

GROUP PROJECT

Lead Score Case Study

By:

- Nguyen Huy Thang
- Bala Subramanyam
- Nairit Dutta

Problem Statement

- X Education sells online courses to industry professionals.
- X Education gets a lot of leads, its lead conversion rate is very poor.
- To make this process more efficient, the company wishes to identify the most potential leads, also known as 'Hot Leads'.
- If they successfully identify this set of leads, the lead conversion rate should go up as the sales team will now be focusing more on communicating with the potential leads rather than making calls to everyone.

Business Objective

- Increase the current lead conversion rate from 30% to a target of 80% by focusing on the most promising leads.
- Build a logistic regression model that assigns a lead score (between 0 and 100) to each lead. This score will help in identifying 'Hot Leads' who are more likely to convert into paying customers.
- Ensure the model is adaptable to accommodate potential changes in the company's requirements.

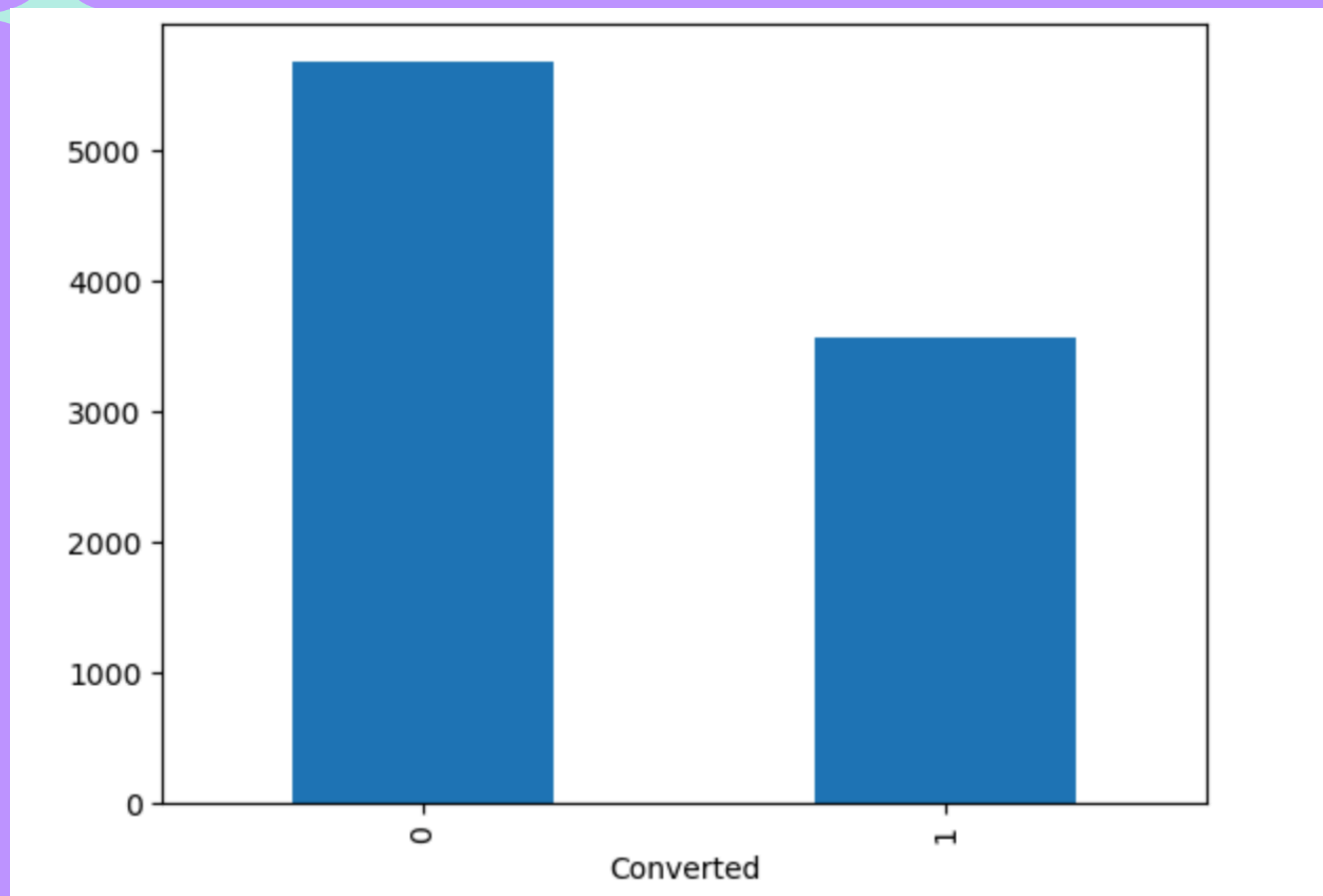
Solution Methodology

- Exploratory Data Analysis and Data Cleaning
- Data Preparation
- Model Building
- Model Evaluation
- Conclusions and Recommendations for the Company Strategy

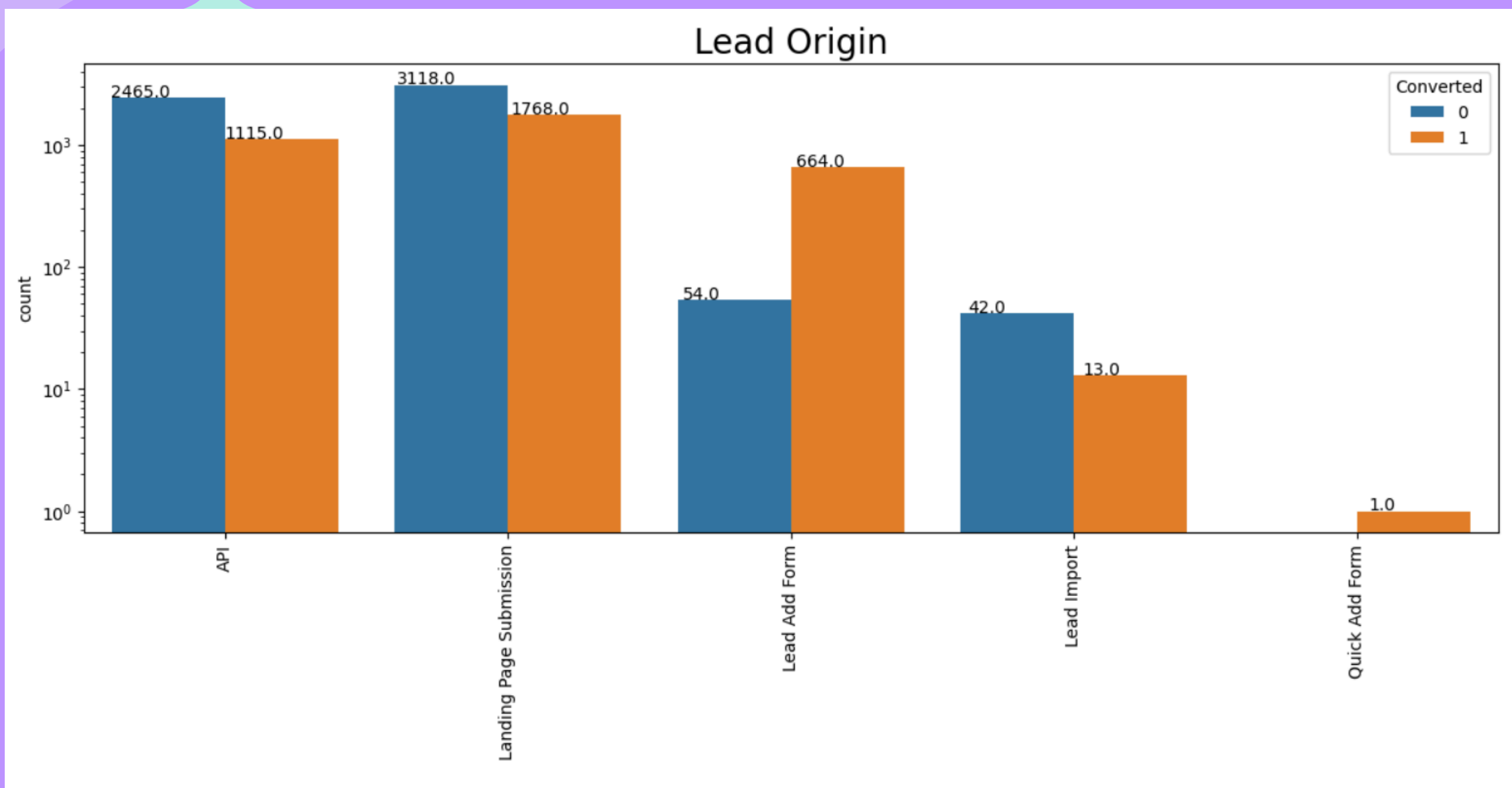
Data Cleaning

- Total number of rows is 37, total number of columns is 9240.
- Eliminate redundant variables.
- Delete columns with null values greater than 40%.
- Remove outlier values.

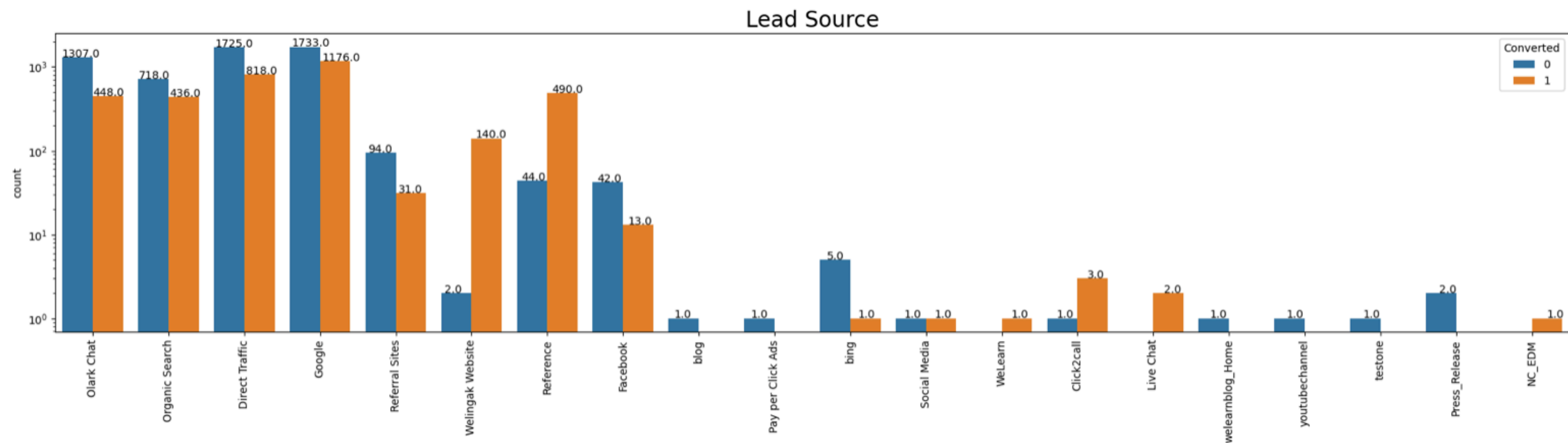
EDA



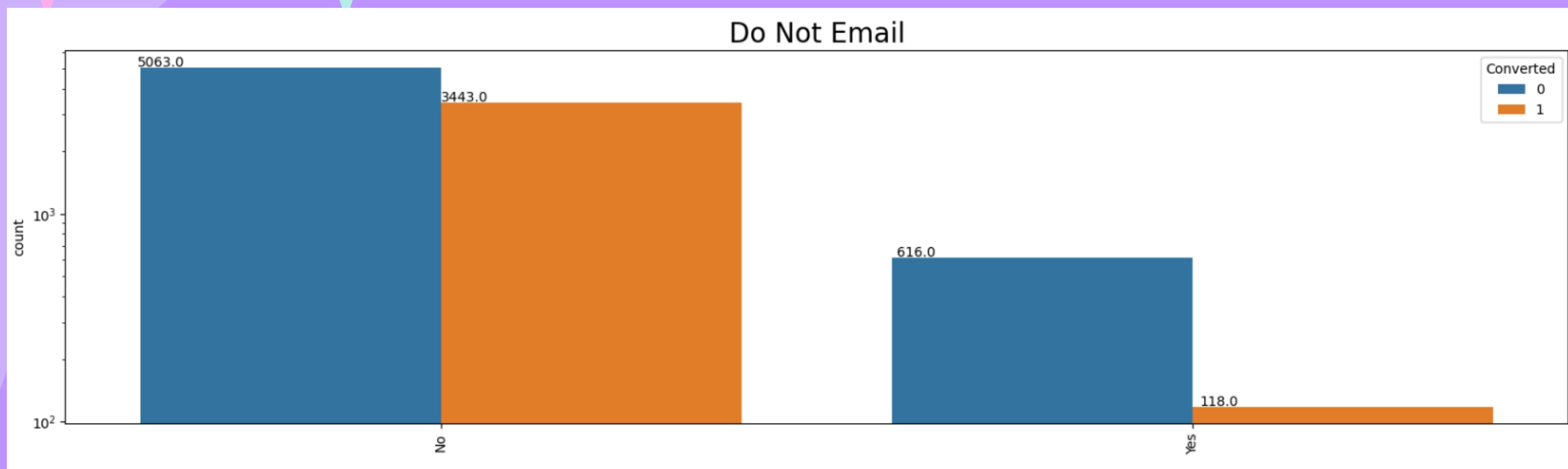
EDA



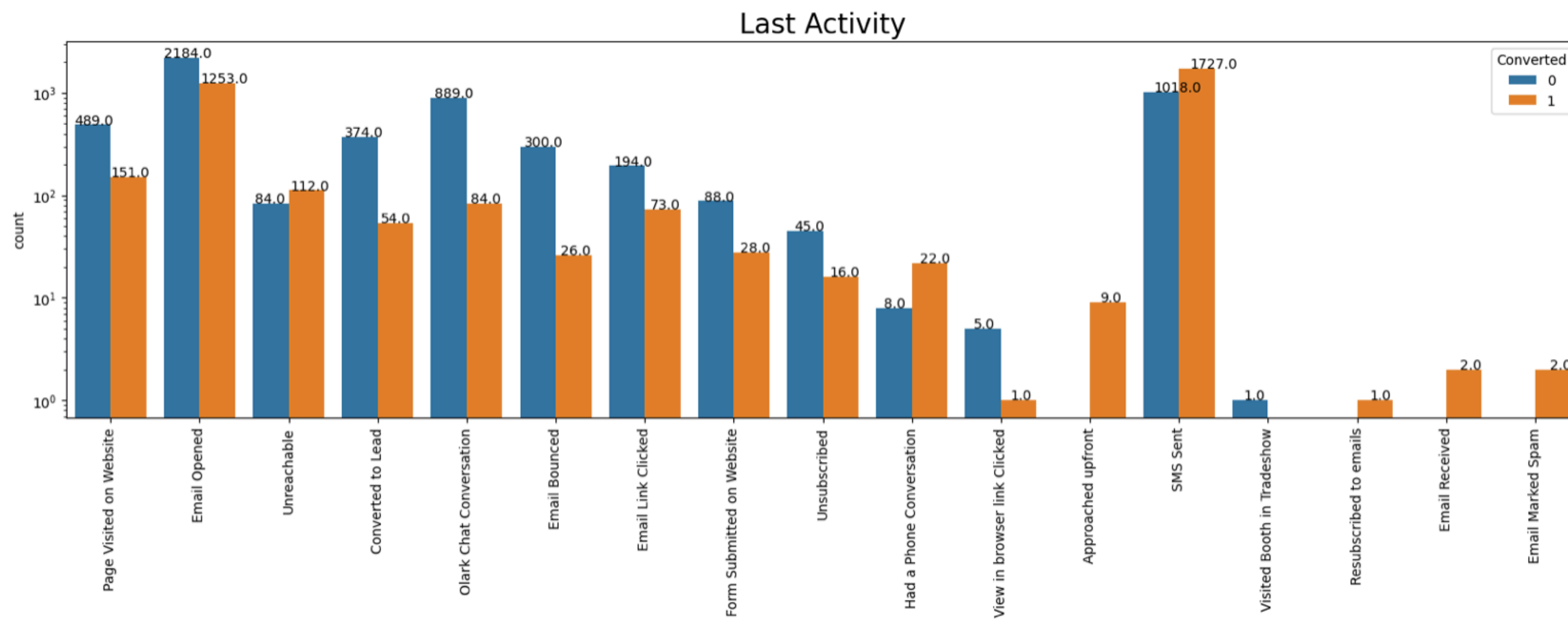
EDA



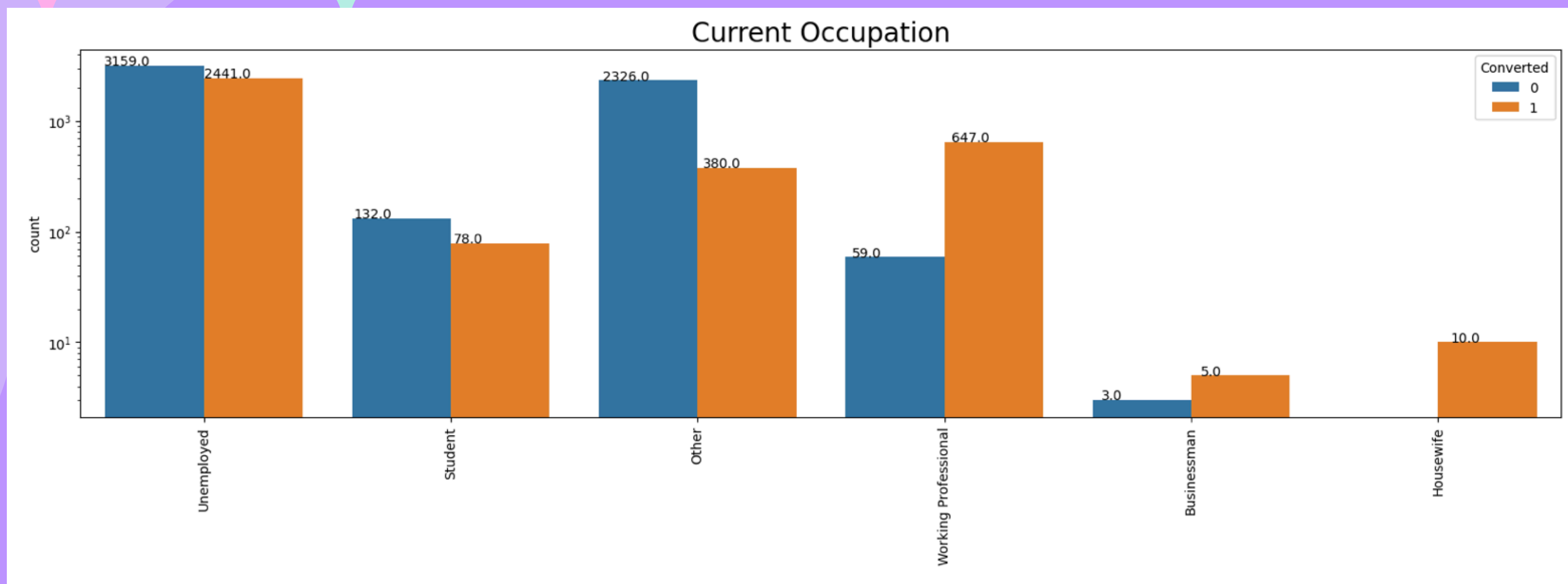
EDA



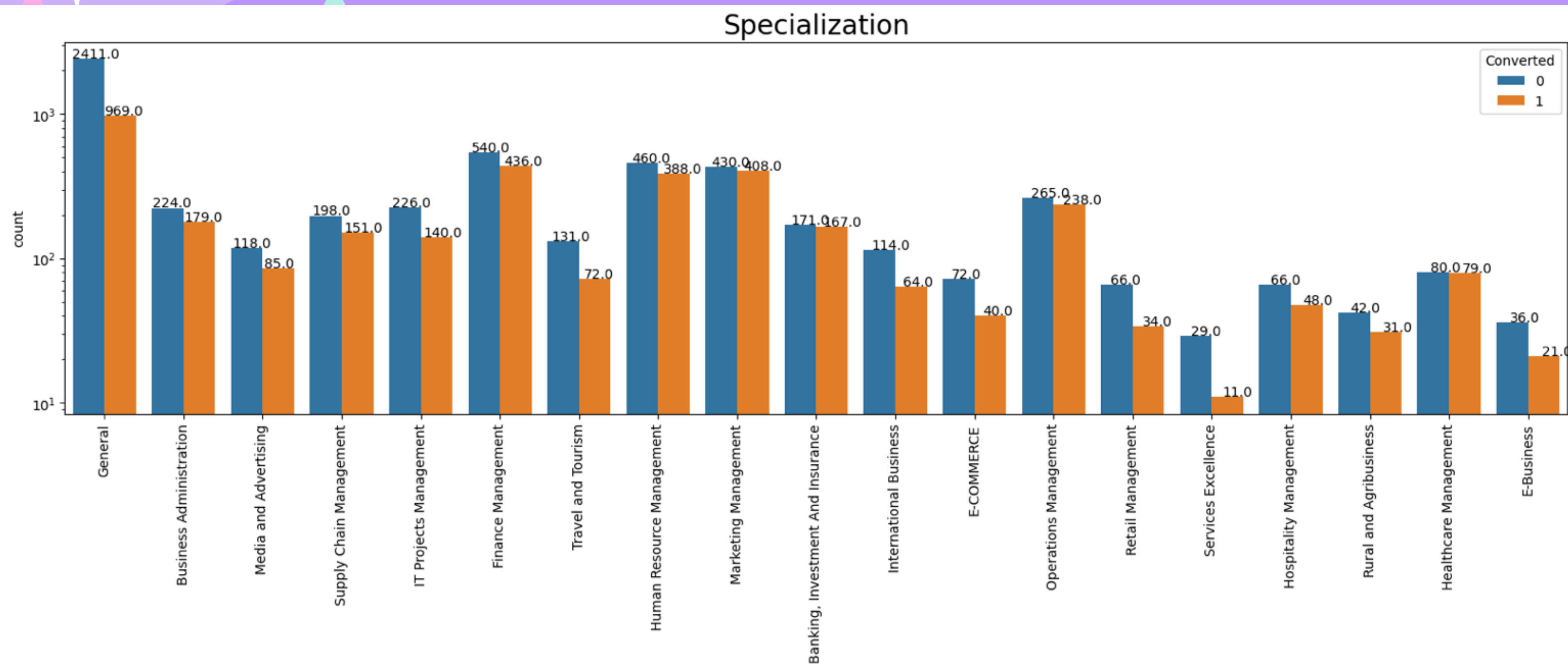
EDA



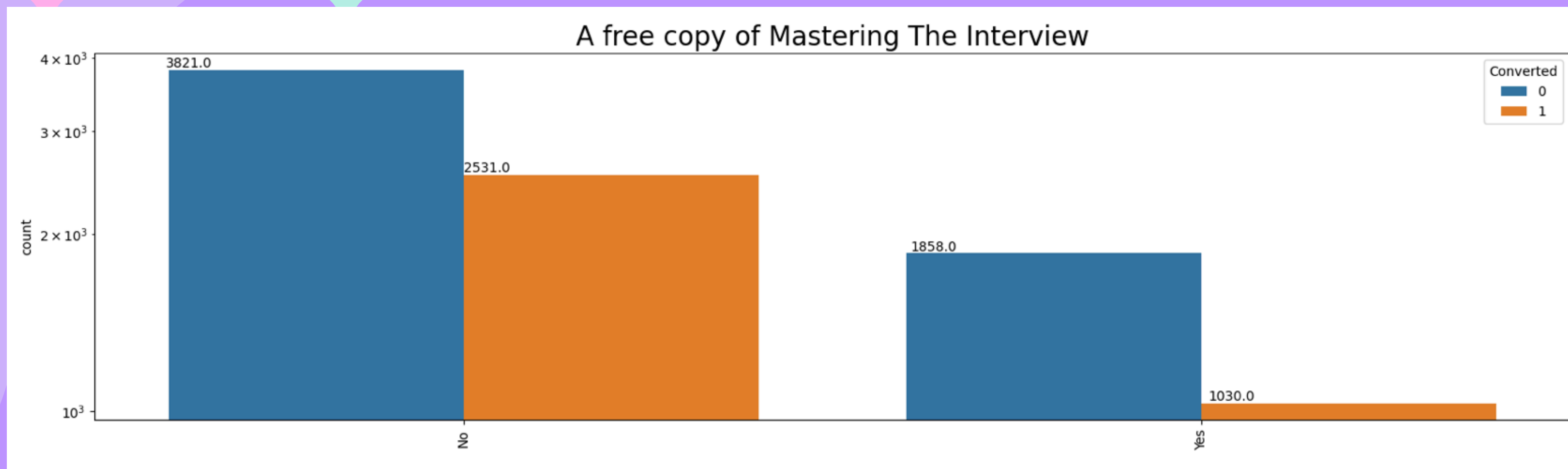
EDA



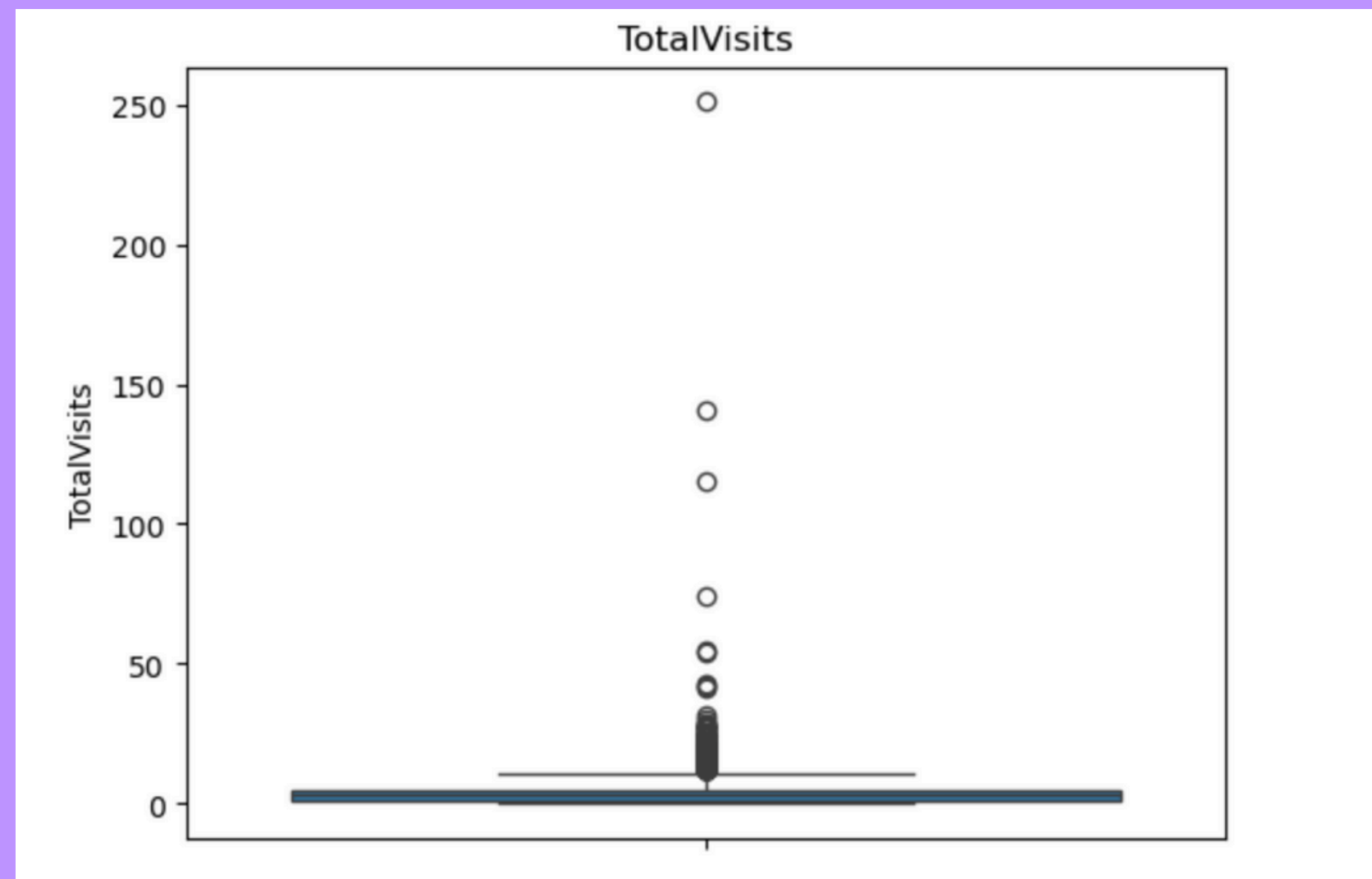
EDA



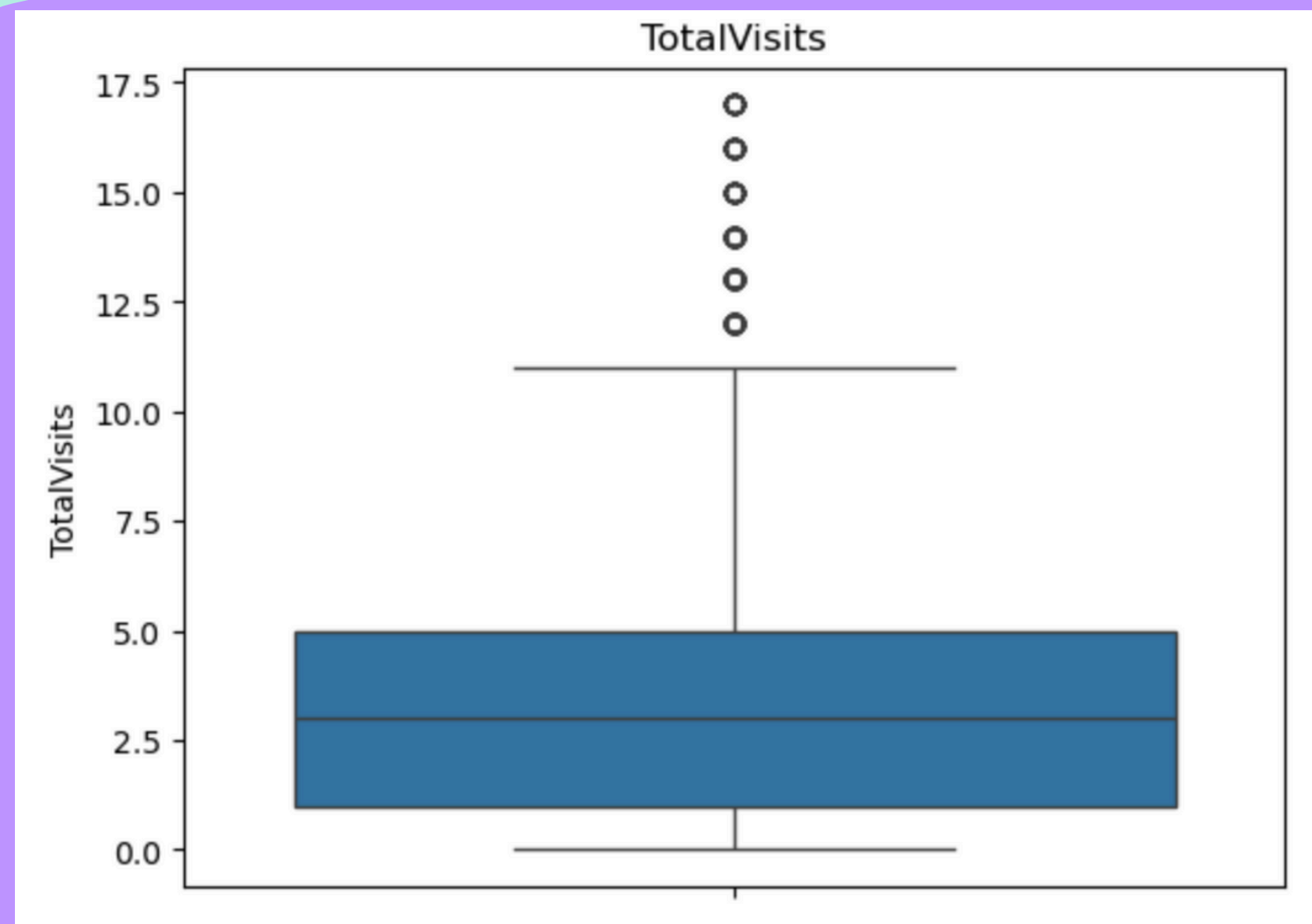
EDA



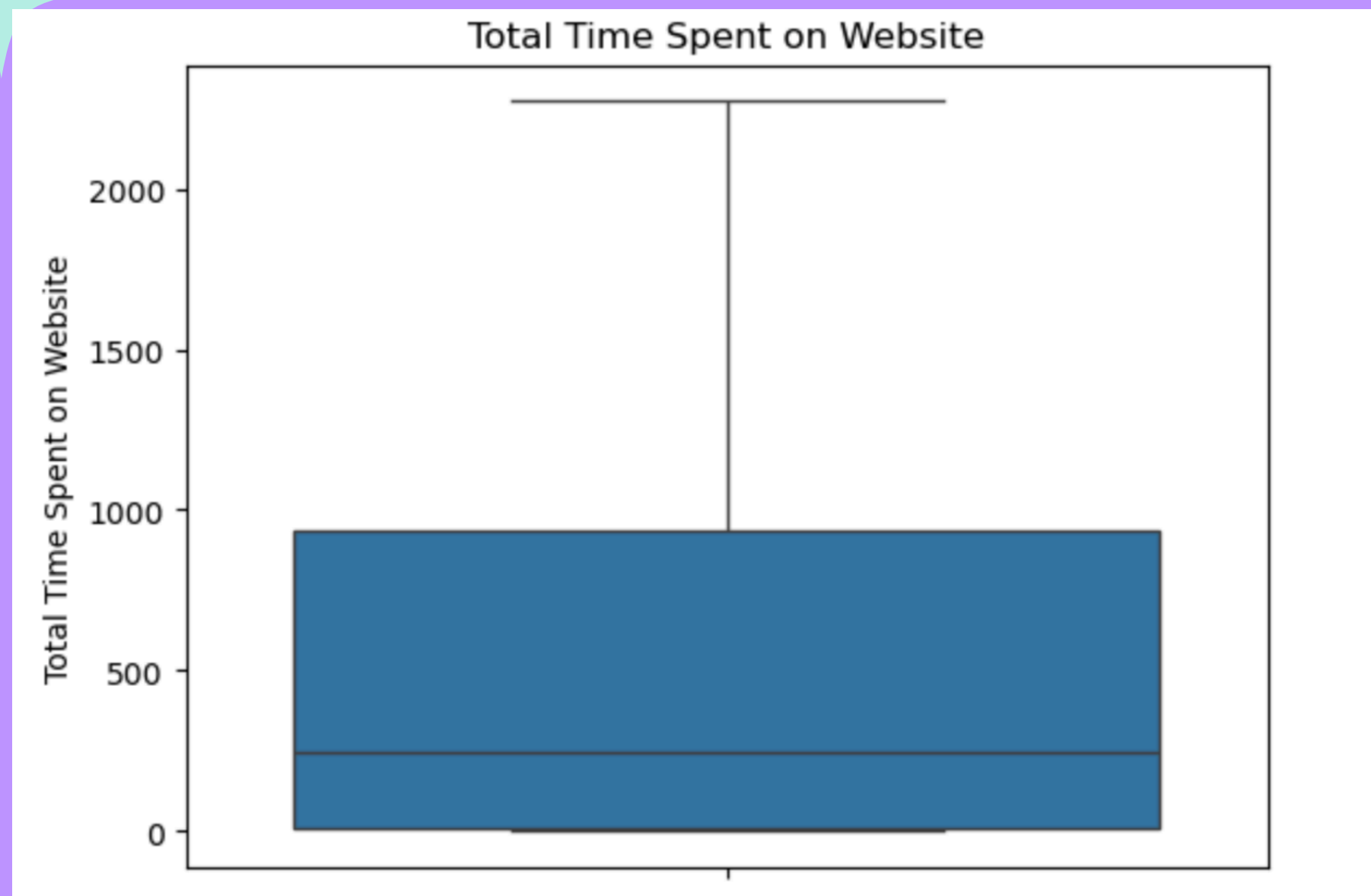
EDA



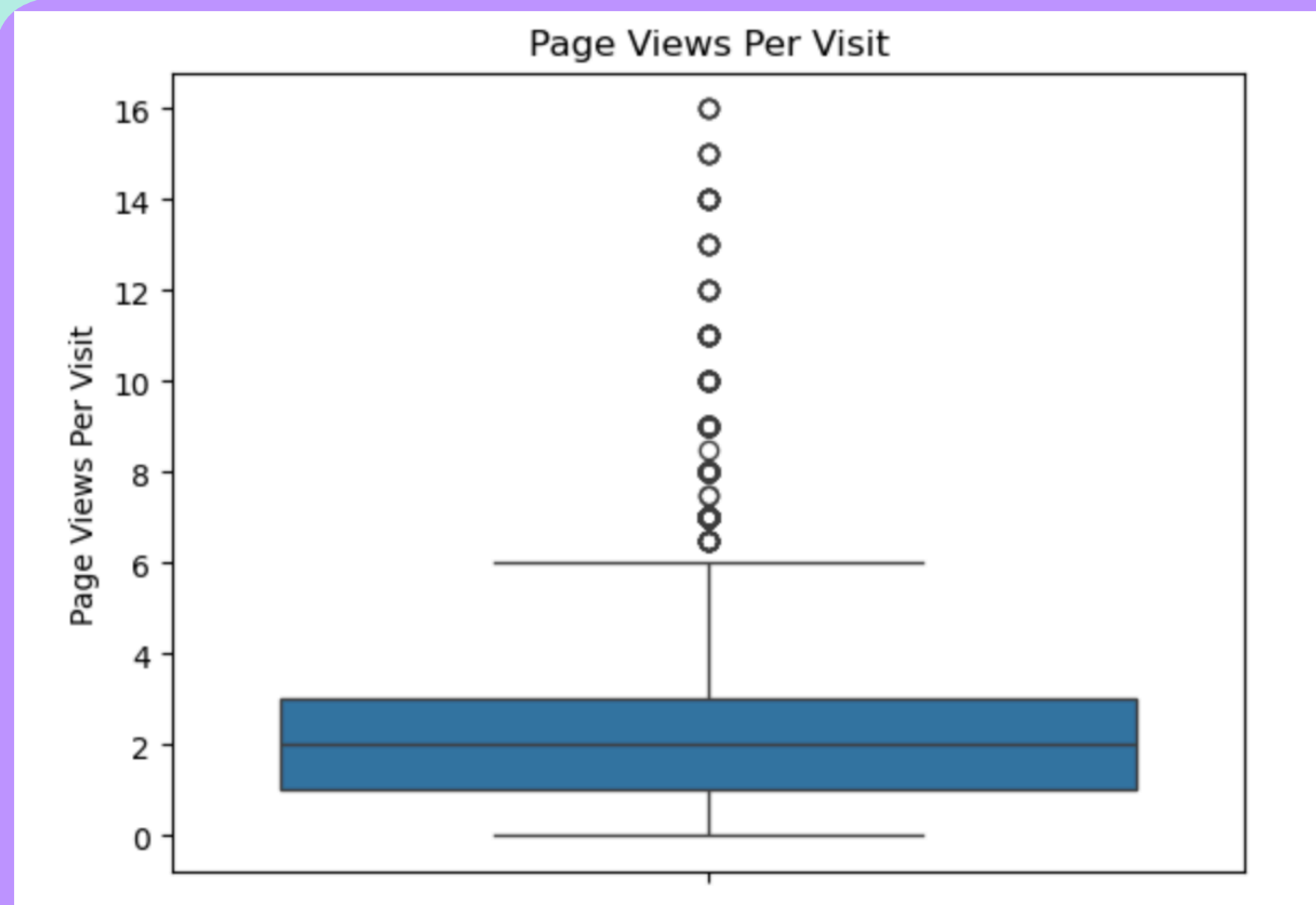
EDA



EDA



EDA



EDA



Data Preparation

- Categorical variables are converted to 0 or 1.
- Split the dataset into train and test sets with a 70:30 ratio.
- Data points are standardized using the StandardScaler for numerical features to avoid bias from variables with higher scales in the model.

Model Building

- The initial model is created by selecting 15 features using the RFE technique.
- Insignificant variables have been identified and removed to optimize the model for better performance. VIF is also checked to assess multicollinearity among the model's features.
- An important step in logistic regression is to find the optimal threshold for probability that aligns with business needs and ensures good accuracy, sensitivity, and specificity. We obtained 0.398 as the cutoff point from the chart that balances accuracy, sensitivity, specificity, and probability range.
- Precision and recall views were then examined to refine thresholds for predictions in the test dataset. The cutoff point obtained from the Recall-Precision chart was 0.445. Model evaluation was completed by checking all performance metrics, including the ROC curve, accuracy, sensitivity, specificity, recall, and precision, all of which were found to be at acceptable levels. Optimal point 0.398 was selected as it had higher sensitivity and recall of 81% .

Model Evaluation

Predicted values are calculated using the model. The lead score is 100 times the predicted log odds, ranging from 0 to 100.

Conclusions and Recommendations

- The model evaluation steps indicate that the metrics for accuracy, precision, and recall are all at acceptable levels. The recall score is also slightly higher than the precision score, which aligns well with future business needs.
- Focusing on leads from the "Add" form is crucial, as they contribute the most to conversions. Prioritizing these leads can enhance your conversion rates significantly.
- Professionals should be the ideal target.
- Leads from the Welingak website have potential for conversion, so advertising revenue should be focused on this website.
- Individuals who do not wish to receive emails should be avoided.
- Avoid individuals without specific expertise.



Thank You