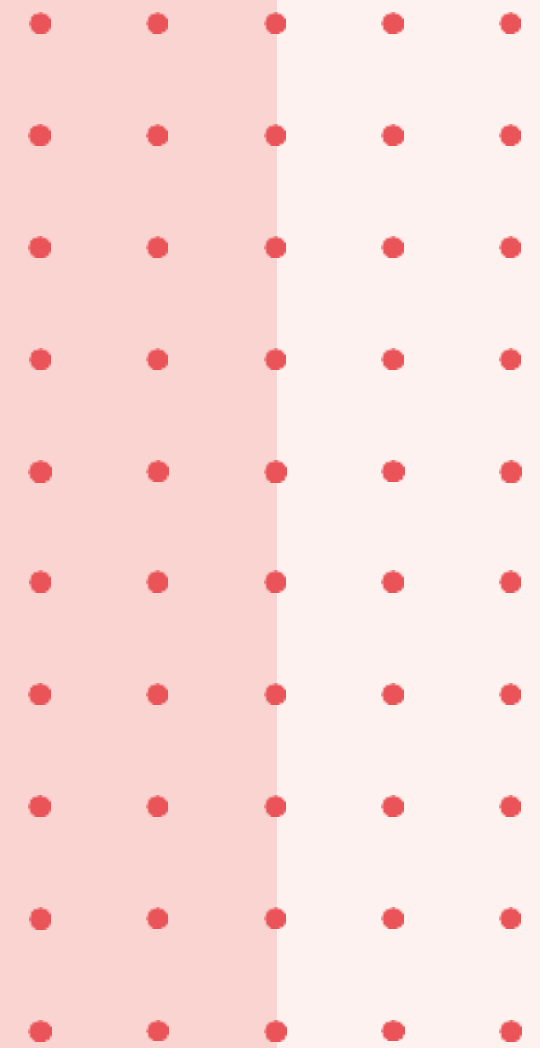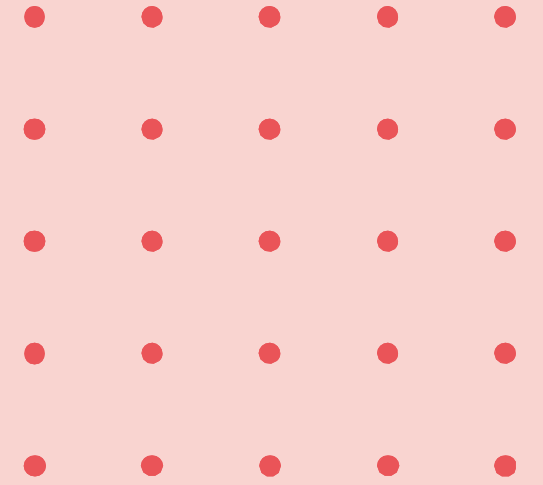# Optimizing Lead Conversion: Developing a Logistic Regression Model

# Data Understanding and Preparation

- Data Loading and Initial Checks
- Loaded dataset from 'Leads.xlsx'.
- Checked for missing values and handled them based on a predefined strategy:

# Missing Values

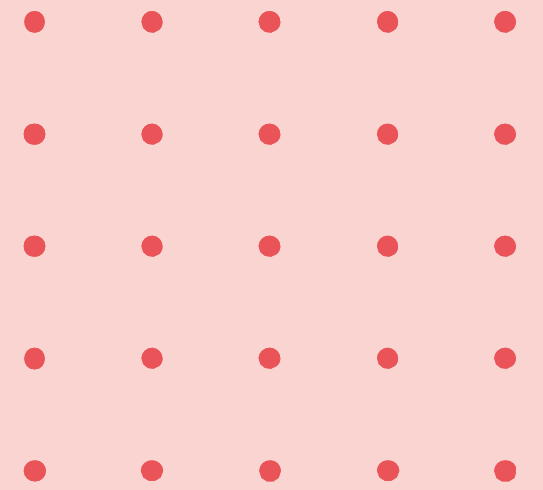| | |
|---|---|
| Prospect ID | 0.000000 |
| Lead Number | 0.000000 |
| Lead Origin | 0.000000 |
| Lead Source | 0.389610 |
| Do Not Email | 0.000000 |
| Do Not Call | 0.000000 |
| Converted | 0.000000 |
| TotalVisits | 1.482684 |
| Total Time Spent on Website | 0.000000 |
| Page Views Per Visit | 1.482684 |
| Last Activity | 1.114719 |
| Country | 26.634199 |
| Specialization | 15.562771 |
| How did you hear about X Education | 23.885281 |
| What is your current occupation | 29.112554 |
| What matters most to you in choosing a course | 29.318182 |
| Search | 0.000000 |
| Magazine | 0.000000 |
| Newspaper Article | 0.000000 |
| X Education Forums | 0.000000 |
| Newspaper | 0.000000 |
| Digital Advertisement | 0.000000 |
| Through Recommendations | 0.000000 |
| Receive More Updates About Our Courses | 0.000000 |
| Tags | 36.287879 |
| Lead Quality | 51.590909 |
| Update me on Supply Chain Content | 0.000000 |
| Get updates on DM Content | 0.000000 |
| Lead Profile | 29.318182 |
| City | 15.367965 |
| Asymmetrique Activity Index | 45.649351 |
| Asymmetrique Profile Index | 45.649351 |
| Asymmetrique Activity Score | 45.649351 |
| Asymmetrique Profile Score | 45.649351 |
| I agree to pay the amount through cheque | 0.000000 |
| A free copy of Mastering The Interview | 0.000000 |
| Last Notable Activity | 0.000000 |

## Initial Missing Values

- Various columns with missing values ranging from 0% to 51%.
- Applied cleaning strategy to handle these missing values effectively.
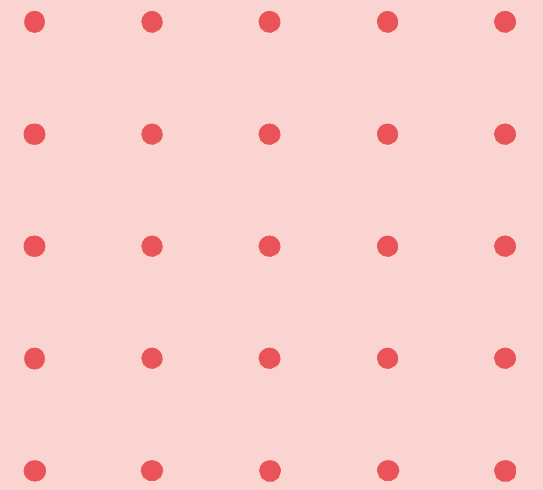
# Data Cleaning

Handling Missing Values and Outliers

- Deleted columns with >40% missing data.
- Dropped rows with missing values in key columns.
- Filled remaining missing values appropriately.
- Removed duplicates and outliers for key numerical columns.

# Standardization and Dummies

Data Standardization and Creating Dummy Variables

- Standardized 'Do Not Email' and 'Do Not Call' columns.
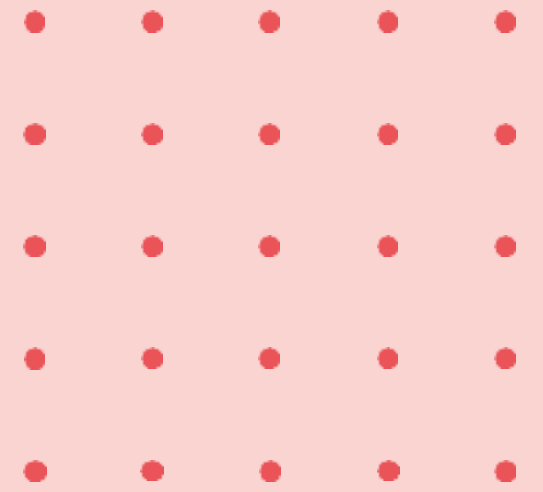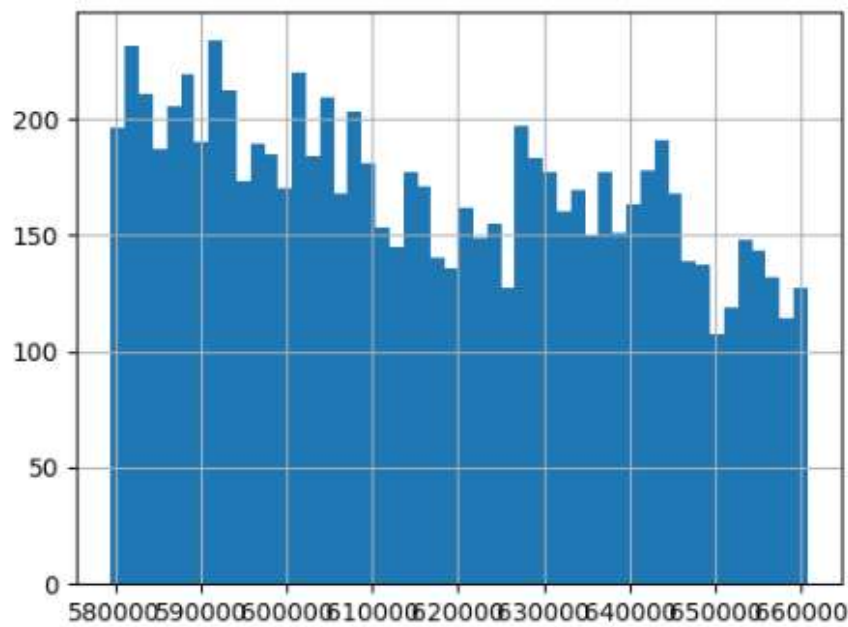- Created dummy variables for categorical columns.
- Derived new metrics like 'Time Per Visit'.

# Univariate Analysis

Plotted histograms for all variables to understand distribution.

# Bivariate Analysis

## Scatter Plots and Box Plots
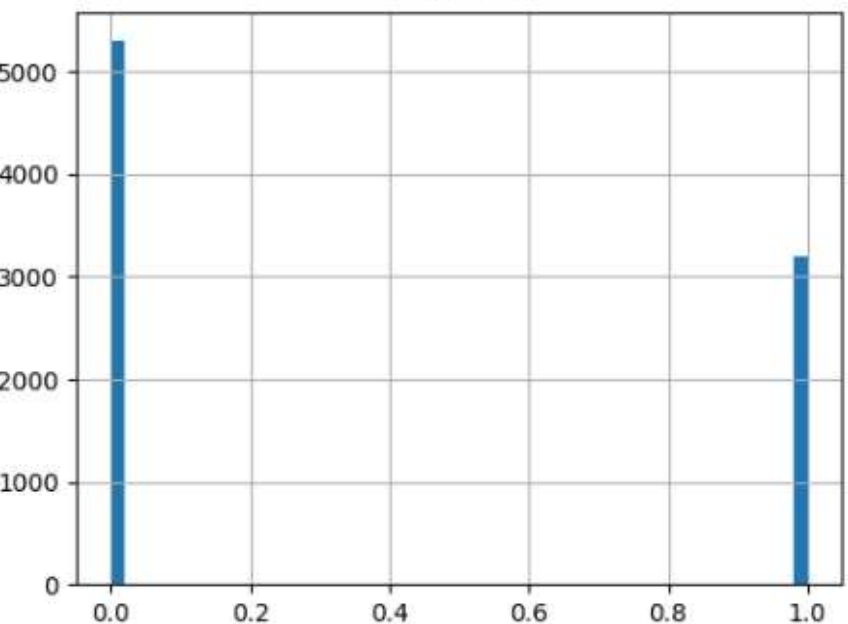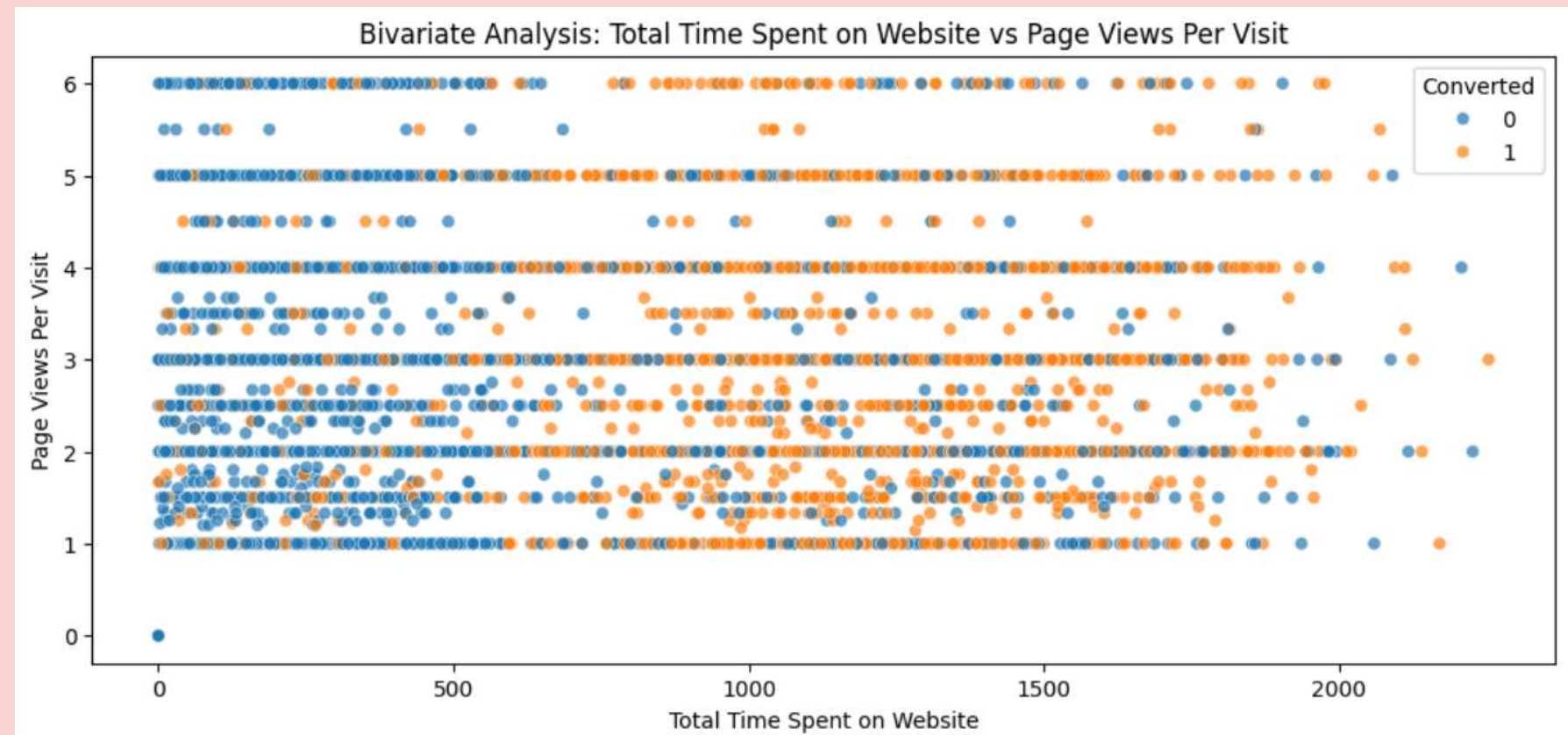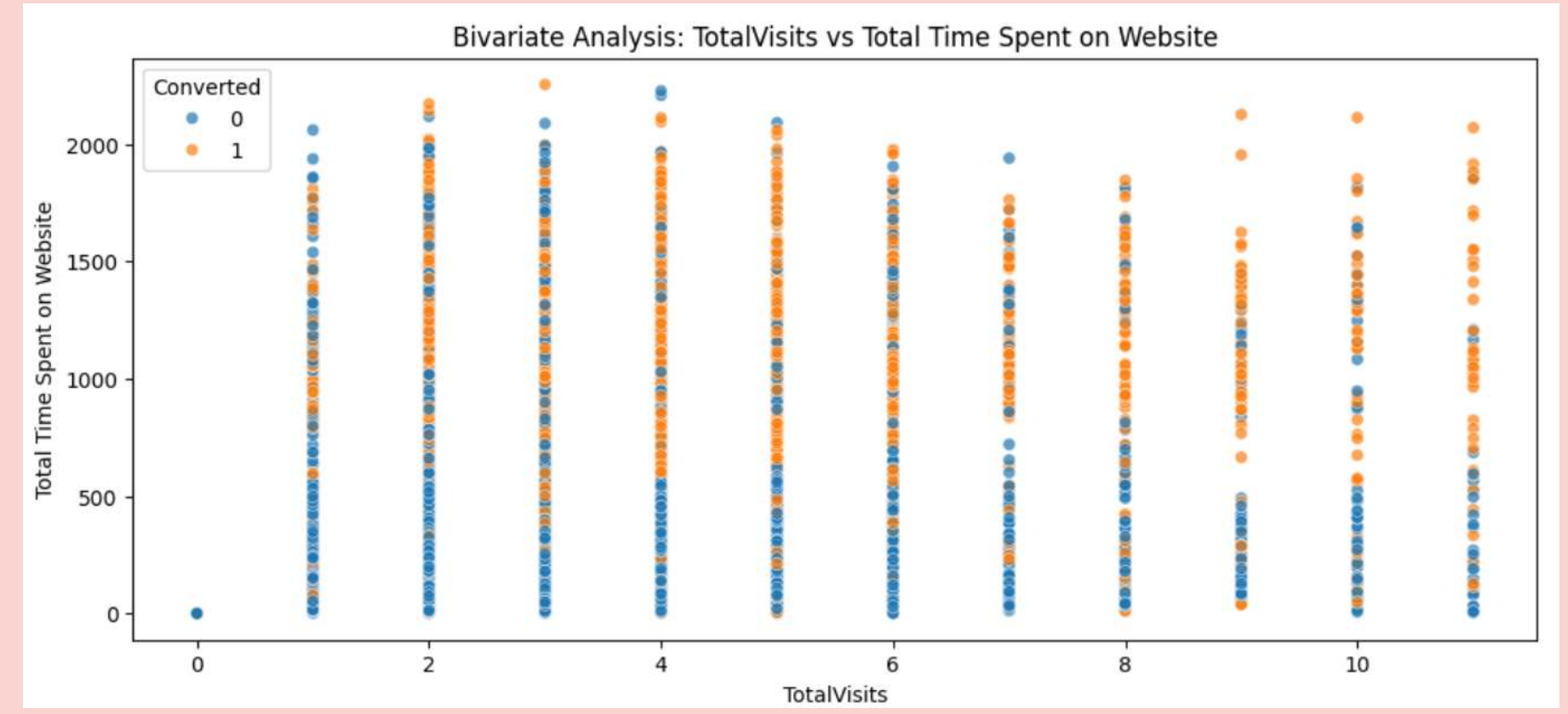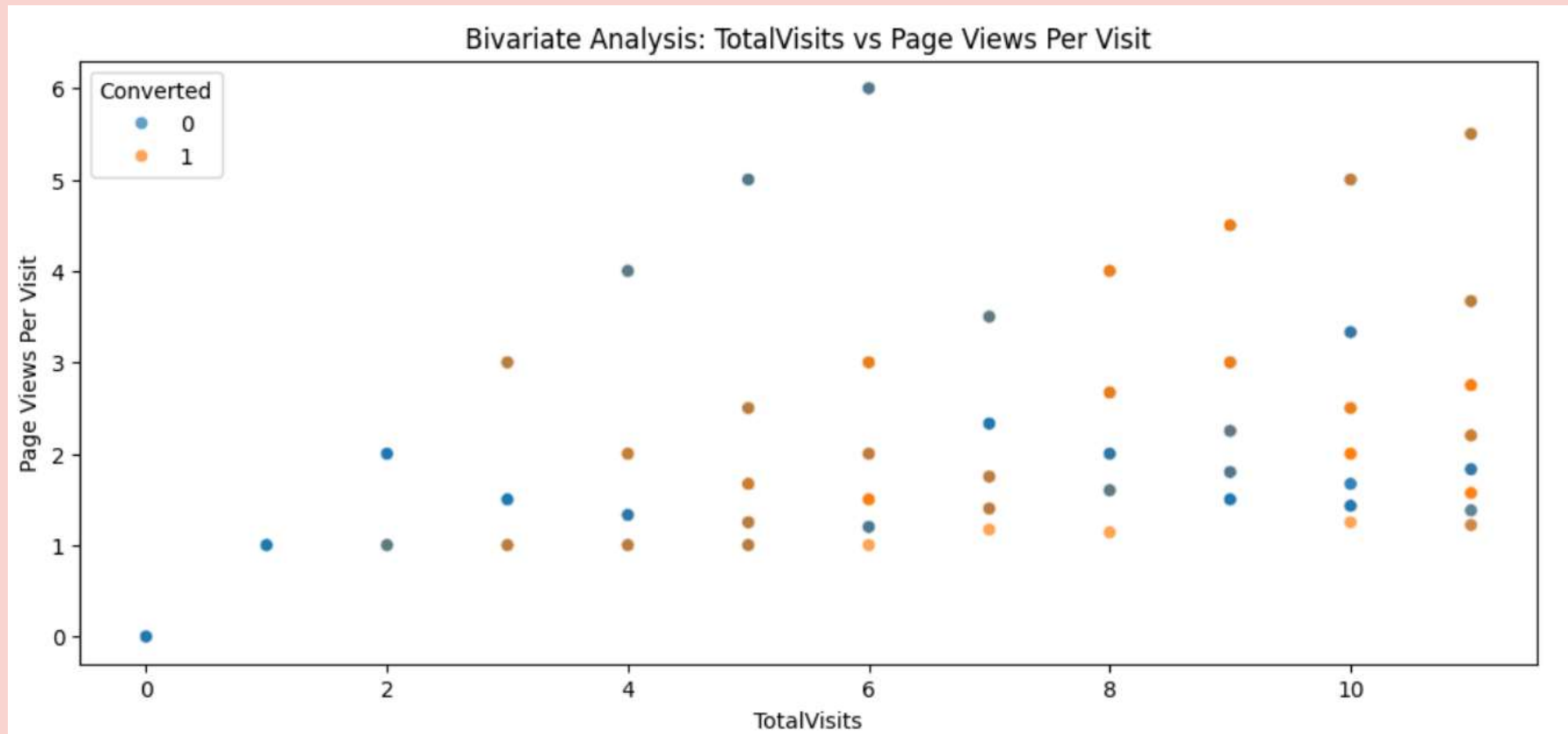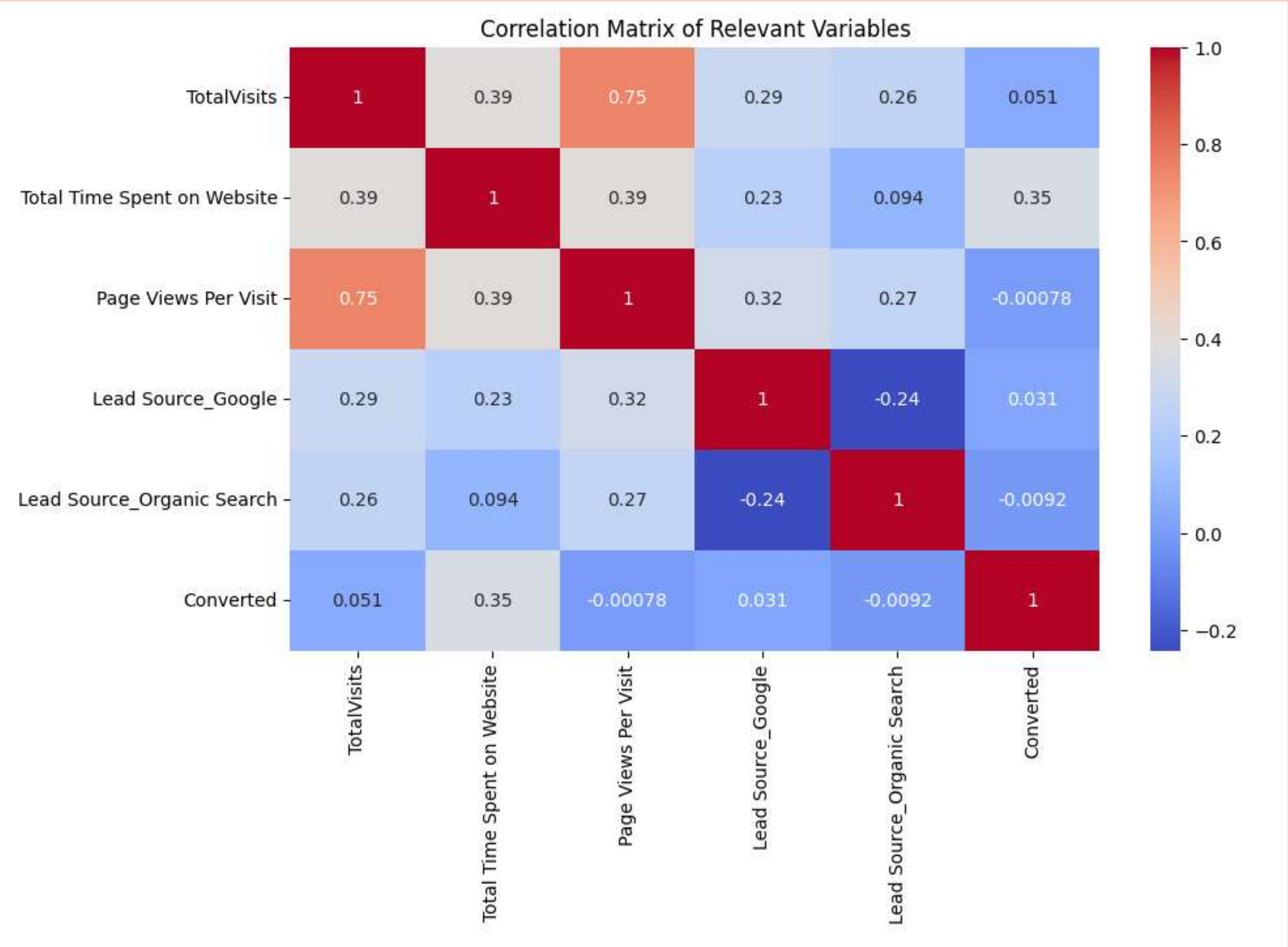
# Correlation Matrix

Heat Map

# Model Building

Logistic Regression Model

o Split data into training and testing sets.
o Handled missing values with SimpleImputer.
o Fitted a logistic regression model with max_iter=1000.

# Model Evaluation

Model Performance Metrics

Evaluated model using:
- o Accuracy: 90%
- o Precision: 91%
- o Recall: 82%
- o F1 Score: 87%
- o ROC-AUC Score: 96%

## Model Interpretation

Top 5 Important Features

o  Identified top 5 features
   contributing to the model:
o  Tags_Will revert after reading the
   email
o  Tags_Lost to EINS
o  Tags_Closed by Horizzon
o  City_Select
o  Lead Origin_Lead Add Form
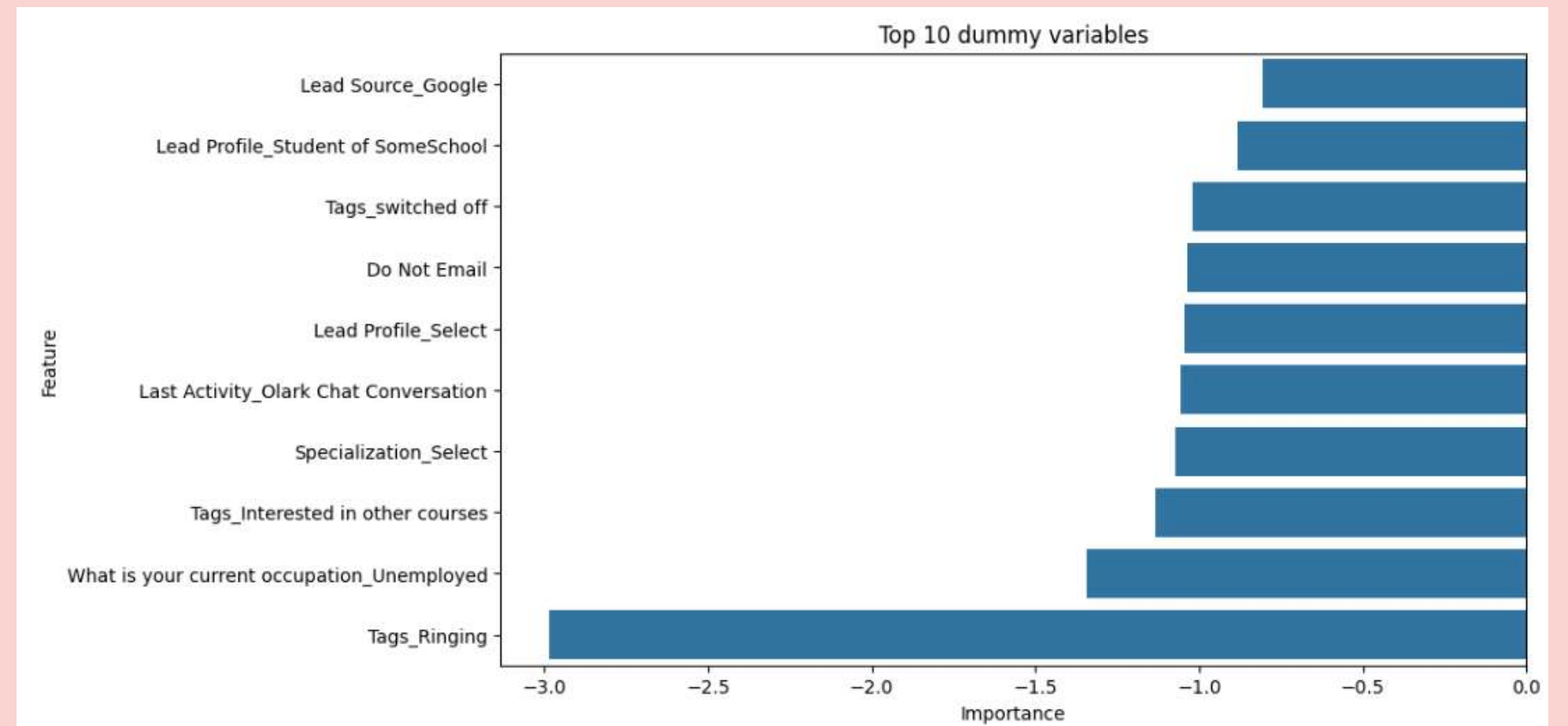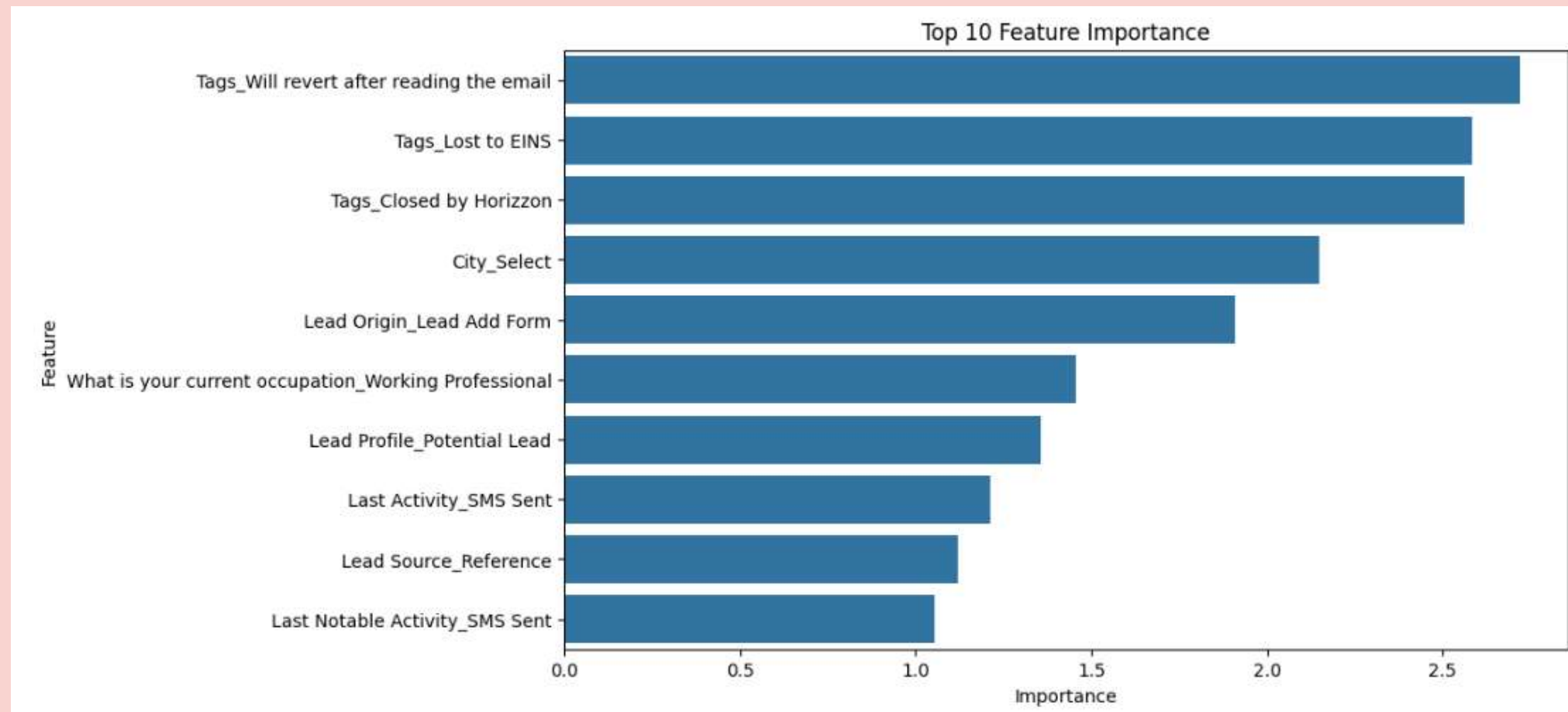
# Recommendations

**Aggressive Strategy with Interns**

o Prioritize leads predicted as 1.
o Increase follow-up calls and emails.
o Offer promotions and provide intern training.

**Minimizing Calls Strategy**

o Focus on leads with high conversion probability (>0.8).
o Reduce follow-up calls and use automated emails.
o Reallocate resources to other tasks.

# Graphs



Top 10 Feature Importance



Top 10 dummy variables

# Conclusion

**Summary of the Analysis and Model**

o Effective data cleaning and preparation.

o Built a robust logistic regression model.

o Provided actionable recommendations based on model insights.

# Thanks!