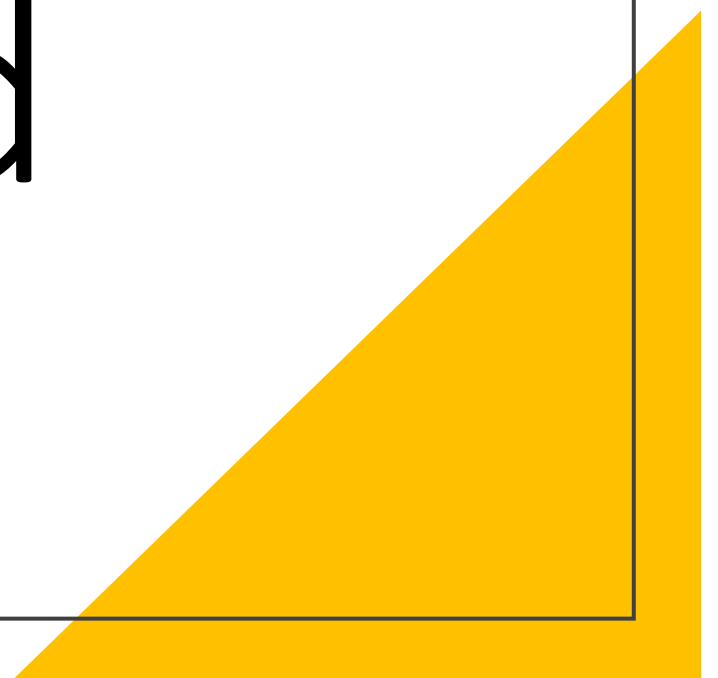


Lead Score Case Study-Upgrad

Group Members:

Avinandita Mahanti

Ashwini Mukati



Approach towards the Solution

1. Check and handle duplicate data:

1. Identify and remove any duplicate records from the dataset to ensure data integrity and avoid bias in analysis.

2. Check and handle NA values and missing values:

1. Identify columns with missing values (NA) and decide on appropriate strategies to handle them (e.g., imputation, removal, or treating them as a separate category).

3. Drop columns with a large amount of missing values and not useful for analysis:

1. Identify columns with a high percentage of missing values and consider dropping them if they don't contribute meaningful information.

4. Imputation of missing values if necessary:

1. For columns with missing values, consider imputing the missing values using techniques like mean, median, mode, or regression imputation, based on the data's characteristics.

5. Check and handle outliers in data:

1. Identify outliers in the data and decide on appropriate strategies for handling them (e.g., removing, transforming, or treating them as special cases).

6. Exploratory Data Analysis (EDA)

1. Perform univariate data analysis: Examine the distribution, value counts, and summary statistics of individual variables to gain insights into their characteristics.
2. Perform bivariate data analysis: Analyze the relationships between pairs of variables, such as correlation coefficients, scatter plots, or other visualizations, to understand their patterns and potential dependencies.

7. Feature Scaling & Dummy Variables and encoding of the data:

1. If needed, apply feature scaling (e.g., normalization or standardization) to ensure all variables are on a similar scale.
2. Create dummy variables or perform encoding to convert categorical variables into numeric format suitable for the model.

8. Classification technique: Logistic Regression

1. Train a logistic regression model on the preprocessed data to predict the binary outcome variable.

9. Validation of the model:

1. Split the dataset into training and testing sets or use cross-validation techniques to evaluate the model's performance on unseen data.
2. Calculate relevant metrics like accuracy, precision, recall, F1-score, and ROC-AUC to assess the model's performance.

10. Model presentation:

- Present the logistic regression model and its results, including coefficients, significance, and model performance metrics.

11. Conclusions and recommendations:

- Summarize the findings from the analysis, highlighting important variables and their impact on the target variable.
- Provide actionable recommendations based on the insights gained from the model to improve decision-making and achieve the desired outcomes.
- It's important to note that each step may require iterations and adjustments based on the nature of the data and the objectives of the analysis. Additionally, proper validation and interpretation of the model results are essential for making informed decisions.

Data Pre-Processing

1. Data Size:

1. Total Number of Rows: 37
2. Total Number of Columns: 9240

2. Dropping Single Value Features:

1. Drop features with a single unique value, such as "Magazine," "Receive More Updates About Our Courses," "Update me on Supply Chain Content," "Get updates on DM Content," "I agree to pay the amount through cheque," etc.

3. Removing Unnecessary Columns:

1. Remove "Prospect ID" and "Lead Number" columns, which are not necessary for the analysis.

4. Dropping Features with Low Variance:

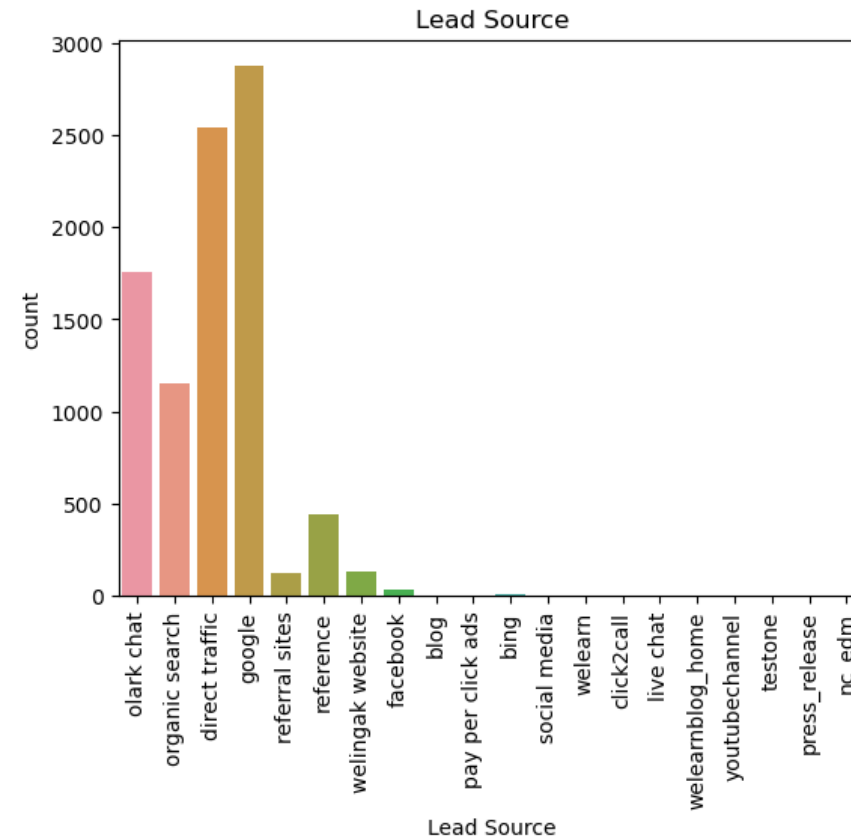
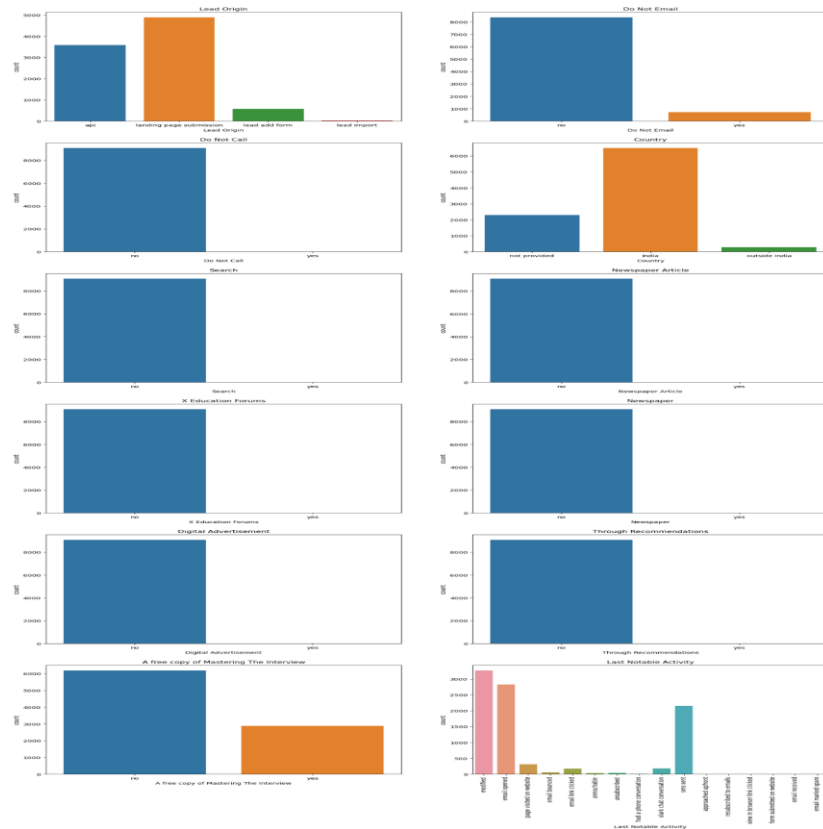
1. Drop features with low variance, such as "Do Not Call," "What matters most to you in choosing course," "Search," "Newspaper Article," "X Education Forums," "Newspaper," "Digital Advertisement," etc. These features may not provide valuable information for analysis due to lack of variation.

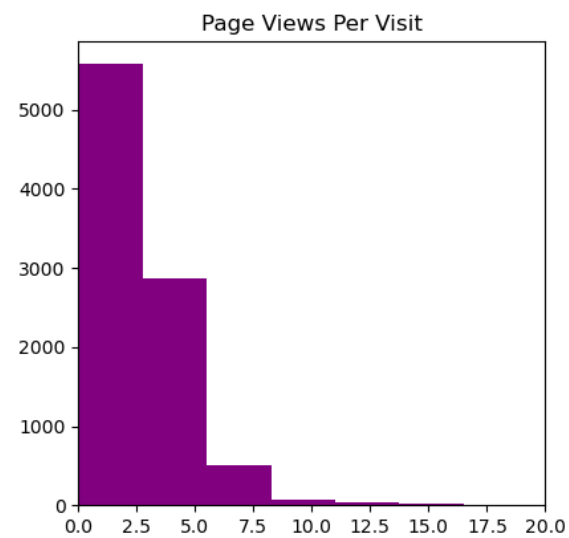
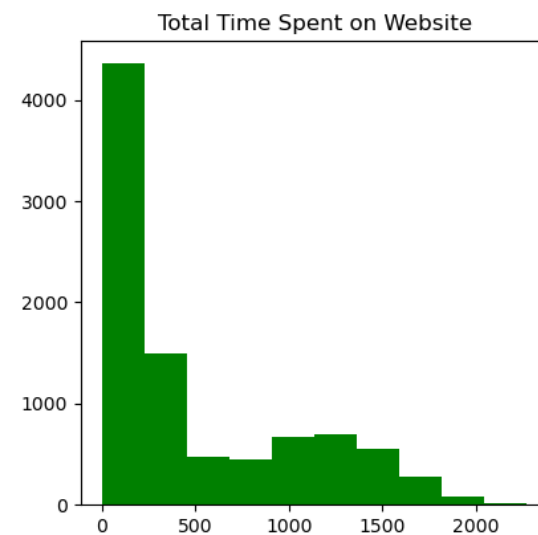
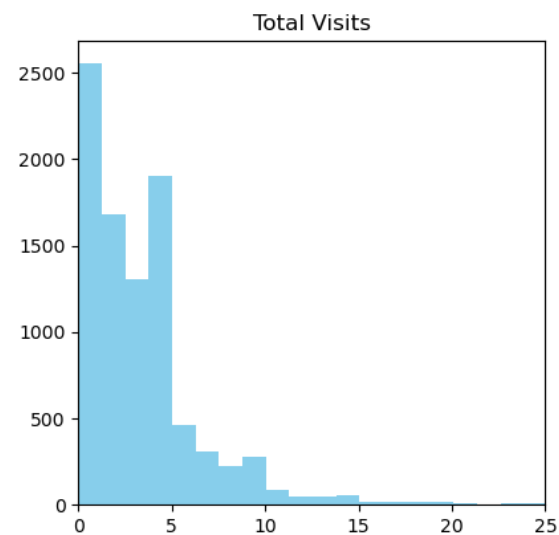
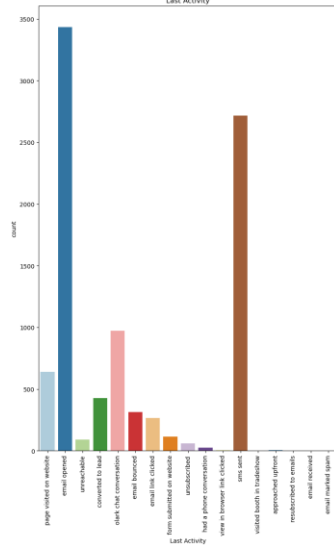
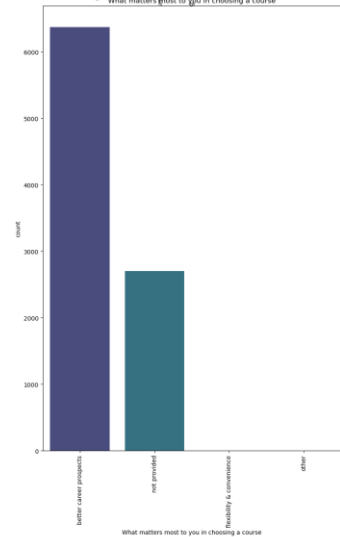
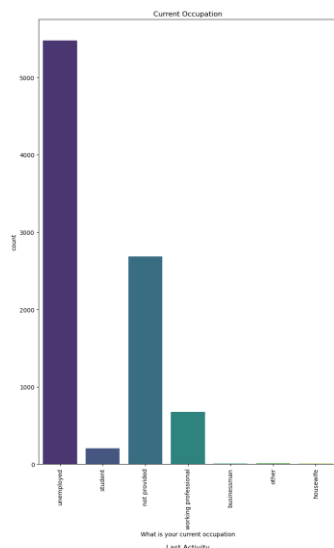
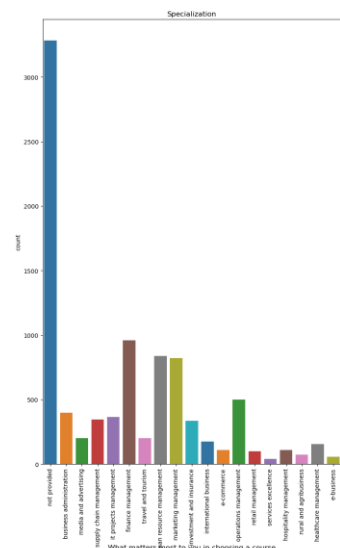
5. Dropping Columns with High Missing Values:

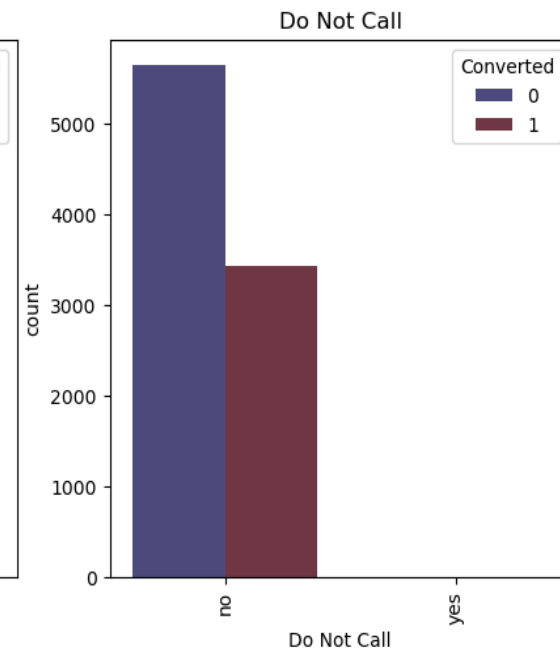
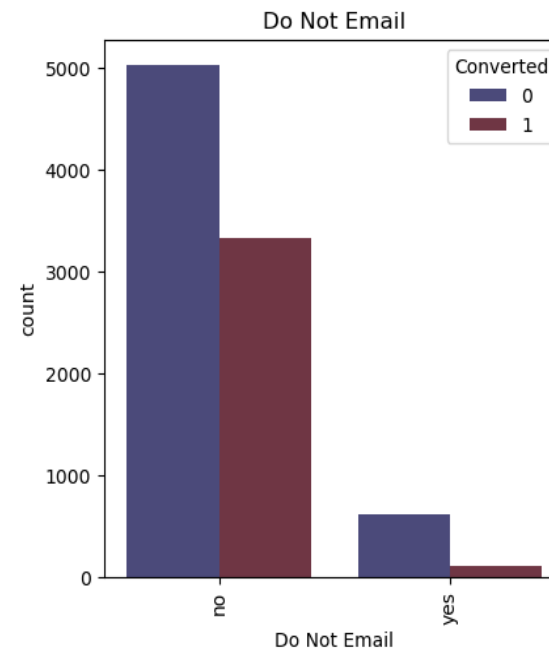
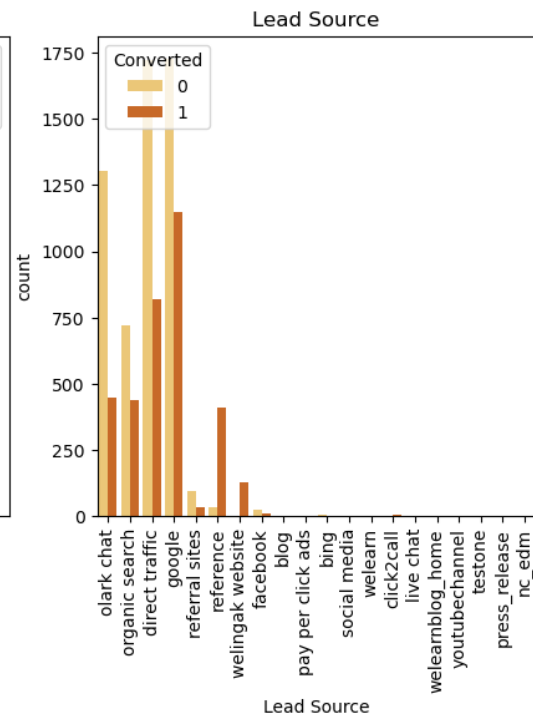
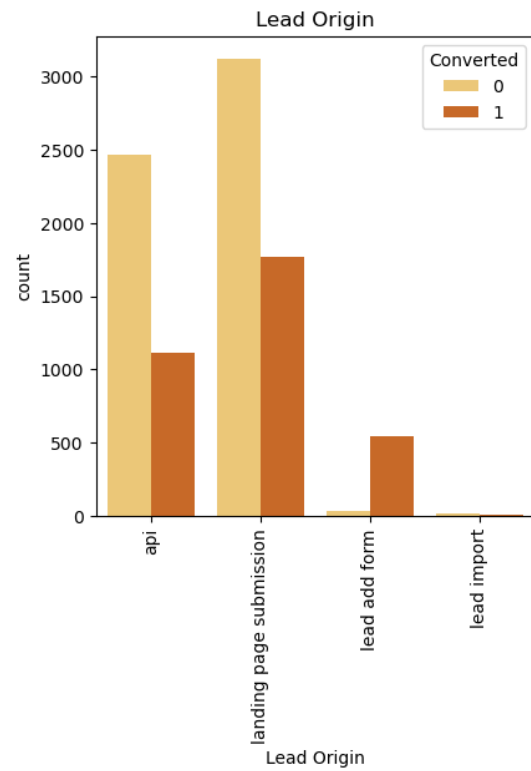
1. Drop columns with more than 35% missing values, such as 'How did you hear about X Education' and 'Lead Profile.' These columns have a significant amount of missing data, making them less informative for the analysis.

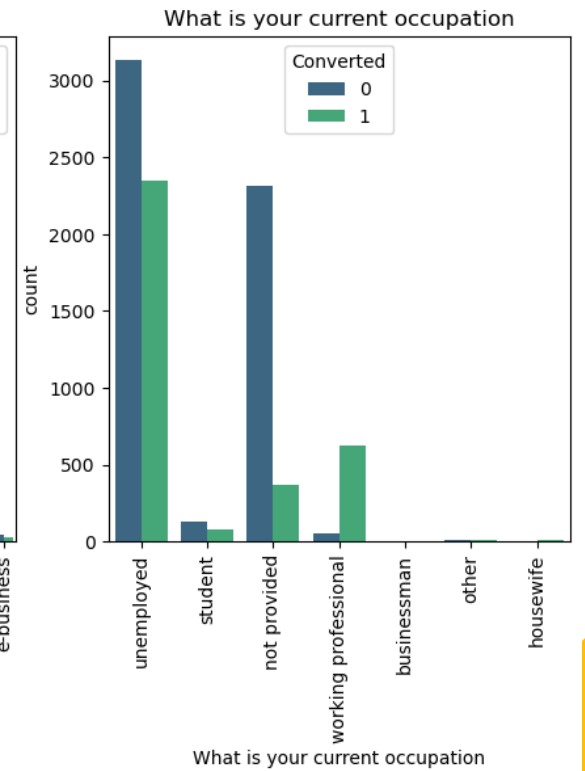
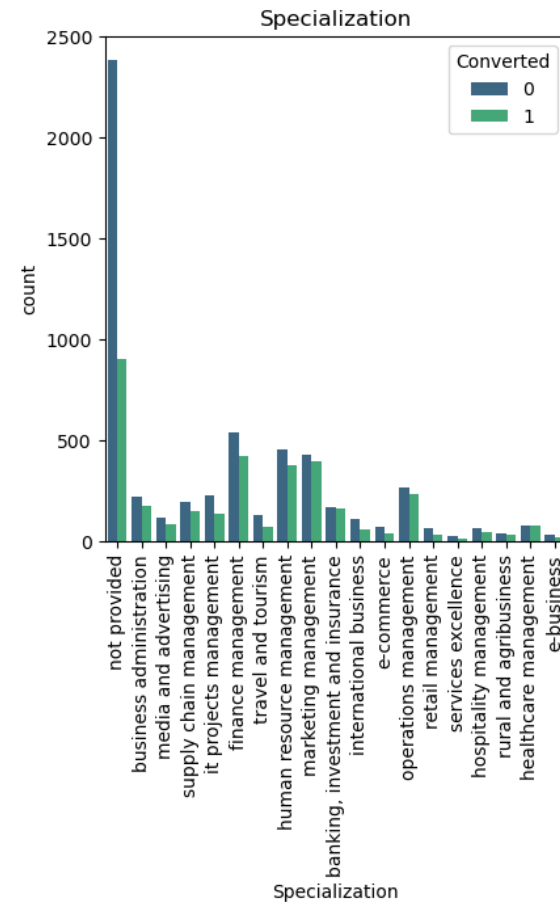
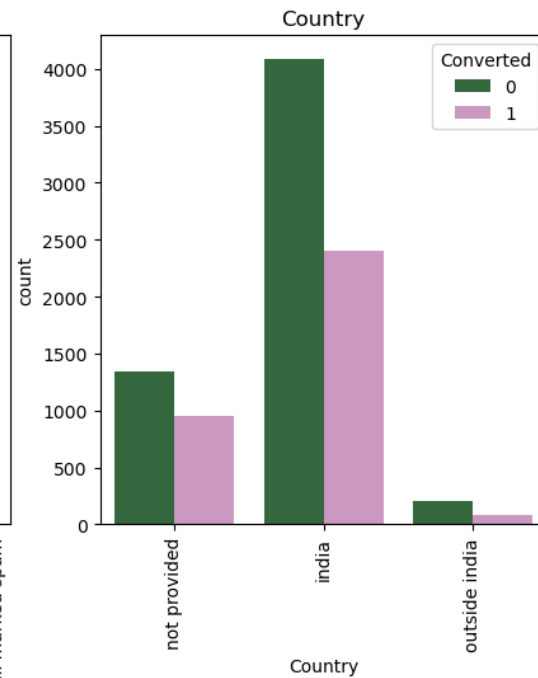
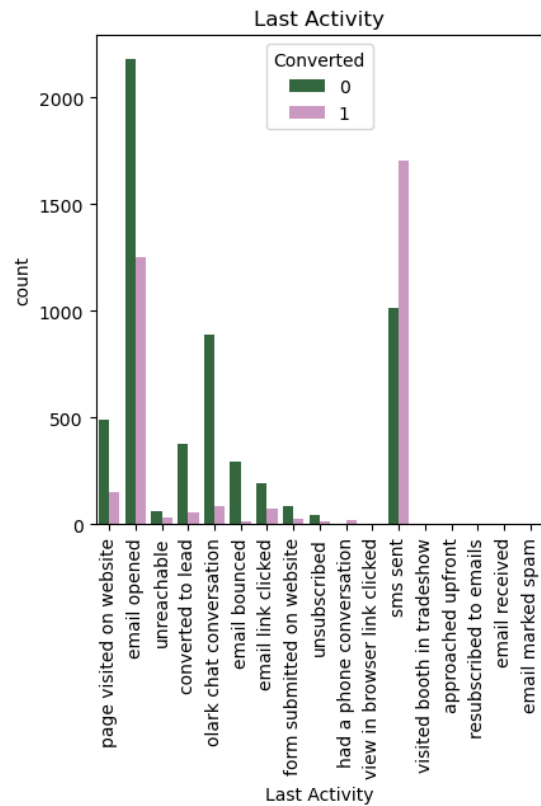
After performing these data cleaning and manipulation steps, you will be left with a more streamlined and relevant dataset that can be used for further analysis, feature engineering, and building the logistic regression model for classification. Remember to keep track of the changes made to the dataset and document the preprocessing steps for transparency and reproducibility.

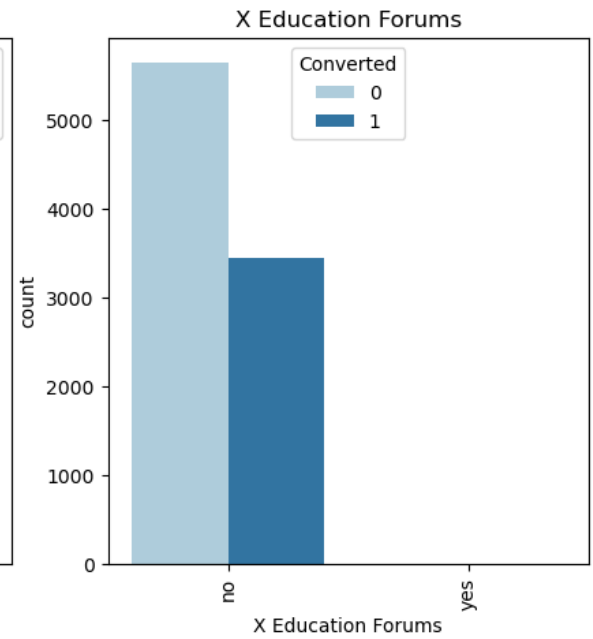
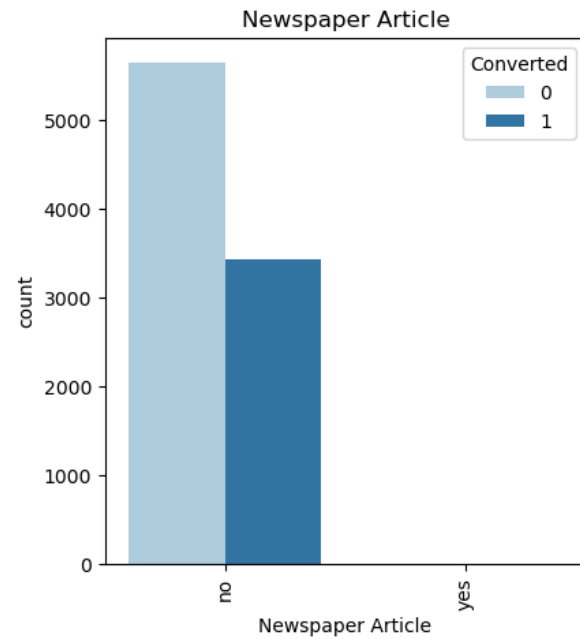
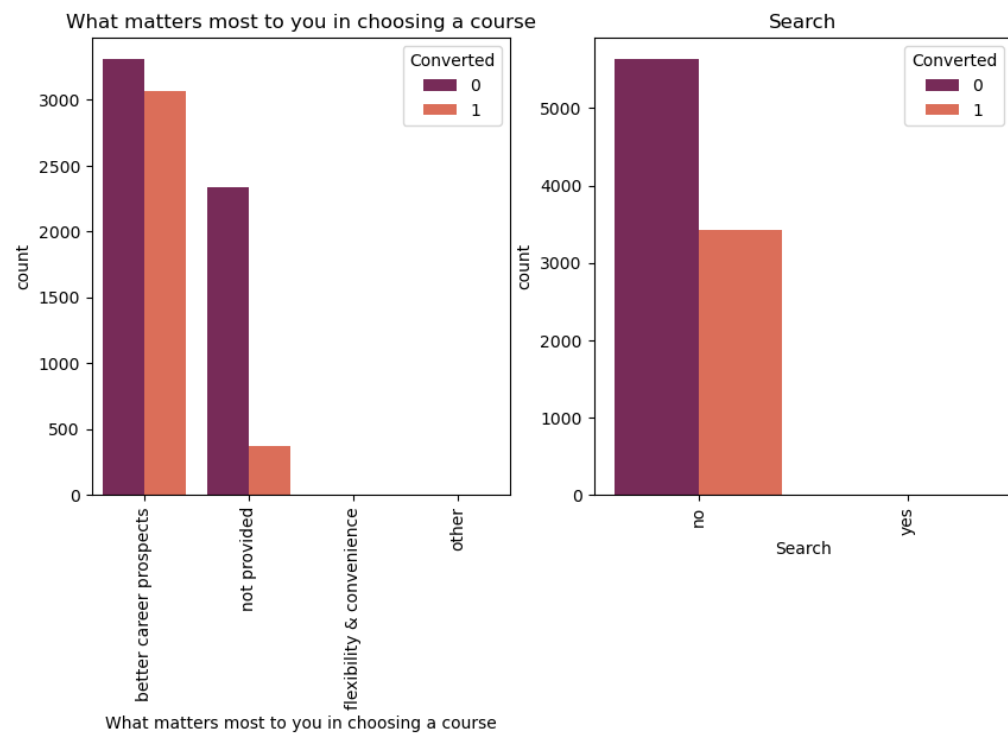
Graphs from Exploratory Data Analysis

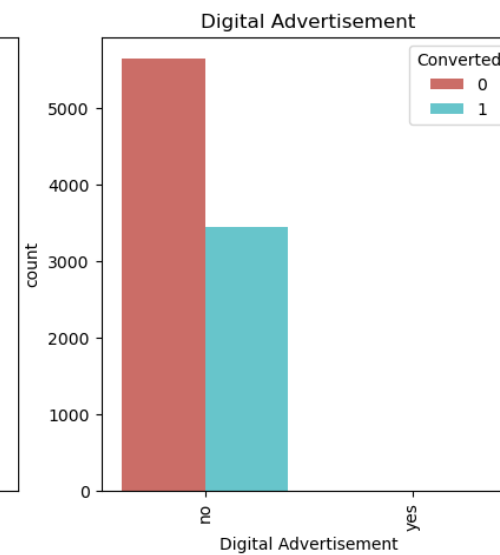
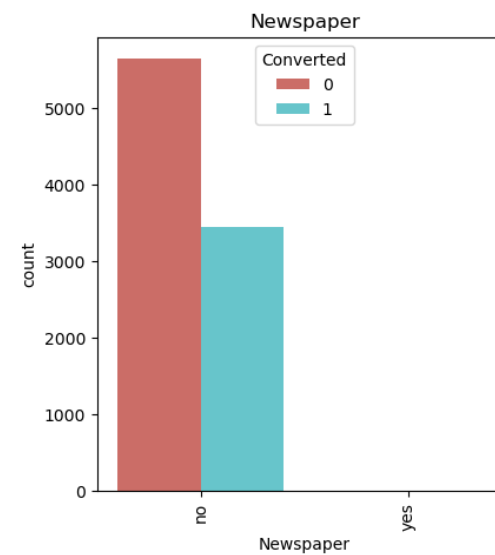
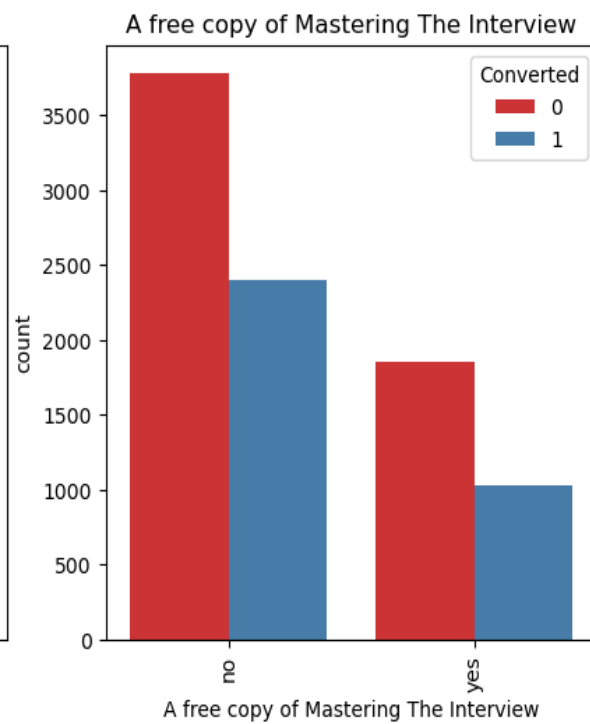
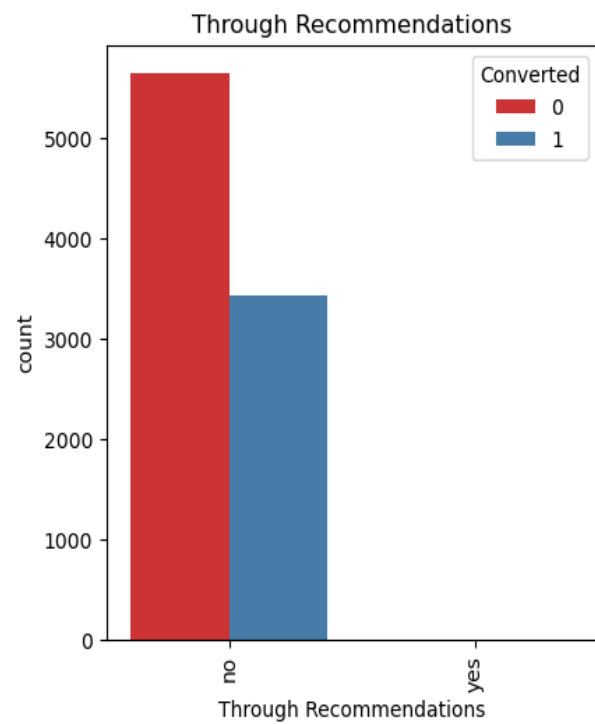


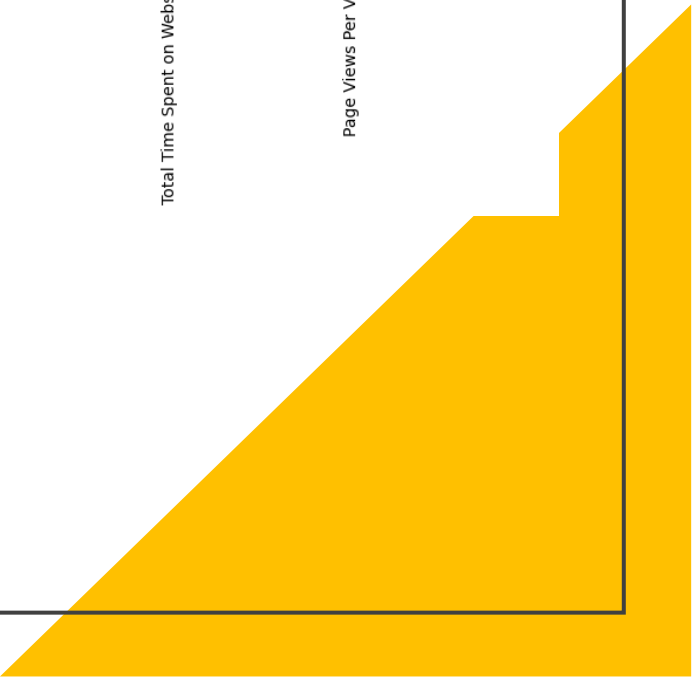
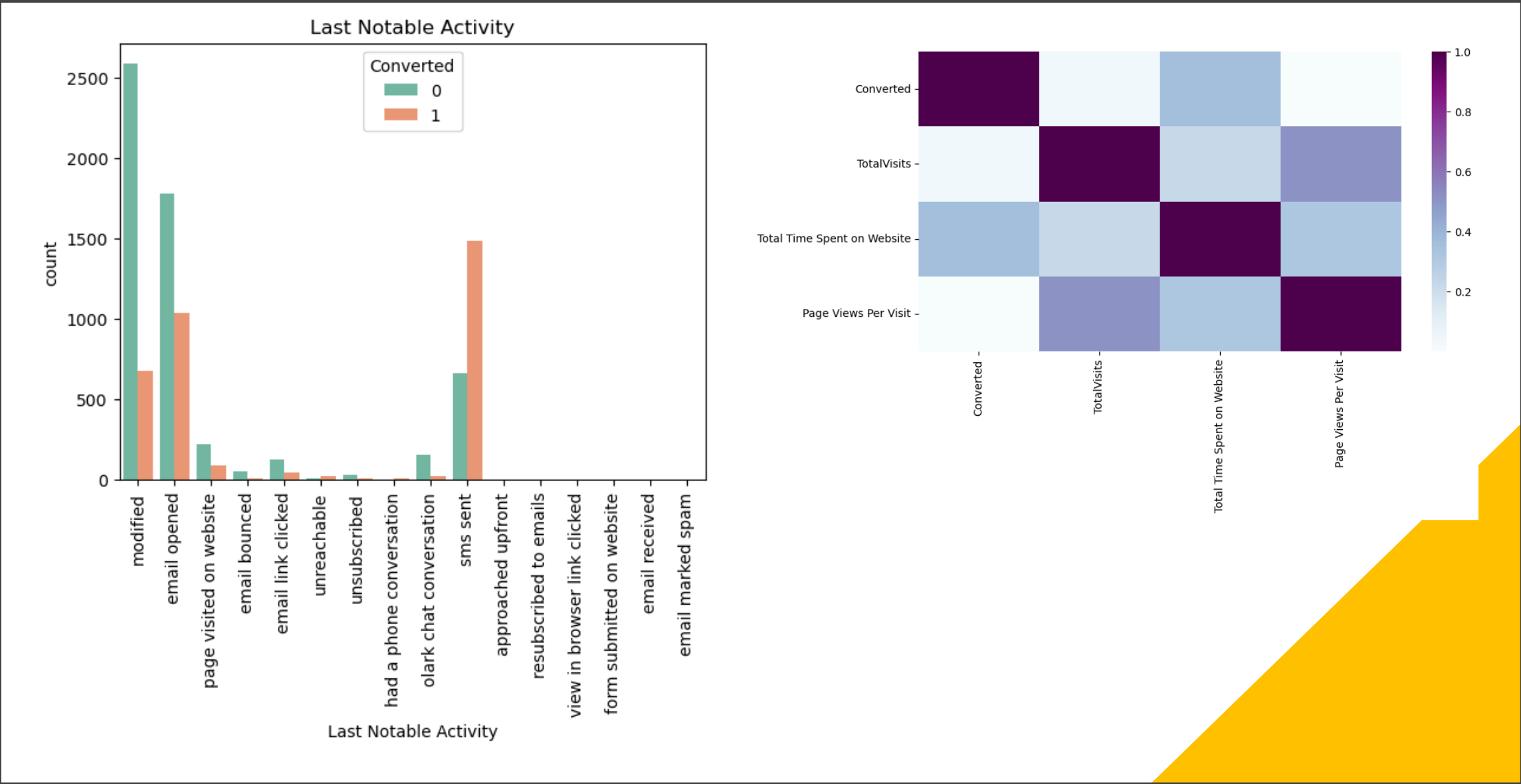


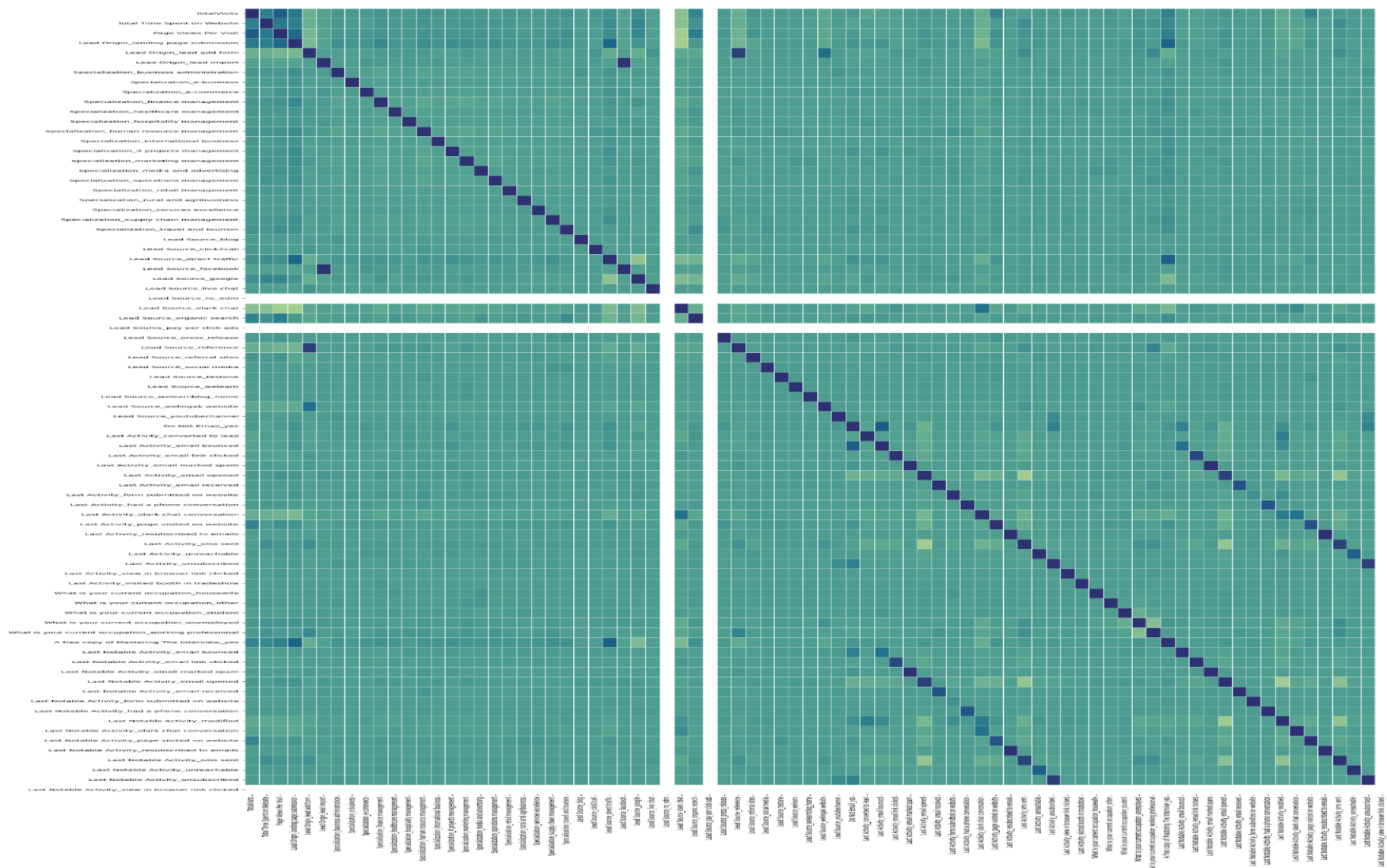












Conversion of Data

1. Numerical Variables Normalization:

1. Numerical variables have been normalized, which means that their values have been transformed to a common scale to ensure that they have similar ranges. Common normalization techniques include Min-Max scaling or Z-score normalization.

2. Creation of Dummy Variables:

1. For object type variables (categorical variables), dummy variables have been created. Dummy variables are binary representations of each category within a categorical variable. They are used to convert categorical data into a numerical format suitable for analysis and modeling.

3. Total Rows for Analysis:

1. After the data conversion, the dataset contains 8792 rows, which are now suitable for analysis.

4. Total Columns for Analysis:

1. The dataset now consists of 43 columns that have been prepared for analysis after the data conversion process.

With numerical variables normalized and object type variables converted into dummy variables, the dataset is ready for further analysis and modeling. You can now proceed with building and validating the logistic regression model or any other classification technique based on your objectives and the characteristics of the data.

Keep in mind to handle class imbalances (if any) and perform thorough feature engineering to maximize the model's performance. Additionally, consider using appropriate evaluation metrics to assess the model's accuracy and performance on unseen data during the validation process.

Model Building

1. Splitting the Data into Training and Testing Sets:

1. The dataset is divided into two subsets: the training set, which is used to train the model, and the testing set, which is used to evaluate the model's performance on unseen data. The chosen split ratio is 70:30, with 70% of the data used for training and 30% for testing.

2. Using RFE for Feature Selection:

1. RFE (Recursive Feature Elimination) is a feature selection technique that recursively removes the least significant features from the dataset until the desired number of features is reached. It helps identify the most relevant features for building the model.

3. Running RFE with 15 Variables as Output:

1. RFE is applied with the aim of selecting 15 most important variables as the model's input features. These 15 variables are considered to have the highest impact on the target variable.

4. Building Model by Removing Insignificant Variables:

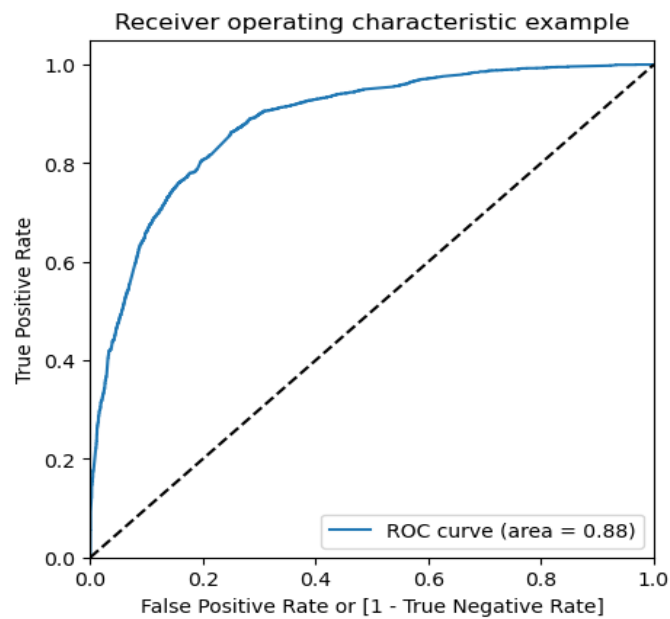
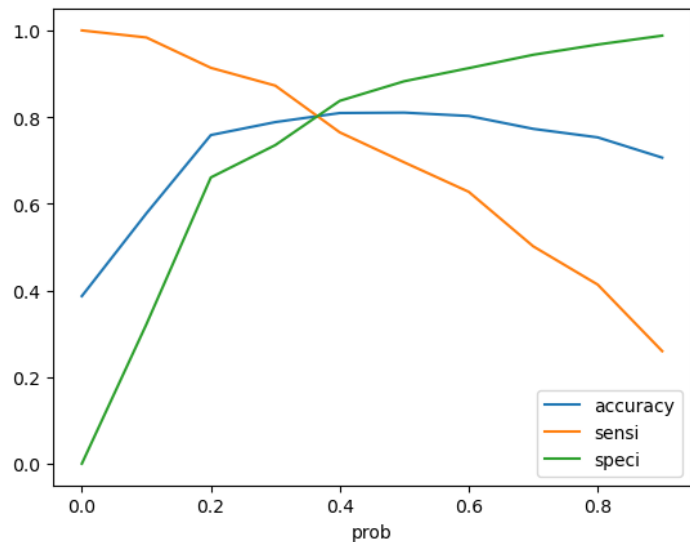
1. After RFE, the model is built using the selected 15 variables. However, some of these variables may have p-values greater than 0.05, indicating they are not statistically significant in predicting the target variable. Additionally, some variables may have high Variance Inflation Factor (VIF) values (greater than 5), suggesting potential multicollinearity issues. To address these concerns, the variables with high p-values and VIF values are removed from the model.

5. Predictions on the Test Dataset:

1. After building the model with the selected significant variables, it is used to make predictions on the test dataset (30% of the original data). This allows for an evaluation of the model's performance on unseen data.

6. Overall Accuracy of 81%:

1. The accuracy of the model, when applied to the test dataset, is found to be 81%. This indicates that approximately 81% of the model's predictions on the test data are correct.
- It's important to remember that model building is an iterative process, and the choice of variables to include or exclude should be based on various statistical tests, domain knowledge, and the specific objectives of the analysis. Additionally, further fine-tuning and validation are crucial to ensure the model's reliability and generalizability to new data.



Finding Optimal Cutoff

- Optimal cut-off probability balances sensitivity and specificity in a classification model.
- The chosen optimal cut-off is 0.35, providing a favorable trade-off between true positive and true negative rates.
- Consider the problem's context and misclassification costs when selecting the cut-off.
- Validate the model's performance on a separate dataset to ensure generalizability.

Case Study Conclusion

Based on the analysis, X Education can significantly improve its chances of converting potential buyers into actual customers by prioritizing certain key variables. These variables, ranked in descending order of importance, include:

1. The total time spent on the Website.
 2. Total number of visits.
 3. Lead source (Google, Direct traffic, Organic search, and Welingak website).
 4. Last activity (SMS and Olark chat conversation).
 5. Lead origin as Lead add format.
 6. Current occupation as a working professional.
- By focusing on these crucial factors, X Education can fine-tune its marketing strategies, cater to potential buyers' preferences, and enhance its chances of success in enticing them to enroll in their courses. This strategic approach will likely lead to increased customer conversion rates and overall growth for the organization.