

LEAD SCORING PROJECT

BY-

Y. SWETHA ASHERVADAM

Problem Statement



- ❖ *X Education, an online education company, faces challenges in converting leads into paying customers despite generating a significant number of leads daily.*
- ❖ *With a conversion rate of only 30%, the company aims to improve its lead identification process to target potential customers more effectively.*
- ❖ *They seek to develop a model that assigns lead scores to prioritize leads based on their likelihood of conversion, ultimately aiming for an 80% conversion rate.*

Data Scenario Categorization

- ❖ The dataset provided contains approximately **9000** data points, including attributes like Lead Source, Total Time Spent on Website, Total Visits, and Last Activity.
- ❖ The target variable, '**Converted**', indicates whether a lead was converted (**1**) or not (**0**).
- ❖ The dataset also includes categorical variables with levels such as '**Select**', which need handling as they are equivalent to null values.



Objectives

- 1) Build a logistic regression model to assign lead scores ranging from 0 to 100 to each lead, aiding the company in targeting potential customers effectively. Higher scores indicate hotter leads with a higher likelihood of conversion, while lower scores represent colder leads less likely to convert.
- 2) Ensure the model's adaptability to accommodate potential changes in the company's requirements in the future, addressing any additional problems or challenges that may arise.



STEPS FOLLOWED

Importing Data and inspecting the Dataframe

- Reading the data from the csv file using pandas library

Data Exploration

- Handling missing values
- Handling the 'Select' value present in four variables
- Checking Data Imbalance

Exploratory Data Analysis

- Univariate Analysis
- Bivariate Analysis
- Segmented Analysis on Target variable 'Converted'
- Multivariate Analysis

Data cleaning and preparation

- Dropping the imbalanced variables
- Replacing lower frequency values if required
- Treatment of Outliers

Feature Engineering

- Converting the binary variables (Yes/No) to 0/1
- Handling categorical variables
 - Mapping categorical variables to integers
 - Dummy variable creation

Test-train split and scaling

- Train data – 70% and Test data – 30%

Model Building

- Feature elimination based on correlations
- Feature selection using RFE (Coarse Tuning)
- Manual feature elimination (using p-values and VIFs)

Model Evaluation

- Accuracy
- Sensitivity and Specificity
- Optimal cut-off using ROC curve
- Precision and Recall

Predictions on the test set

- Predictions using the test data

Data Manipulation

Handling Missing Values

Dropping single value variables or imbalanced variables

Capping the outliers



Handling Missing Values



This dataset has 9240 rows and 37 columns.

Identified and removed columns with null values exceeding 40% threshold for analysis.

After removal, 12 columns still had null values.

Imputed missing values in 'Tags', 'What matters most to you in choosing a course' and 'What is your current occupation' column with 'Not Specified'.

Imputed missing values in 'Country' column with mode and clubbed lower frequency values.

Addressed the variables containing 'Select' values, requiring handling as follows:

- Dropped 'Lead Profile' column due to combined presence of 'Select' and null values exceeding 70%.
- Dropped 'City' column due to combined presence of 'Select' and null values exceeding 39%.
- Dropped 'How did you hear about X Education' column due to combined presence of 'Select' and null values constituting 74%.
- Imputed 'Specialization' column's 'Select' and null values with 'Not Specified' due to constituting 36% of the values.

Eliminated rows with negligible null values (ranging from approximately 1% to 2%) in remaining four key variables.

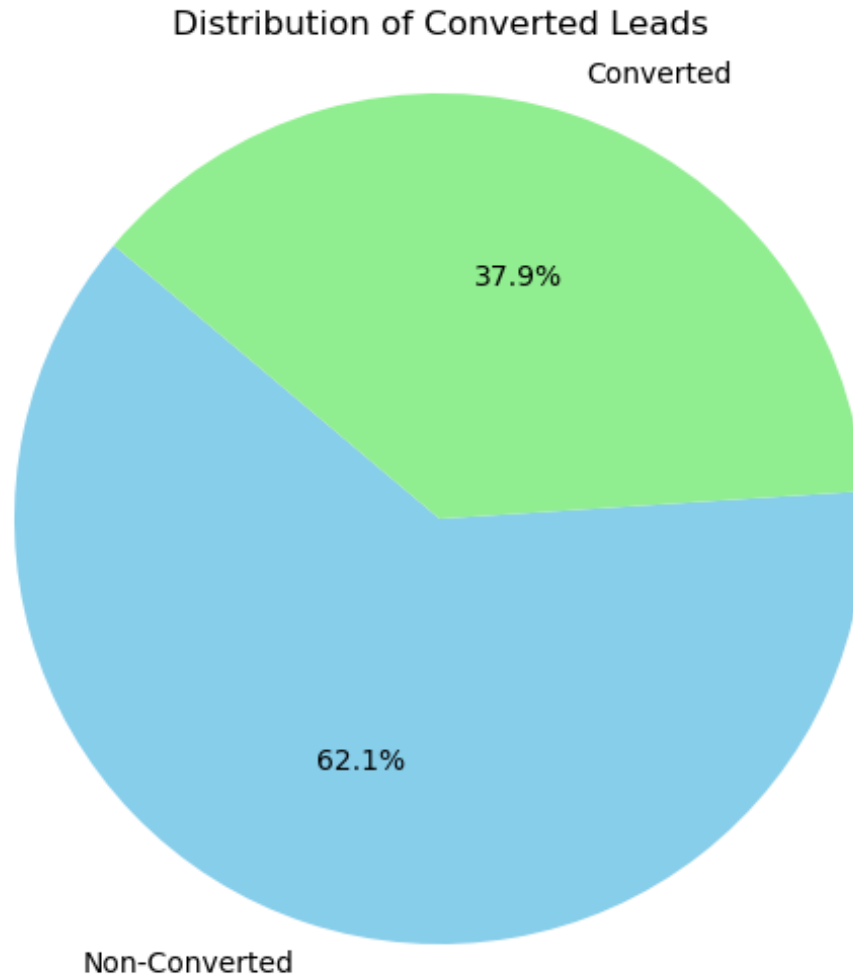
Dropped columns deemed imbalanced:

- Prospect ID, Lead Number, Do Not Call, Magazine, Search, Newspaper Article, X Education Forums, Newspaper, Through Recommendations, Digital Advertisement, I agree to pay the amount through cheque, Update me on Supply Chain Content, Receive More Updates About Our Courses, Get updates on DM Content.

Capped outliers in 'TotalVisits' and 'Page Views Per Visit' to the 95th percentile.

After data cleaning, the dataset contains 9074 rows and 15 columns.

Data Imbalance



The number of converted leads are (3435, 29)

The number of non converted leads are (5639, 29)

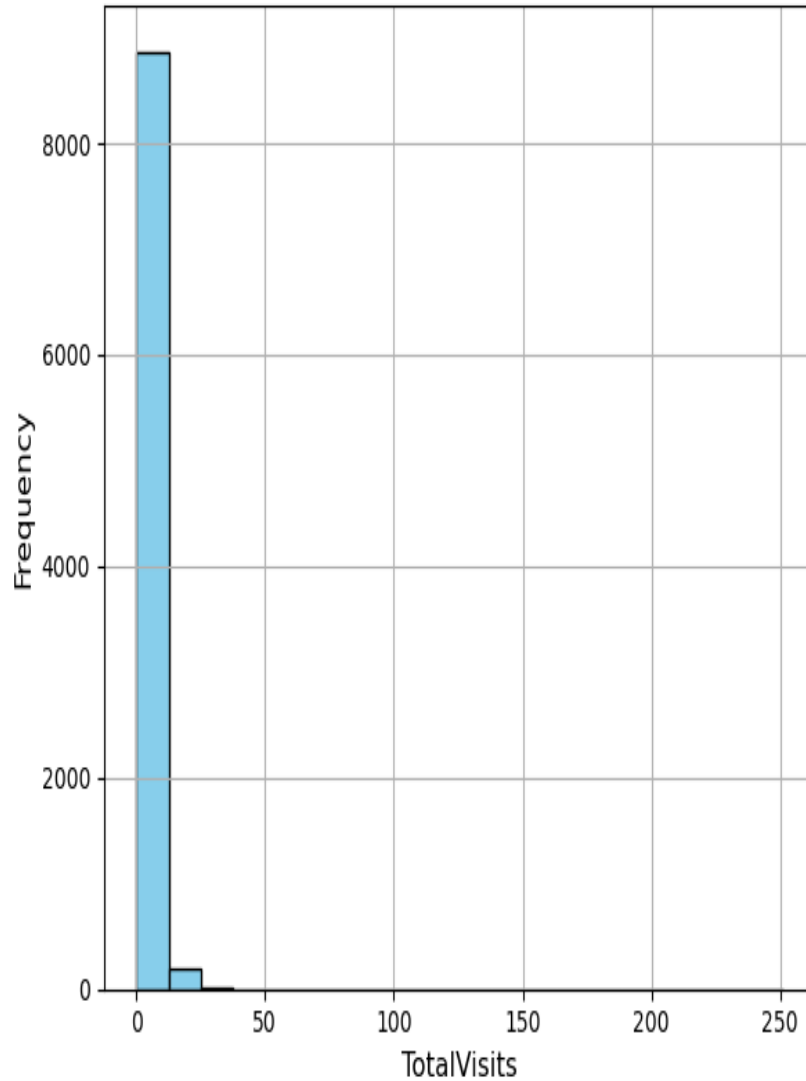
Exploratory Data Analysis



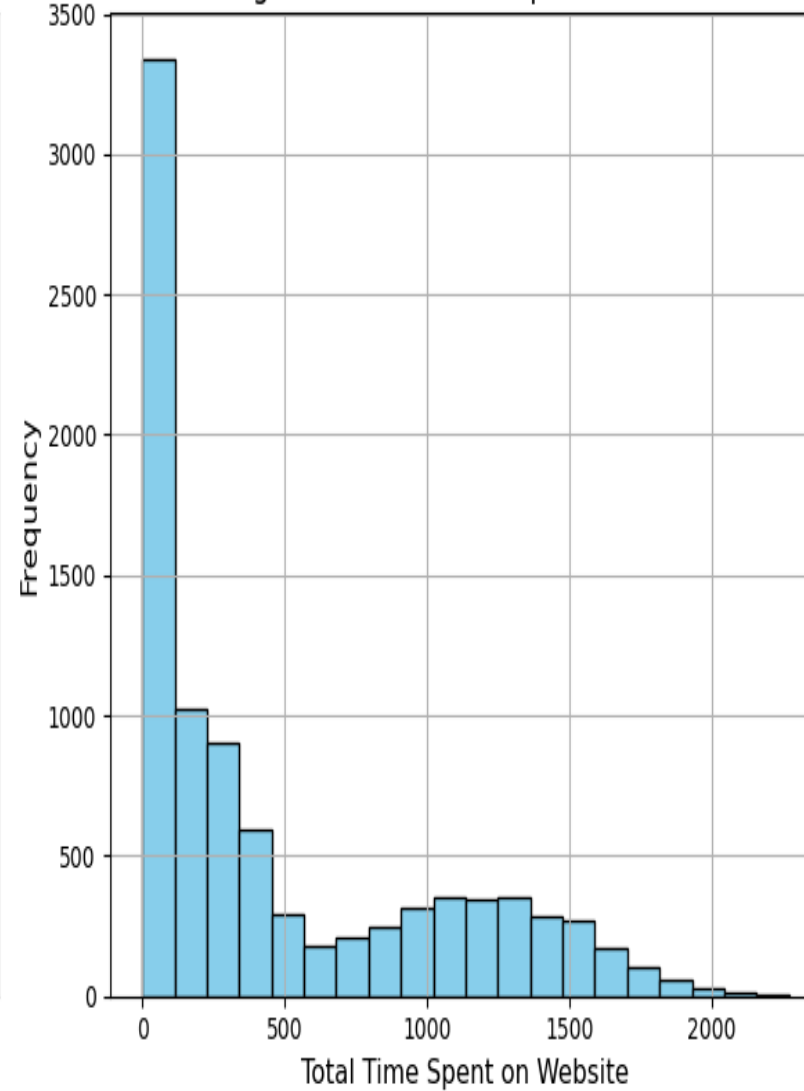
UNIVARIATE ANALYSIS

Histograms of Numeric Variables

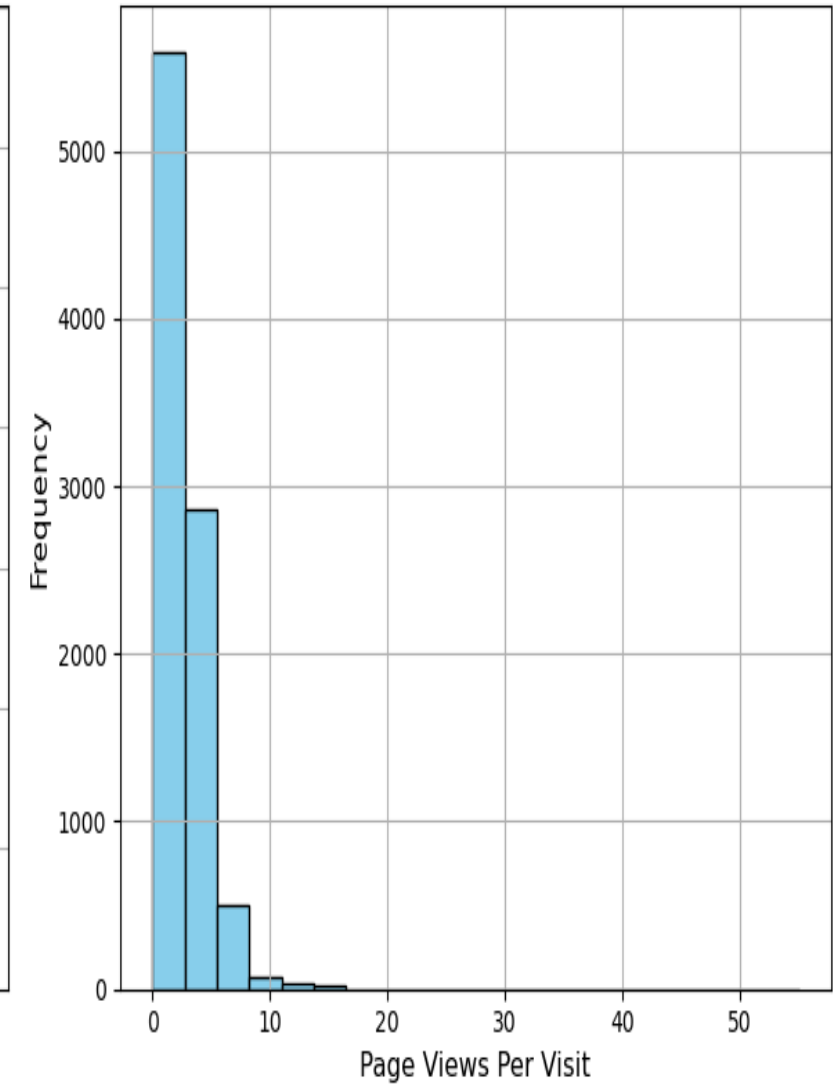
Histogram of TotalVisits



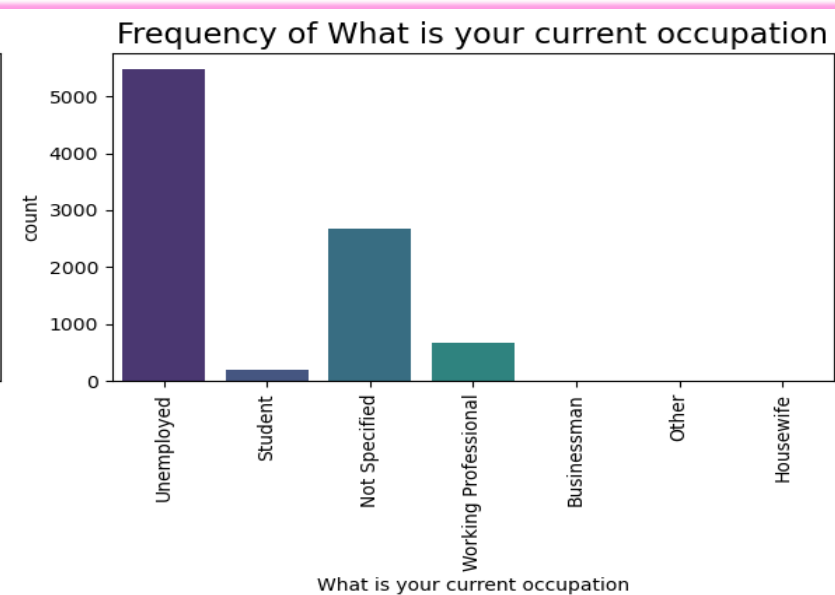
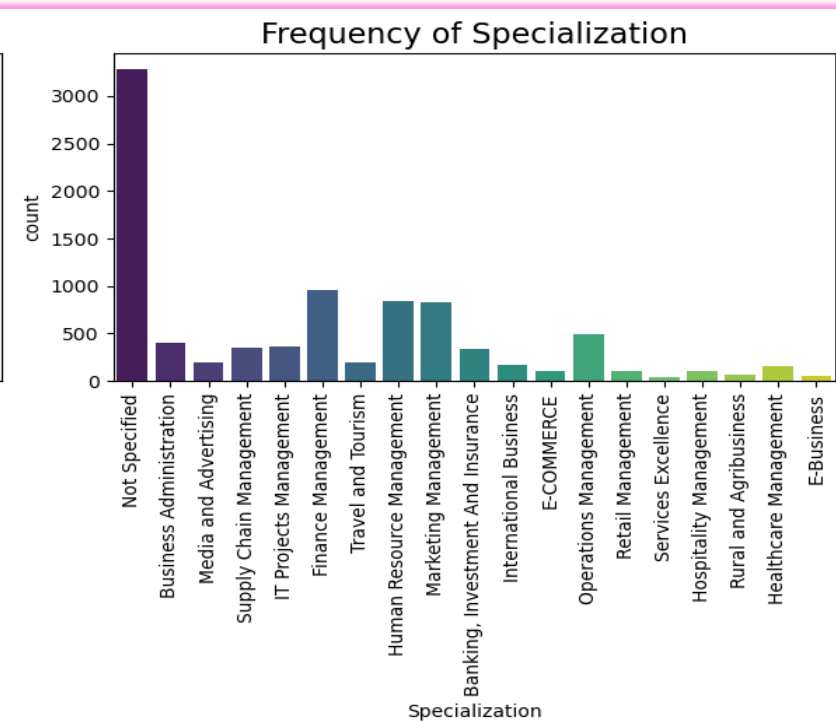
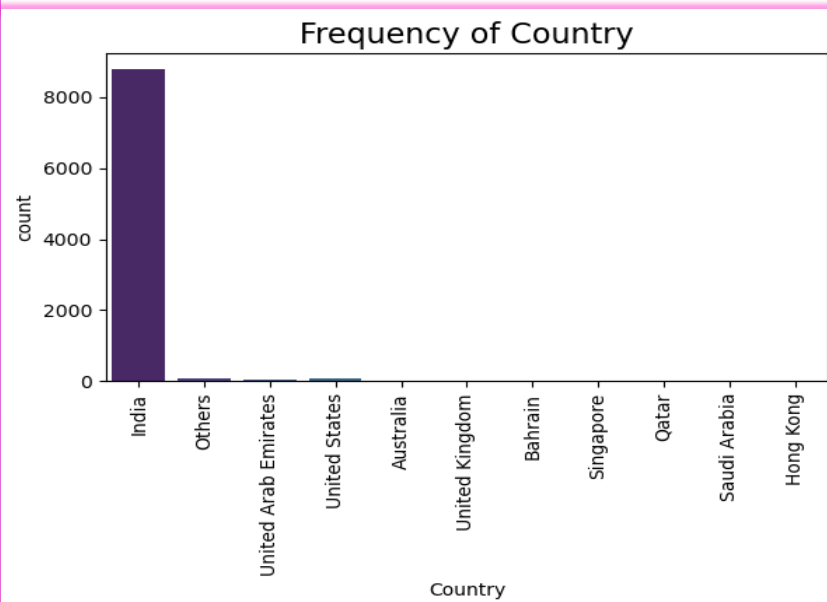
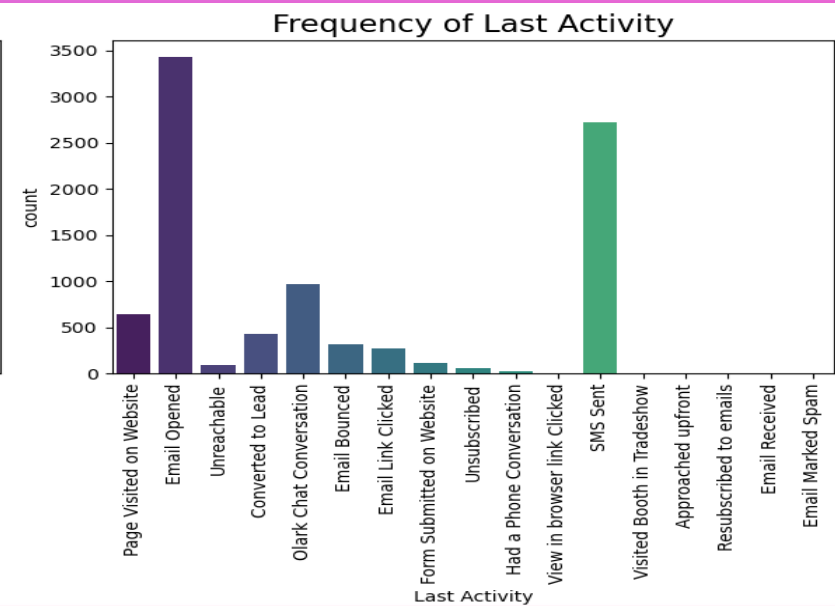
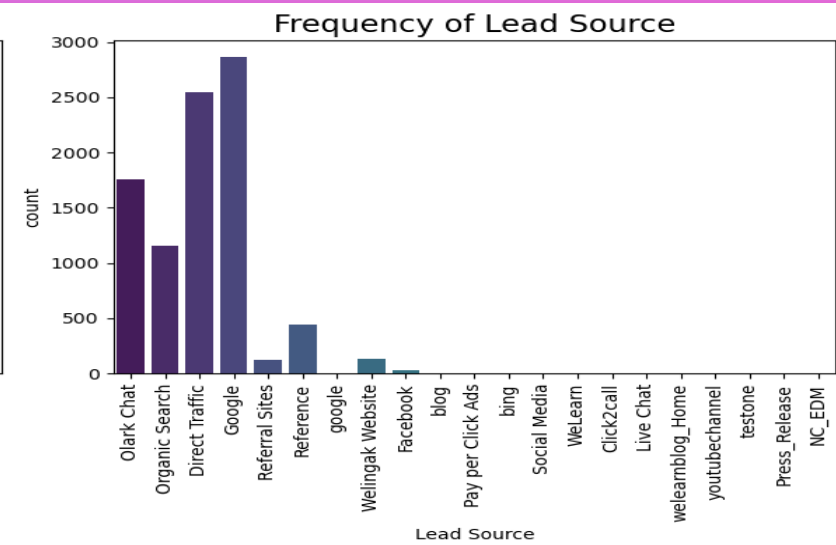
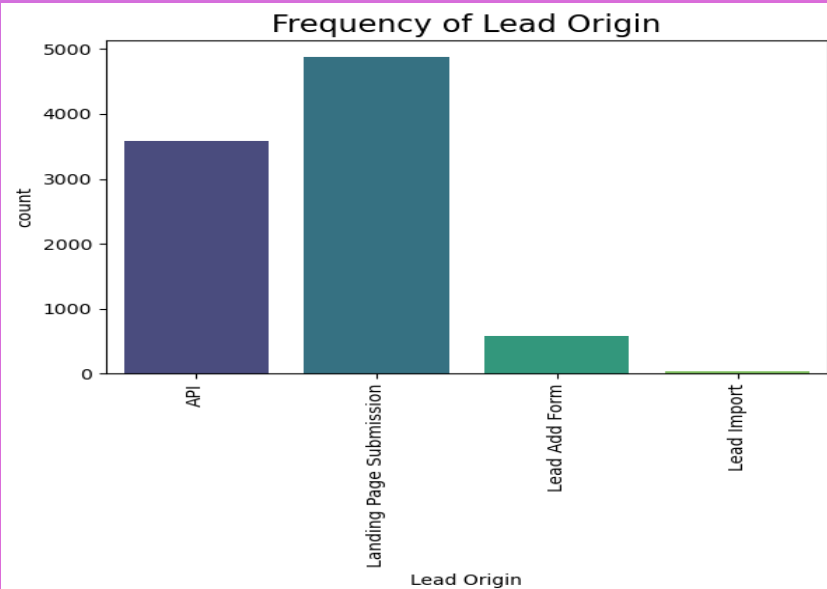
Histogram of Total Time Spent on Website



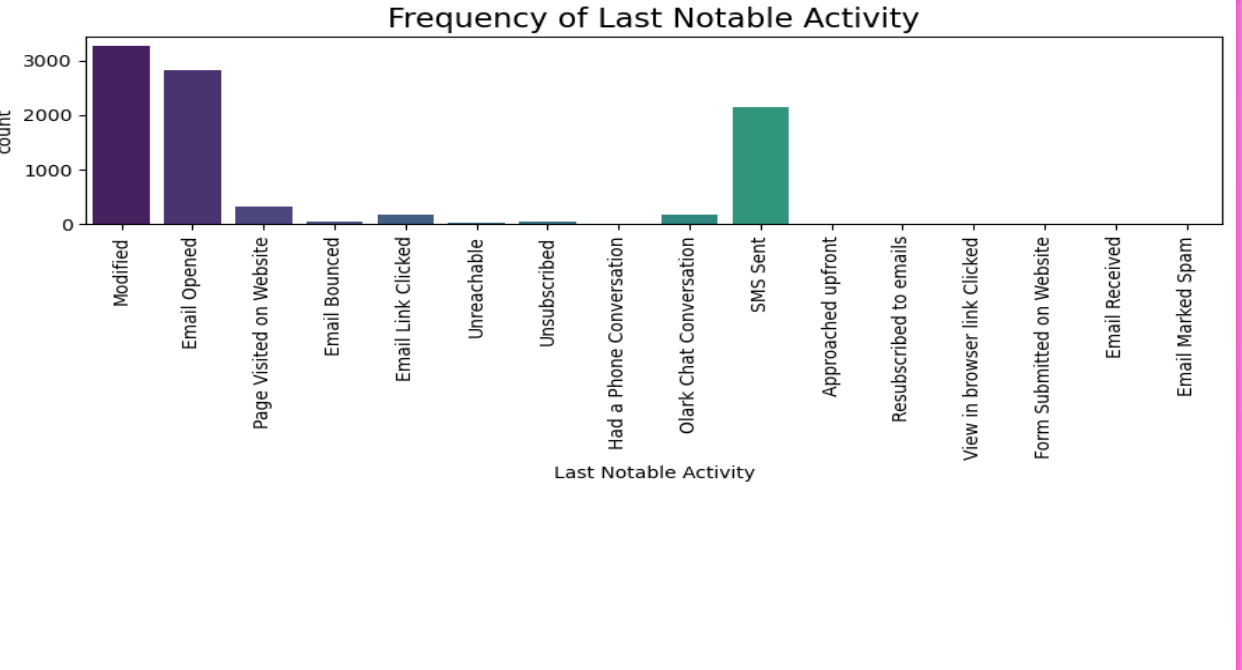
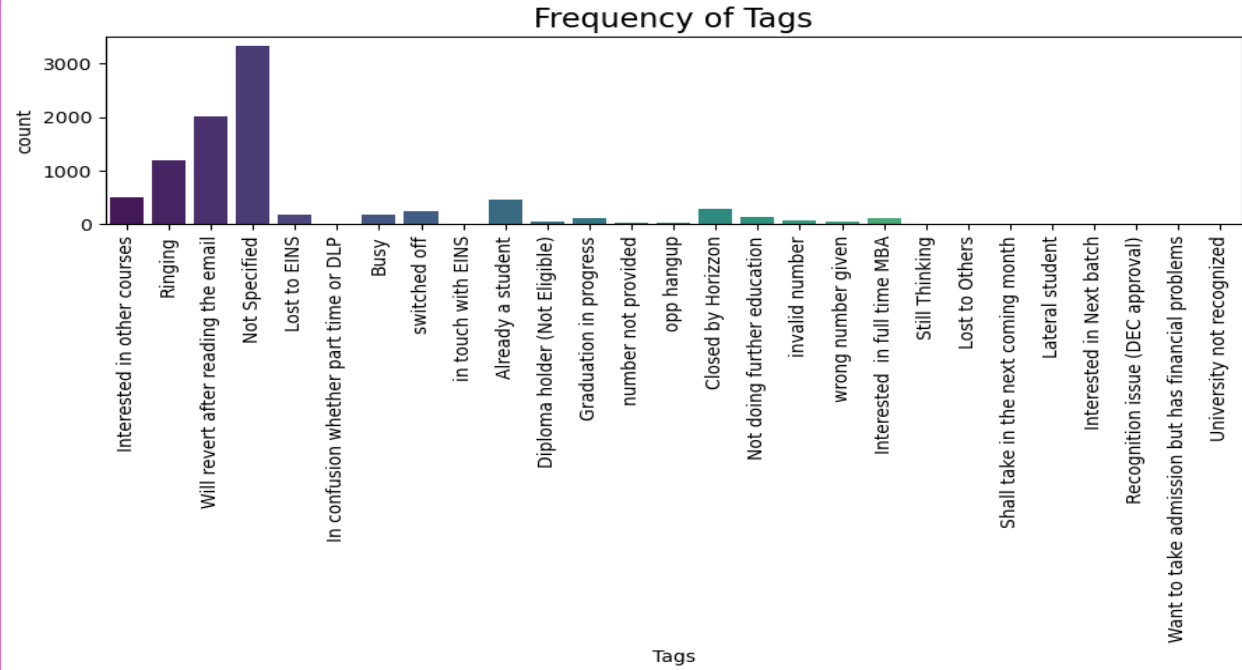
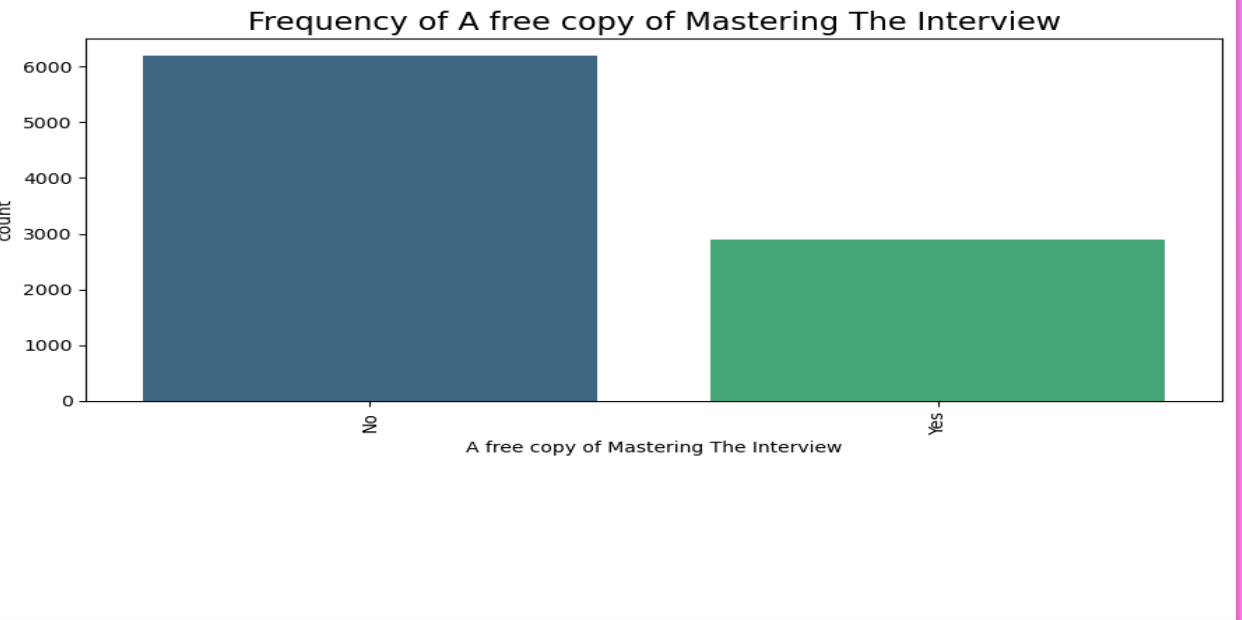
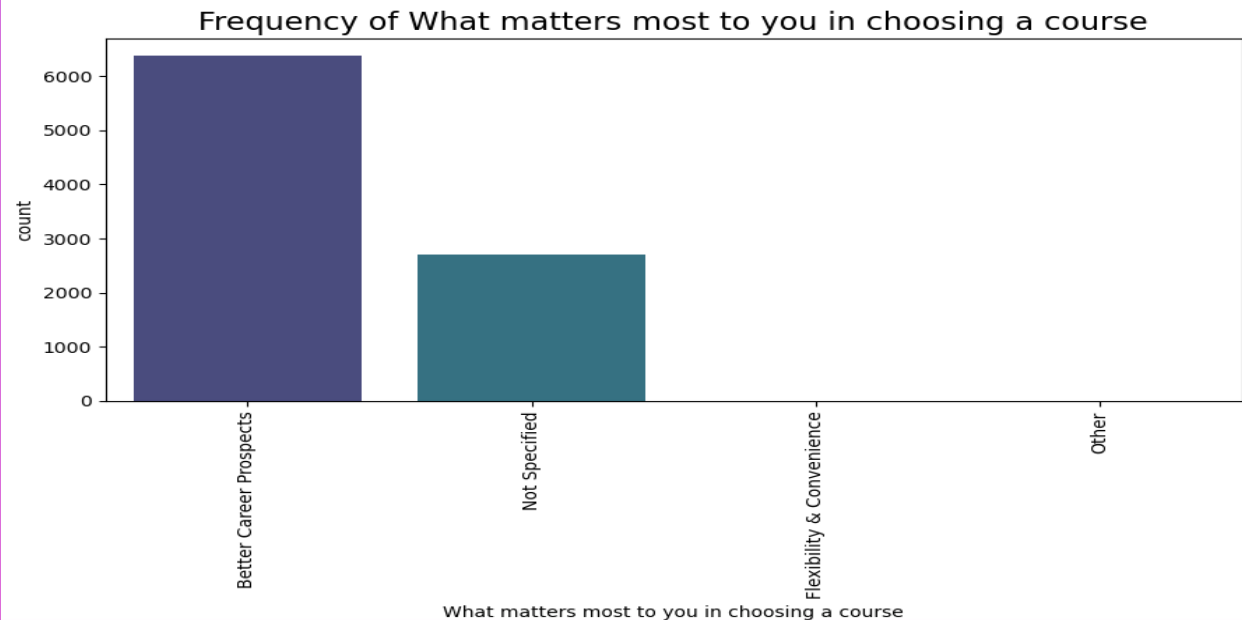
Histogram of Page Views Per Visit



Bar charts for categorical variables

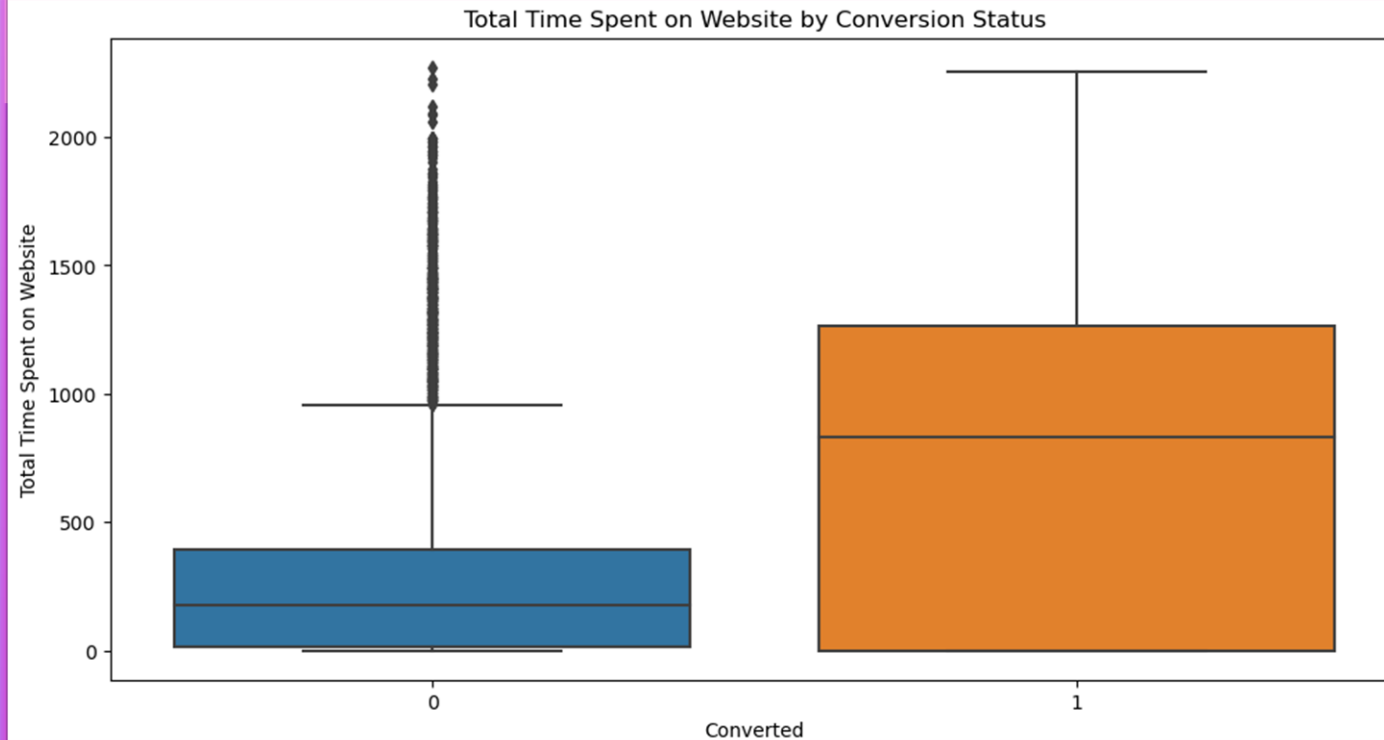


Bar charts for categorical variables

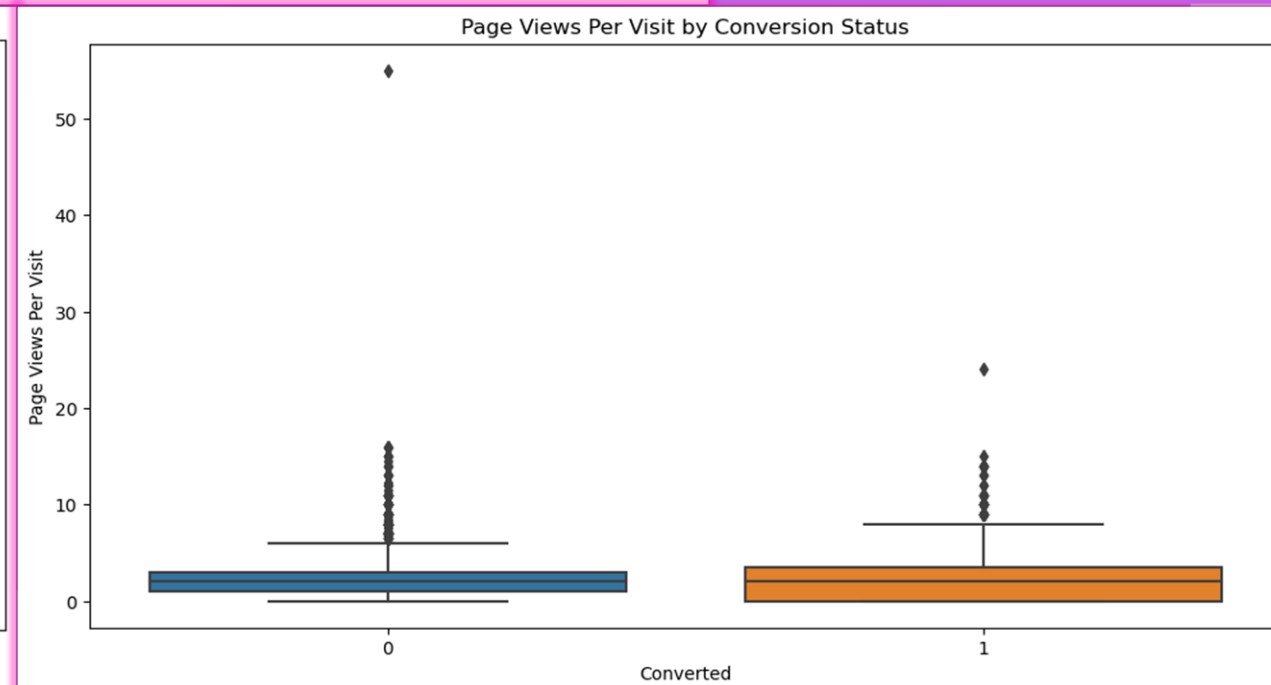
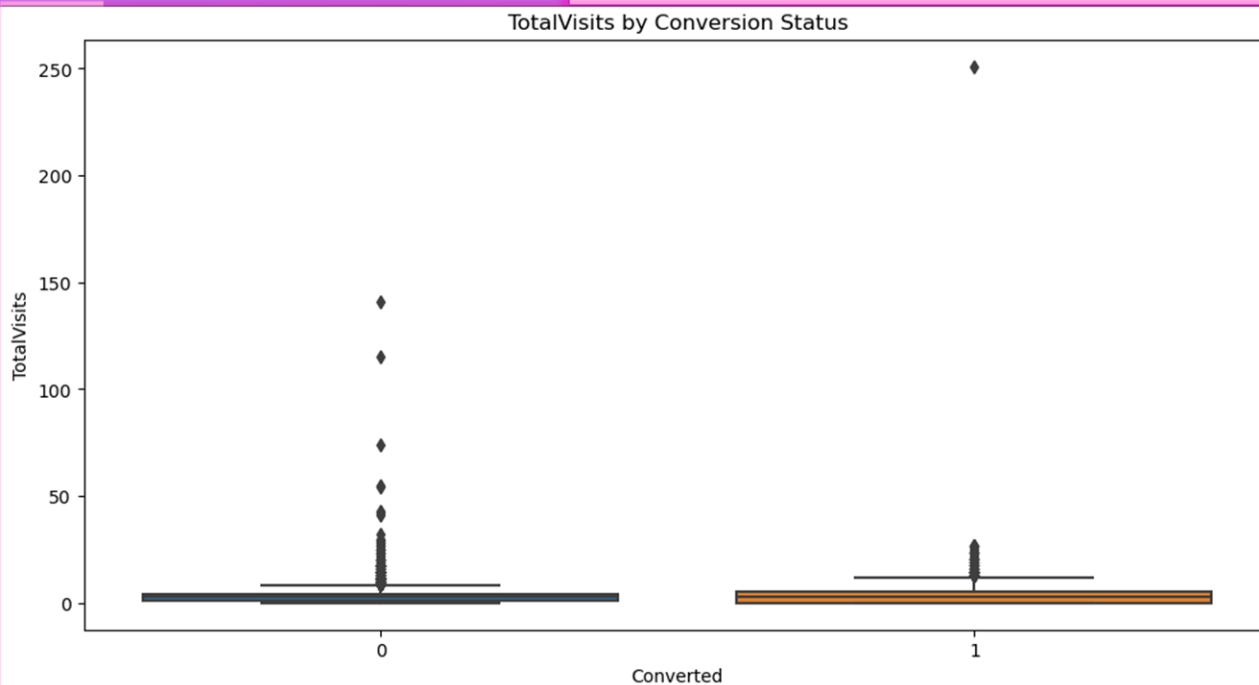


BIVARIATE ANALYSIS

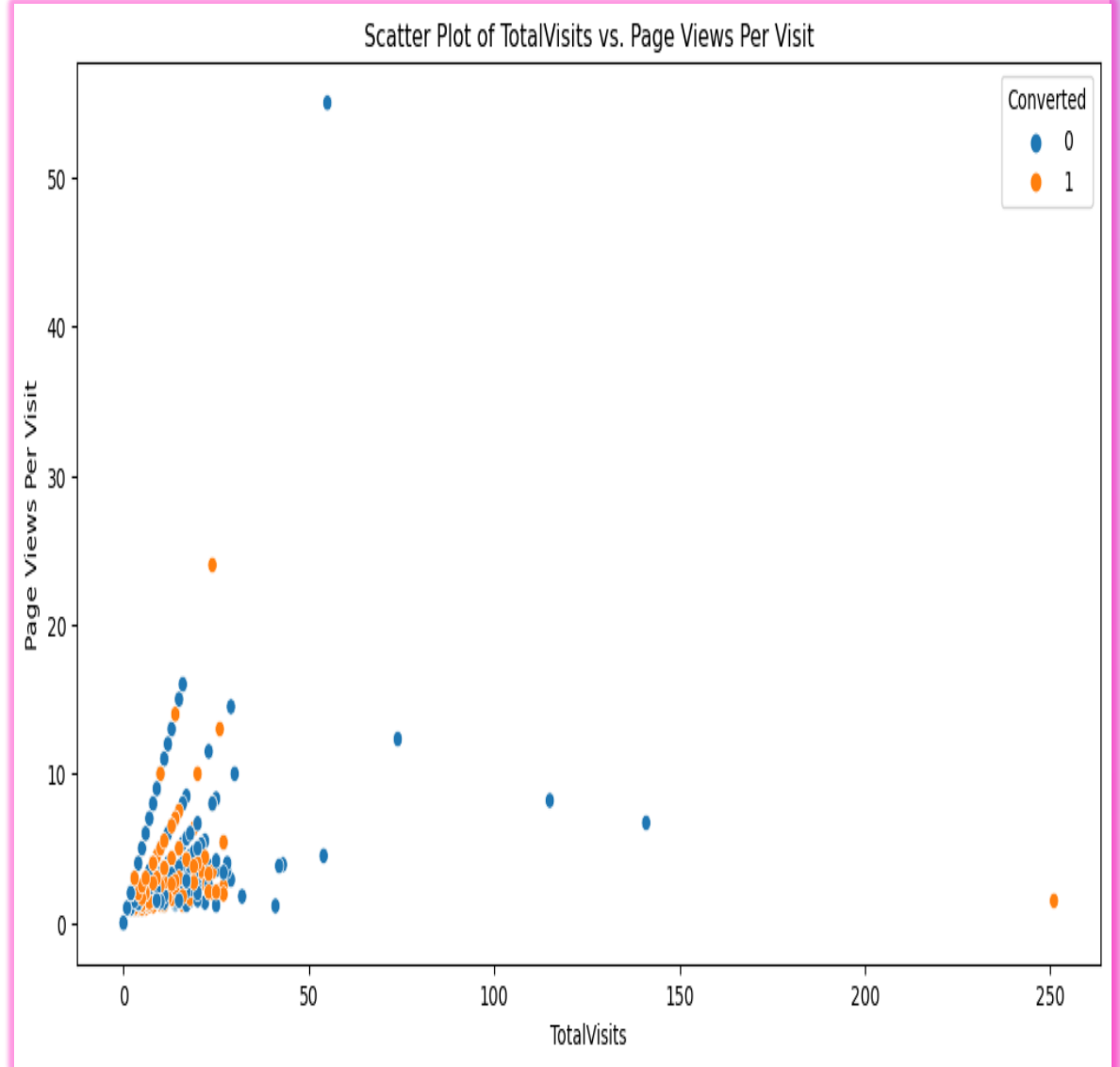
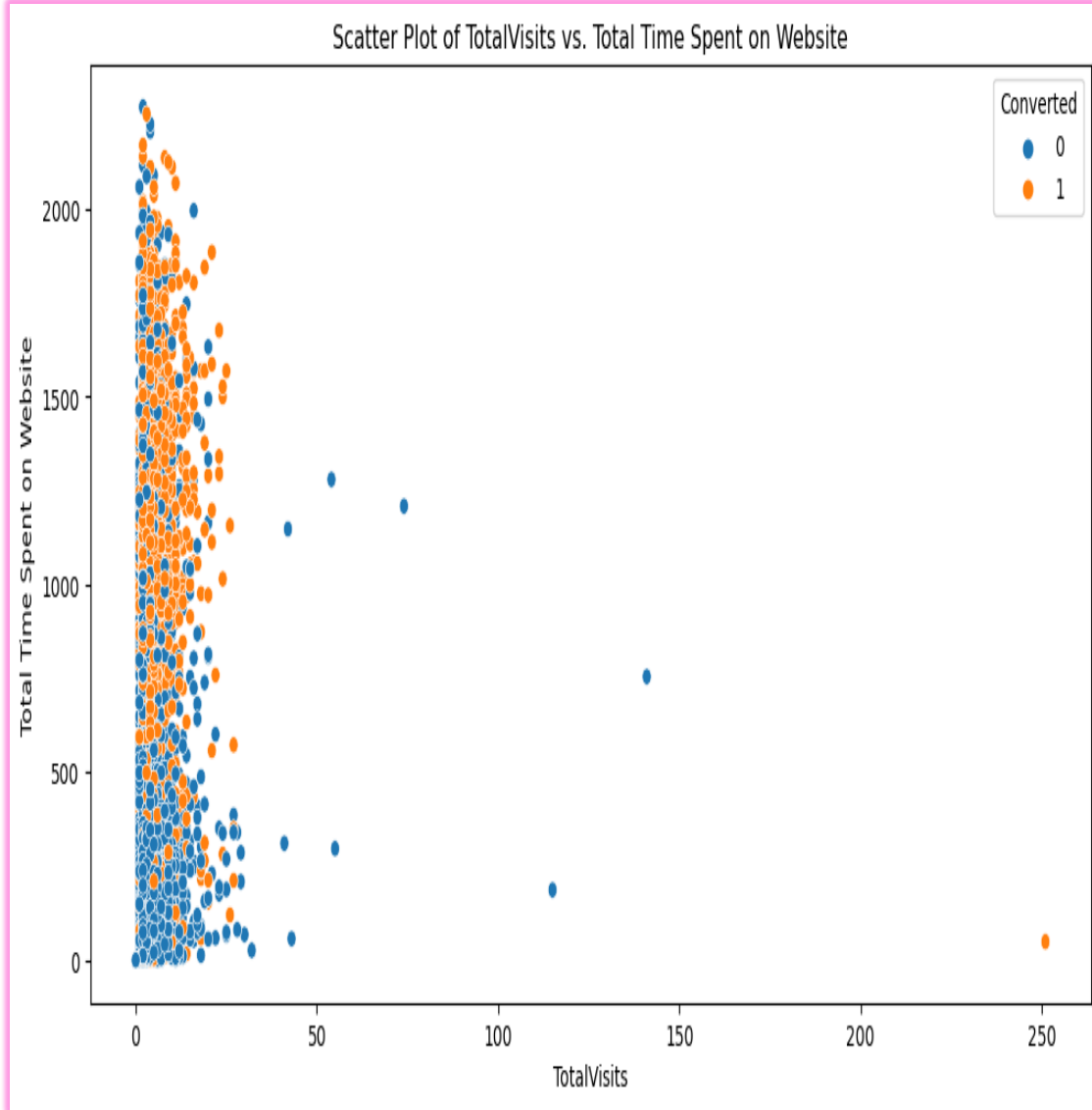
*Box plots
for
numerical
variables*



*Against
Conversion
status*

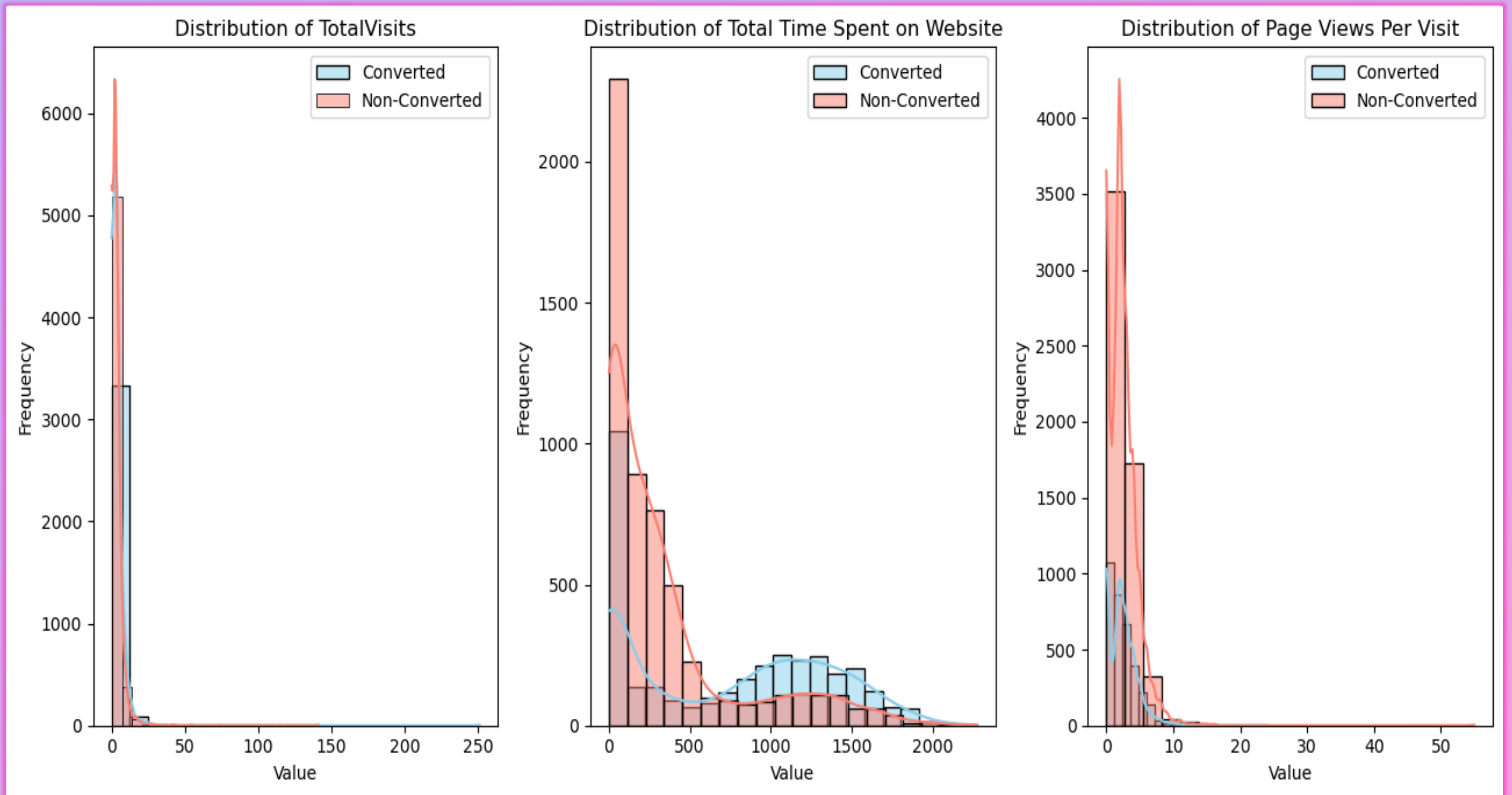


Scatter plots for numerical variables against the target variable 'Converted'

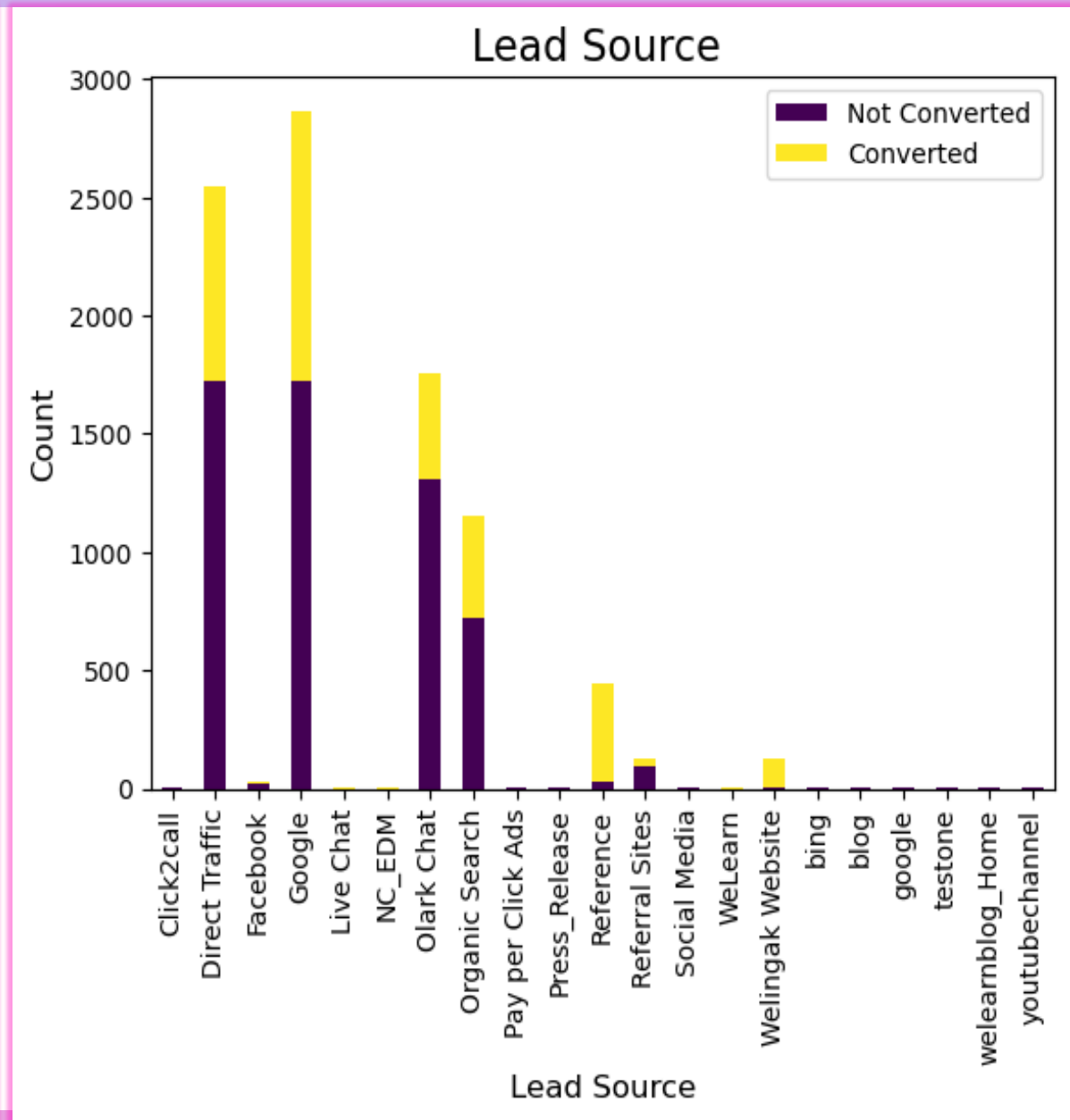
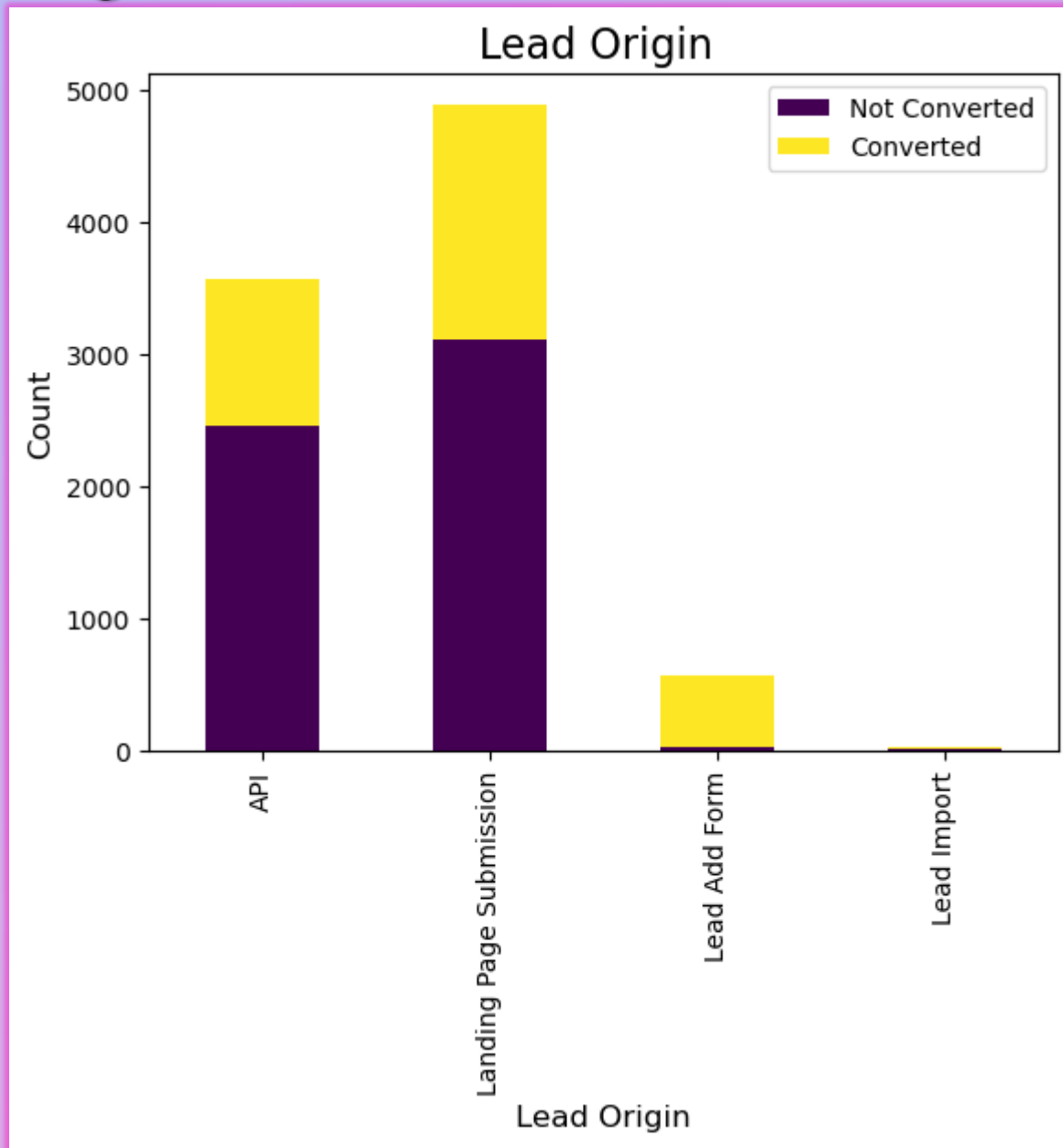


SEGMENTED ANALYSIS

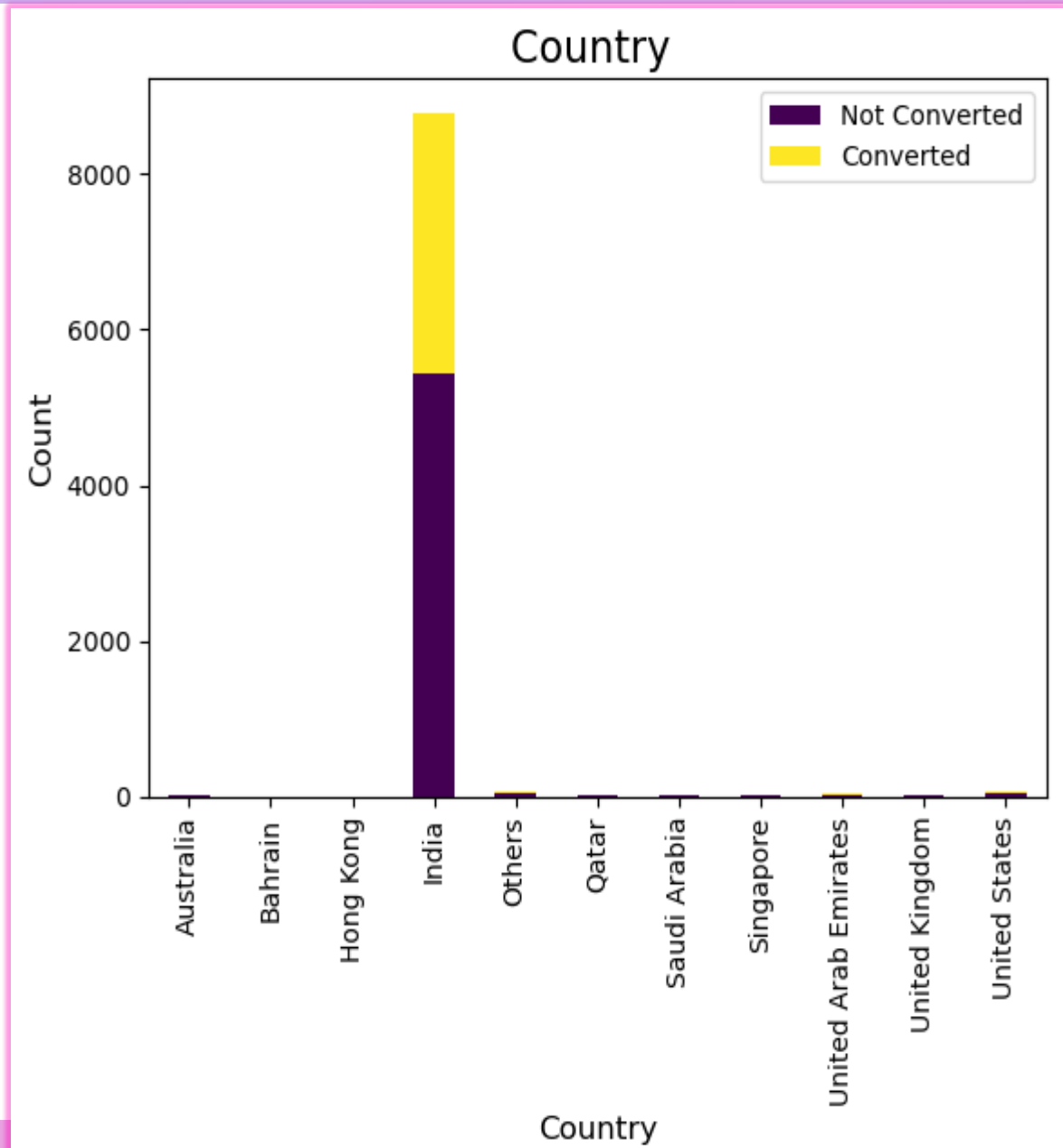
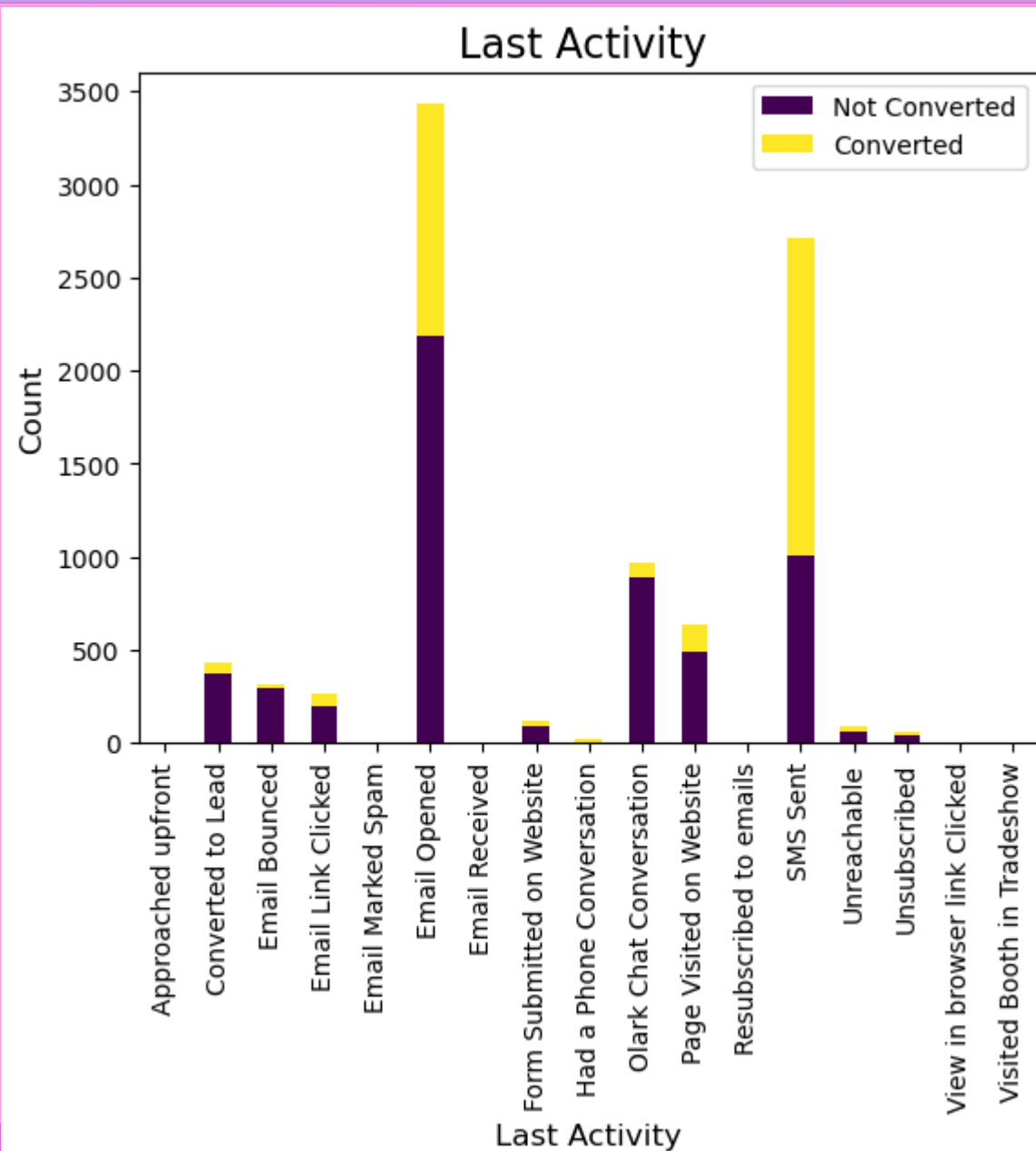
Numeric variables



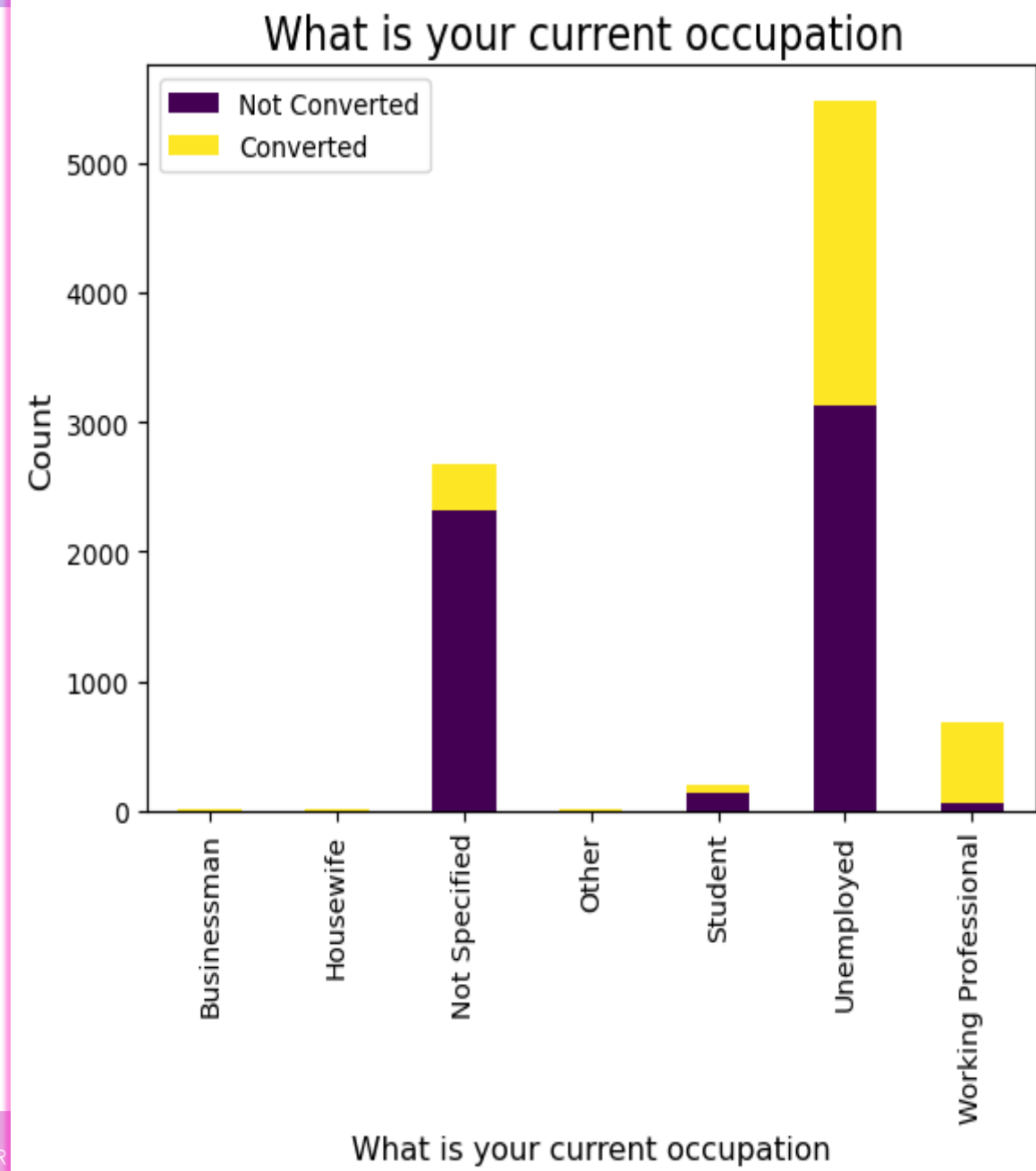
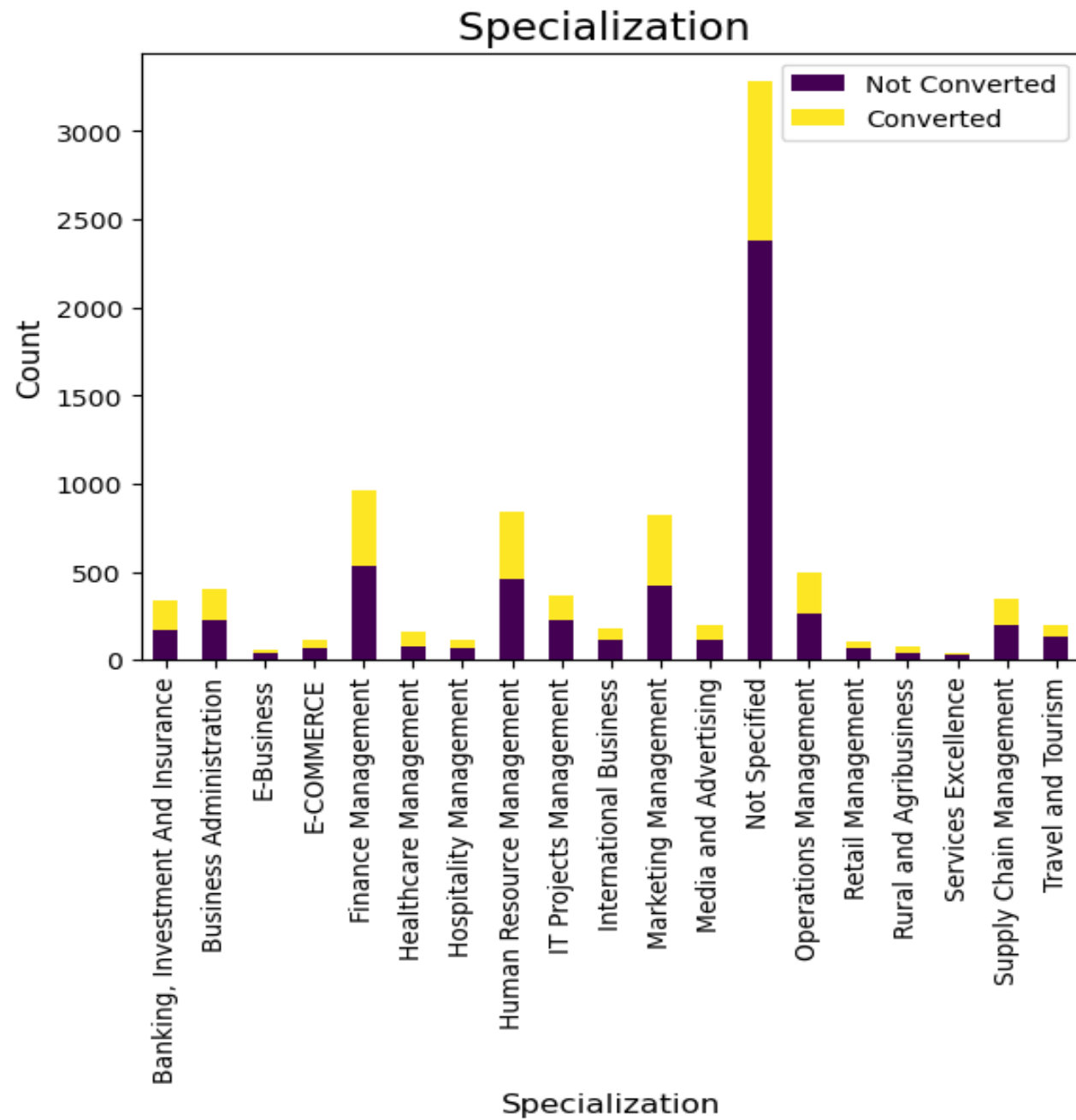
Categorical variables



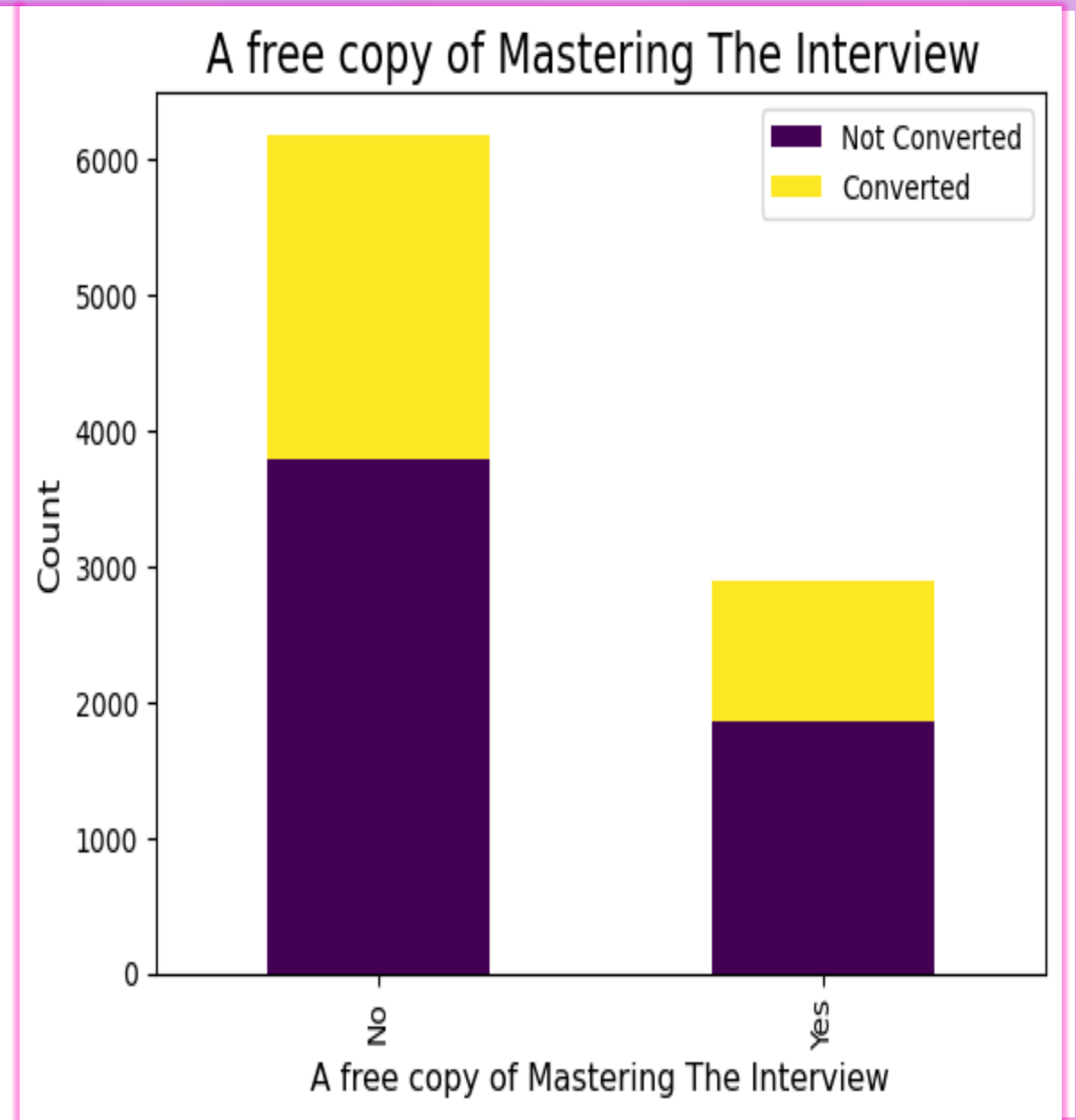
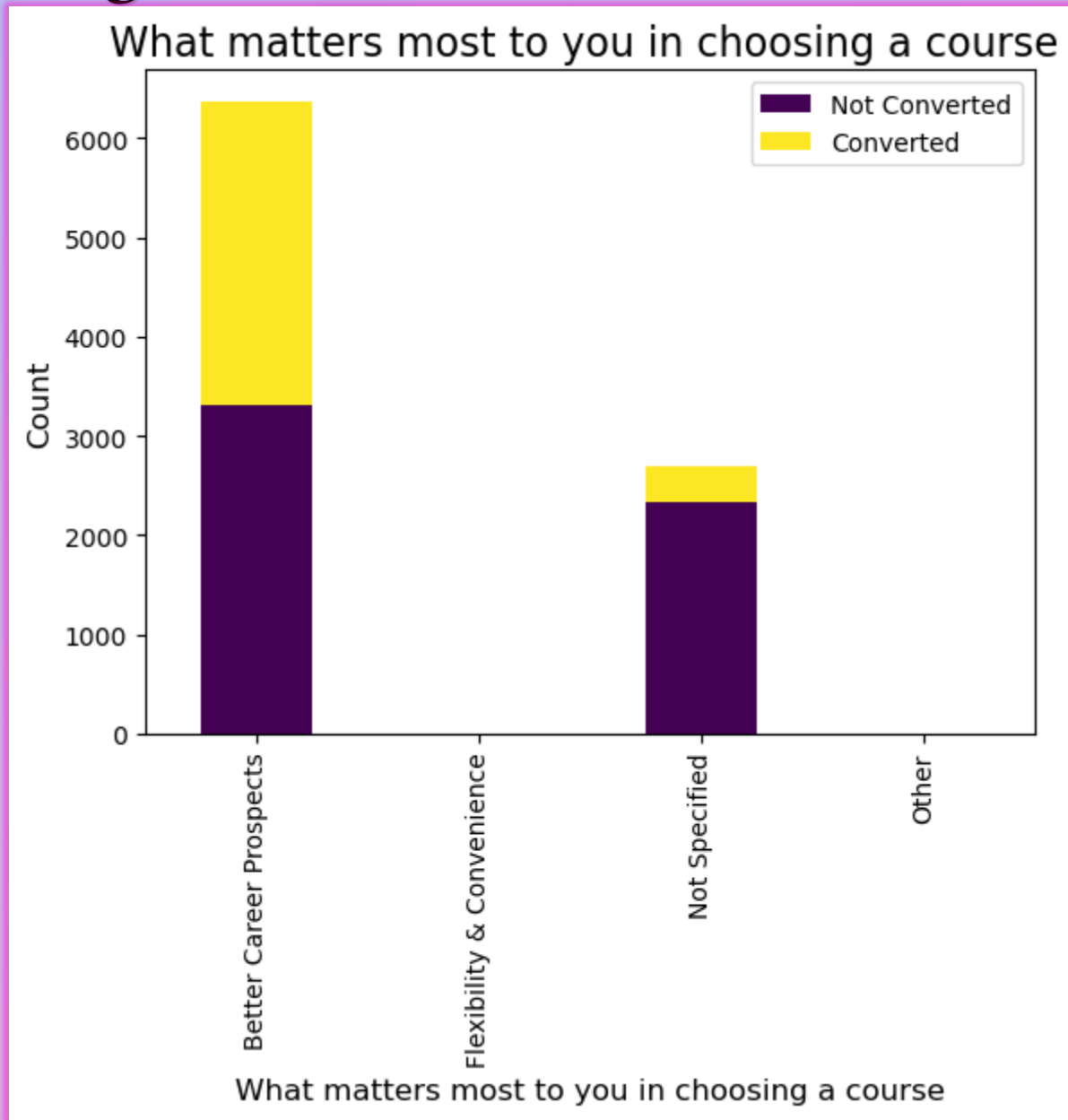
Categorical variables



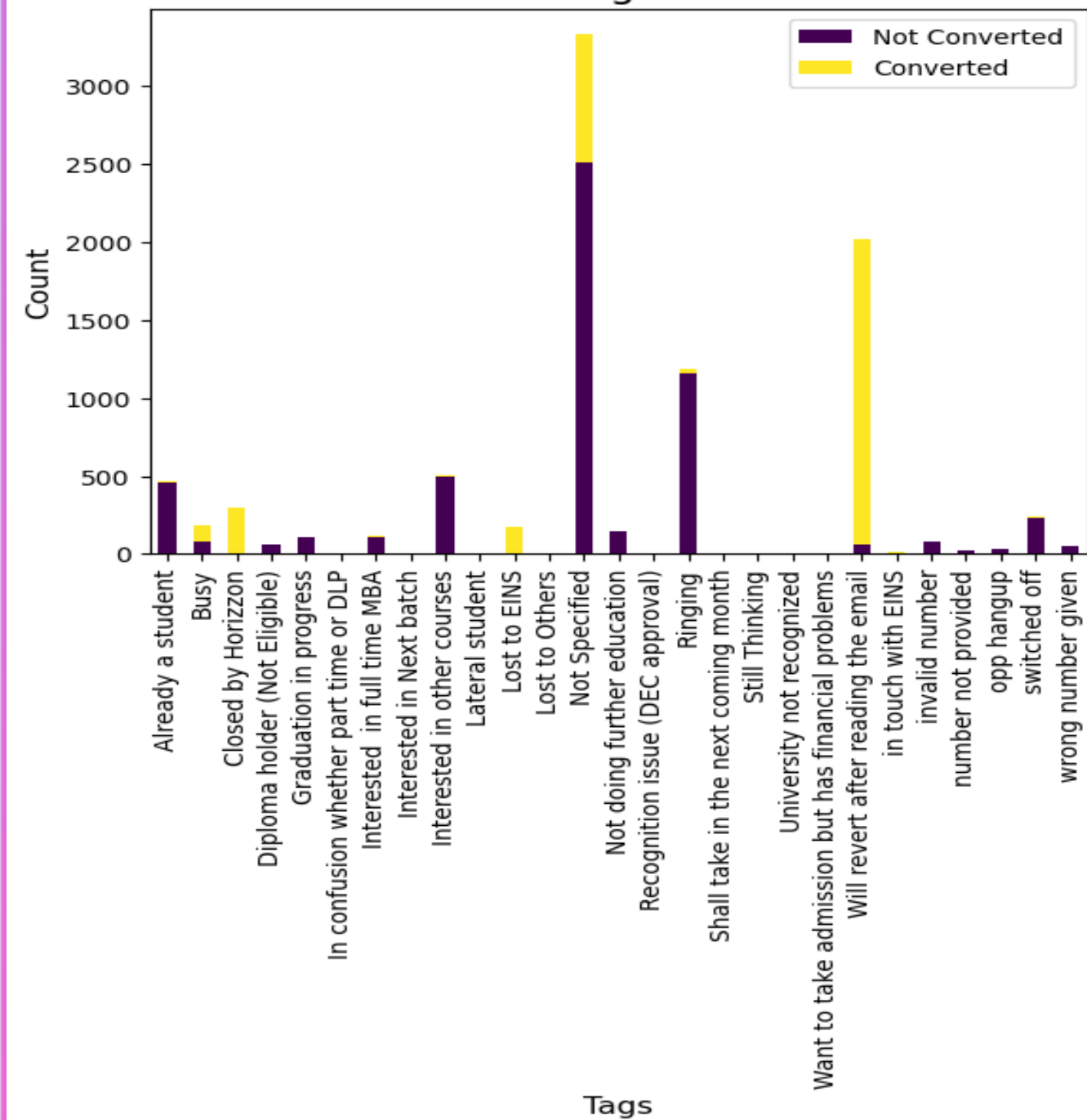
Categorical variables



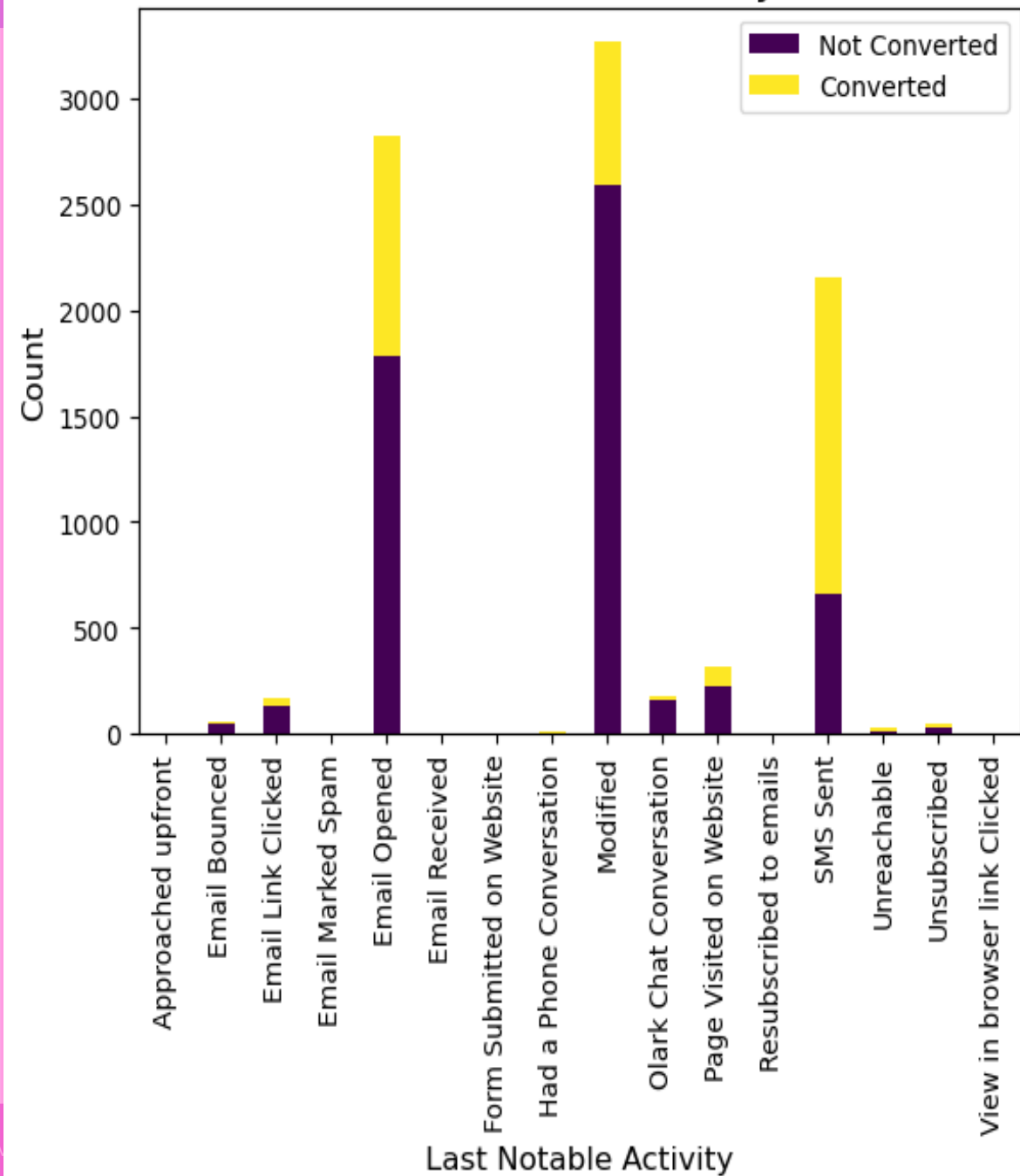
Categorical variables



Tags

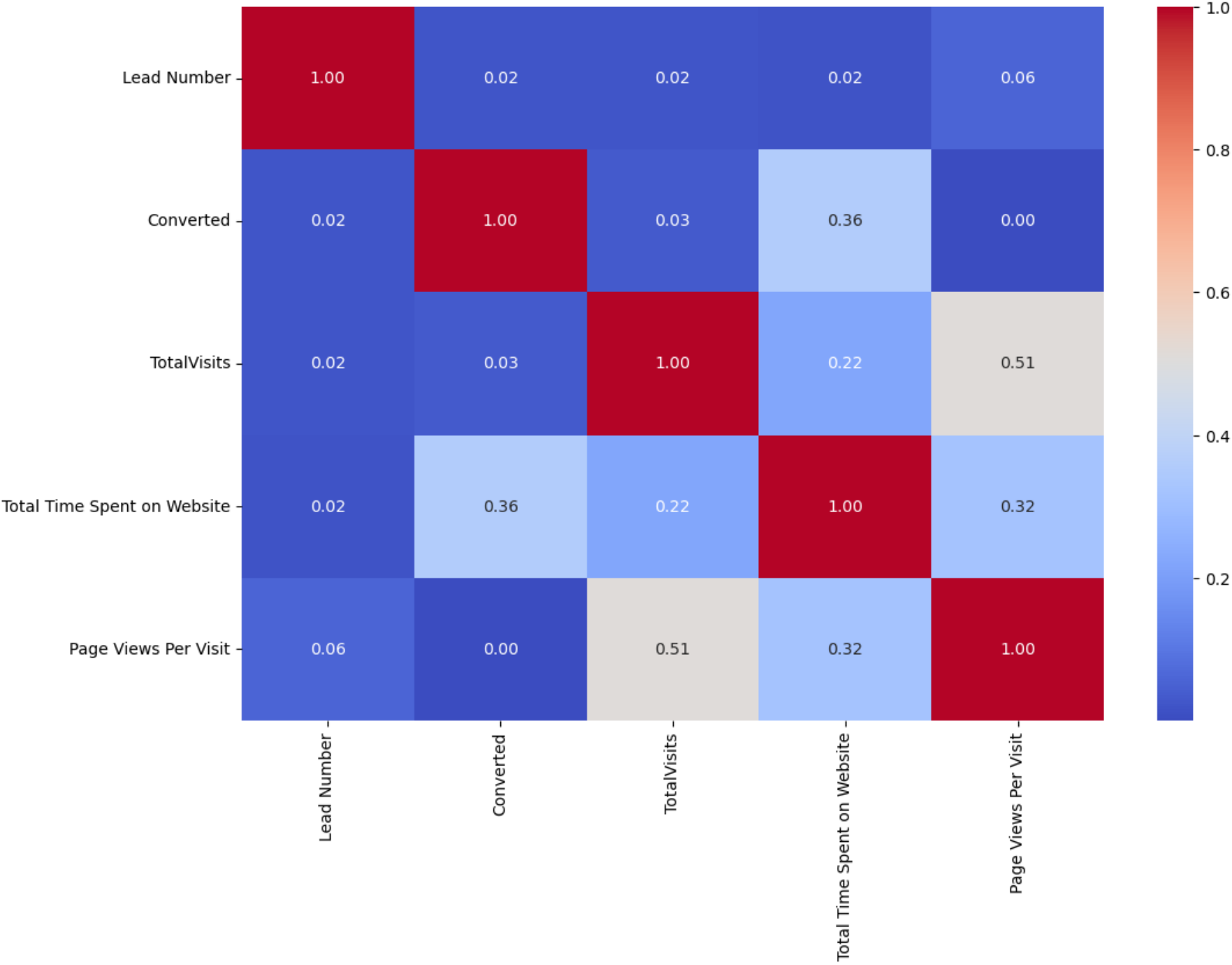


Last Notable Activity

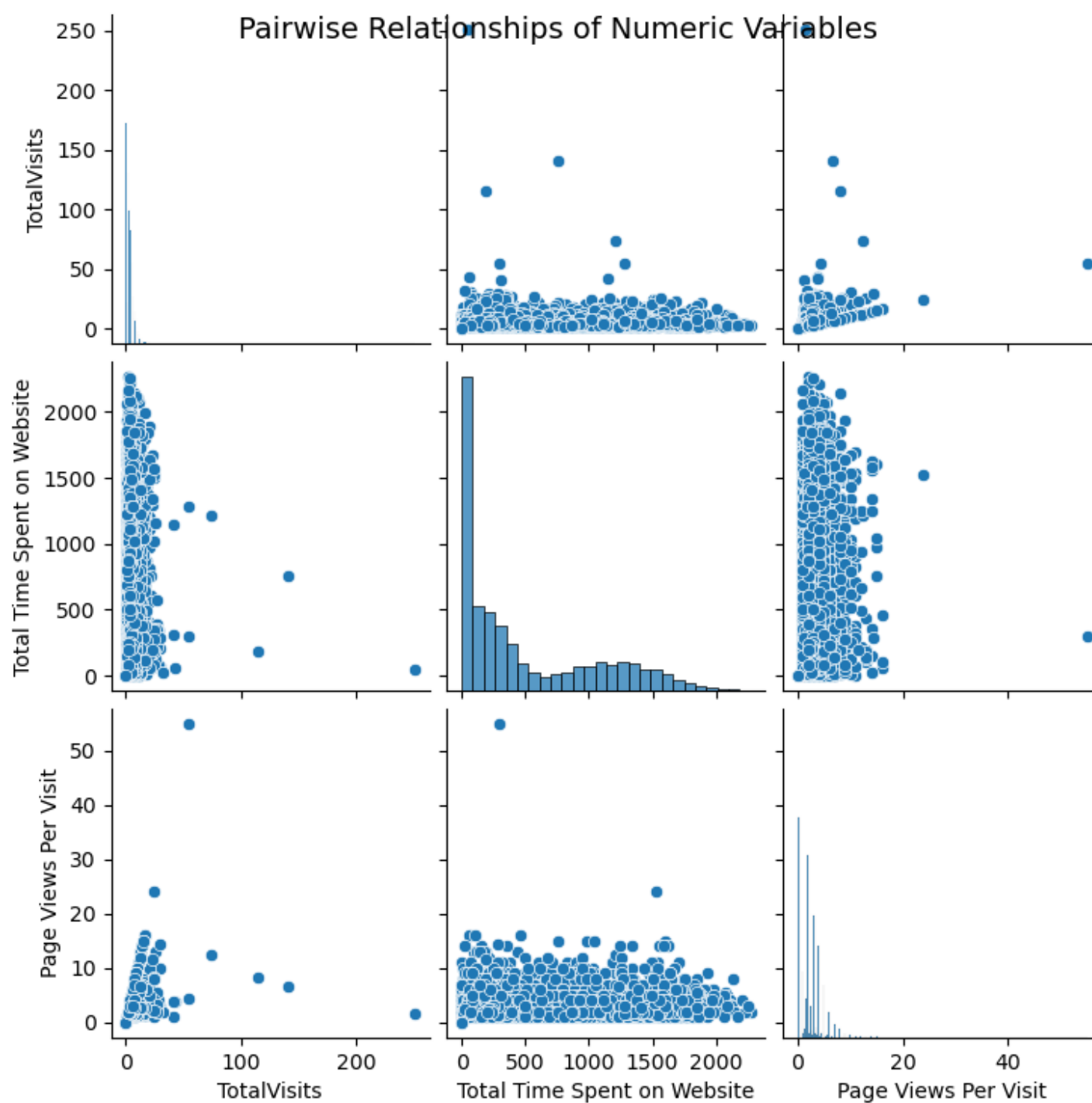


MULTIVARIATE ANALYSIS

Correlation Matrix



Pairwise Relationships of Numeric Variables



DATA CONVERSION

Converted the following columns to binary variables (0,1)

Do Not Email

A free copy of Mastering The Interview'

Created dummy variables- one hot encoded for the following columns-

Lead Origin

What matters most to you in choosing a course

Lead Source

Last Activity

Country

Specialization

Tags

What is your current occupation

'Last Notable Activity'

MODEL BUILDING

Initial dataset consisted of 9074 rows and 118 columns.

Performed feature scaling on the dataset.

Observed a conversion rate of approximately 38%.

Conducted train-test split with a ratio of 70% for training data and 30% for testing data.

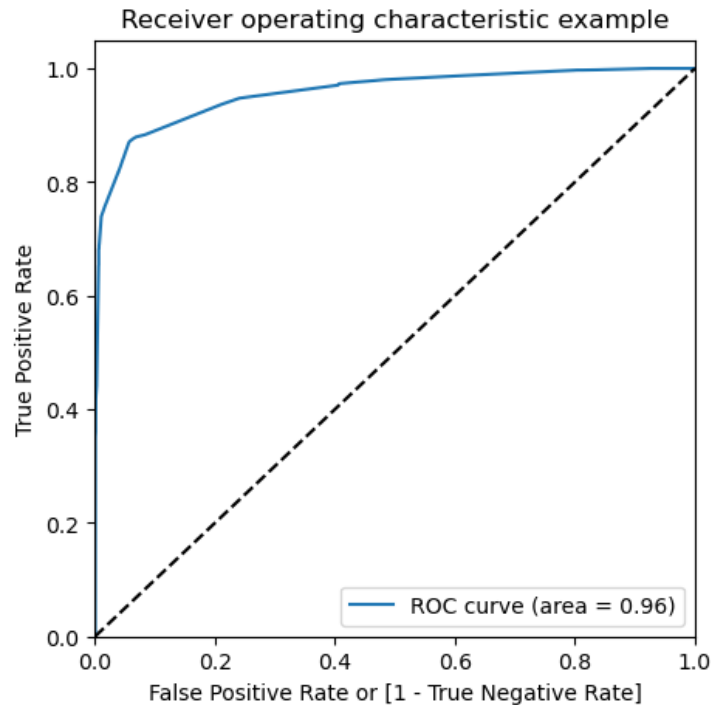
Utilized Logistic Regression model for prediction.

Employed Recursive Feature Elimination (RFE) technique with 15 variables as output for feature selection.

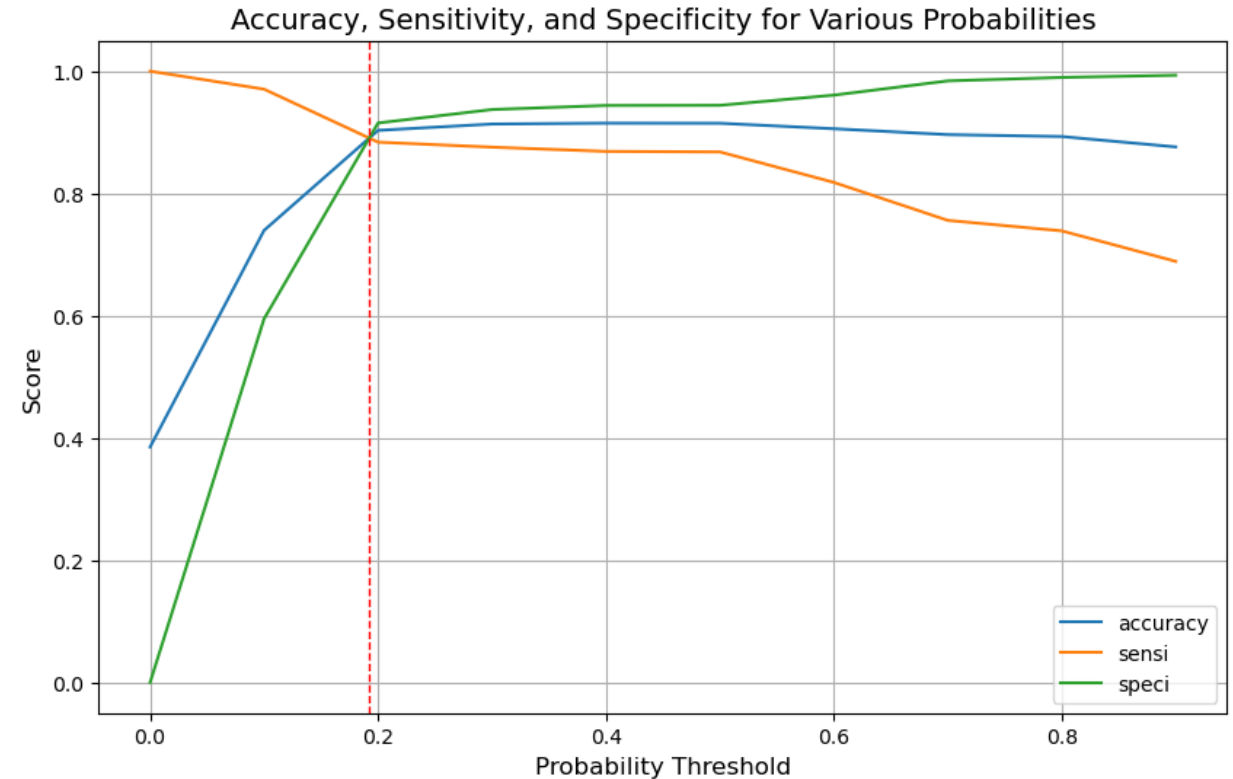
Developed models iteratively, removing features with p-values greater than 0.05 and Variance Inflation Factor (VIF) greater than 5.

The 7th model iteration was finalized for predictions on the test data.

ROC CURVE



Accuracy, Sensitivity, And Specificity Curves



- With the current cut off as seen in the graph as **0.192**, the **accuracy** is **90.32%**, **sensitivity** is **88.39%** and **specificity** is **91.52%**
- The **precision** is **86.72%** , **Recall** is **88.39%** and **F1** score is **88.39%**.

Predictions on the Test data

Performed scaling on the required variables as part of preprocessing.

Trained the test set using the final model derived from the feature selection process.

Utilized the trained model to make predictions on the test data.

Employed a cutoff value of 0.192 for classification.

Achieved an **Accuracy** of **89.94%** on the test data.

Sensitivity, representing the true positive rate, stood at **87.46%**.

Specificity, representing the true negative rate, was measured at **91.35%**.

Conclusion

The most important variables that mattered the most in the potential buyers are –

- The Total time spent on the website and the Page views per visit
- The 'Lead Add Form' is the lead origin with the highest proportion of converted leads among all the leads.
- When the Lead source is
 - Google
 - Direct Traffic
 - Organic Search
 - Weliangk website
- When the Last Activity was –
 - Reading the 'SMS Sent'
 - Opening the 'Emails' received
- When their current occupation is a working professional.
- When the tag is 'Will revert after reading the email' , 'Closed by Horizon', and 'Lost to EINS'



THANK YOU!