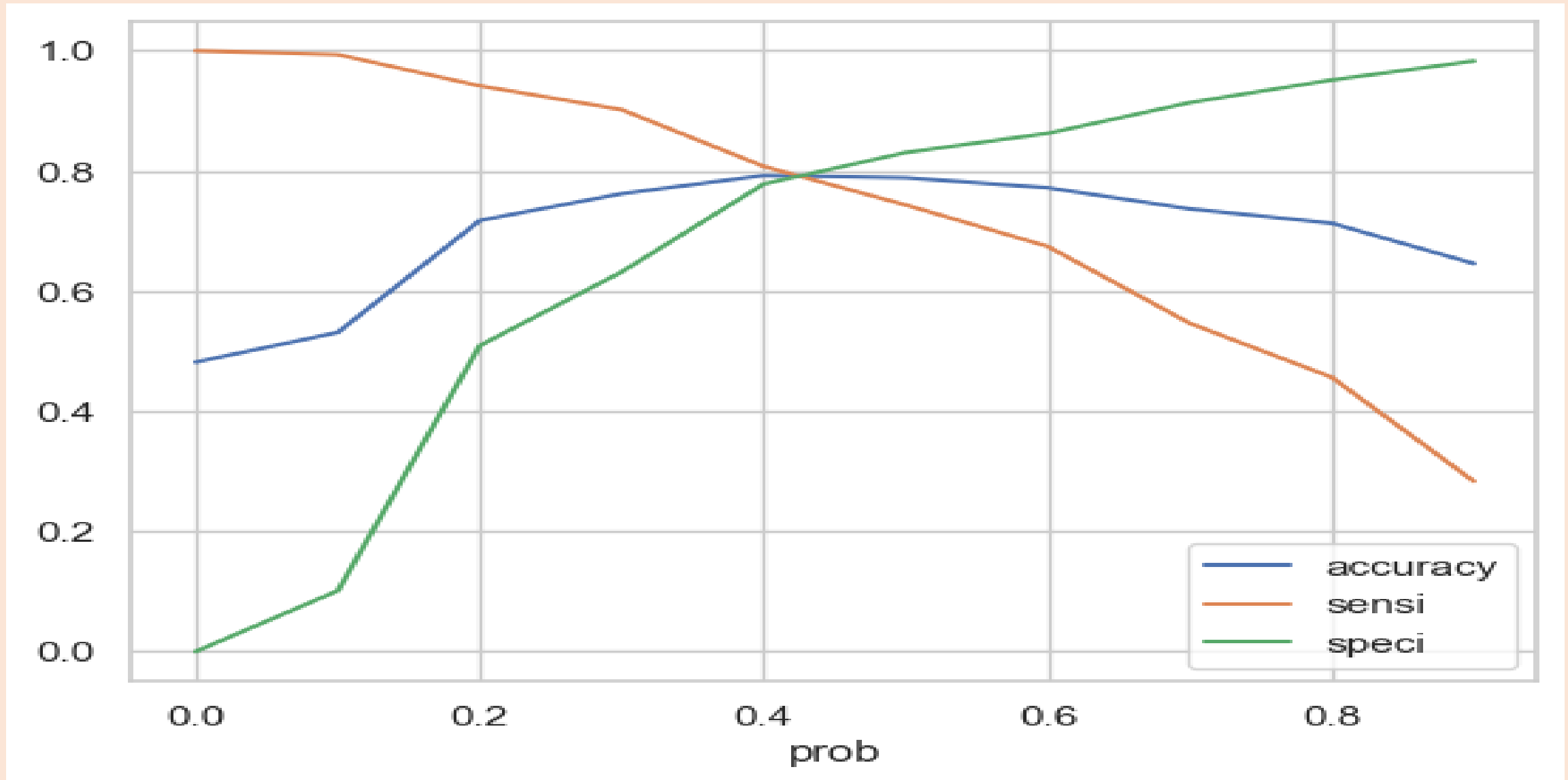
Correlation



ROC Curve Graphs



Accuracy sensitivity and specificity Graphs



Data Conversion

- Numerical Variables are Normalised
- Dummy Variables are created for object type variables
- Total Rows for Analysis: 6373
- Total Columns for Analysis: 75

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

Conclusion

- The most numbers of leads are from India and in terms of city highest number are from Mumbai.
- First, sort out the best prospects from the leads you have generated. 'Total Visits' , 'Total Time Spent on Website' , 'Page Views Per Visit' which contribute most towards the probability of a lead getting converted.
- The leads are joined course for Better Career Prospects, most of having Specialization from Finance Management. Leads from HR, Finance & marketing management specializations are high probability to convert.
- Talking to last notable Activity, making improvement in customer engagement through email & calls will help to convert leads. As the leads which are opening email have high probability to convert.