



Lead Scoring Case Study

Group Members

- Abhishek Kumar Singh
- Vivek Abhyankar
- Vicky Kumar



Problem Statement

- X Education sells online courses to industry professionals.
- X Education gets a lot of leads, its lead conversion rate is very poor. For example, if, say, they acquire 100 leads in a day, only about 30 of them are converted.
- 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

- X education wants to know most promising leads.
- For that they want to build a Model which identifies the hot leads.
- Deployment of the model for the future use.



Solution Methodology

- Data cleaning and data manipulation
- EDA (Univariate and Bivariate data analysis)
- Feature Scaling, Dummy Variables and encoding of the data
- Classification technique: Logistic Regression
- Validation of the model
- Model presentation
- Conclusions and recommendations



Data cleaning

- Check and handle duplicate data.
- Check and handle NA values and missing values.
- Drop columns, if it contains large amount of missing values and not useful for the analysis.
- Imputation of the values, if necessary.
- Check and handle outliers in data.

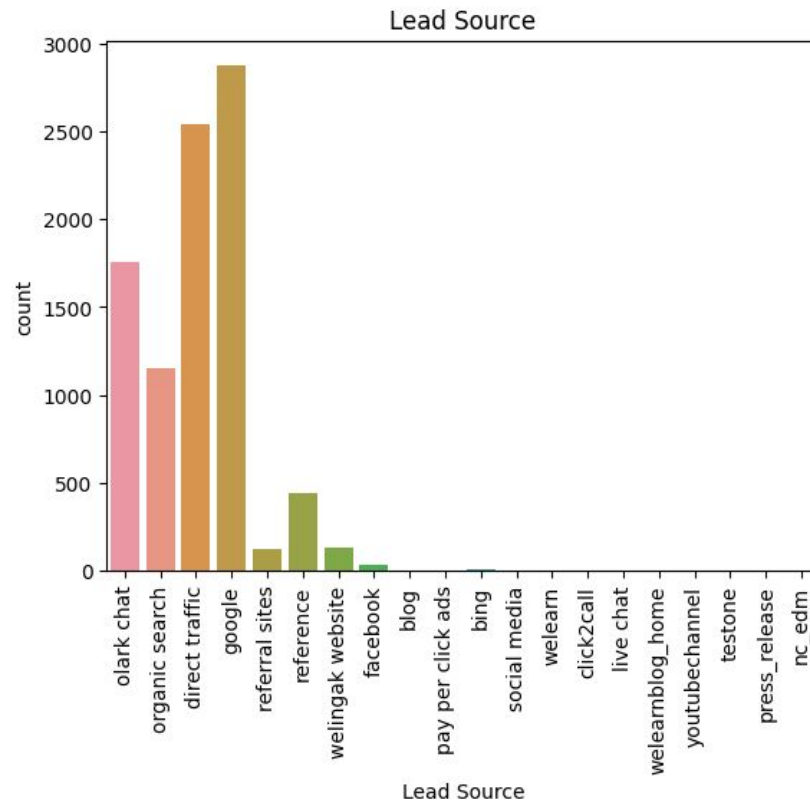


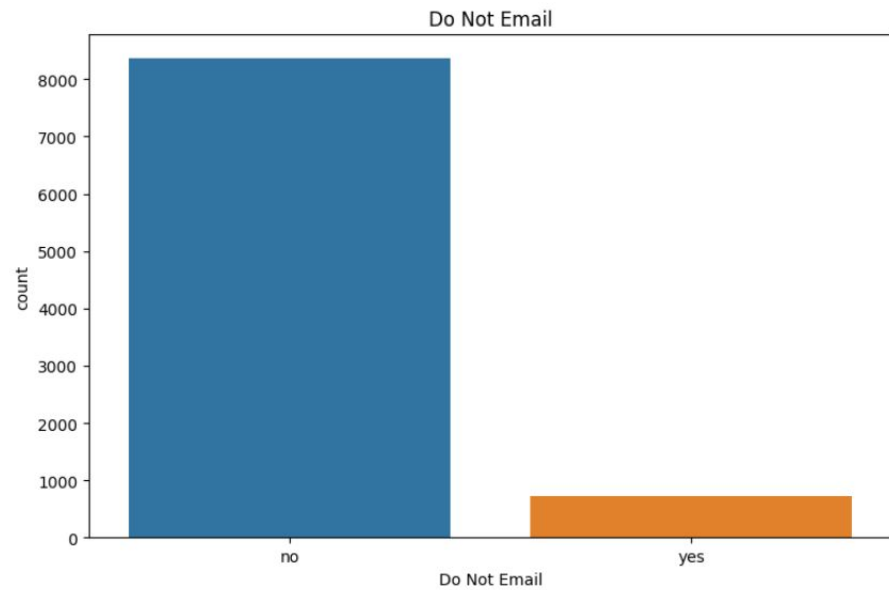
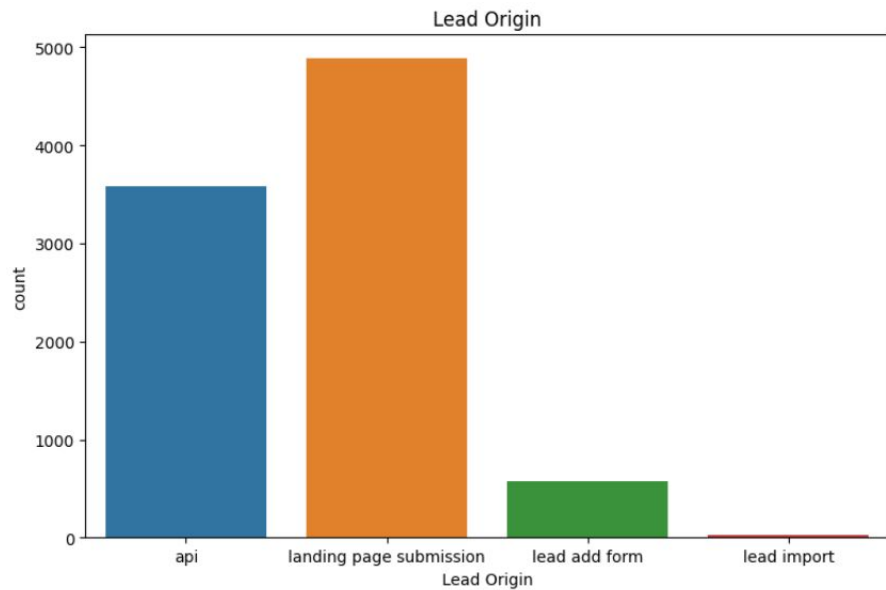
Data Manipulation

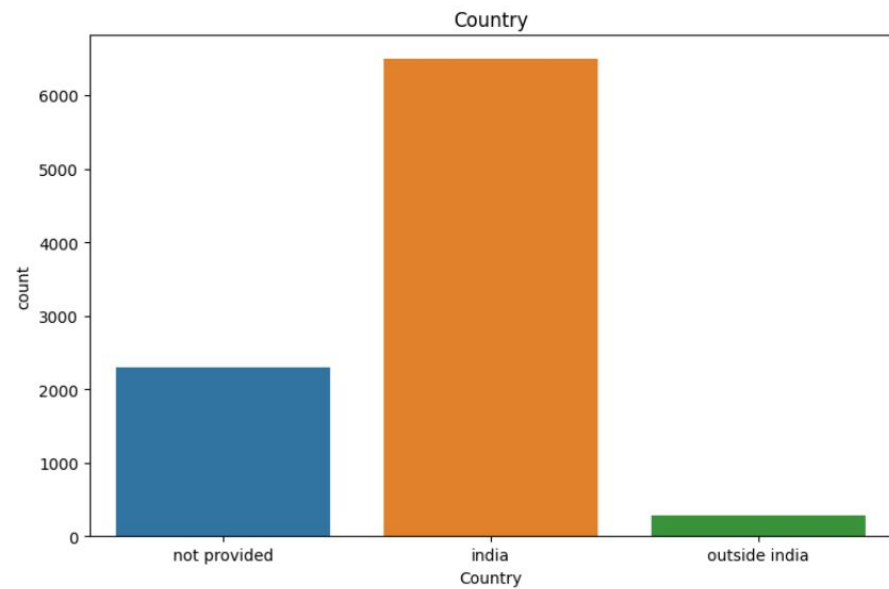
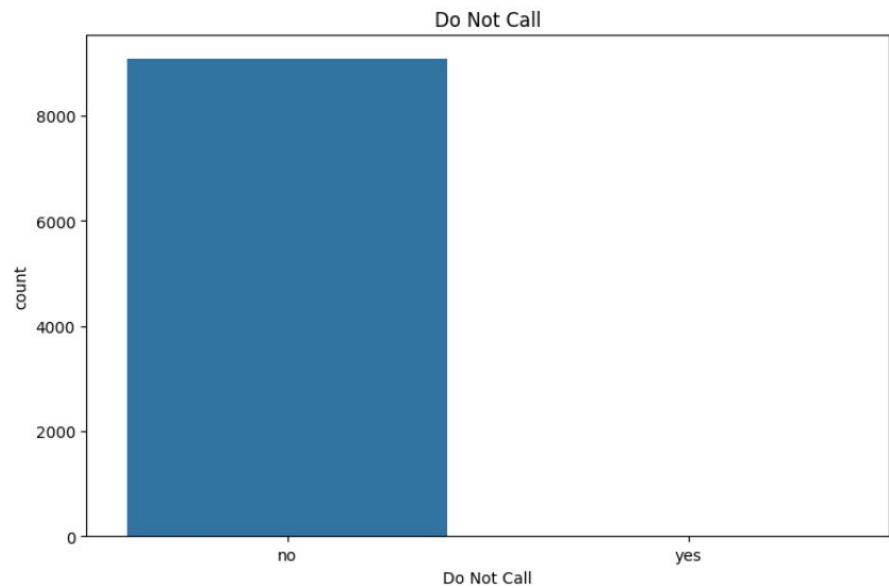
- Total Number of Columns =35, Total Number of Rows =9240
- Single value features like “Magazine”, “Receive More Updates About Our Courses”, “Update me on Supply”
- Chain Content”, “Get updates on DM Content”, “I agree to pay the amount through cheque” etc. have been dropped.
- Removing the “Prospect ID” and “Lead Number” which is not necessary for the analysis.
- After checking for the value counts for some of the object type variables, we find some of the features which has no enough variance, which we have dropped, the features are:
 - “Do Not Call”, “What matters most to you in choosing course”, “Search”, “Newspaper Article”, “X Education Forums”, “Newspaper”, “Digital Advertisement
- Dropping the columns having more than 35% as missing value such as ‘How did you hear about X Education’ and ‘Lead Profile’.

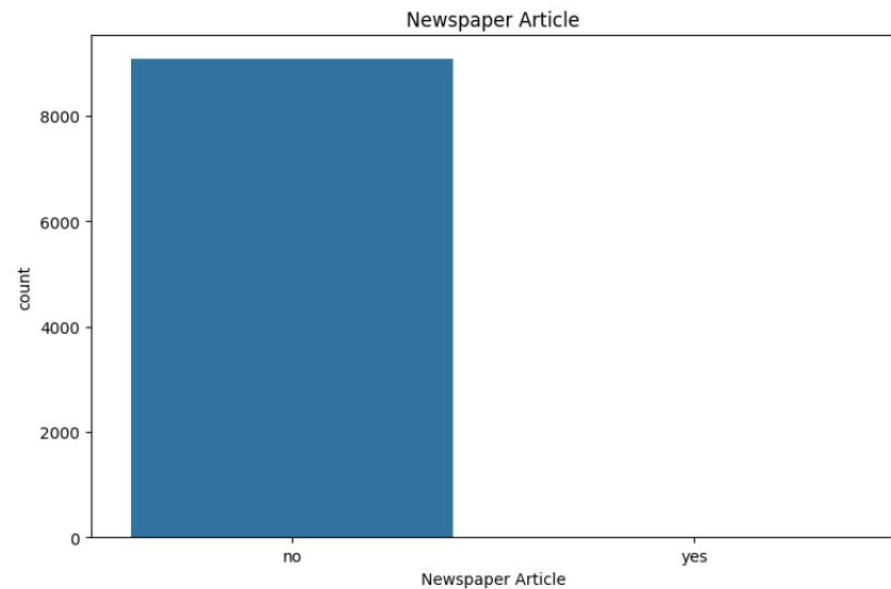
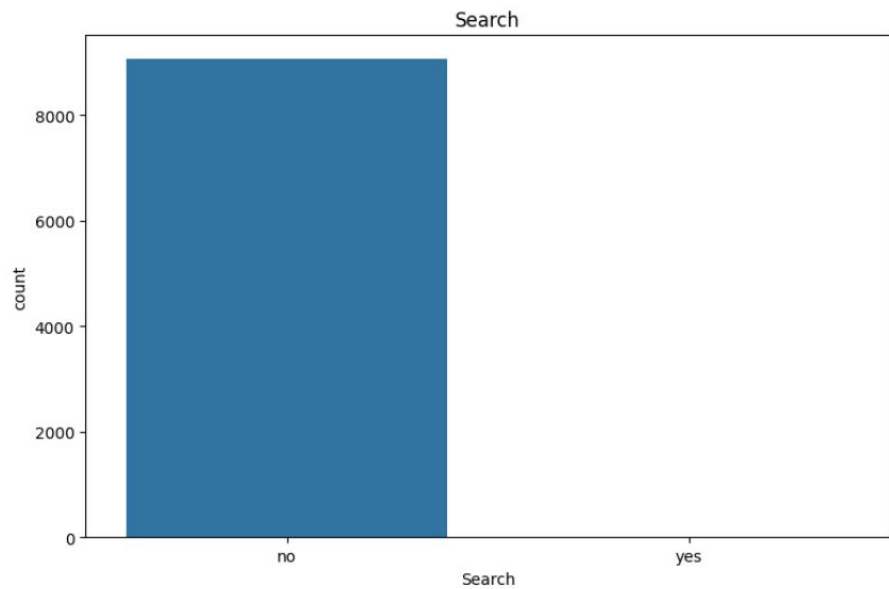
EDA

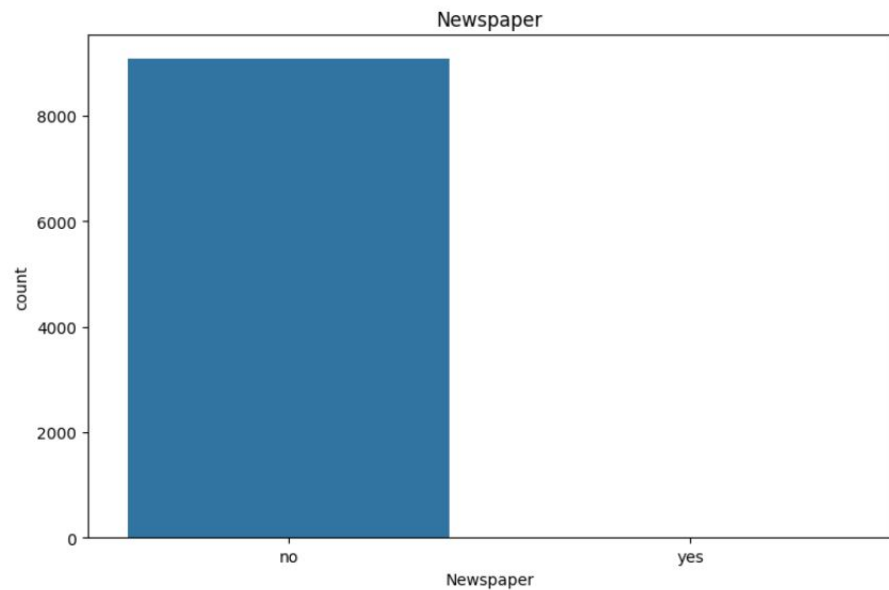
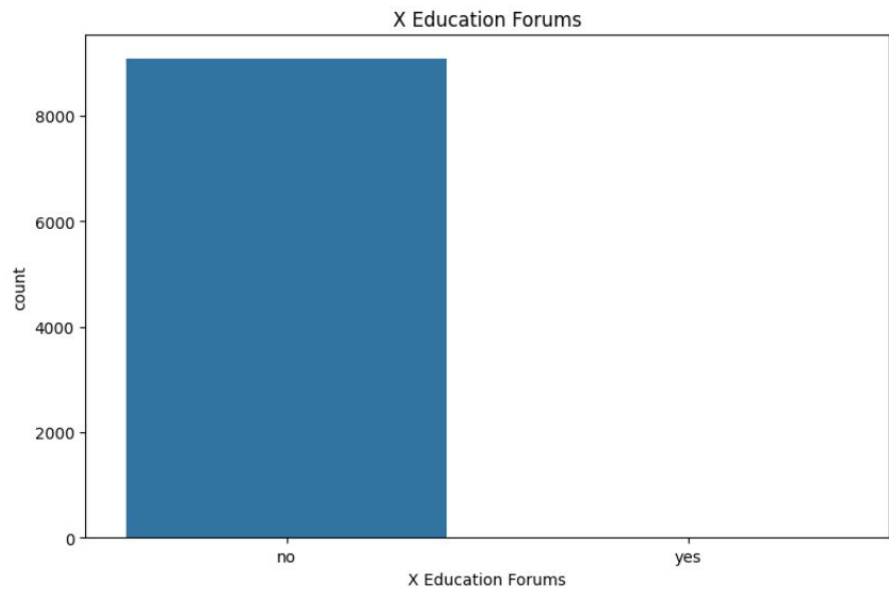
- Univariate data analysis
 - value count
 - distribution of variable
- Bivariate data analysis:
 - correlation coefficients
 - pattern between the variables

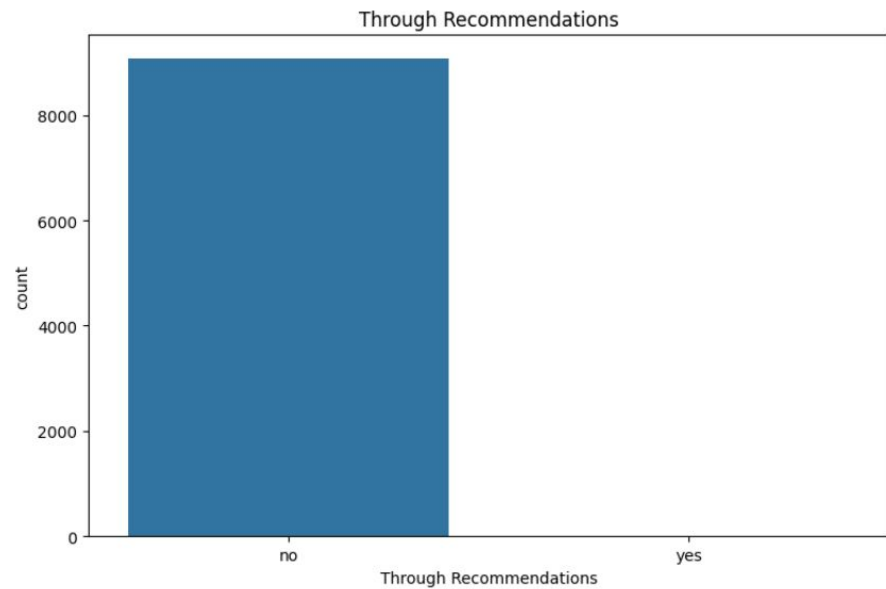
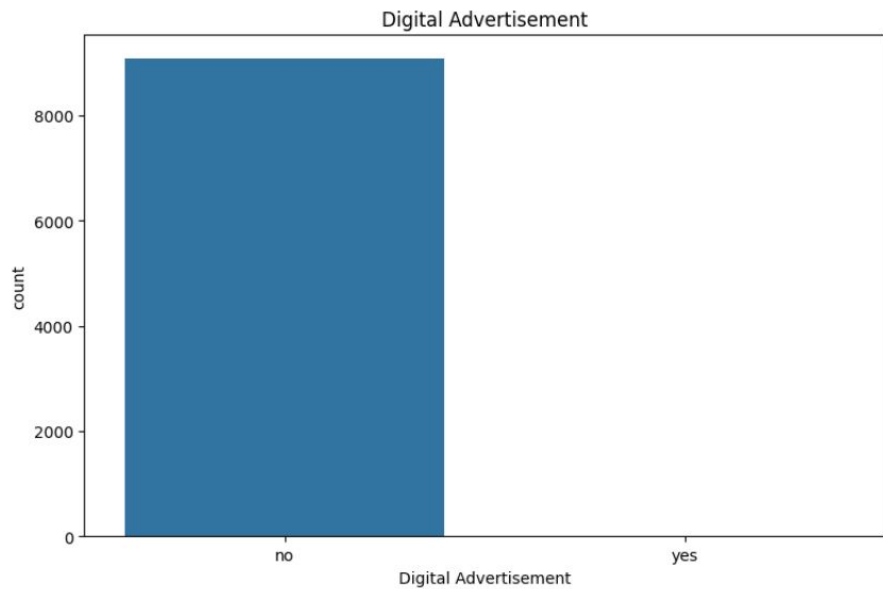




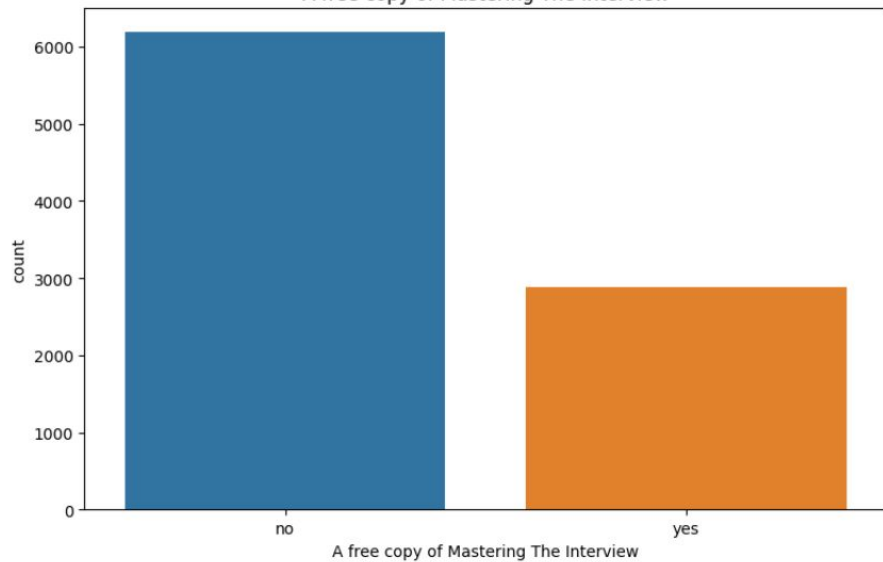




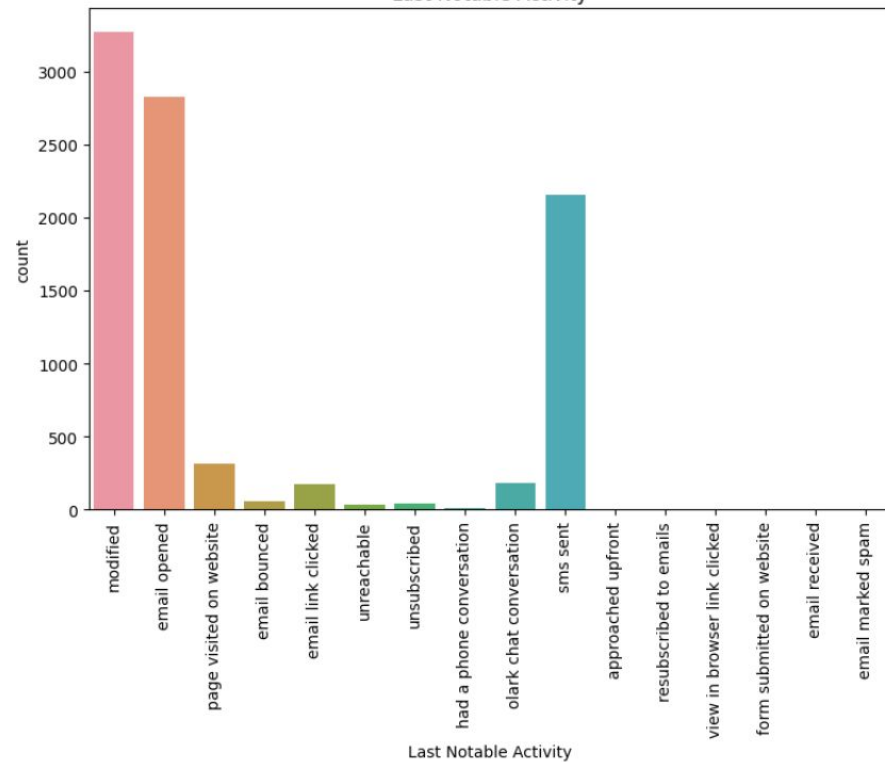


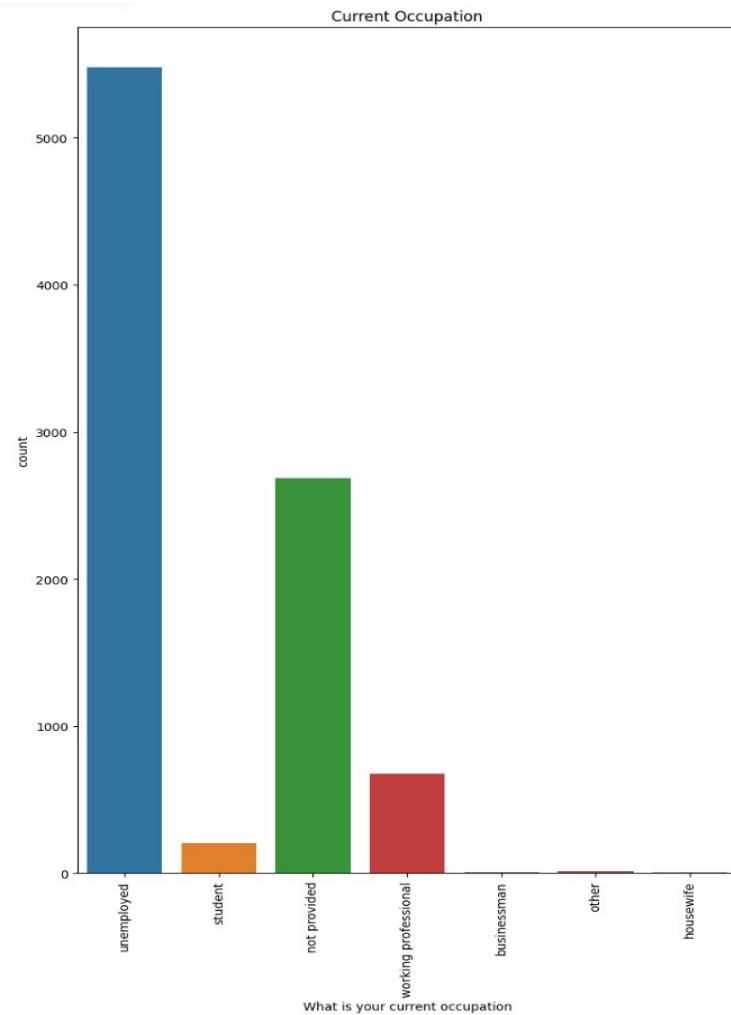
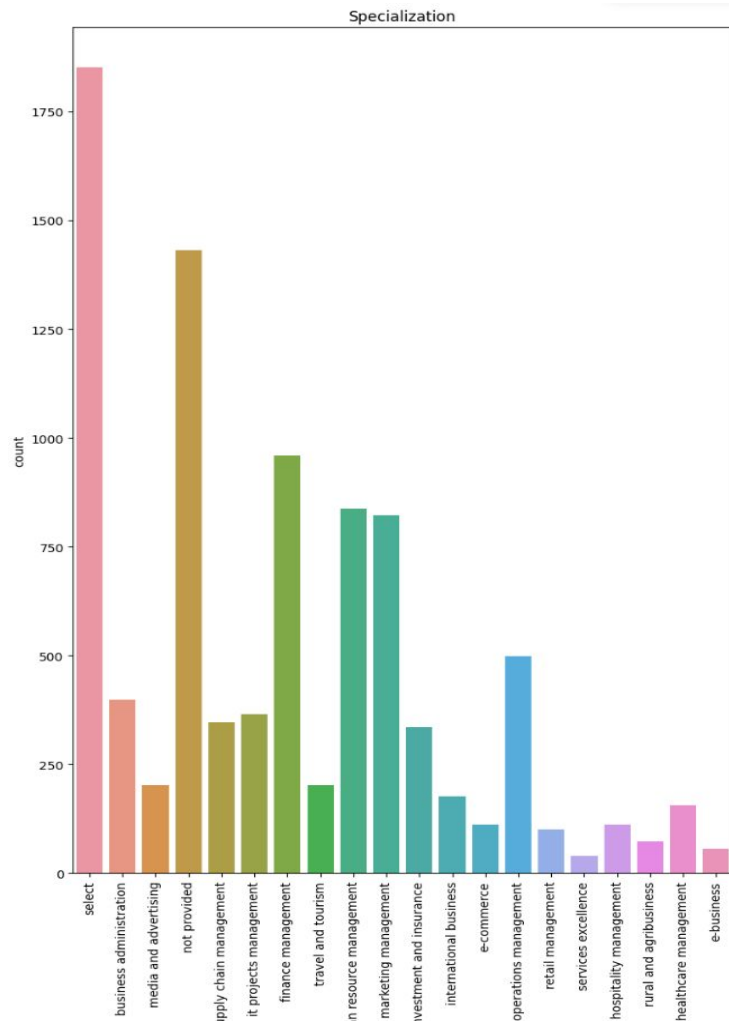


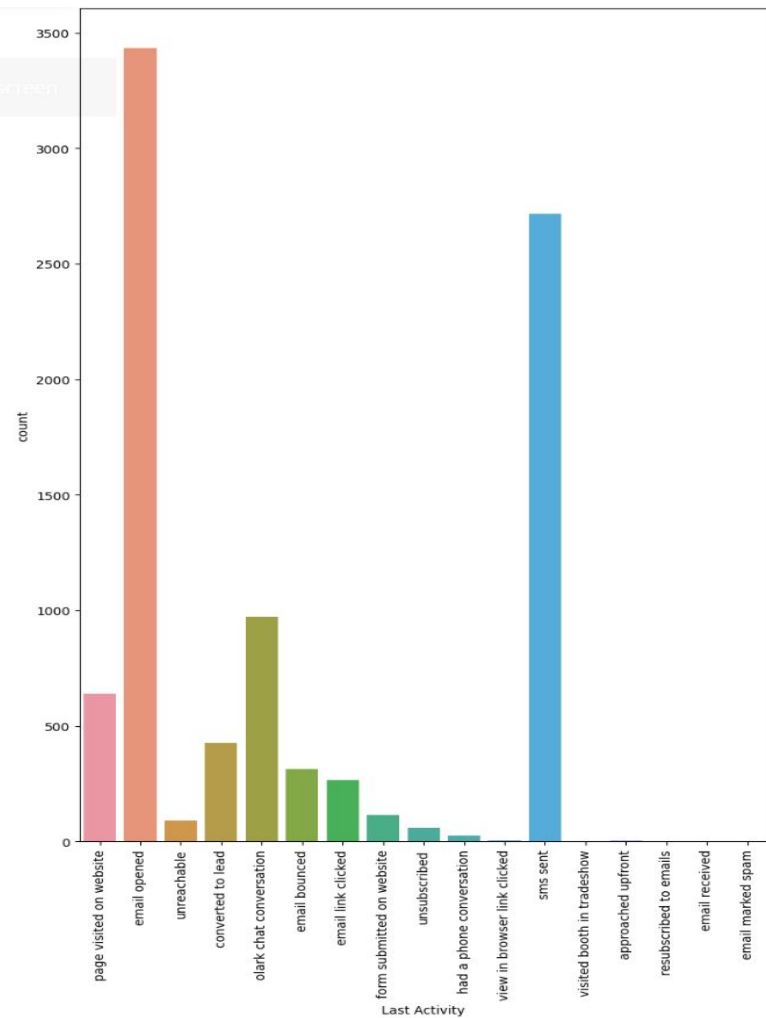
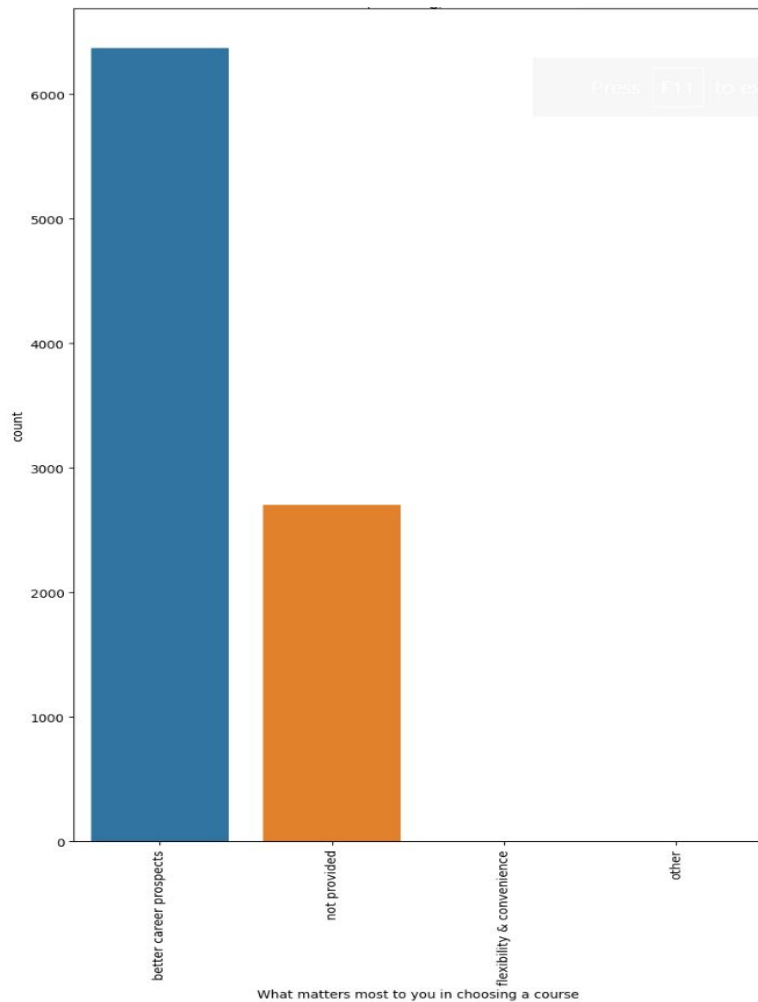
A free copy of Mastering The Interview



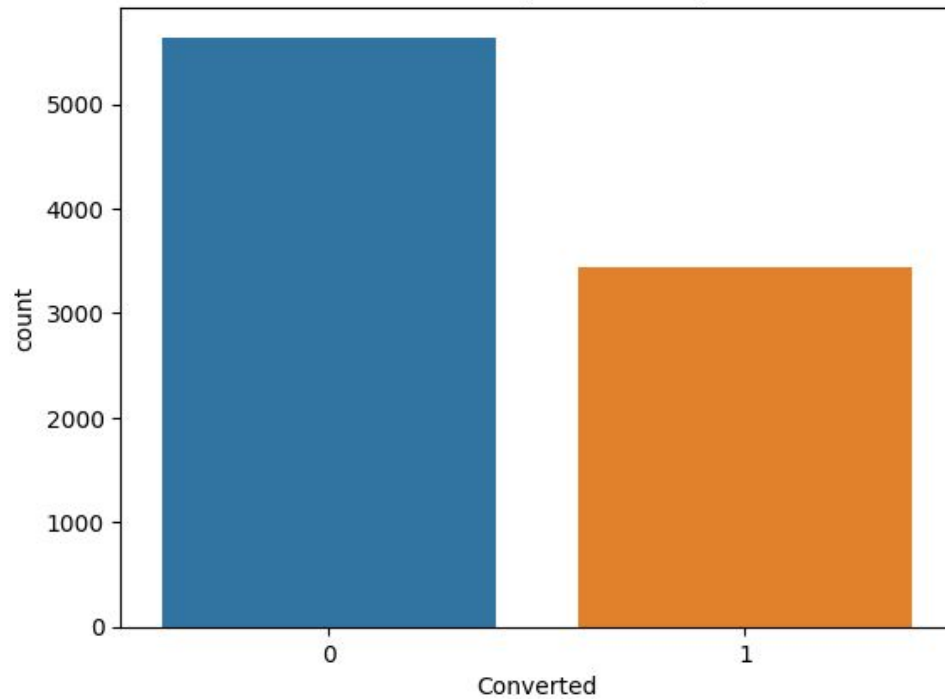
Last Notable Activity



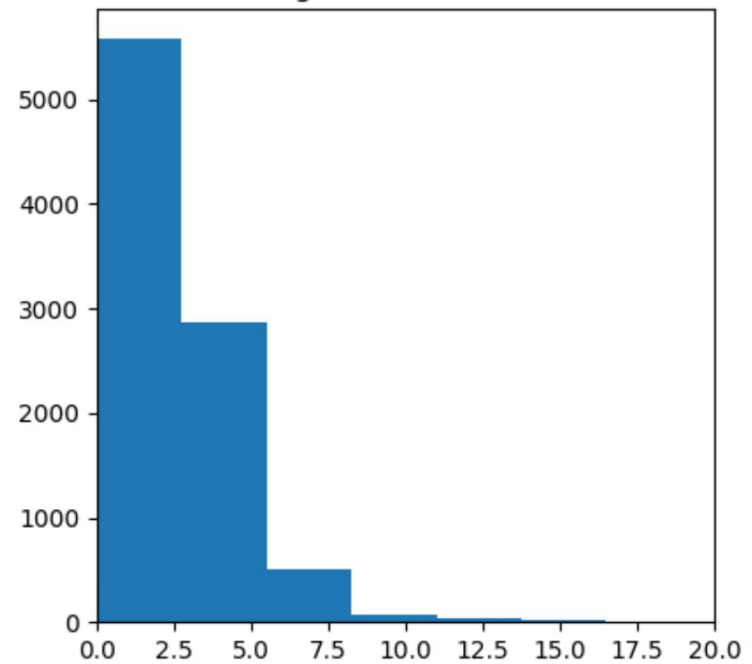




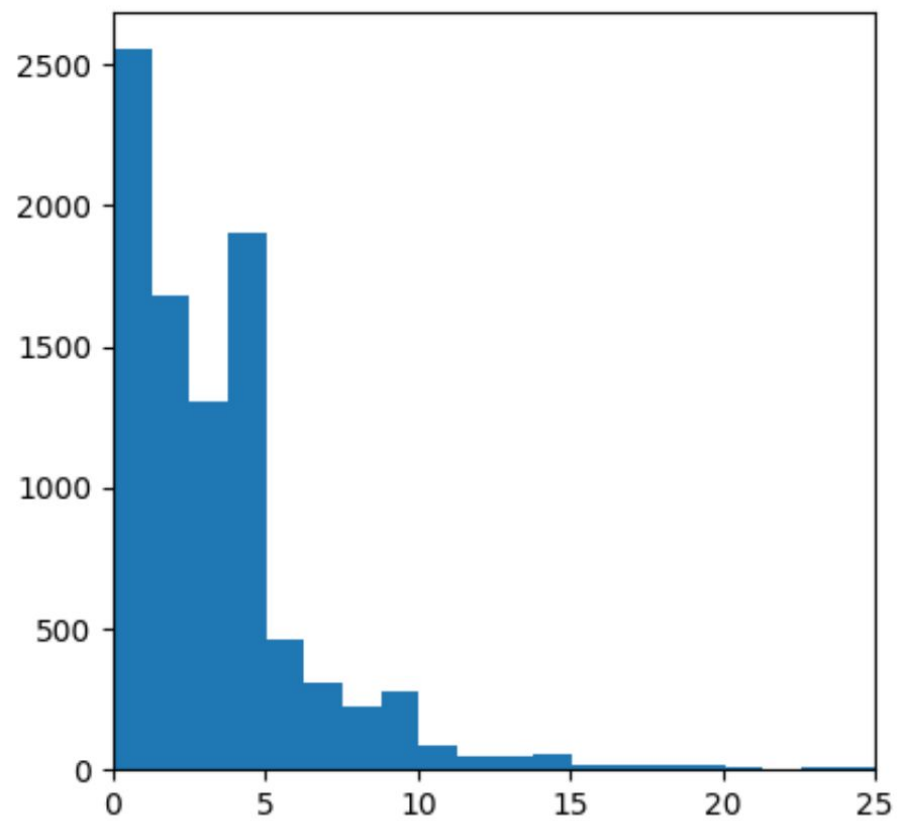
Converted("Y variable")



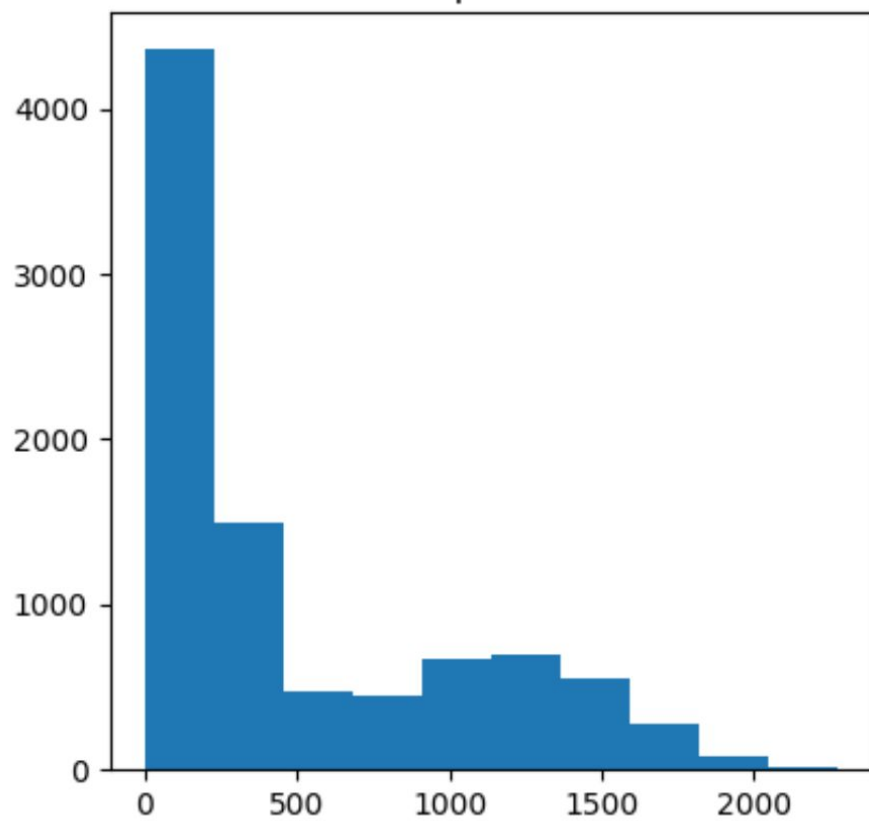
Page Views Per Visit



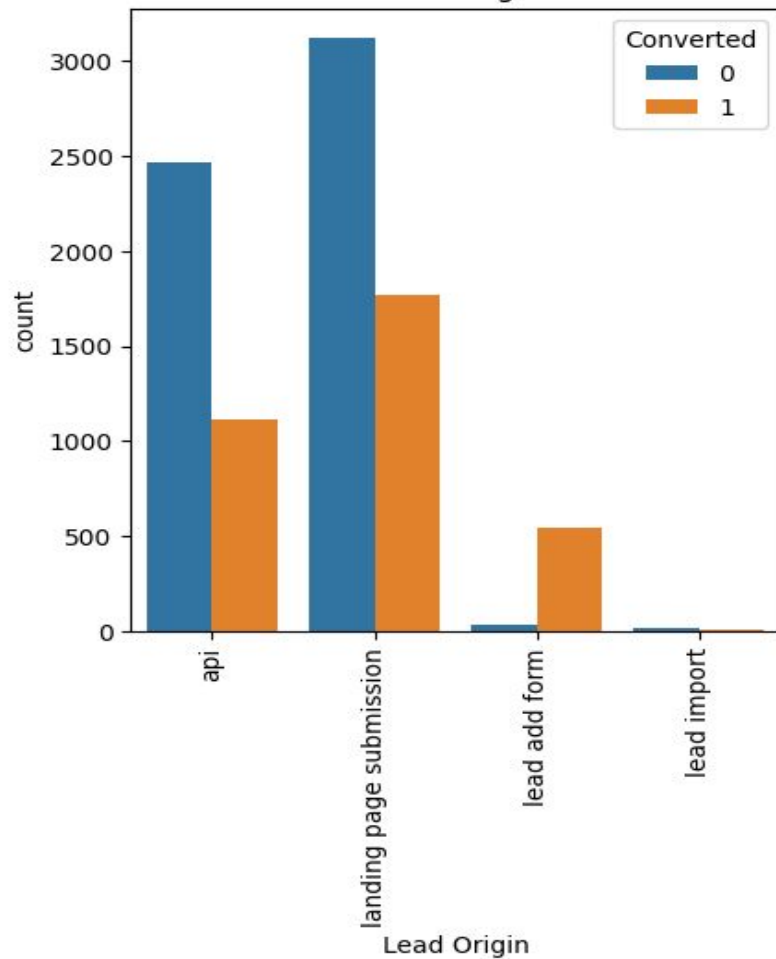
Total Visits



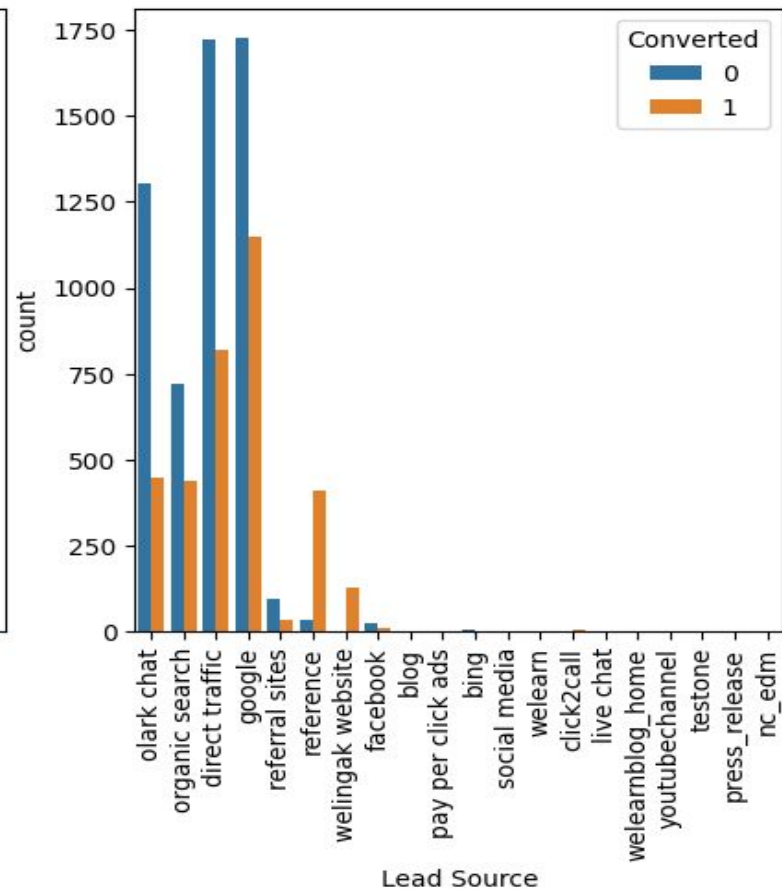
Total Time Spent on Website



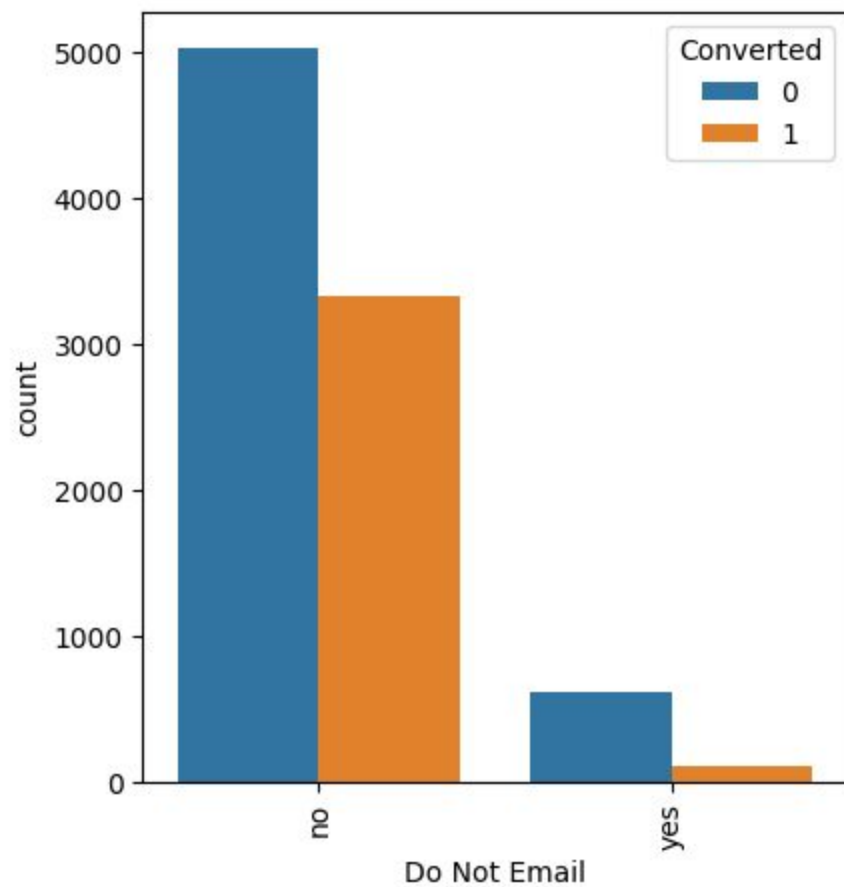
Lead Origin



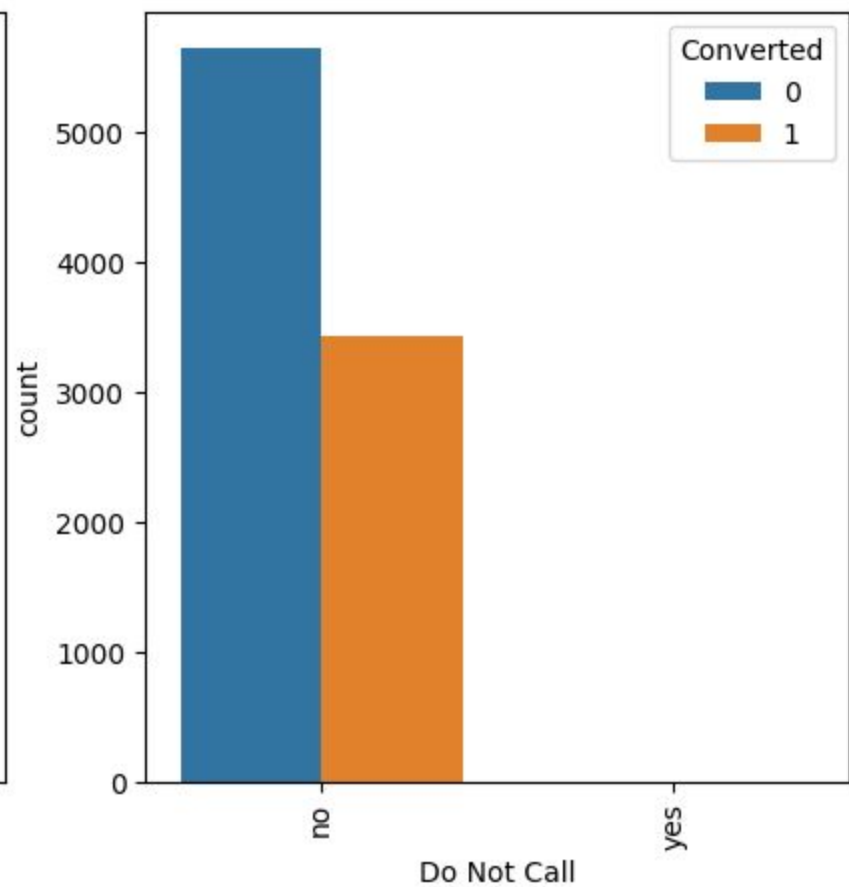
Lead Source



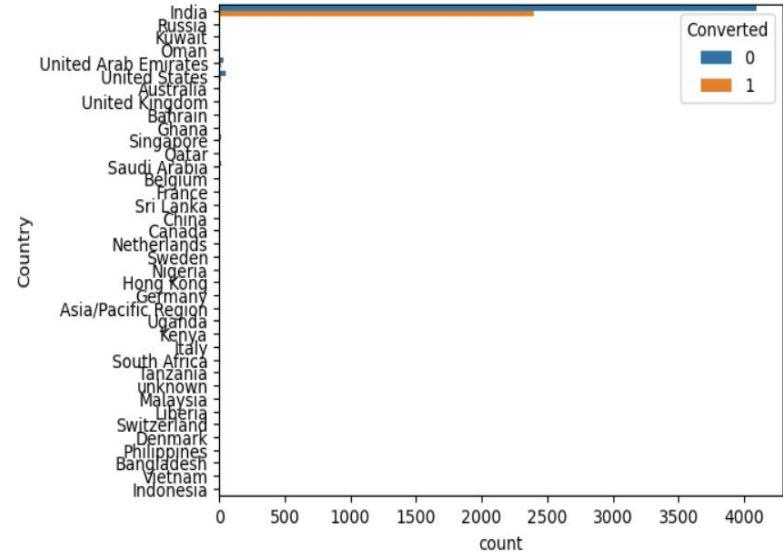
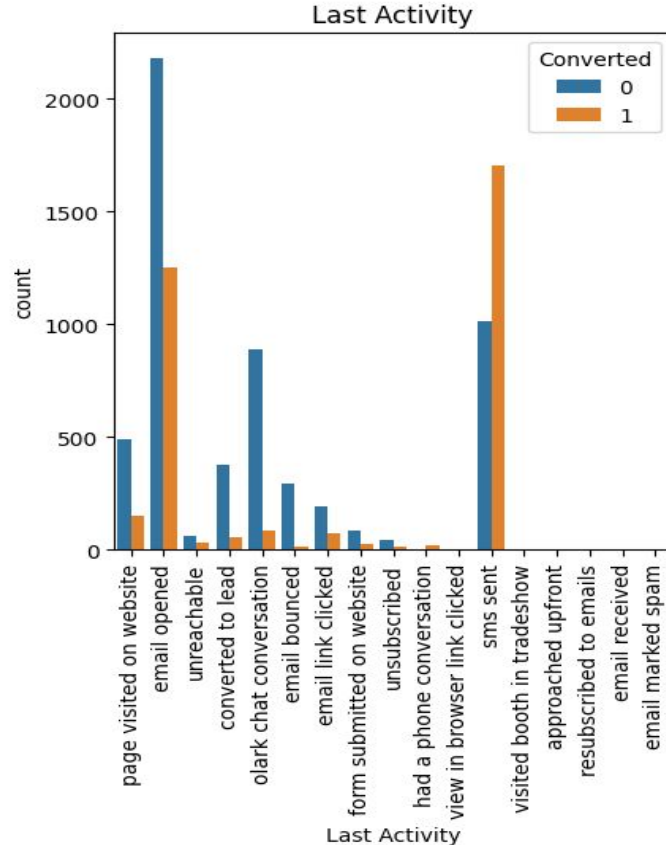
Do Not Email

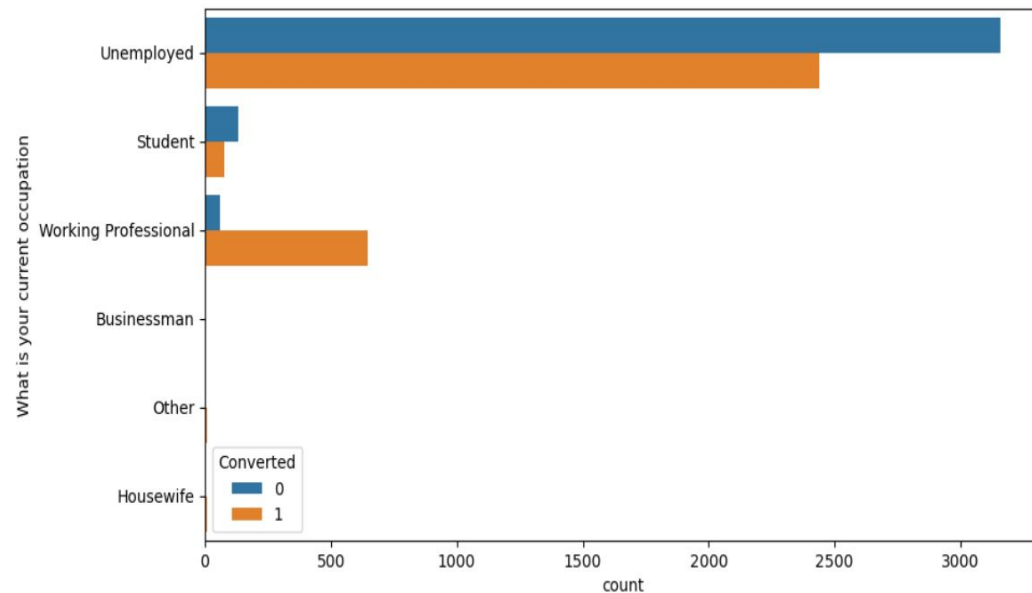
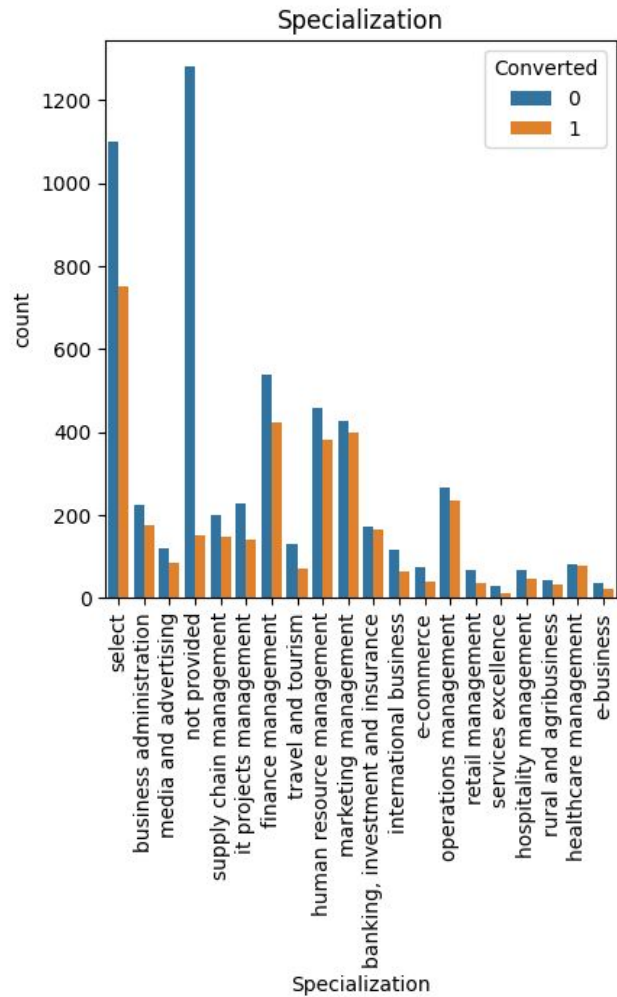


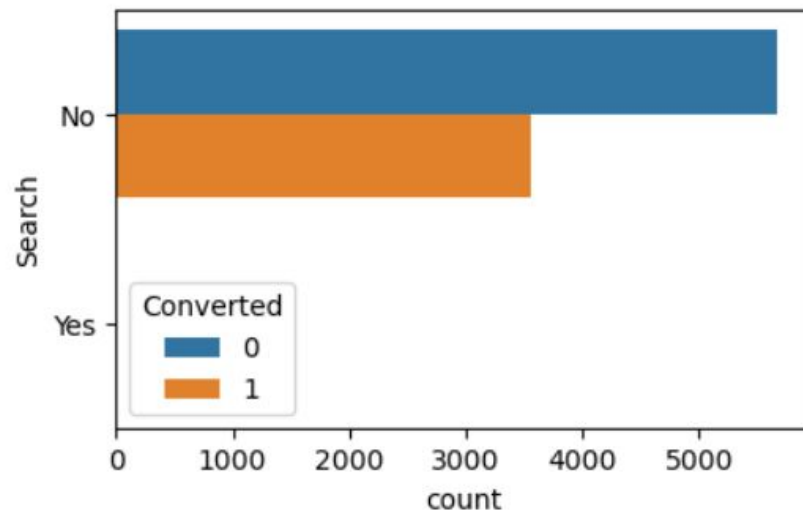
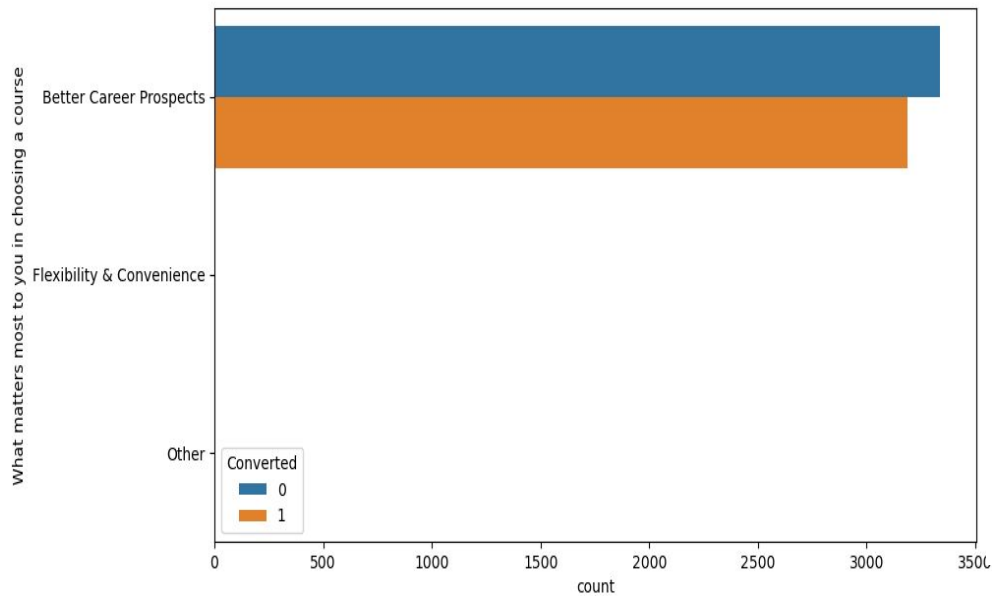
Do Not Call



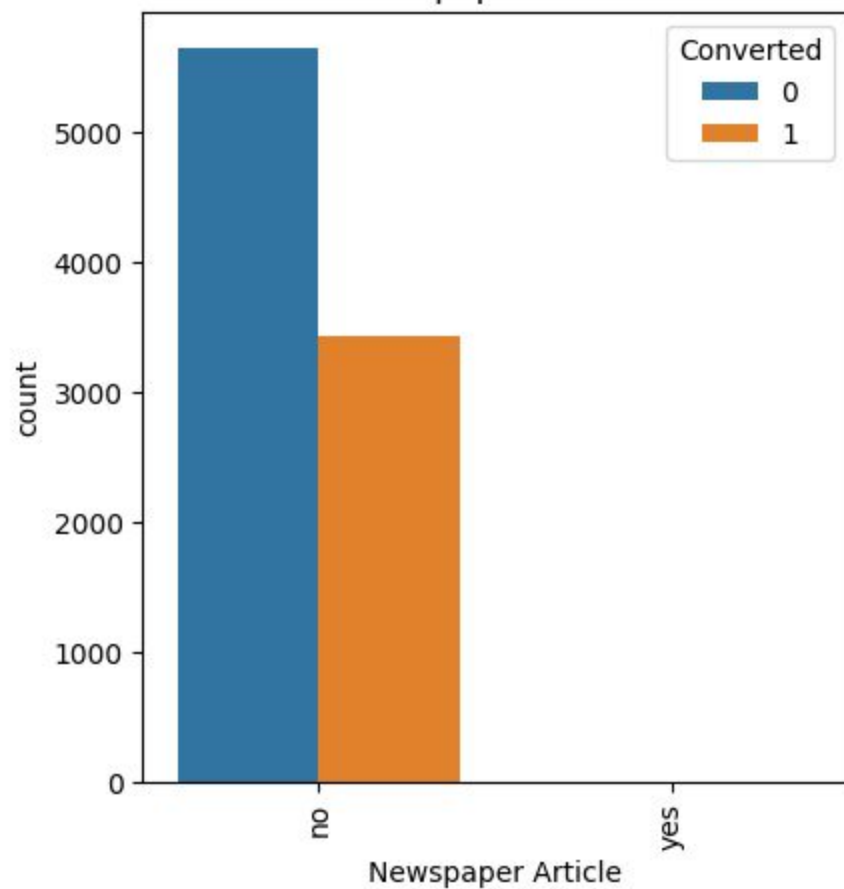
Categorical Variable Relation



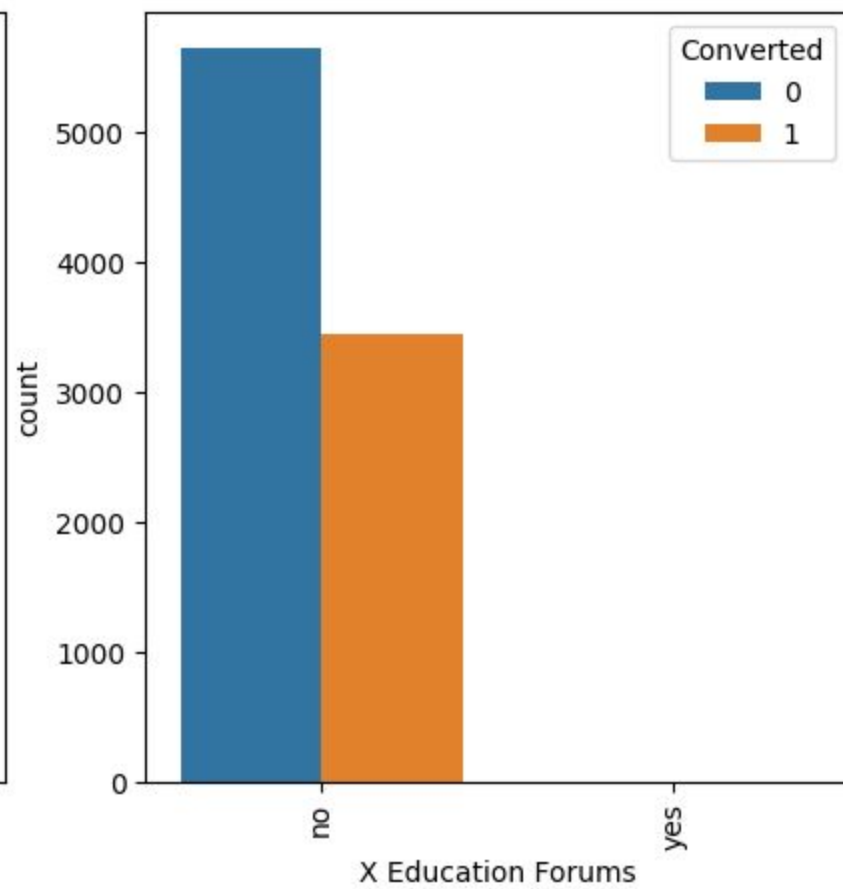




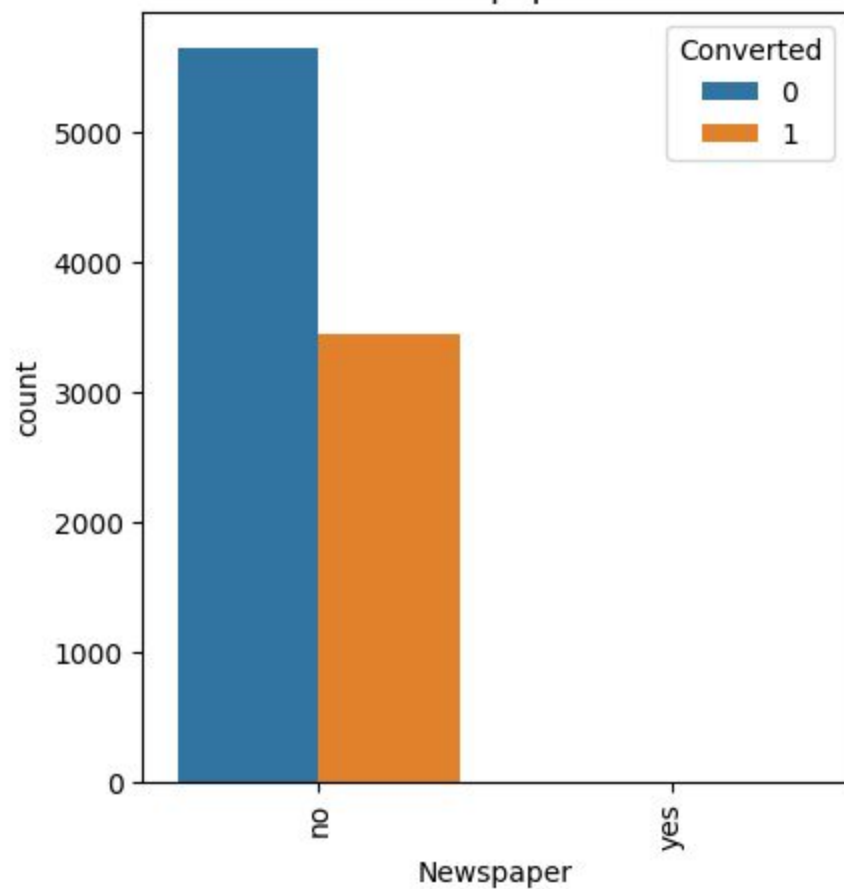
Newspaper Article



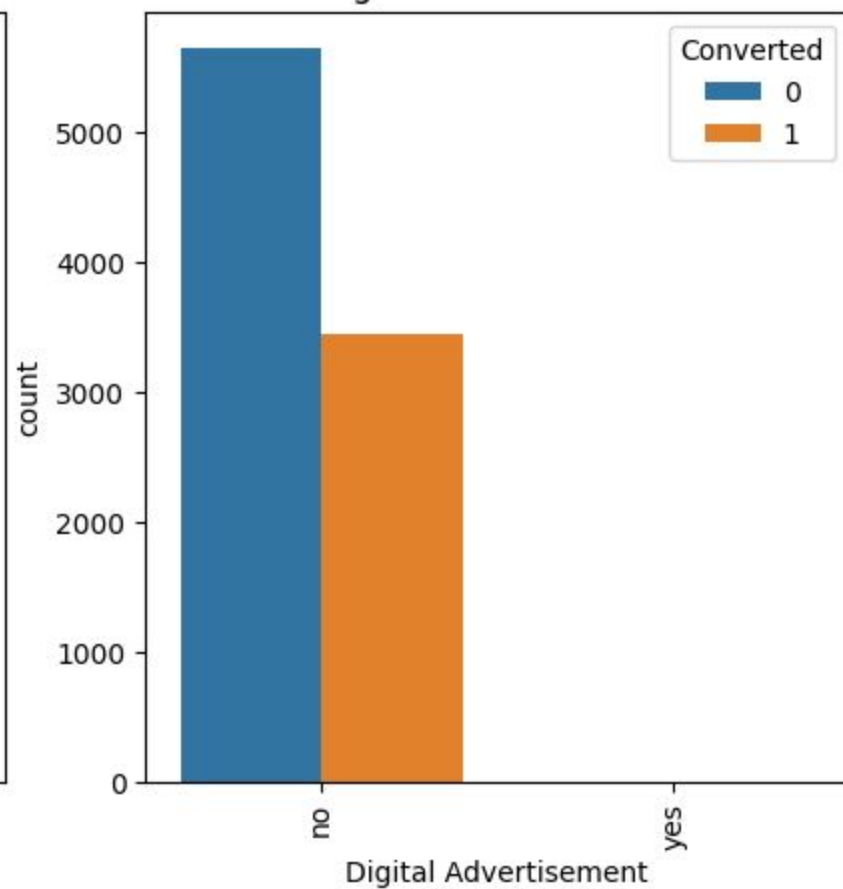
X Education Forums



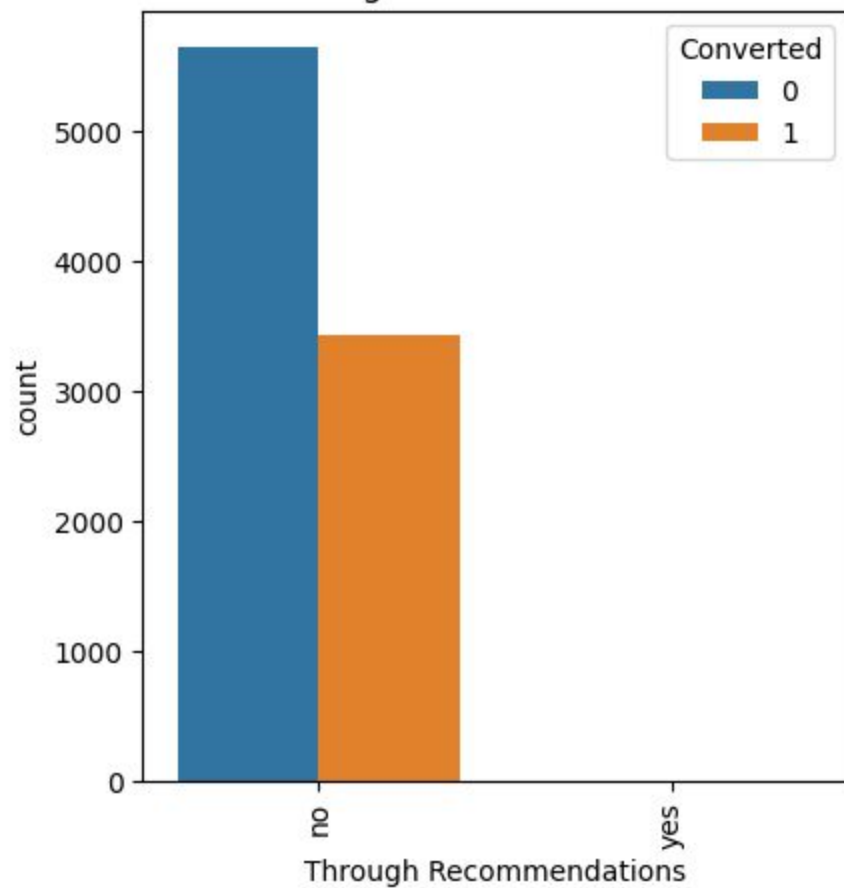
Newspaper



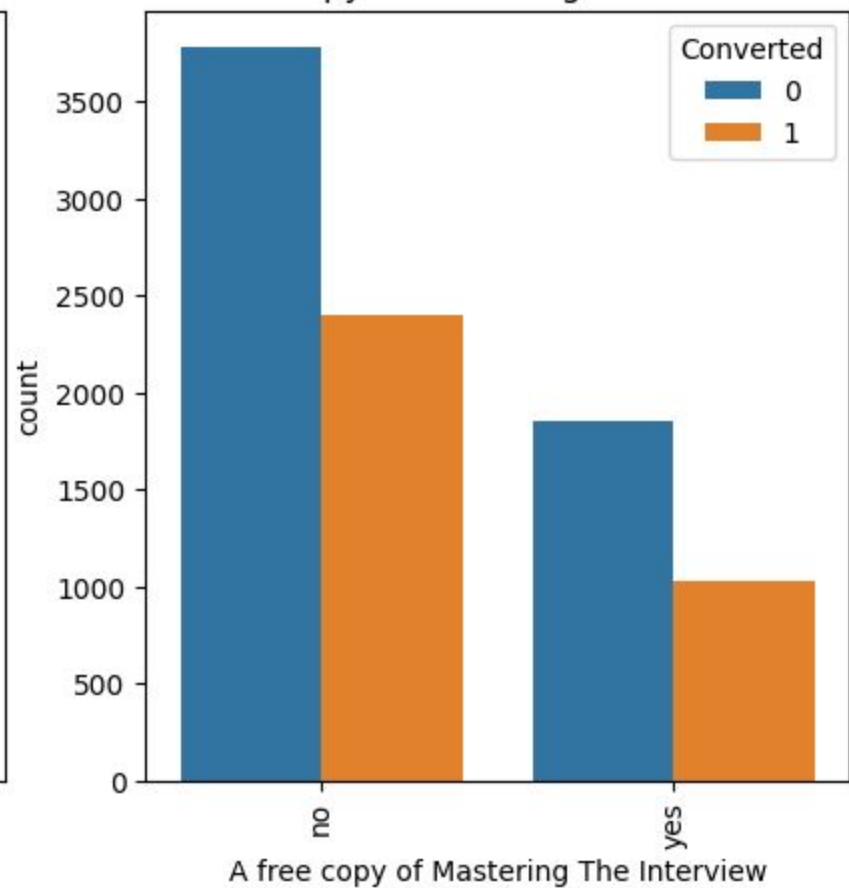
Digital Advertisement



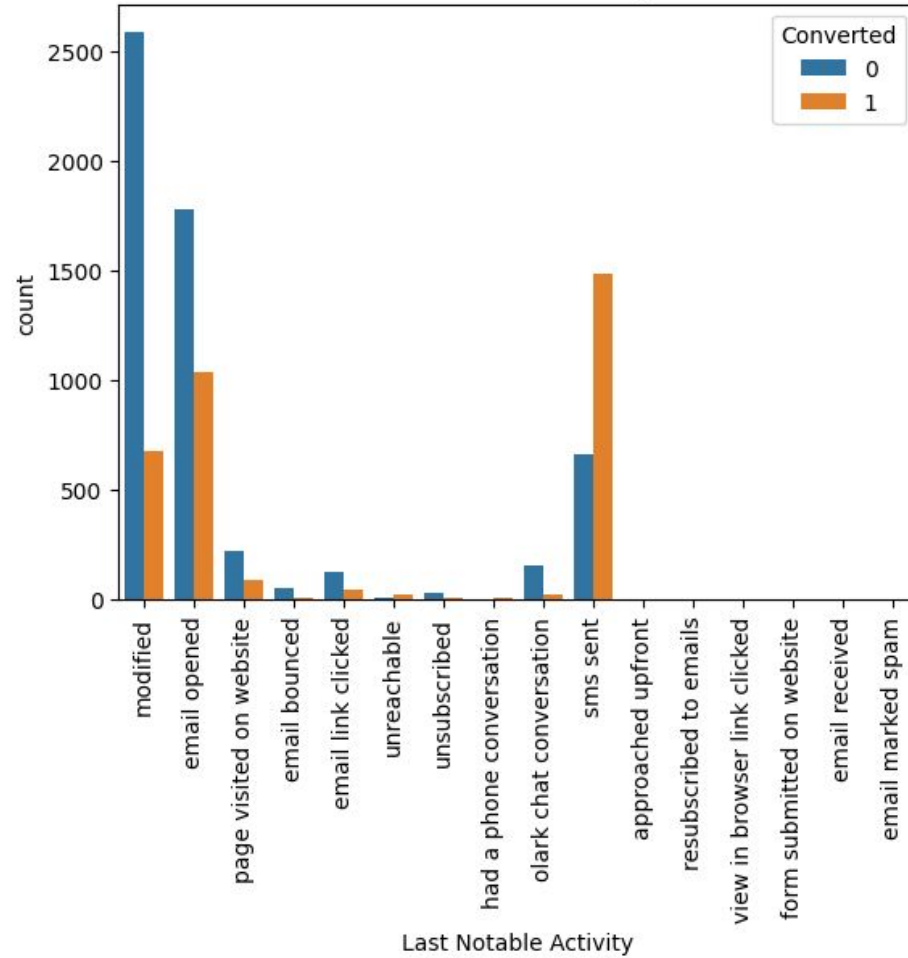
Through Recommendations

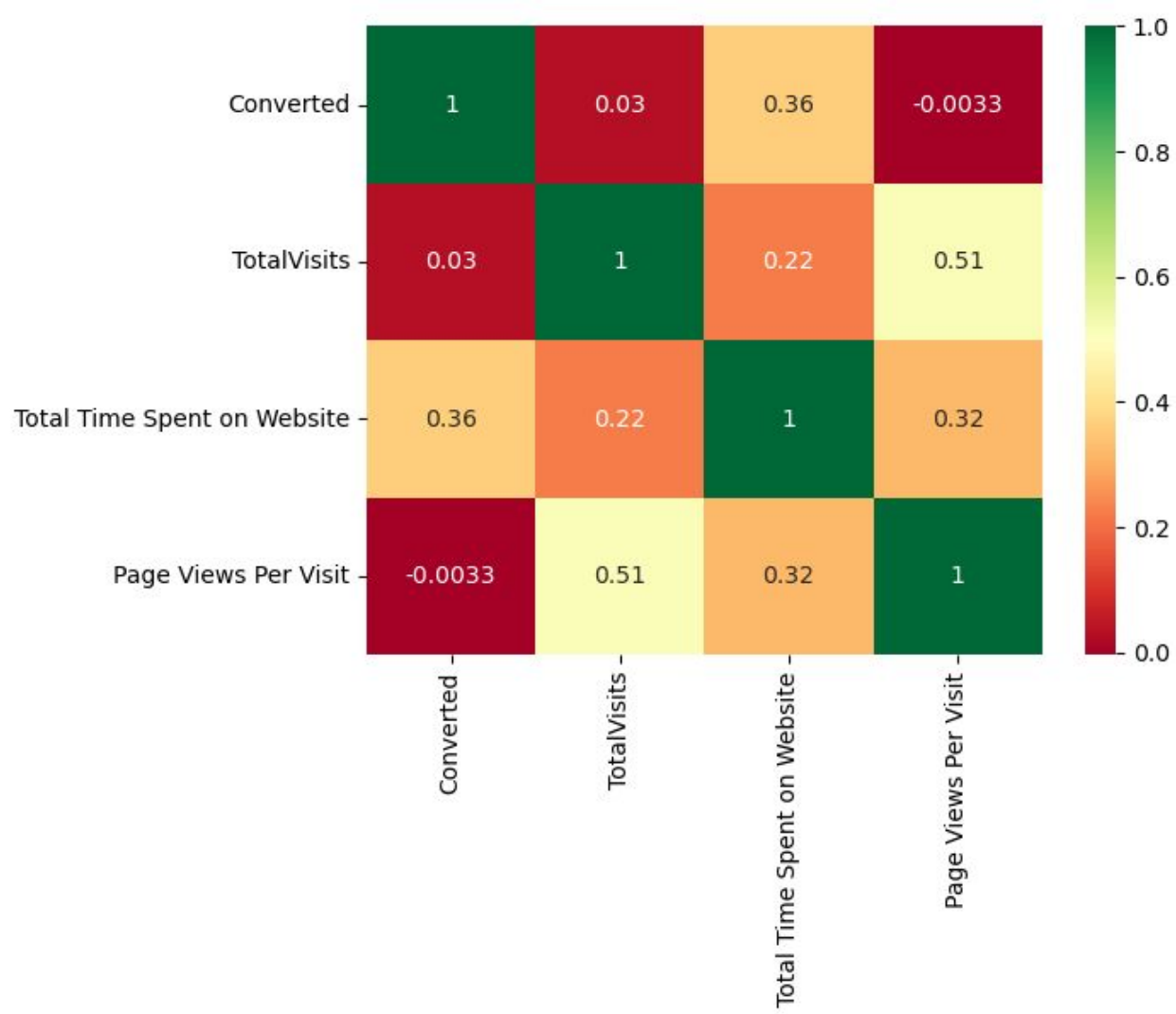


A free copy of Mastering The Interview



Last Notable Activity







Data Conversion

- Numerical Variables are Normalised
- Dummy Variables are created for object type variables
- Total Rows for Analysis: 9074
- Total Columns for Analysis: 13

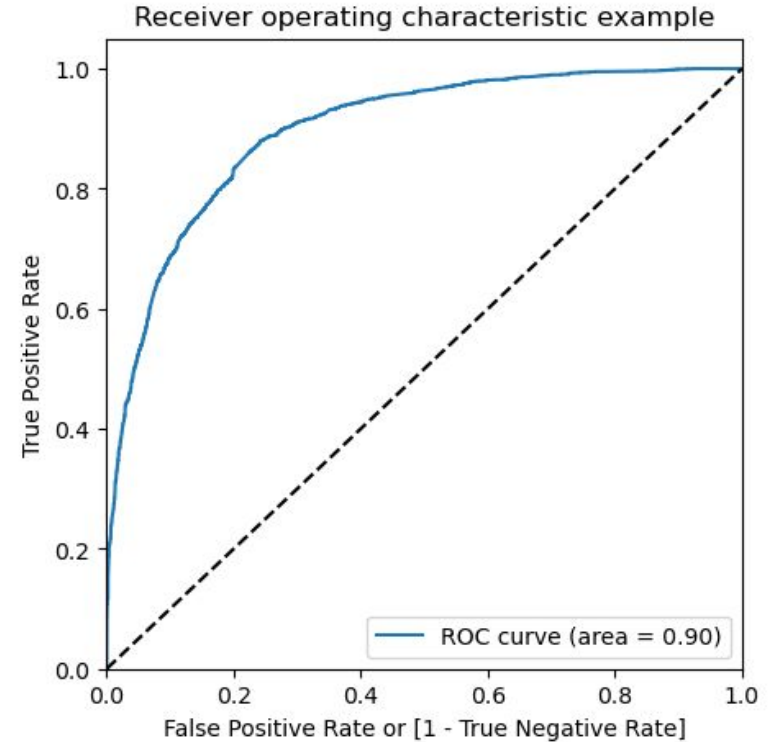


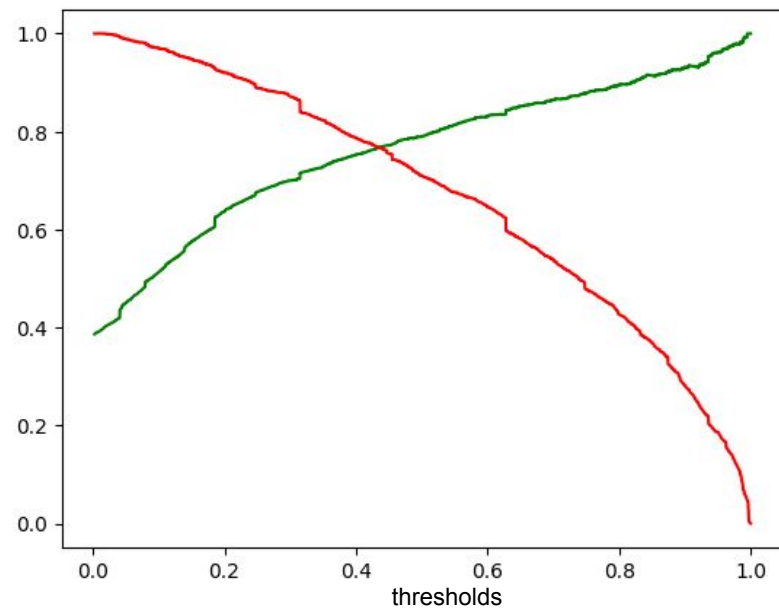
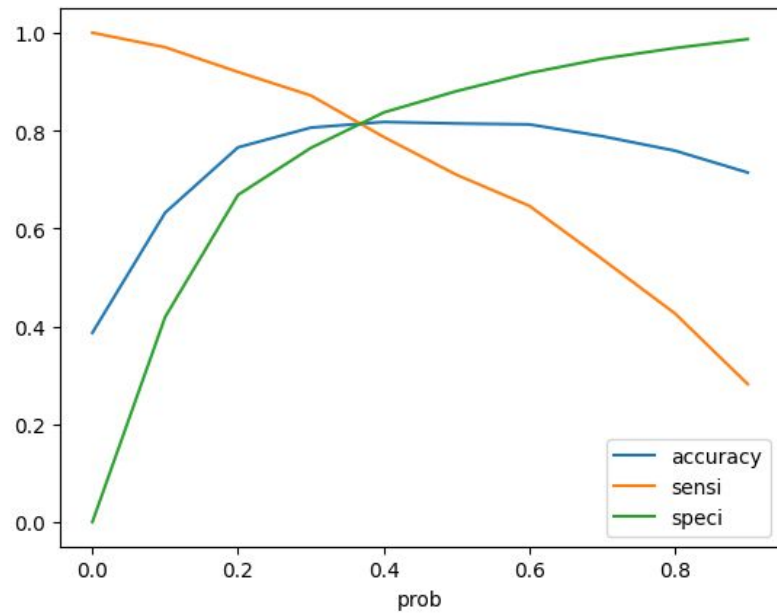
Model Building

- Splitting the Data into Training and Testing Sets
- The first basic step for regression is performing a train-test split, we have chosen 70:30 ratio.
- Use RFE for Feature Selection
- Running RFE with 20 variables as output
- Building Model by removing the variable whose p- value is greater than 0.05 and vif value is greater than 5
- Predictions on test data set
- Overall accuracy 81%

ROC Curve

- Finding Optimal Cut off Point
- Optimal cut off probability is that
- probability where we get balanced sensitivity and specificity.
- From the second graph it is visible that the optimal cut off is at 0.35.





Metric for Train Set

$$\text{Precision} = \frac{\text{TP}}{(\text{TP} + \text{FP})} = 79 \%$$

$$\text{Recall} = \frac{\text{TP}}{(\text{TP} + \text{FN})} = 70 \%$$

$$\text{Sensitivity} = \frac{\text{TP}}{(\text{TP} + \text{FN})} = 81 \%$$

$$\text{Specificity} = \frac{\text{TN}}{(\text{TN} + \text{FP})} = 80 \%$$

		Actual Values	
		Positive	Negative
Predicted Values	Positive	TP	FP
	Negative	FN	TN

Confusion Matrix of Train Set	
3148	757
453	1993

Metric for Test Set



$$\text{Precision} = \frac{\text{TP}}{(\text{TP} + \text{FP})} = 70 \%$$

$$\text{Recall} = \frac{\text{TP}}{(\text{TP} + \text{FN})} = 80 \%$$

$$\text{Sensitivity} = \frac{\text{TP}}{(\text{TP} + \text{FN})} = 80 \%$$

$$\text{Specificity} = \frac{\text{TN}}{(\text{TN} + \text{FP})} = 81 \%$$

		Actual Values	
		Positive	Negative
Predicted Values	Positive	TP	FP
	Negative	FN	TN

Confusion Matrix of Test Set	
1403	331
194	795

Conclusion



It was found that the variables that mattered the most in the potential buyers are:

1. The total time spend on the Website.
2. Total number of visits.
3. When the lead source was:
 - a. Reference
 - b. Olarc Chat
 - c. Welingak website
4. When the last activity was:
 - a. SMS
 - b. Phone conversation
 - c. Unsubscribed
 - d. Olarc chat
5. Total time spent on the website.
6. When their current occupation is as a working professional.

Keeping these in mind the X Education can flourish as they have a very high chance to get almost all the potential buyers to change their mind and buy their courses.



Thankyou