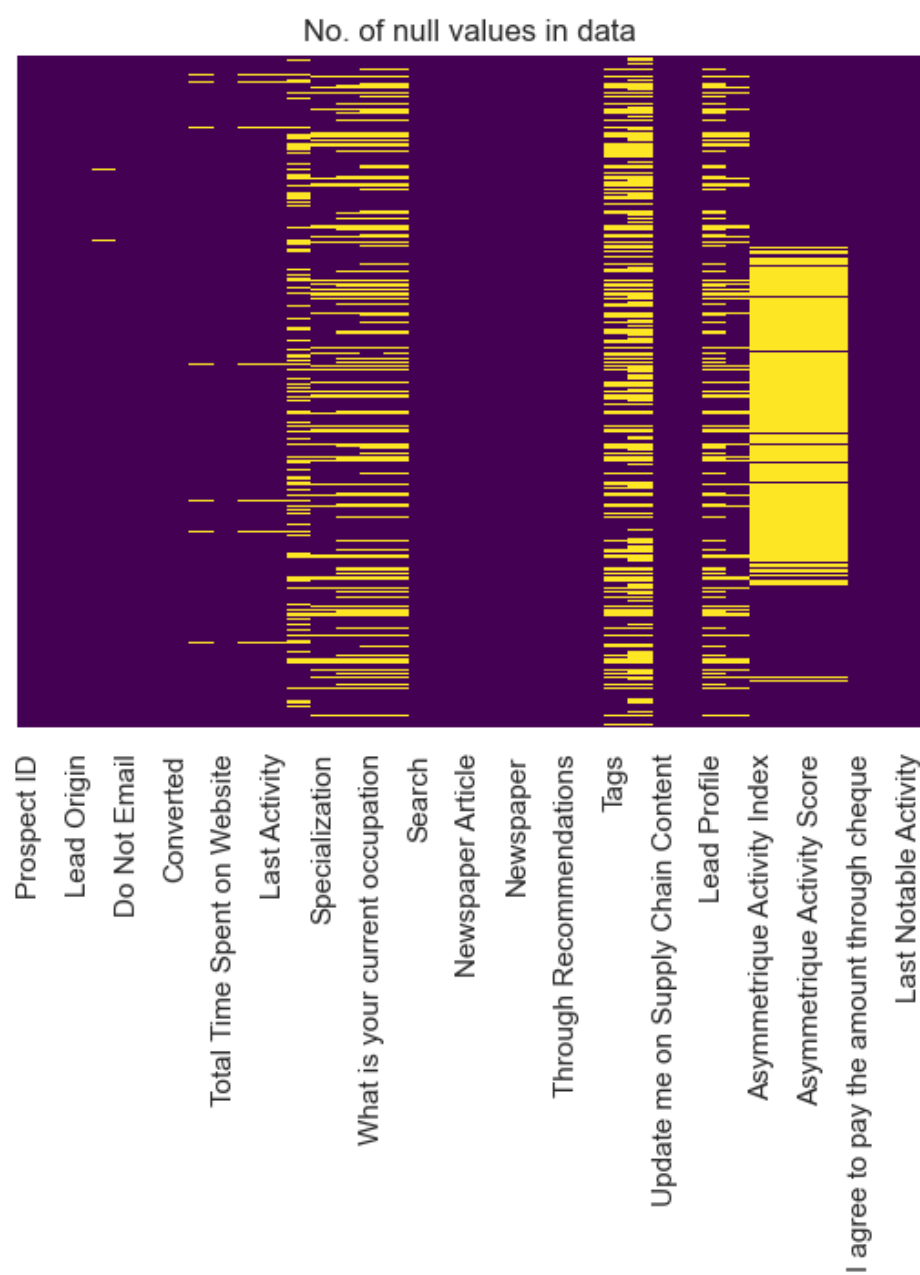# Leads Conversion Using Logistic Regression
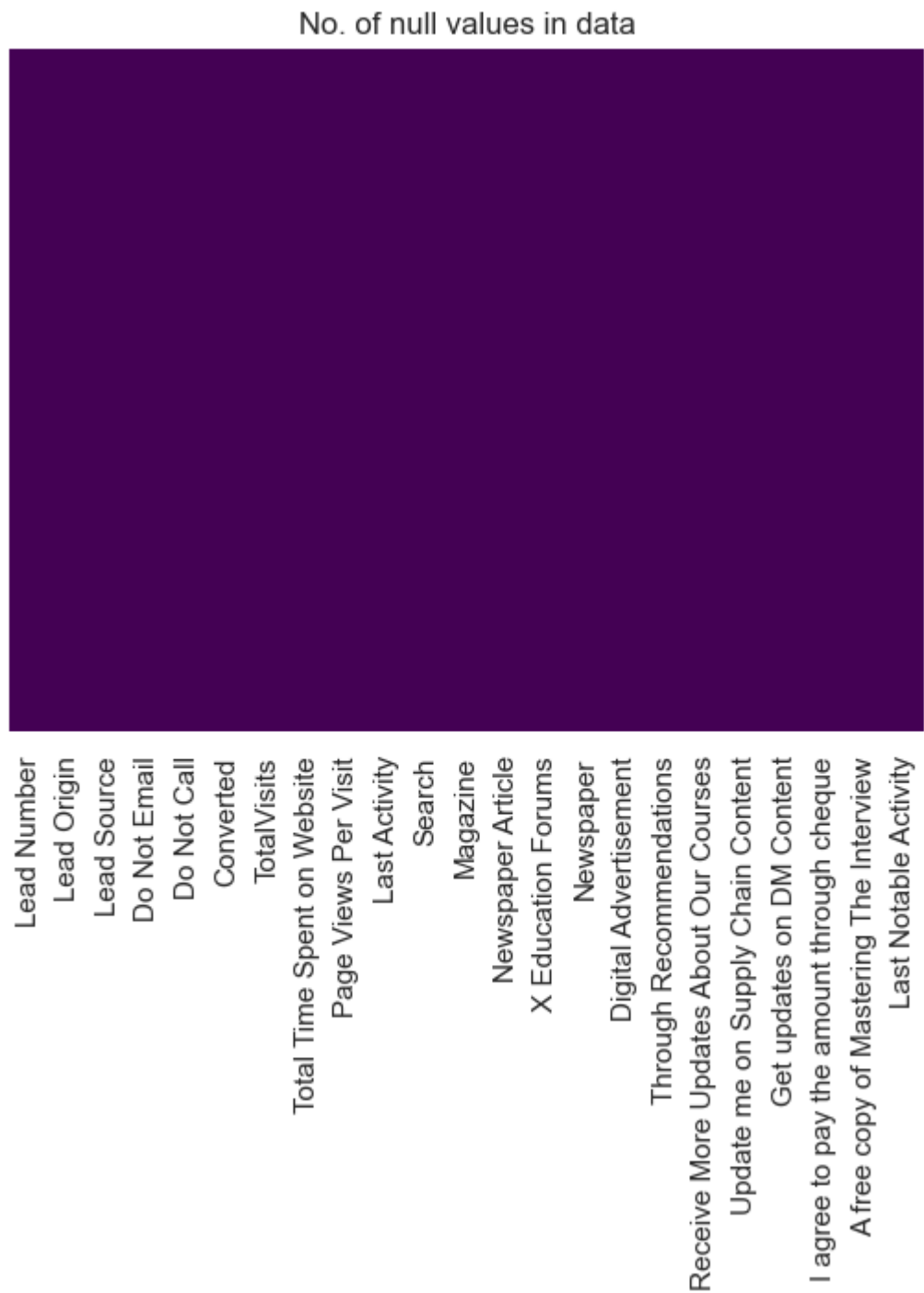
# Steps to Perform.

1. 2. Basic Information.
2. Handling Missing Values
3. EDA : Exploratory Data Analysis
4. Outlier Detection and Capping.
5. Features Selection Based on correlation and Select K method.
6. Model Selection and Training.

# Basic Data Manipulation.



No. of null values in data

Here Yellow lines show the missing values higher the yello lines, higher the missing values.

No. of null values in data
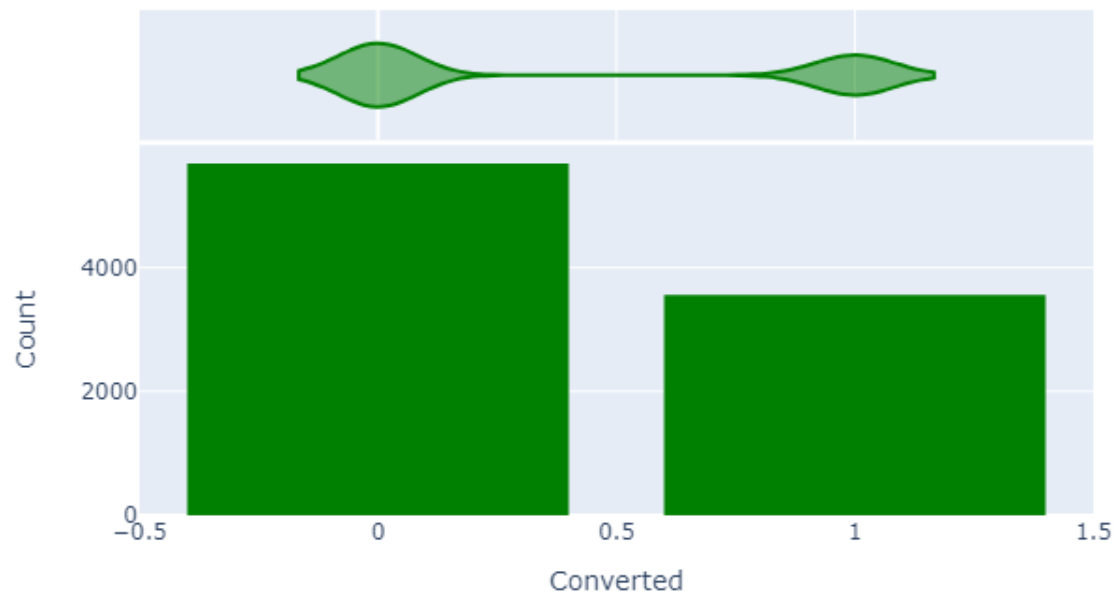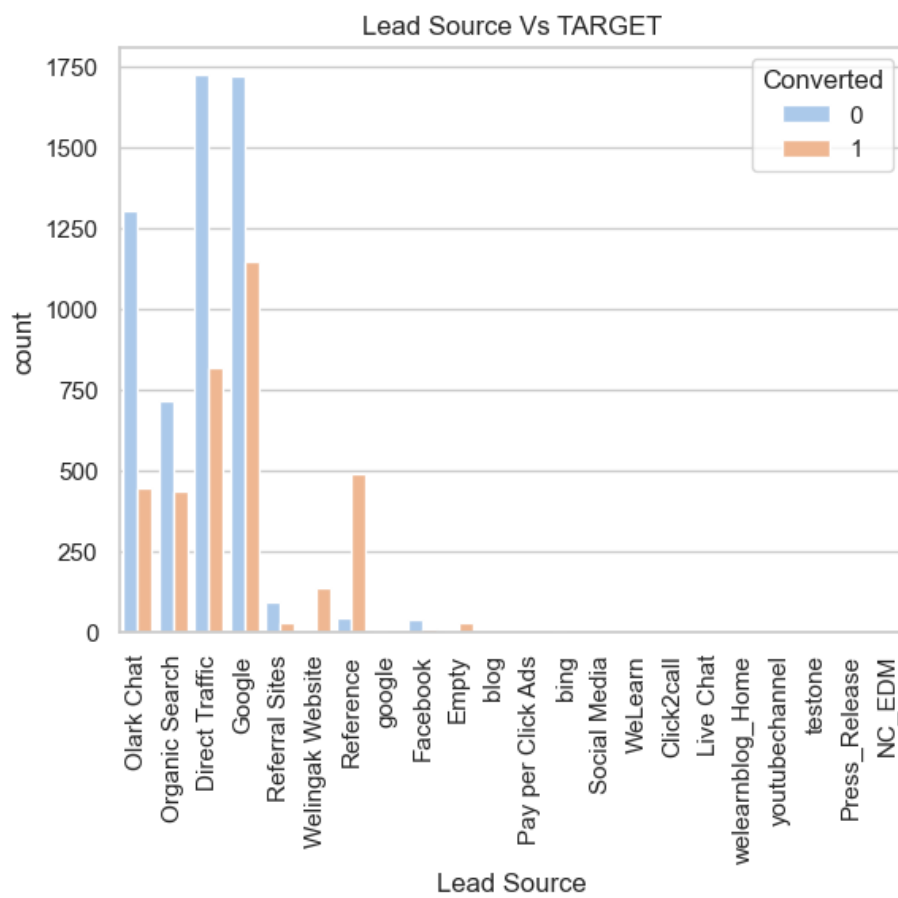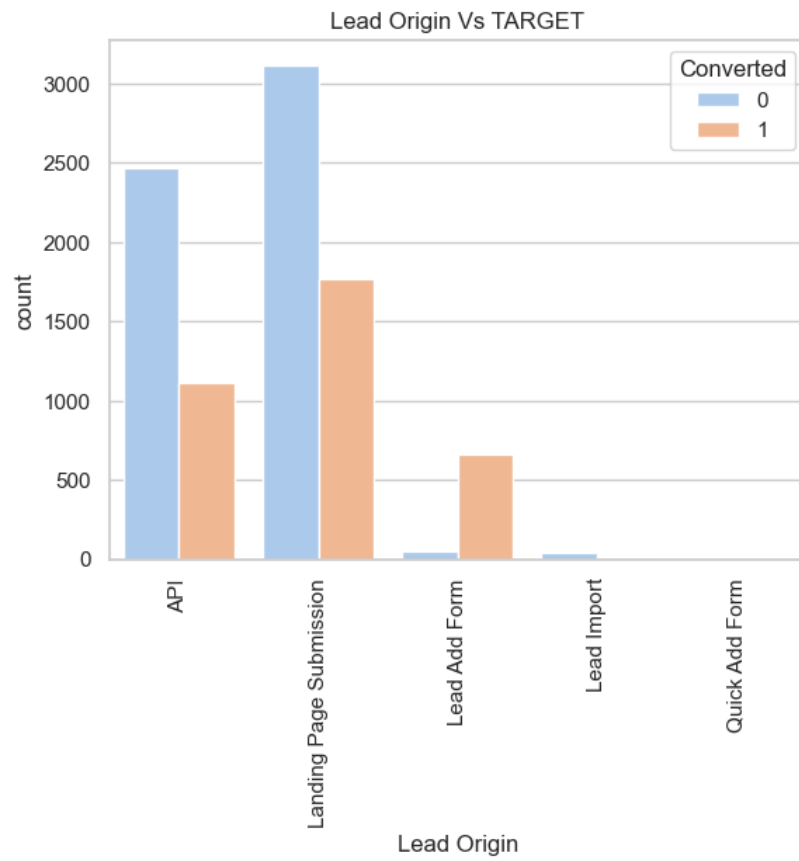


All the null null imputed. for categorical col we replace the values with new category and for numerical col we replaced with mean value.
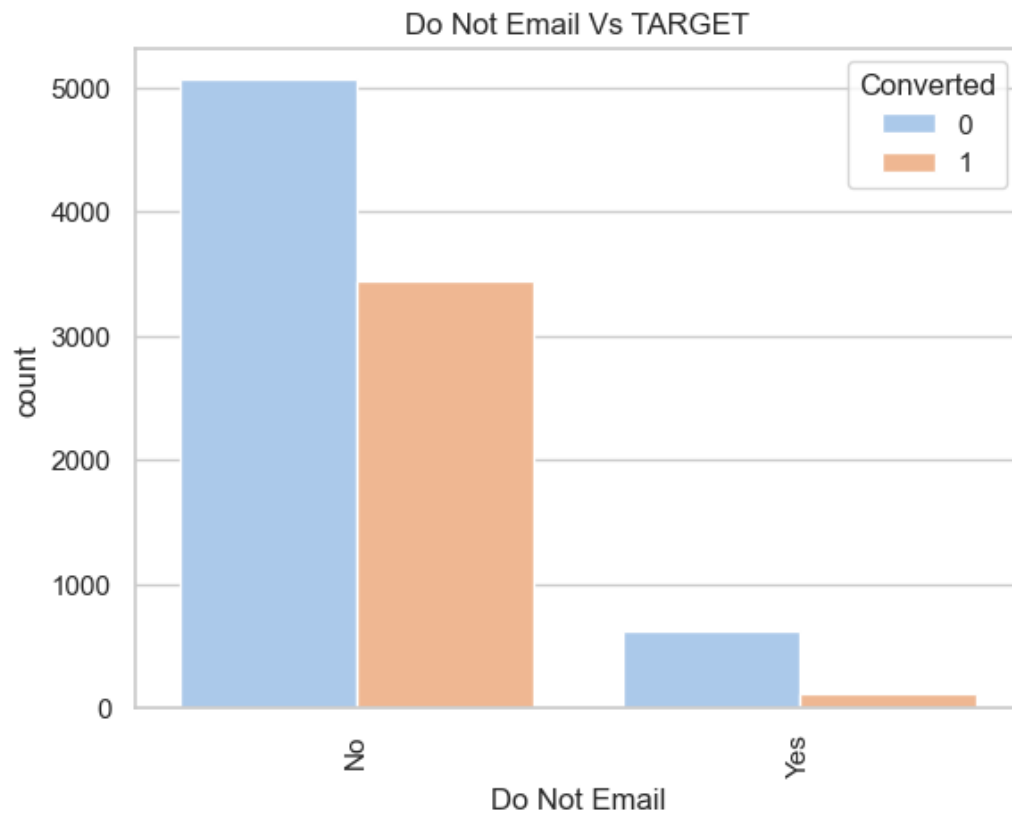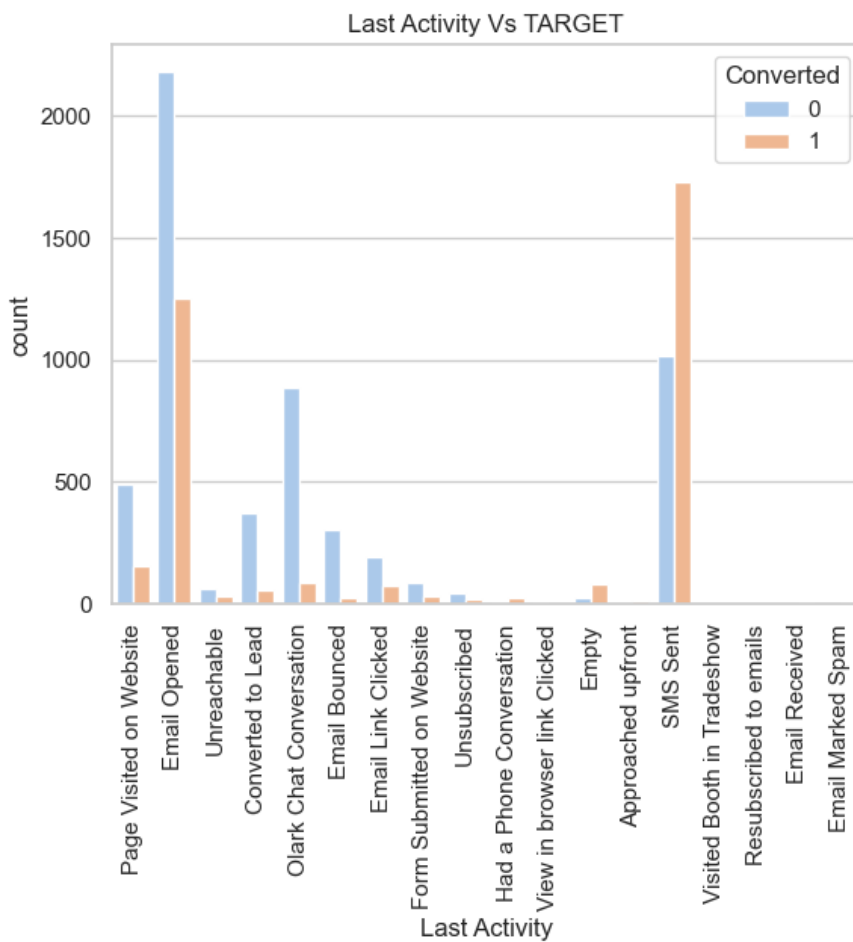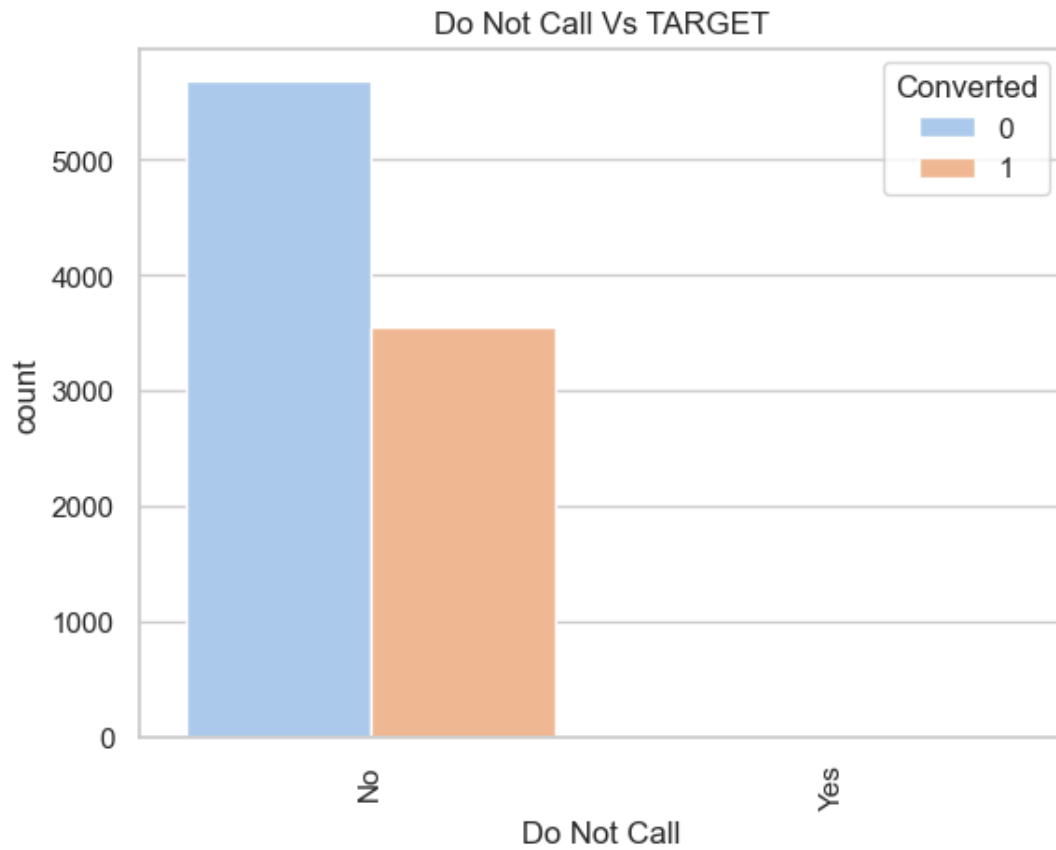
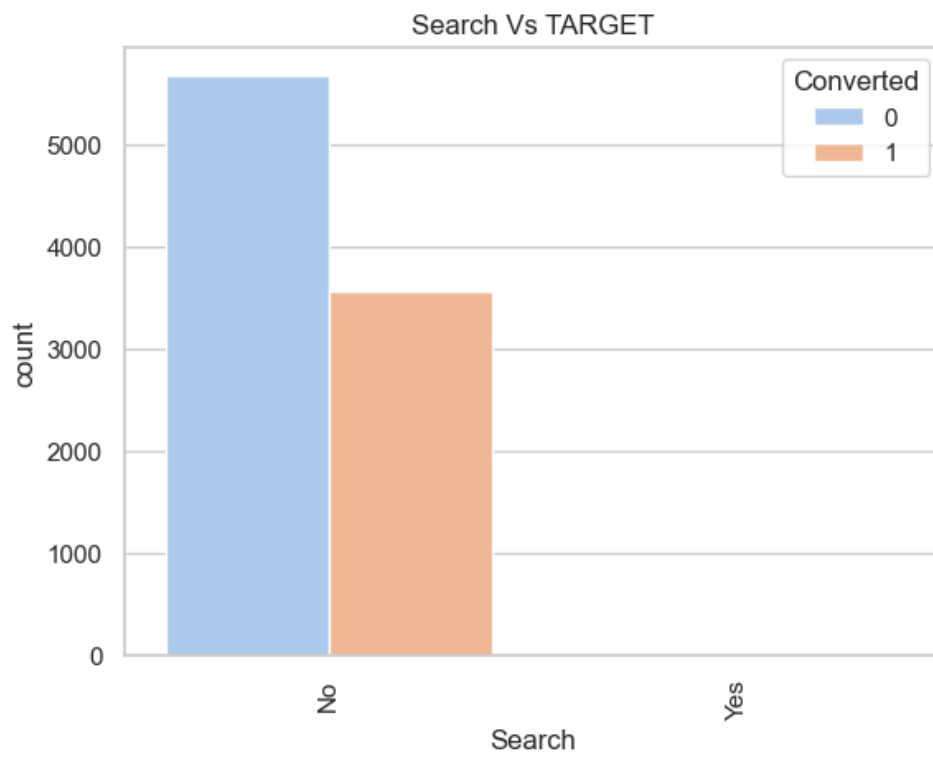Distribution of Converted



As we earlier predicted that Target column has higher number of 0 values then 1. We have more negative examples in our training then positive.

Lead Origin Vs TARGET



Lead Source Vs TARGET

Do Not Call Vs TARGET



Last Activity Vs TARGET

Search Vs TARGET

Magazine Vs TARGET



Newspaper Article Vs TARGET

## X Education Forums Vs TARGET



## Newspaper Vs TARGET

Digital Advertisement Vs TARGET


Through Recommendations Vs TARGET

Receive More Updates About Our Courses Vs TARGET



Update me on Supply Chain Content Vs TARGET

## Get updates on DM Content Vs TARGET



## I agree to pay the amount through cheque Vs TARGET

A free copy of Mastering The Interview Vs TARGET

Last Notable Activity Vs TARGET

## Observations

1. Here some categories have higher number of value then others in a single column.

2. higher number categories supports to not conversion and a smaller number of categories support conversion.

3. Like in Last Notable Activity if we sent the SMS then chance of conversion is higher.

4. Also if a source is references then there is higher chance of conversion of that person.

5. Olark chat conversion has higher negative impact on the conversion. So, we need the inchance the bot response.

## Lead Source Vs Converted



| Segment | Percentage |
|---------|-----------|
| WeLearn | 11.154751% |
| Social Media | 5.577376% |
| Referral Sites | 2.766378% |
| Reference | 10.235634% |
| Pay per Click Ads | 0.000000% |
| Welingak Website | 10.997643% |
| welearnblog_Home | 1.859125% |
| bing | 0.000000% |
| Organic Search | 4.214447% |
| Olark Chat | 2.847481% |
| NC_EDM | 11.154751% |
| Click2call | 8.366064% |
| Direct Traffic | 3.588119% |
| Empty | 8.985773% |
| Facebook | 2.636578% |
| Google | 4.461123% |
| Live Chat | 11.154751% |

## Last Activity Vs Converted



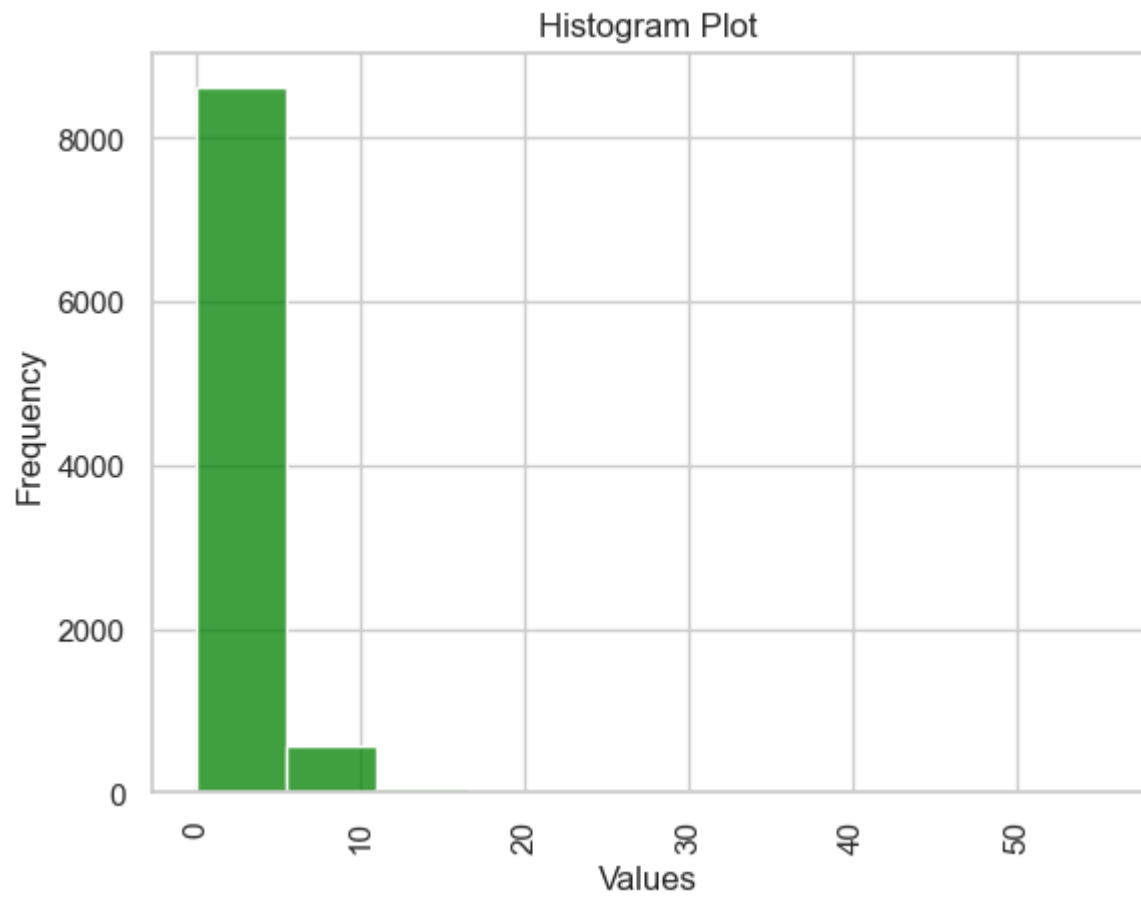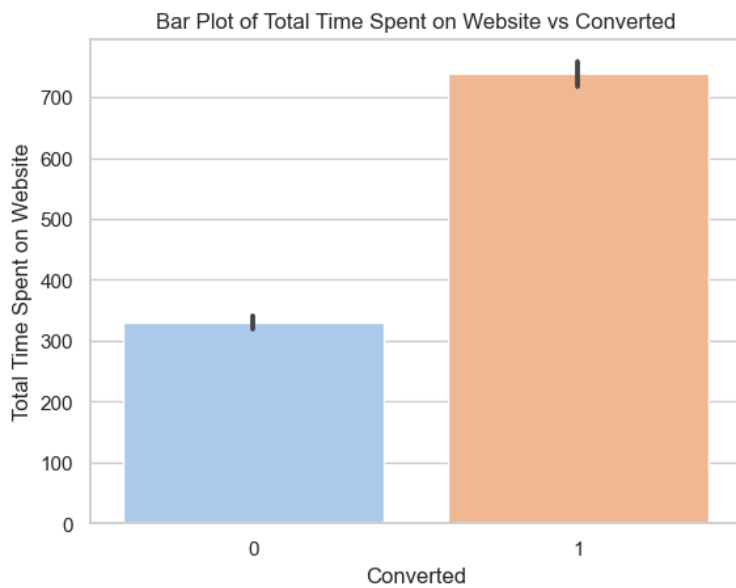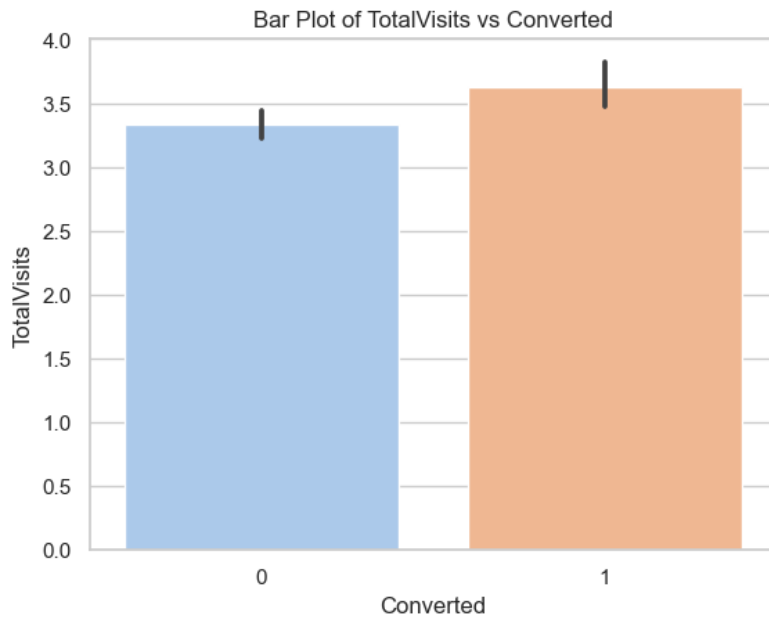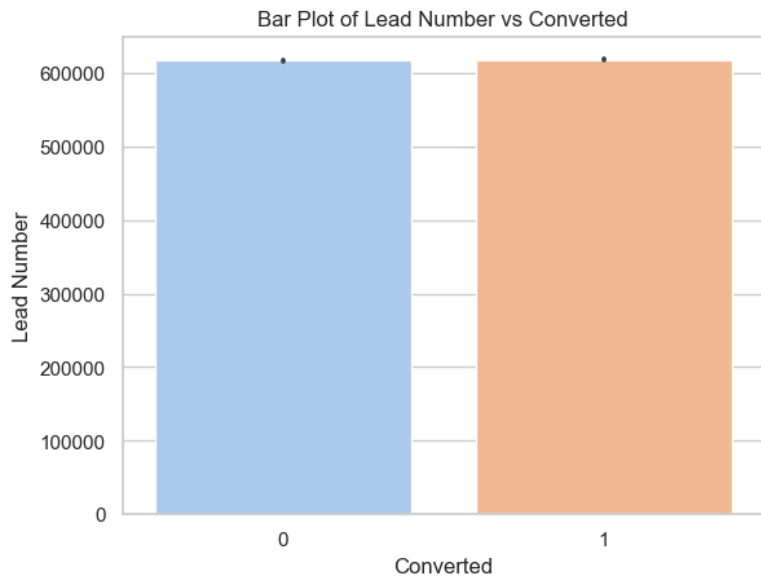| Segment | Percentage |
|---------|-----------|
| SMS Sent | 7.562988% |
| Resubscribed to emails | 12.021079% |
| Page Visited on Website | 2.836223% |
| Olark Chat Conversation | 1.037791% |
| Had a Phone Conversation | 8.815458% |
| Unreachable | 4.007027% |
| Unsubscribed | 3.153070% |
| View in browser link Clicked | 2.003513% |
| Visited Booth in Tradeshow | 0.000000% |
| Form Submitted on Website | 2.901640% |
| Empty | 9.453470% |
| Approached upfront | 12.021079% |
| Converted to Lead | 1.516678% |
| Email Bounced | 0.958736% |
| Email Link Clicked | 3.286662% |
| Email Marked Spam | 12.021079% |
| Email Opened | 4.382430% |
| Email Received | 12.021079% |

Last Notable Activity Vs Converted
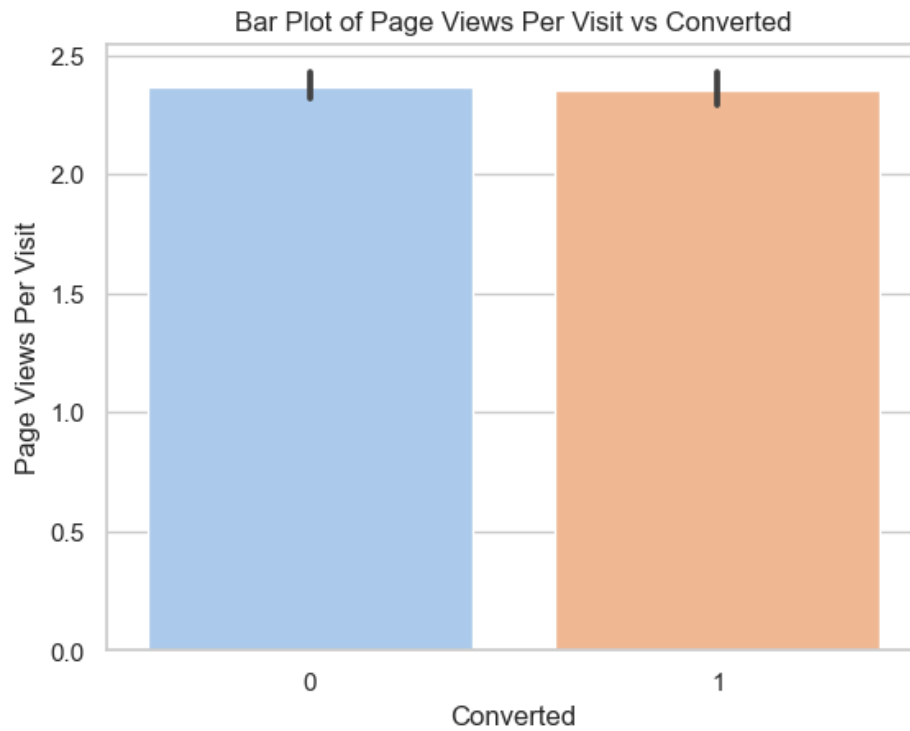
## Observation

1. As we can see that some the categorical columns contribute 100% to not conversion.

2. Also some of the cat cols like Do not email. Do not call has higher number of negative corelation to conversion.

3. Also references or new paper articles are the effective sources till now.

Histogram Plot



Histogram Plot



Histogram Plot

Histogram Plot

Distribution of these continuous numerical columns are not normalized these are left skewed. Also Some of them have outliers because of these we can see the bar at one side and one single value on other side.
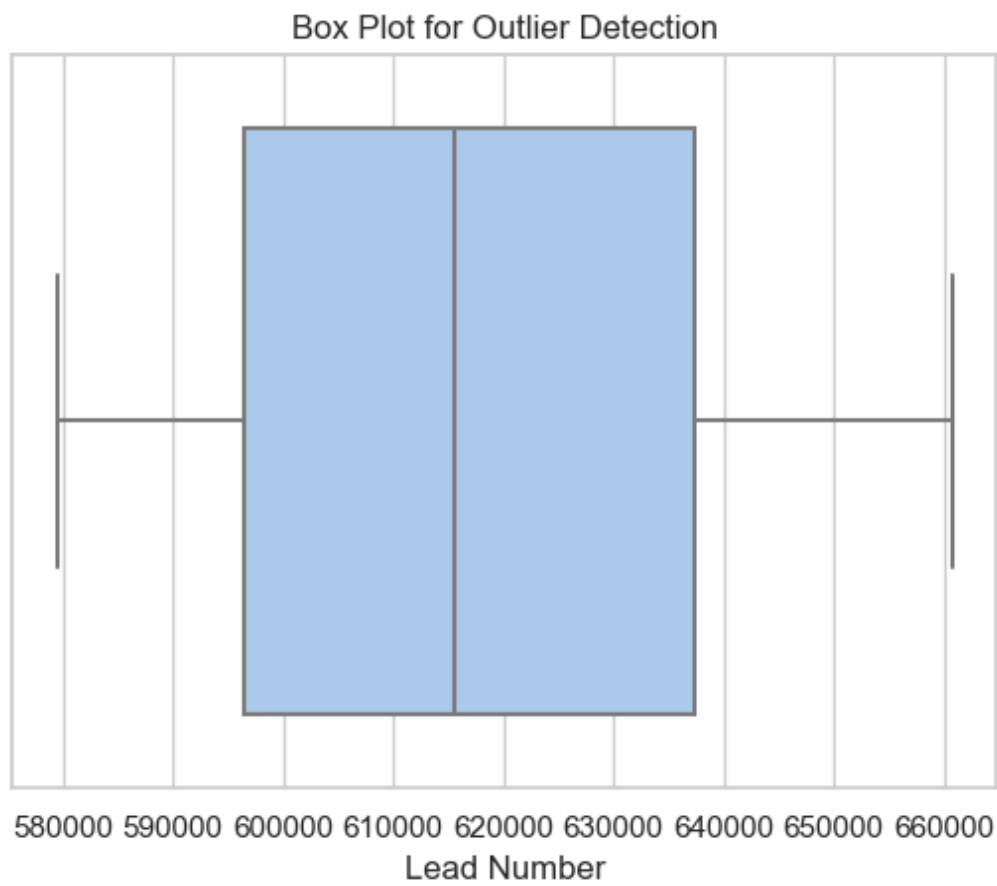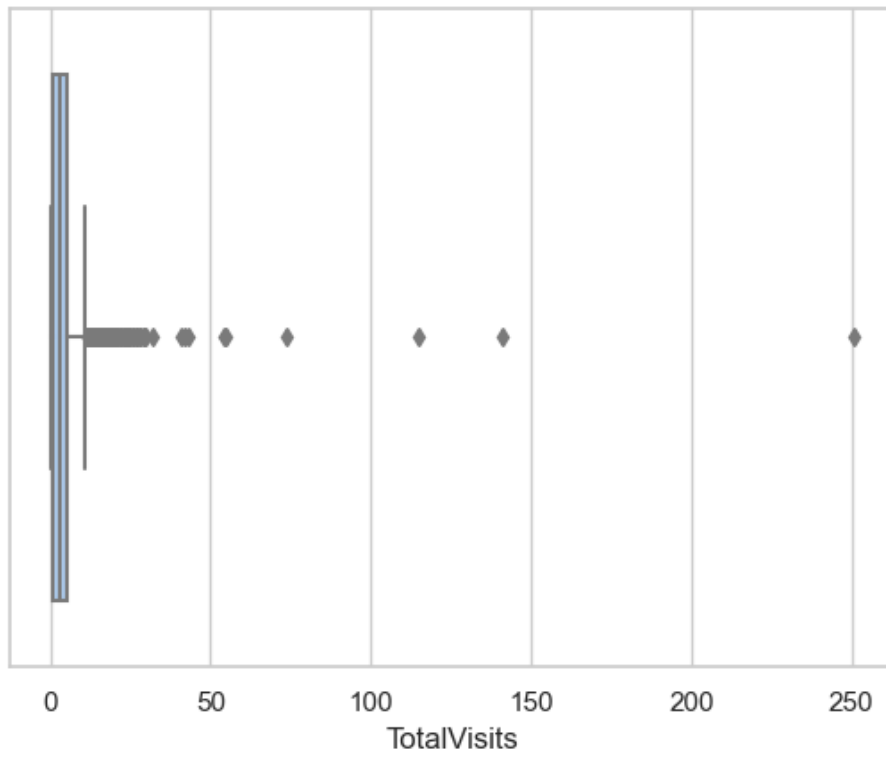
## Bar Plot of Lead Number vs Converted



## Bar Plot of TotalVisits vs Converted



## Bar Plot of Total Time Spent on Website vs Converted

Bar Plot of Page Views Per Visit vs Converted

## Observations

1. Here we can see that time spent and TotalVisit on website has higher potitive impact on conversion

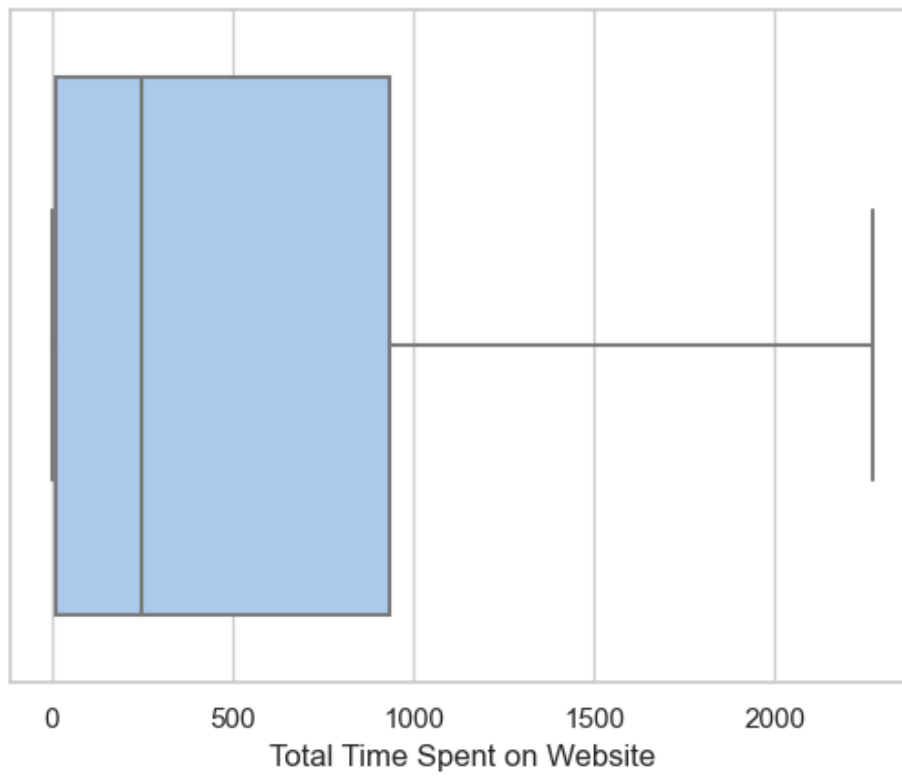2. And page views has netural impact on the conversion.
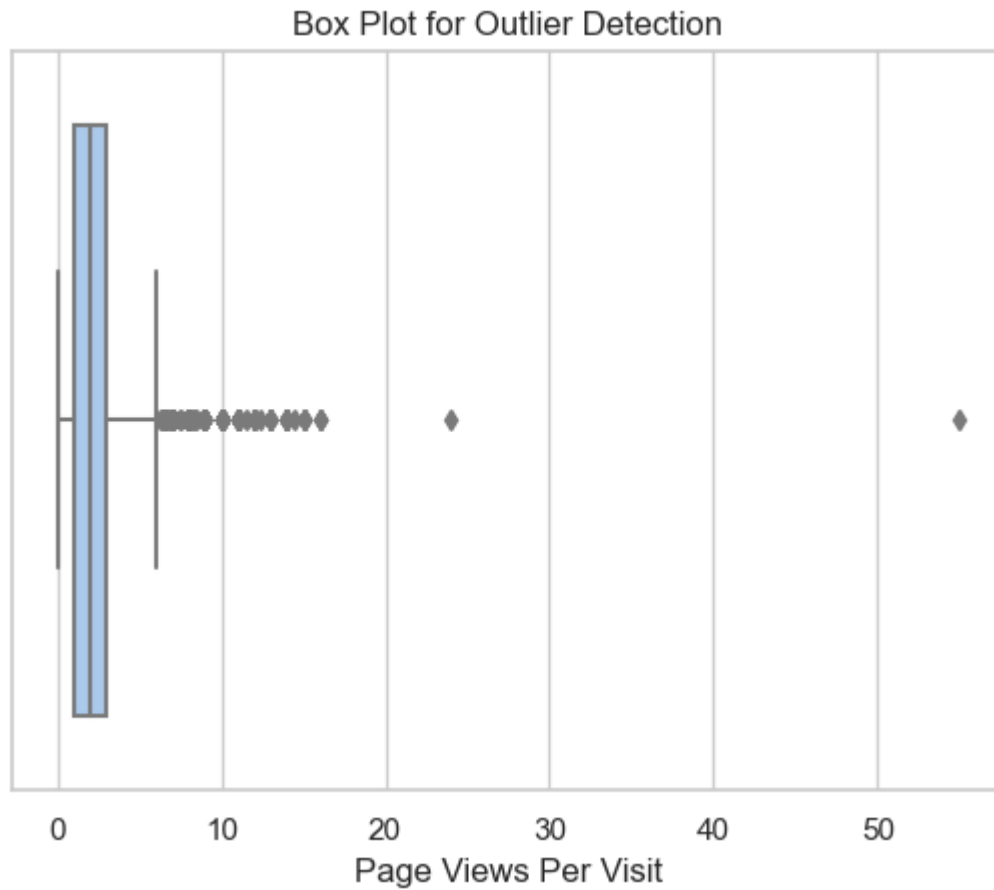
# Outlier detection



Box Plot for Outlier Detection

Box Plot for Outlier Detection

TotalVisits



Box Plot for Outlier Detection

Total Time Spent on Website
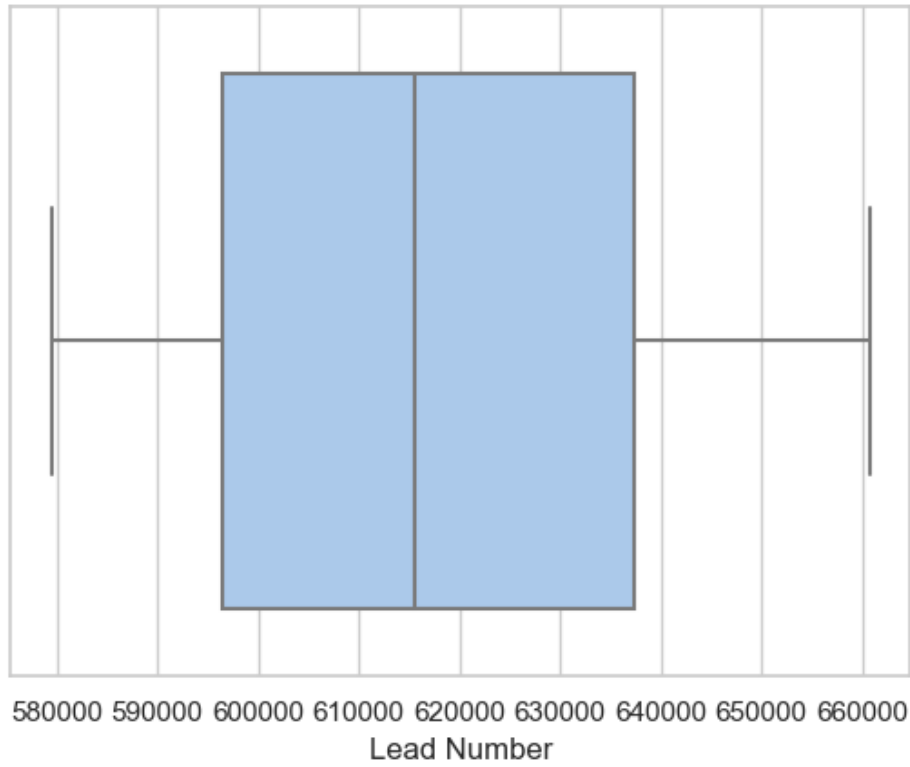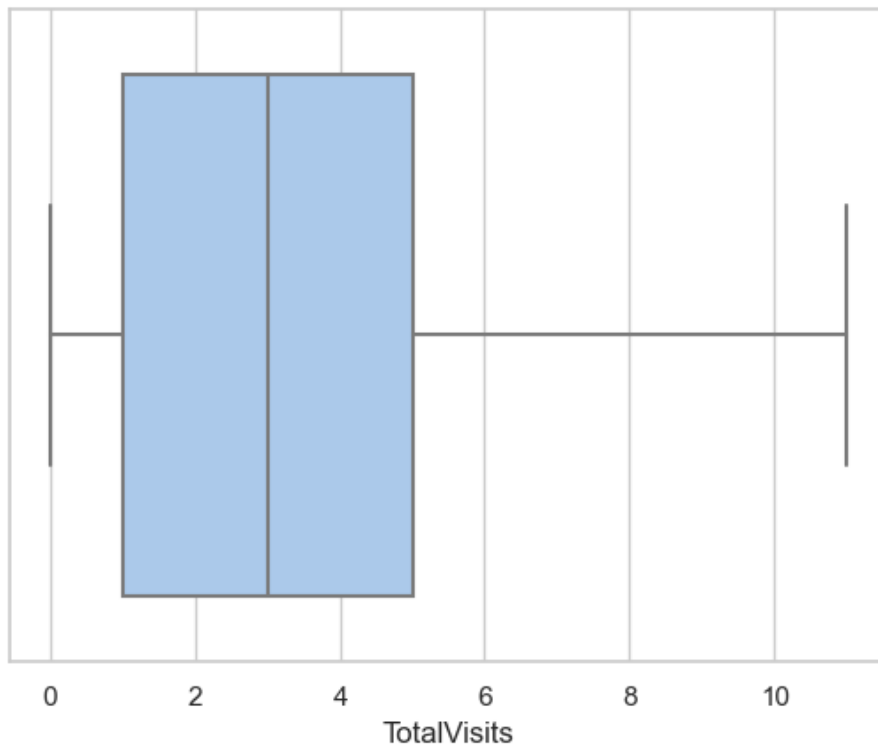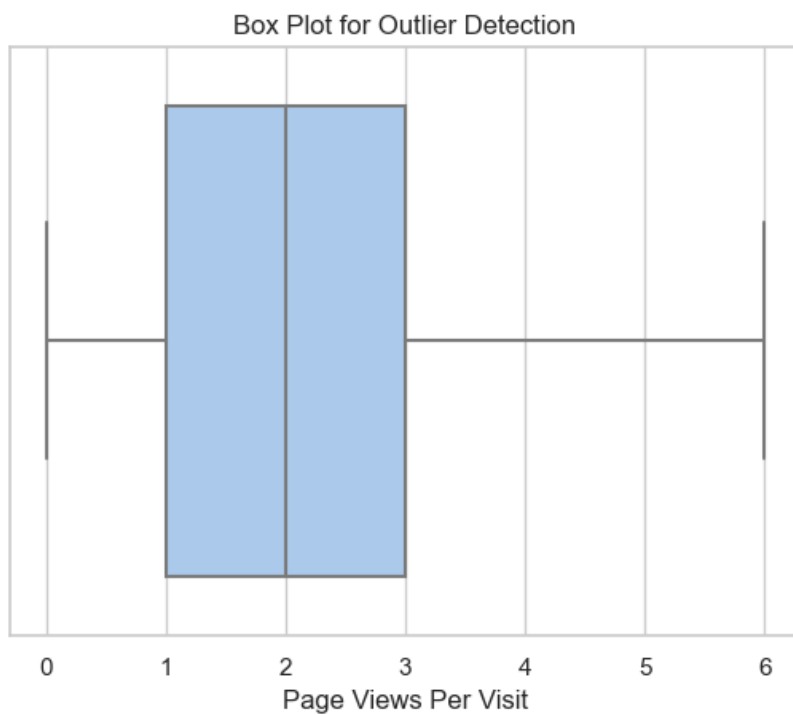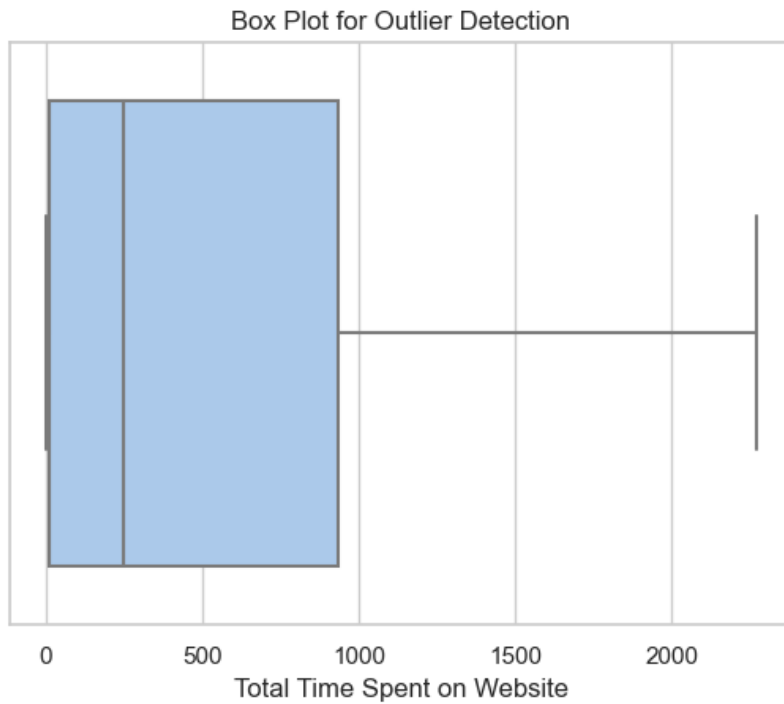
Box Plot for Outlier Detection

## Observations

1. In total visits and page views per visit has outliers.

2. Removing outliers does not solve the problem we need to cap the outlier that means if some day total visit per page goes higher than a threshold then we will cap this with thresh value for better results from the model.
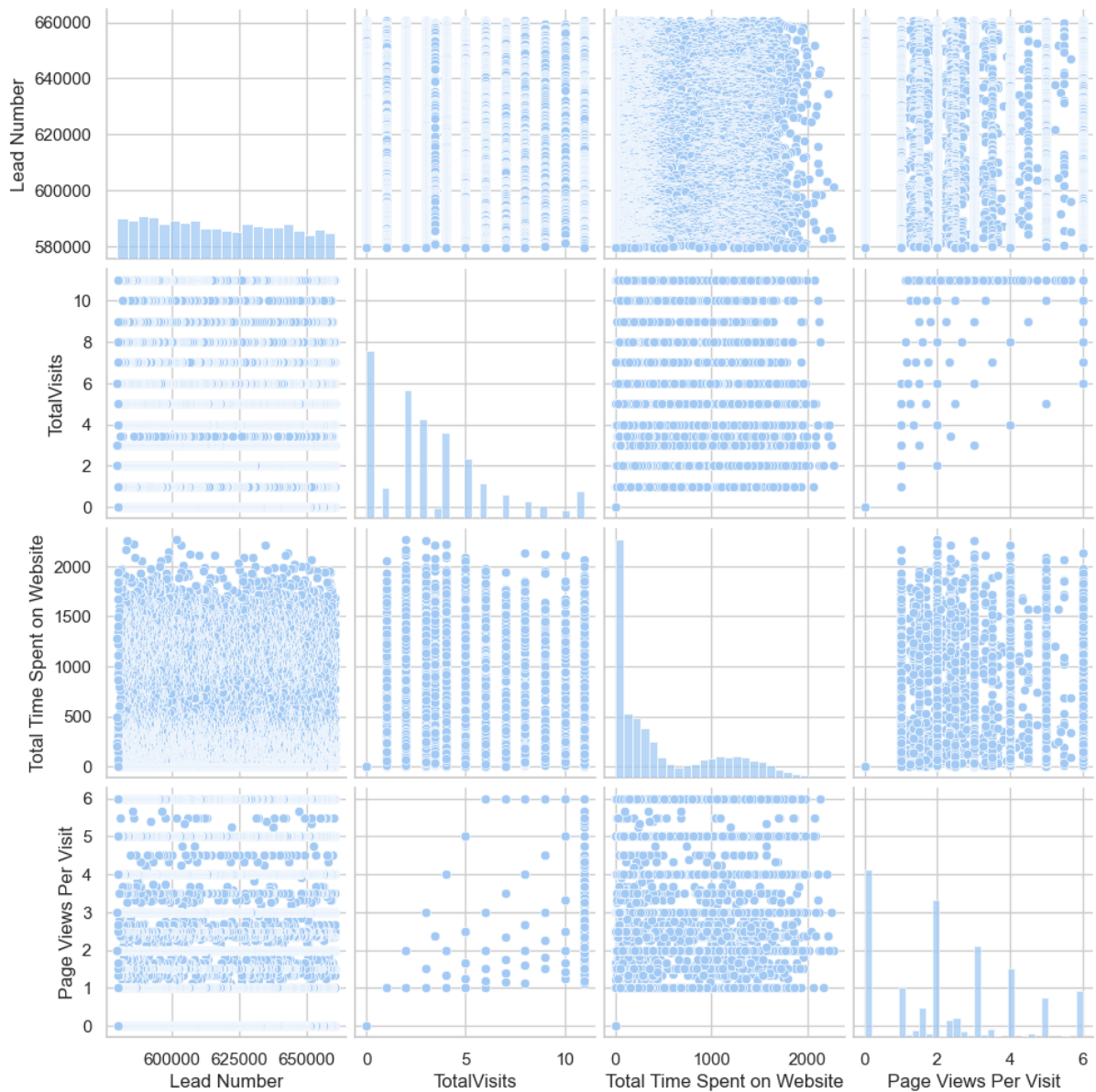
Box Plot for Outlier Detection

Lead Number

Box Plot for Outlier Detection

TotalVisits

Box Plot for Outlier Detection

Total Time Spent on Website



Box Plot for Outlier Detection
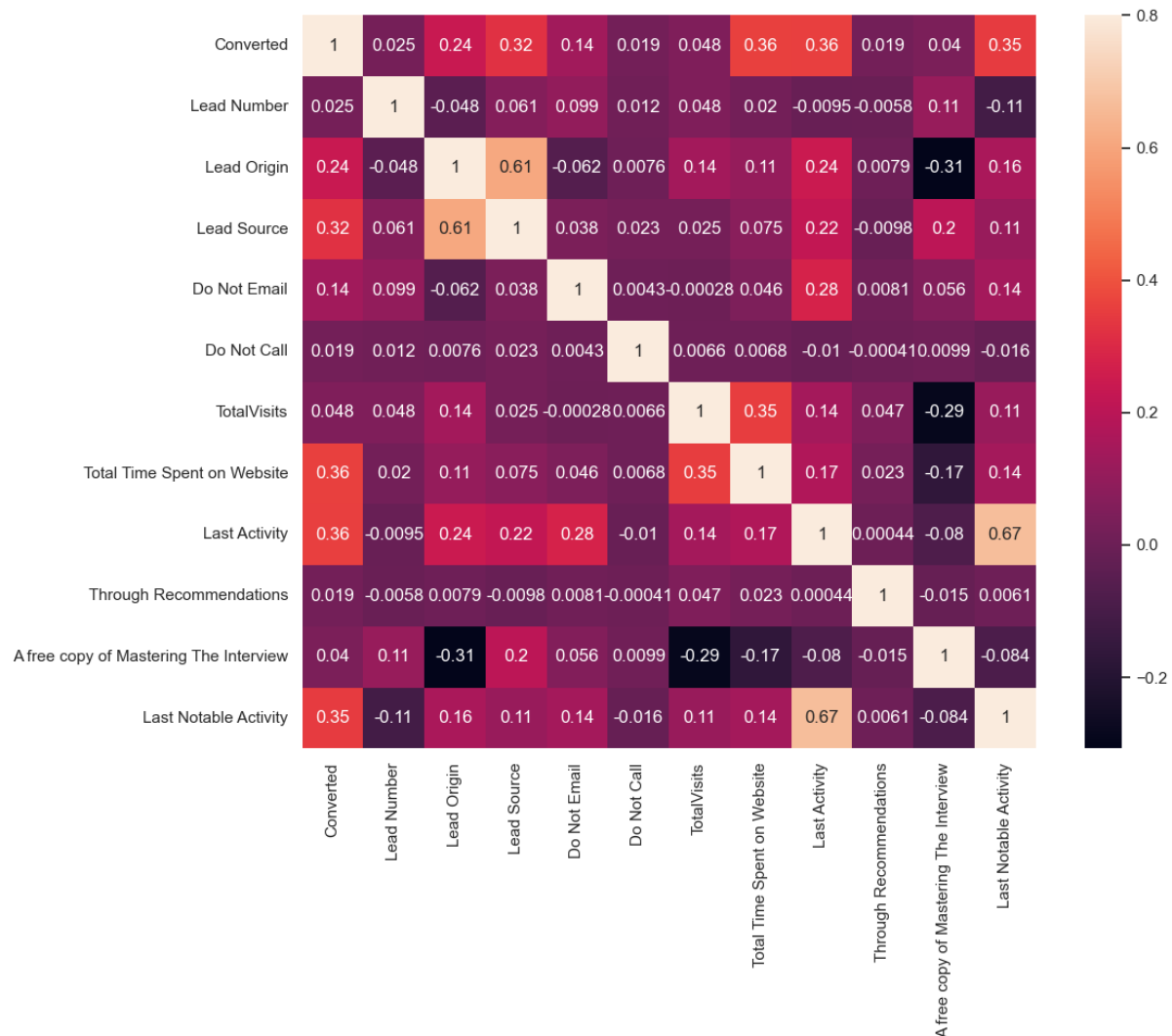
Page Views Per Visit

## Observations

Now there is not outliers, we cap the outliers using z-index method. where we cap the outlier between IQR range.

## Observation

In the pair plot we can see the scatter plots between multiple independent numerical variables that follows some trend or has some correlation between them.
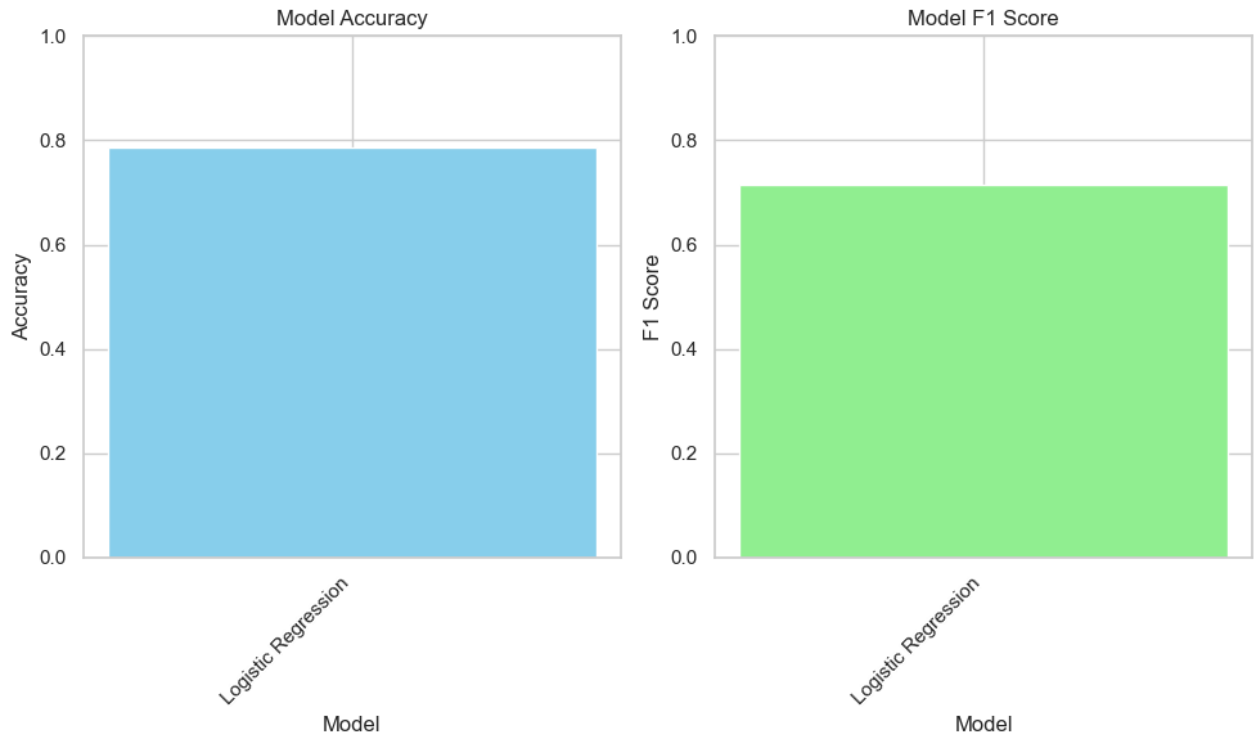
# Features Selections or Features Engineering.



## Observations

Heat maps showcase the same that we have earlier discussed.

# Modelling



## Observations

1. Model achived accuracy of 72% and it successfully classify 1445 in positive and 730 in negative side.

2. F-Score is 71 which is good.

3. In the next step we will implemnt hyper parameter tunning for further increasing the accuracy of the model