# Lead Scoring Case Study Using Logistic Regression

Submitted by:

1.Soma Prasad

2.Tarun Shivkumar Ramanathan

3.Satweek Bandi

# Contents:-

☐ Statement of the Problem

☐ Objective of the Business

☐ Solution Methodology

☐ EDA

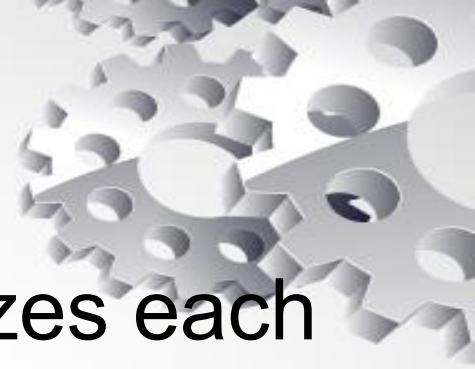☐ Building of the Model

☐ Evaluation of Model

☐ Conclusion

# Statement Of Problem:-

- X Education is selling the online courses to the industry professionals.
- XEducation is having ample amount of leads.those leads conversion rates are lower.If we consider an example,they had 100 leads,only 30 of them joins the organisation.
- This process can be made optimistic is surfing for the most potential leads,called "Hot Leads"
- If everything works great then those set of leads are more into the conversion rate,through which the sales leads will follow up the "Hot Leads".

# Objective of the Business :-

- The Lead X wants us to build a model that analyzes each lead and assigns them a score between 0 and 100. The higher the score, the "hotter" the lead, meaning they're more interested and ready to buy.

- CEO wish us to achieve a lead conversion rate of 80%.

- They wanted a Model to adapt to peak times, optimize manpower, and outline post-target strategies for effective implementation and sustained success.
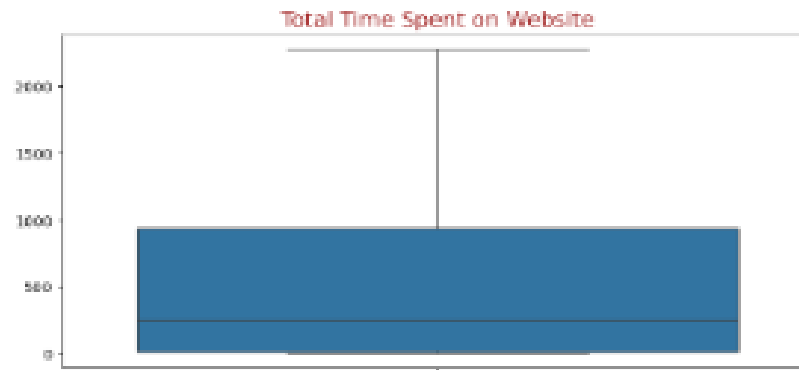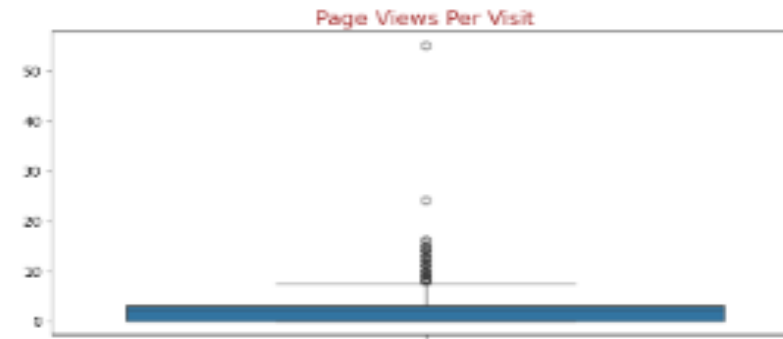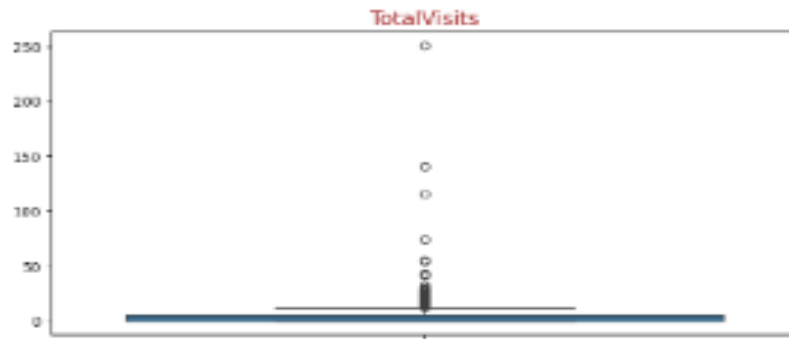
# Solution Methodology:-

- o Data cleaning and data Inspection.

- o EDA

- o Dummy Variable Creation

- o Test-Train split

- o Feature scaling

- o Dropping highly correlated dummy variables

- o Model Building (RFE Rsquared VIF and pvalues)

- o Model Evaluation

- o Checking Accuracy

- o Finding Optimal Cutoff Point

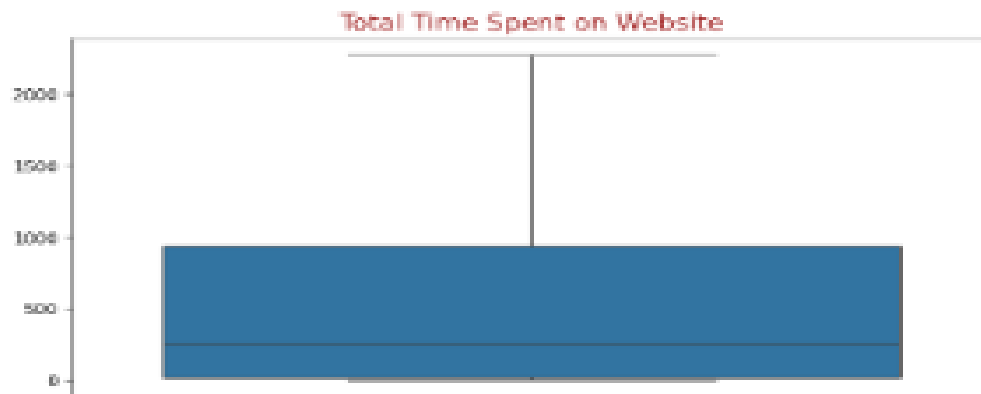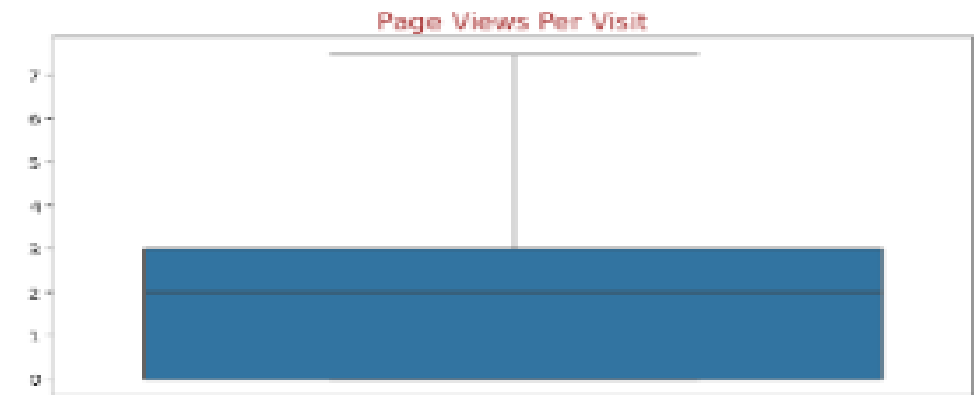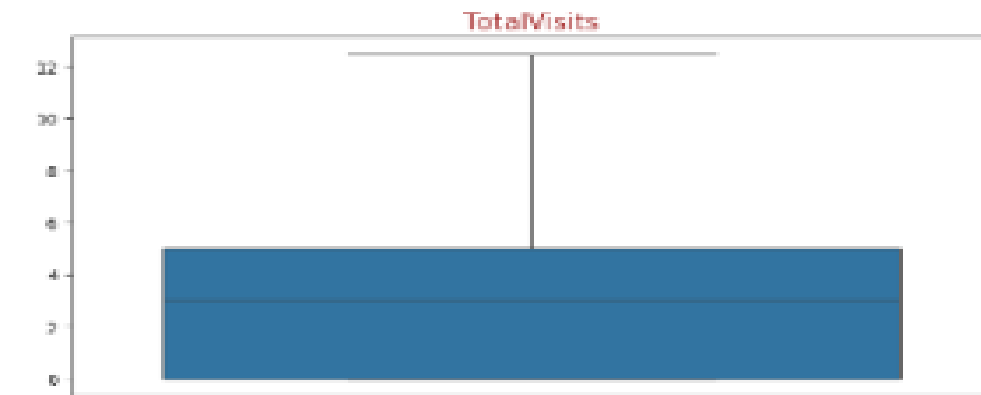- o Making predictions on test set

# Outlier Analysis

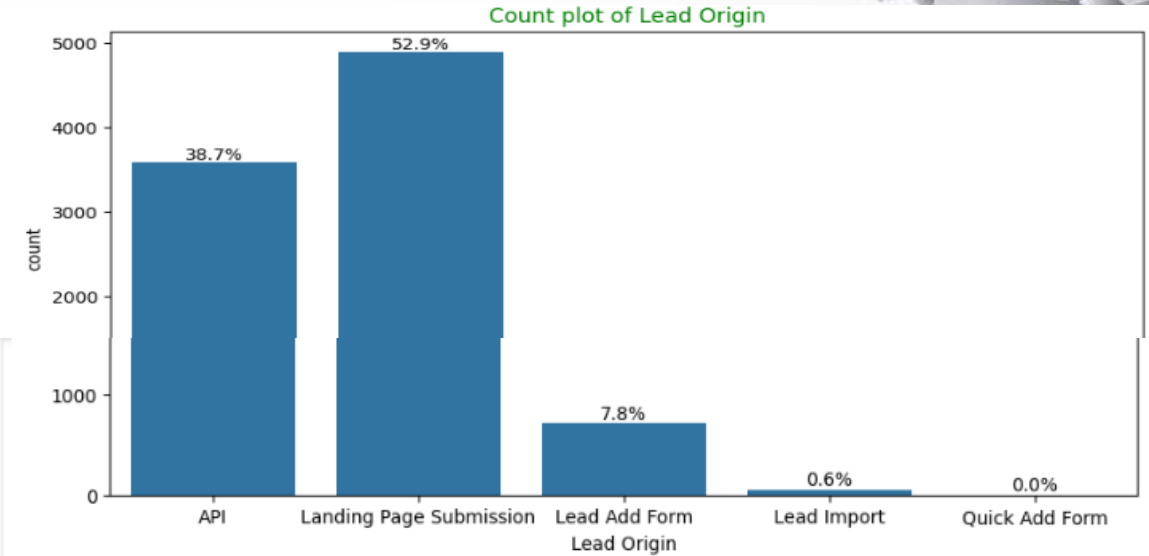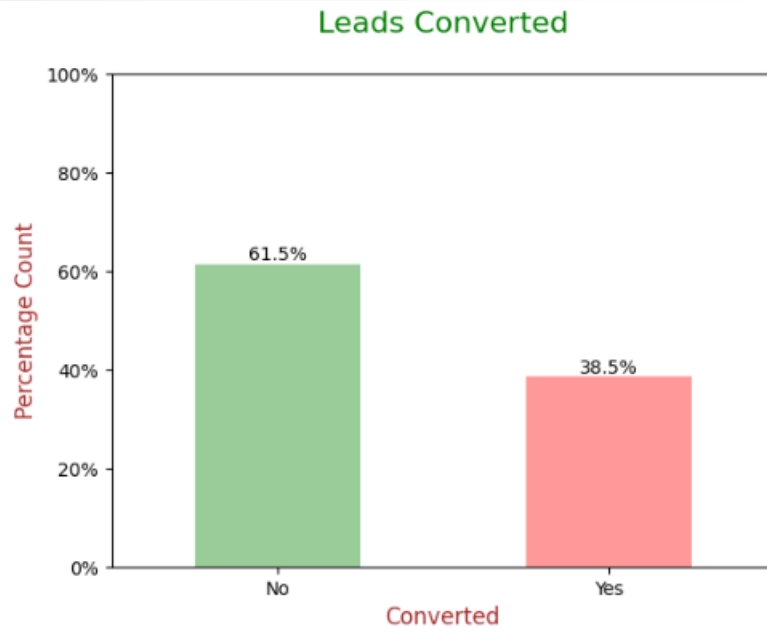# Checking Outliers with Boxplot :-

# Exploratory Data Analysis

# Univariate Analysis

- Most people find us through Google or directly visiting our site, making these the primary lead sources.

- The majority of site visitors are unemployed individuals, highlighting a potential demographic focus for outreach and support.

- Email opening marks the final step for most leads, emphasizing the importance of email communication in our engagement strategy.

- Most leads originate from submitting landing pages, indicating the effectiveness of our landing page design in capturing interest.

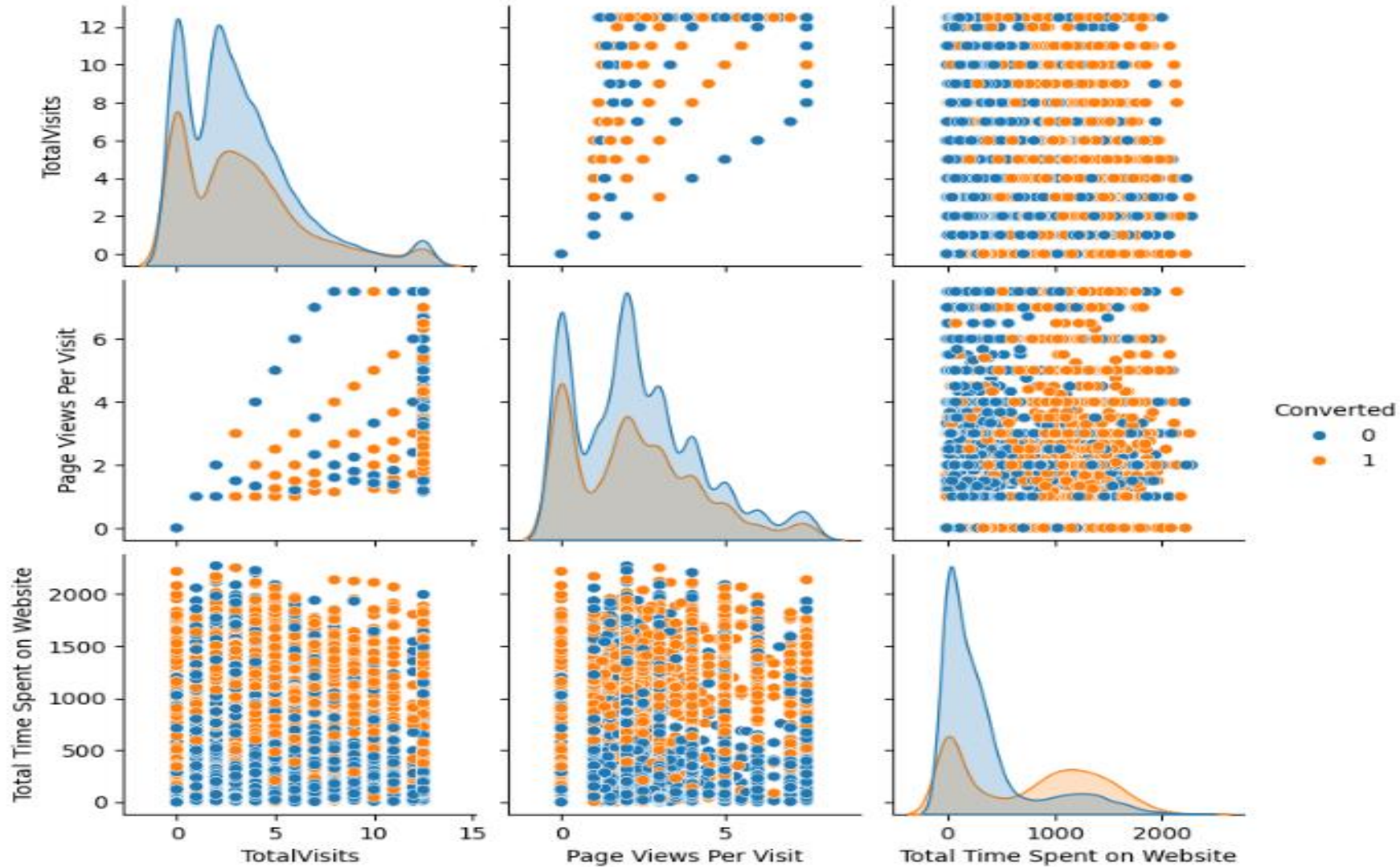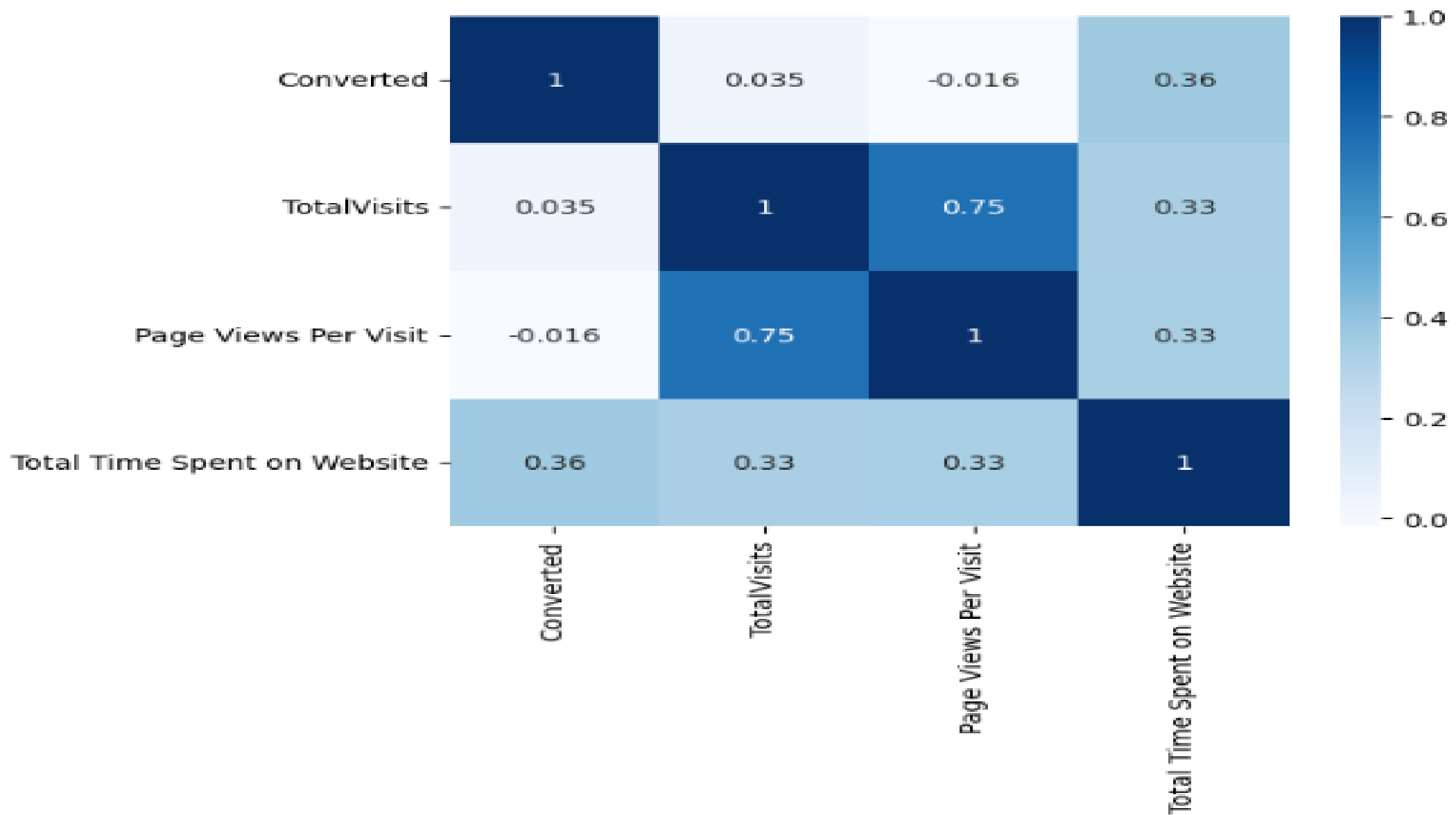# Leads Converted ,Univariate Analysis for Categorical---Variables:-

# Bivariate Analysis

- "Lead forms & landing pages yield most conversions, prioritize Google/Direct Traffic sources."

- "Email/SMS activity drives highest conversions, prioritize outreach to these leads."

- "India leads in conversions, focus on finance, marketing, HR. Professionals prioritize over housewives."

- "Career prospects drive course opt-ins, prioritize clients with this motivation."

- "Tag 'revert after email' has highest conversion rate among others."

- "Target Mumbai for lead conversion first, then Tier II cities."

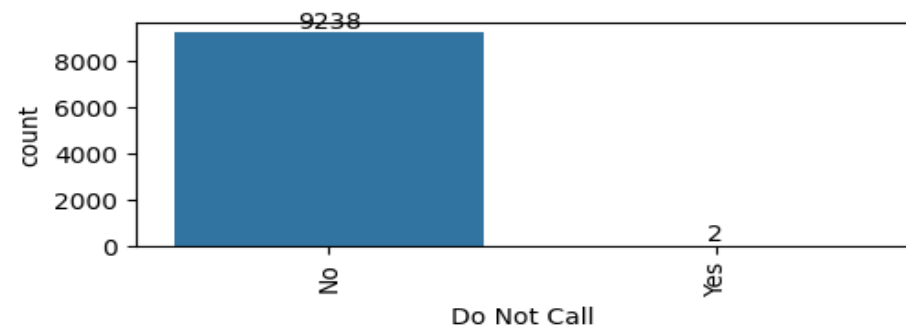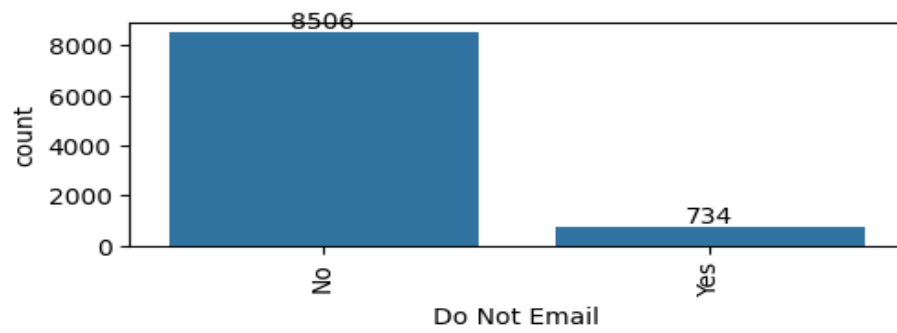- "SMS leads in conversions, followed by email opens."
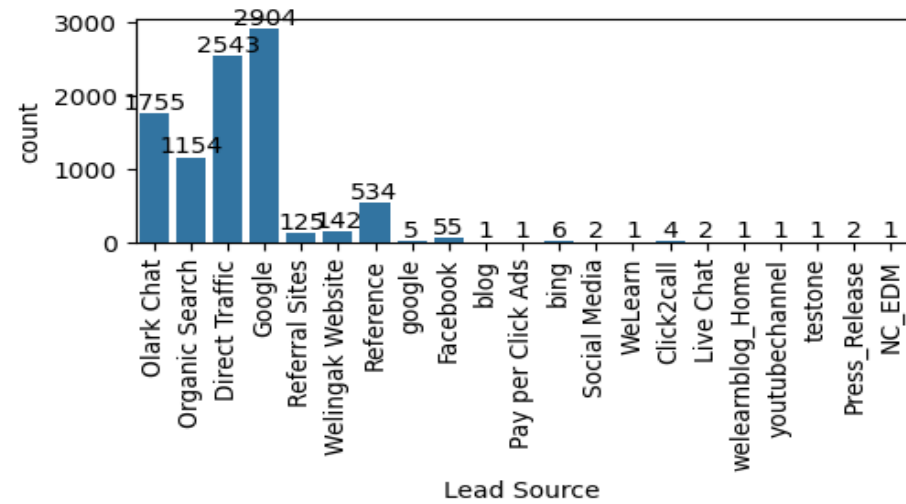
# EDA BIVARIATE ANALYSIS:-

# HEATMAP :-

# Boxplot with Converted as hue

# MODEL BUILDING :-

- "Split data into training and testing sets."
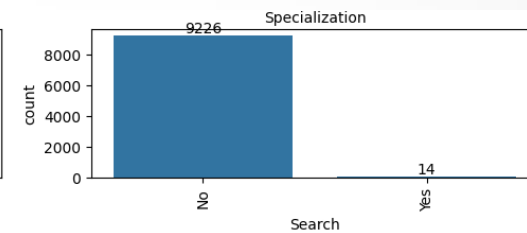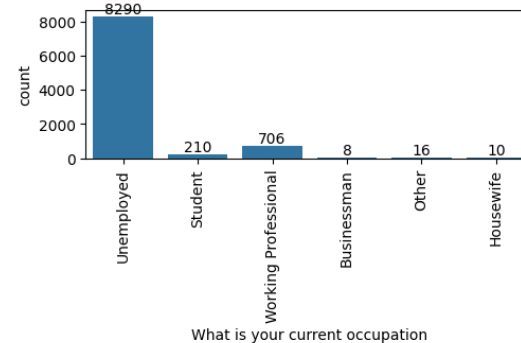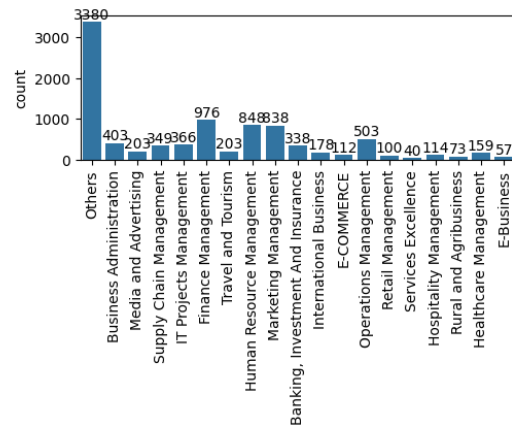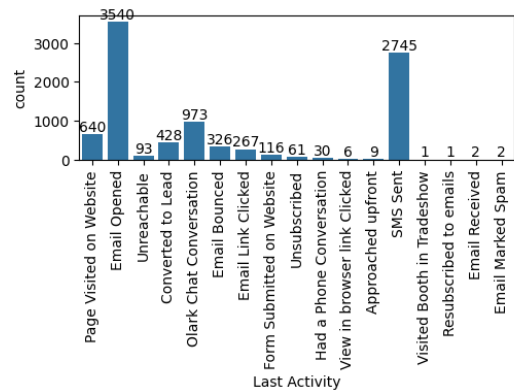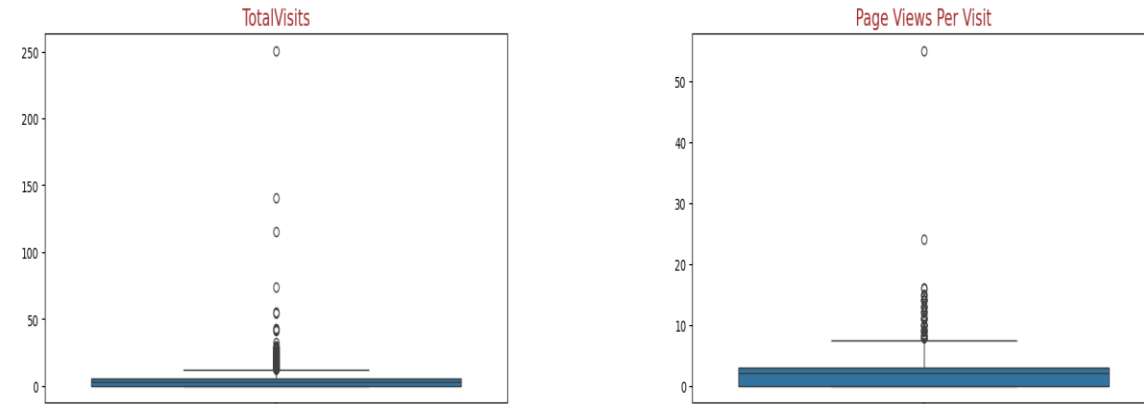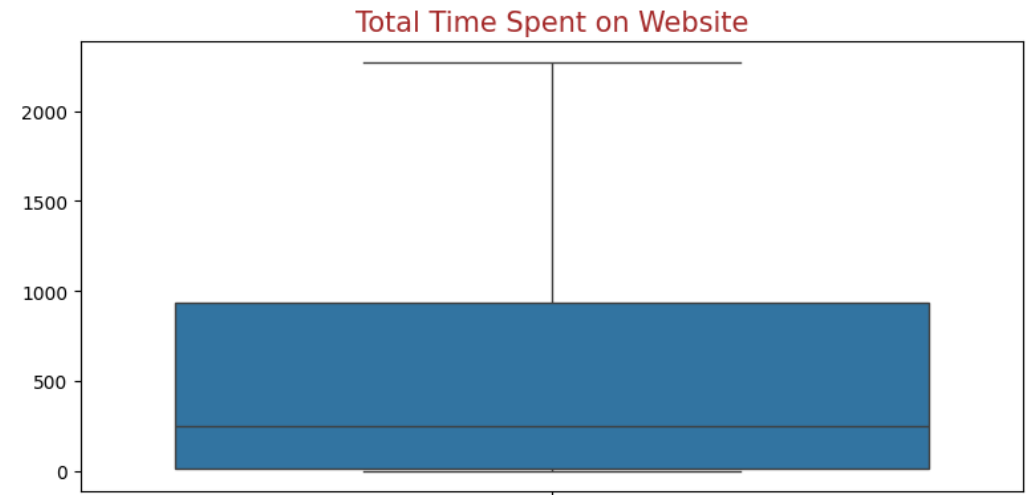- "Scale variables in the training set for consistency."
- "Construct initial model."
- "Use Recursive Feature Elimination (RFE) to remove less important variables."
- "Refine model by eliminating variables with high p-values."
- "Assess multicollinearity with Variance Inflation Factor (VIF)."
- "Make predictions using the training set."
- "Evaluate model accuracy and other metrics."
- "Apply model to predict outcomes using the test set."
- "Analyze precision and recall of test predicting future insights."

# MODEL EVALUATION

Receiver operating characteristic example

(0.345 , 0.8)

ROC curve (area = 0.88)

Receiver operating characteristic example

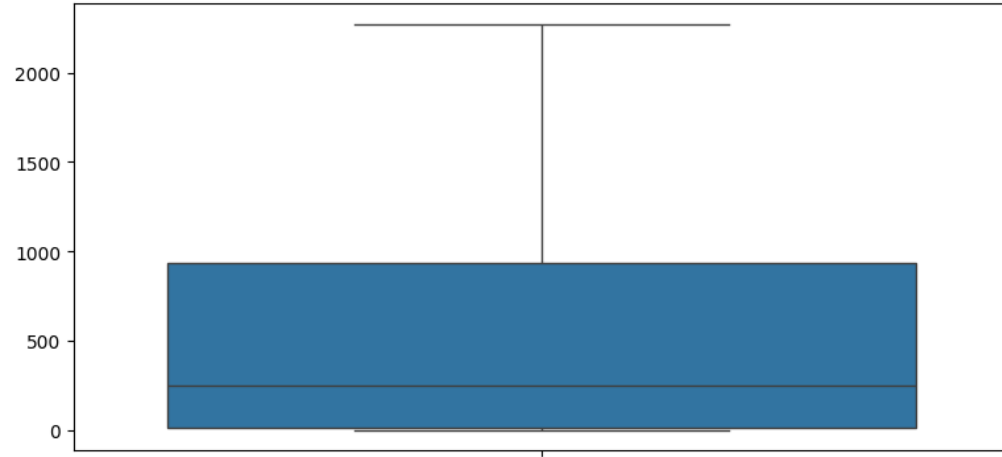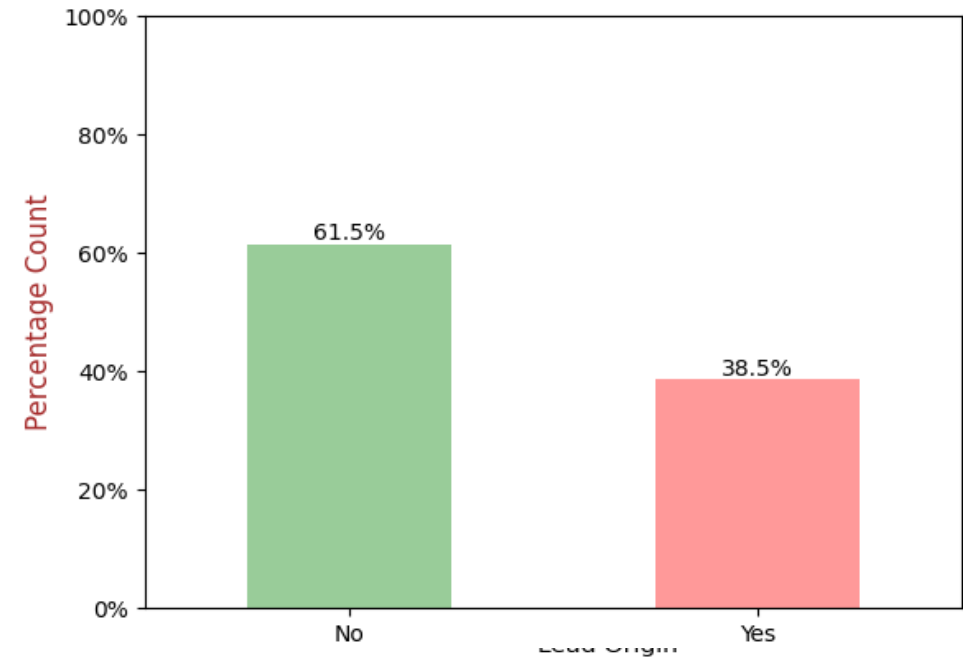```
[ ]  # features and their coefficient from final model
     parameters=logm4.params.sort_values(ascending=False)
     parameters

     Lead Source_Welingak Website                5.388662
     Lead Source_Reference                       2.925326
     Current_occupation_Working Professional     2.669665
     Last Activity_SMS Sent                      2.051879
     Last Activity_Others                        1.253061
     Total Time Spent on Website                 1.049789
     Last Activity_Email Opened                  0.942099
     Lead Source_Olark Chat                      0.907184
     Last Activity_Olark Chat Conversation      -0.555605
     const                                      -1.023594
     Specialization_Hospitality Management      -1.094445
     Specialization_Others                      -1.203333
     Lead Origin_Landing Page Submission        -1.258954
     dtype: float64
```

```
[ ]   # Lets add Lead Score

      y_pred_final['Lead_Score'] = y_pred_final['Converted_Prob'].map( lambda x: round(x*100))
      y_pred_final.head()
```

| | Prospect ID | Converted | Converted_Prob | final_predicted | Lead_Score |
|---|---|---|---|---|---|
| **0** | 4269 | 1 | 0.697934 | 1 | 70 |
| **1** | 2376 | 1 | 0.860665 | 1 | 86 |
| **2** | 7766 | 1 | 0.889241 | 1 | 89 |
| **3** | 9199 | 0 | 0.057065 | 0 | 6 |
| **4** | 4359 | 1 | 0.871510 | 1 | 87 |

**Lead Score:** Lead Score is assigned to the customers

- The customers with a higher lead score have a higher conversion chance
- The customers with a lower lead score have a lower conversion chance.

# Evaluation Metrics :-

- Train set :
- o Accuracy -> 81.7%
- o Sensitivity-> 79.9%
- o Specificity-> 82.7%
- For Test set :
- o Accuracy : 79.8%
- o Sensitivity : 75.99%
- o Specificity : 82.15%
- Evaluation metrics in both test and train dataset are consistent. Therefore

# Evaluation Metrics(Contd.) :-

- final model is performing good.
- Top 3 features contributing to predicting hot leads are:
- o Lead Origin_Lead Add Form
- o Current_occupation_Working Professional
- o Last Activity_SMS Sent

# Recommendations:-

- Lead Origin, Current Occupation, and Last Activity are top contributors to lead conversion probability.

- Focus on Lead Add Form origin, Working Professional occupation, and Customer SMS activity for conversion.

- Improve Specialization-Others, Olark Chat Last Activity, and address issues with bounced emails.