**Problem Statement** : We have been provided with a leads dataset from the past with around 9000 data points. This dataset consists of various attributes such as Lead Source, Total Time Spent on Website, Total Visits, Last Activity, etc. which may or may not be useful in ultimately deciding whether a lead will be converted or not.

**The target variable**, in this case, is the column 'Converted' which tells whether a past lead was converted or not wherein 1 means it was converted and 0 means it wasn't converted.

Another thing that we need to check out for are the levels present in the categorical variables. Many of the categorical variables have a level called 'Select' which needs to be handled because it is as good as a null value

Approach

1  Reading and Understanding Data
2  Data Cleaning
- 2.1  Rename column names
- 2.2  Drop prospect_id column
- 2.3  Replace "Select" category with null values
- 2.4  Handle null values and sales generated columns
  - 2.4.1  Drop columns that have null values > 40% or Sales generated columns
  - 2.4.2  country column
  - 2.4.3  course_selection_reason column
  - 2.4.4  occupation column
  - 2.4.5  specialization column
  - 2.4.6  city column
- 2.5  Handle categorical columns with low number of missing values and low representation of categories
  - 2.5.1  lead_origin column
  - 2.5.2  lead_source column
- 2.6  Handle Binary columns
- 2.7  Handle Numerical columns
  - 2.7.1  lead_number column: change datatype
  - 2.7.2  total_visits column
  - 2.7.3  page_views_per_visit column

3 Exploratory Data Analysis
- 3.1  Numerical columns
  - 3.1.1  Heatmap
  - 3.1.2  Check for outliers
- 3.2  Categorical columns
  - 3.2.1  Lead Origin
  - 3.2.2  Lead Source
  - 3.2.3  Specialization
  - 3.2.4  Occupation
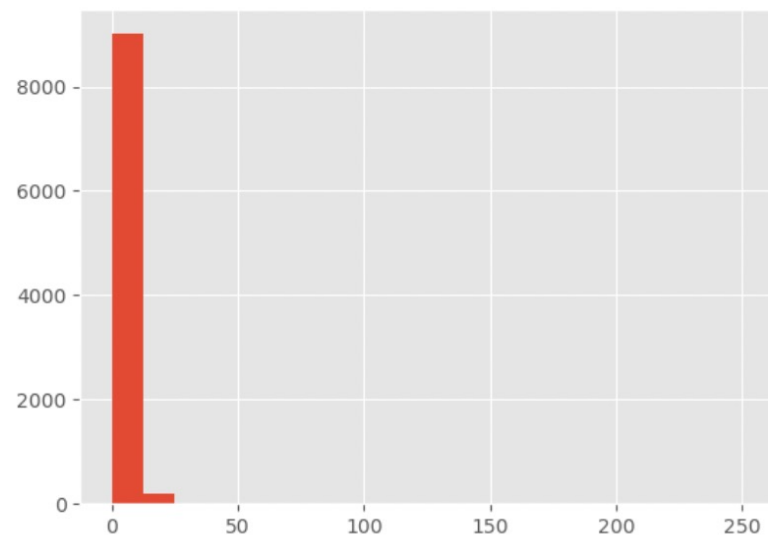  - 3.2.5  City

4  Data Preparation
- 4.1  Converting Binary (Yes/No) to 0/1
- 4.2  Creating dummy variable for categorical columns
- 4.3  Outliers Treatment
- 4.4  Test-Train Split
- 4.5  Feature Scaling
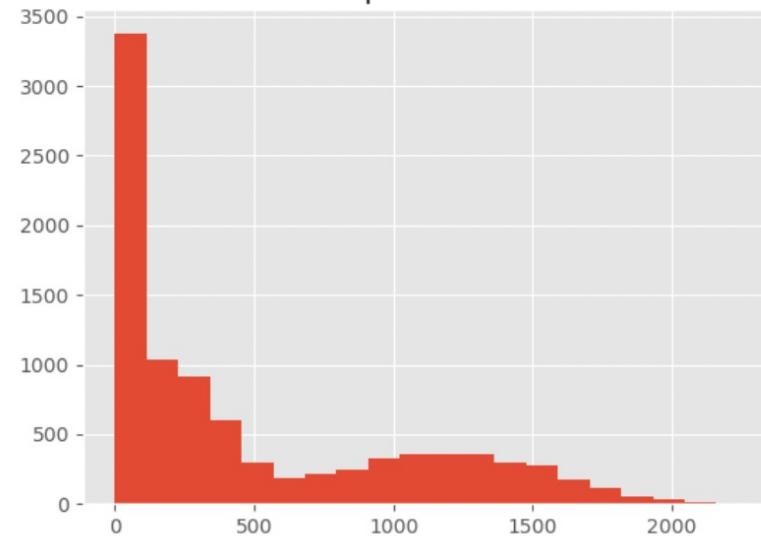- 4.6  Looking at correlations
  - 4.6.1  Drop highly correlated dummy variables

5  Model Building
- 5.1  Model 1: All variables
- 5.2  Feature selection using RFE
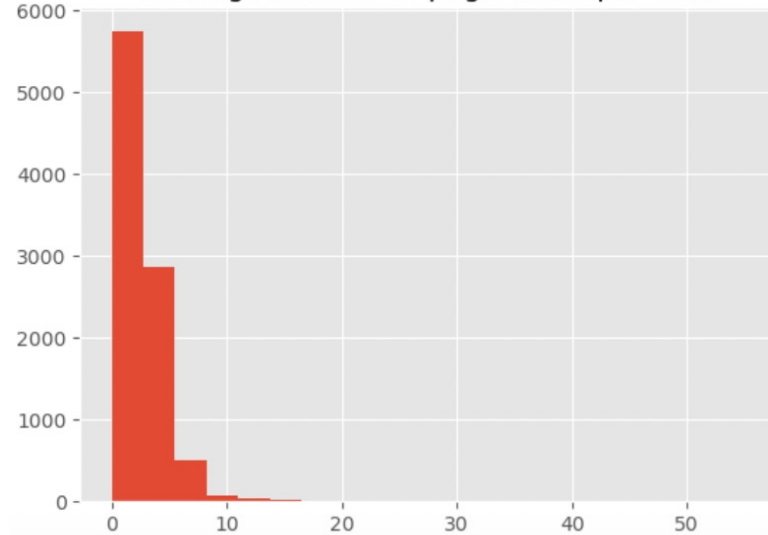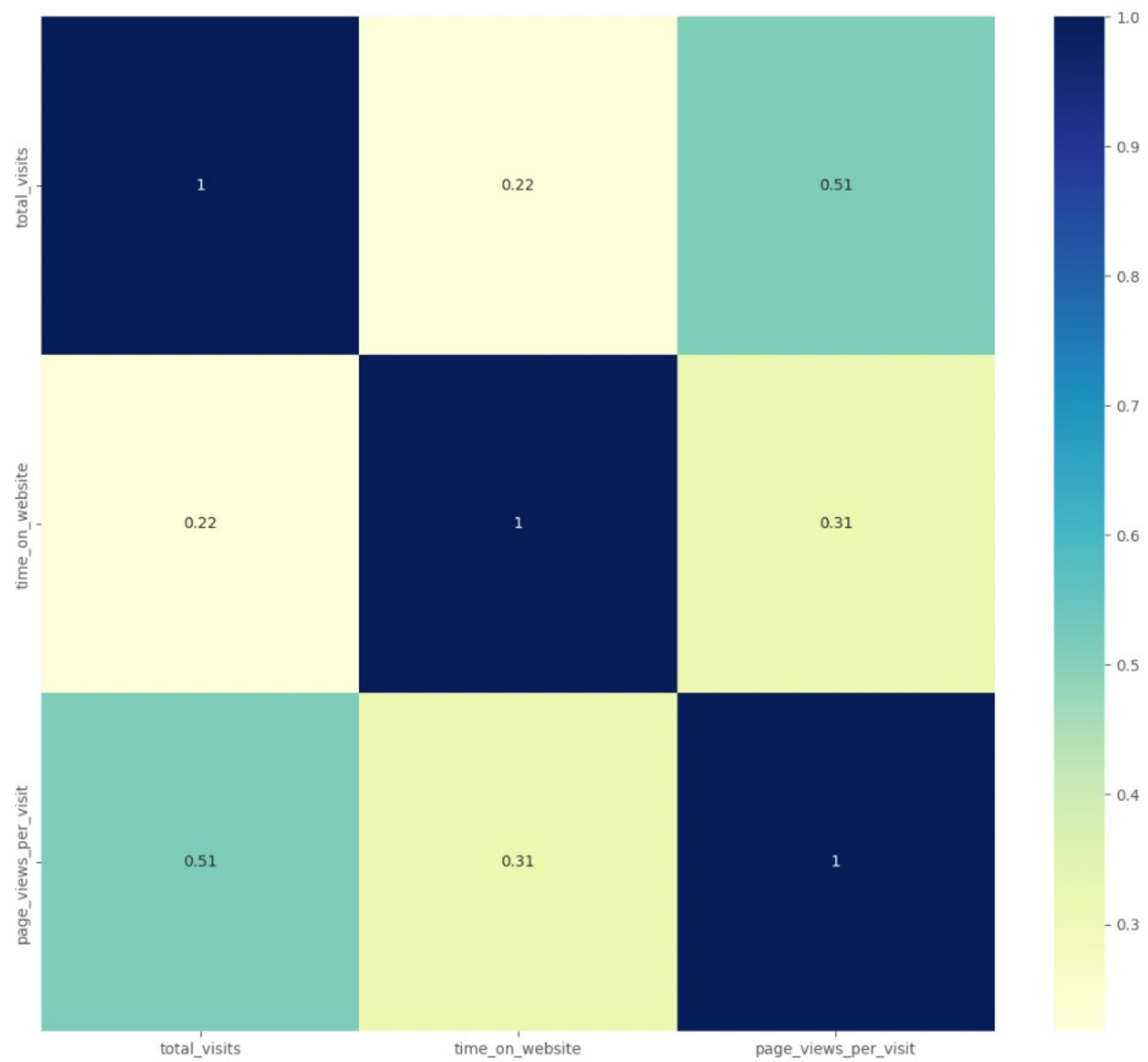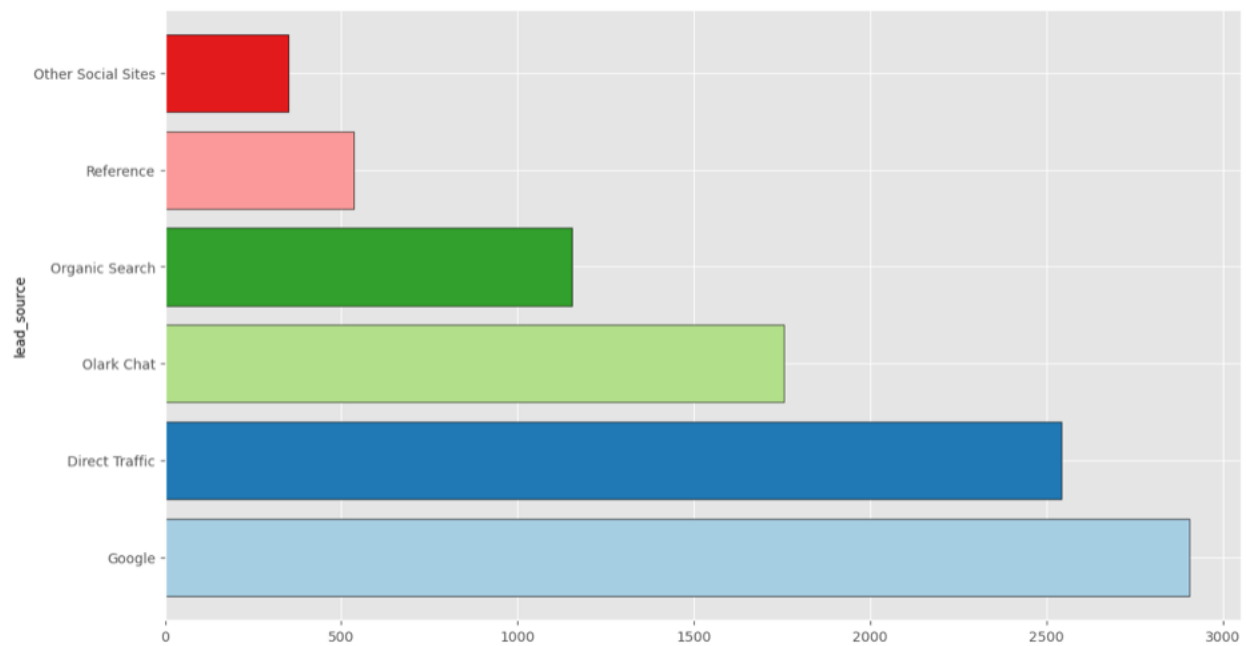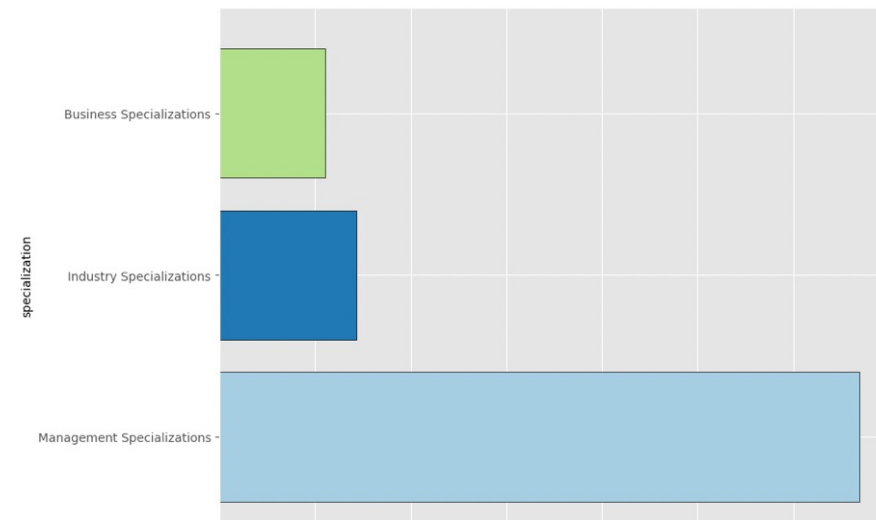- 5.3  Model 2: Assessing the model with statsmodel

**Conclusion-**

The Sensitivity and Specificity, Accuracy, Precision and Recall score we got from test set are almost accurate.
We have high recall score than precision score which is a sign of good model.
In business terms, this model has an ability to adjust with the company's requirements in coming future.
This concludes that the model is in stable state.
Important features responsible for good conversion rate or the ones' which contributes more towards the probability of a lead getting converted are :

- Lead Origin Lead Add Form
- Total Time Spent on Website
- What is your current occupation Working Professional