# LEAD SCORING CASE STUDY FOR X EDUCATION

BY:

Yashika Singh

Yash Jindal

Yogesh

# Problem Statement

- An education company named X Education sells online courses to industry professionals. On any given day, many professionals who are interested in the courses land on their website and browse for courses.

- The company wants its leads generation rate to be increased. For which they need to identify what are the potential leads aka Hot Leads.

- For this company wants to create a model where we can assign lead score to each of the identified leads for higher chances of conversion. Target for the company is 80 percent.

## GOALS TO ACHIEVE

- Build a logistic regression model to assign a lead score between 0 and 100 to each of the leads which can be used by the company to target potential leads.

- A higher score would mean that the lead is hot, i.e. is most likely to convert whereas a lower score would mean that the lead is cold and will mostly not get converted.

## Problem solving methodology

- Reading and Understanding the data

- Data Cleaning

- Data Visualization using EDA

- Data Preparation for Modelling

- Model Building: Building the model with features selected by RFE. Eliminate all features with high p values and VIF values and finalize the model

- Model Evaluation: with various metrics like sensitivity, specificity, precision, recall, etc.
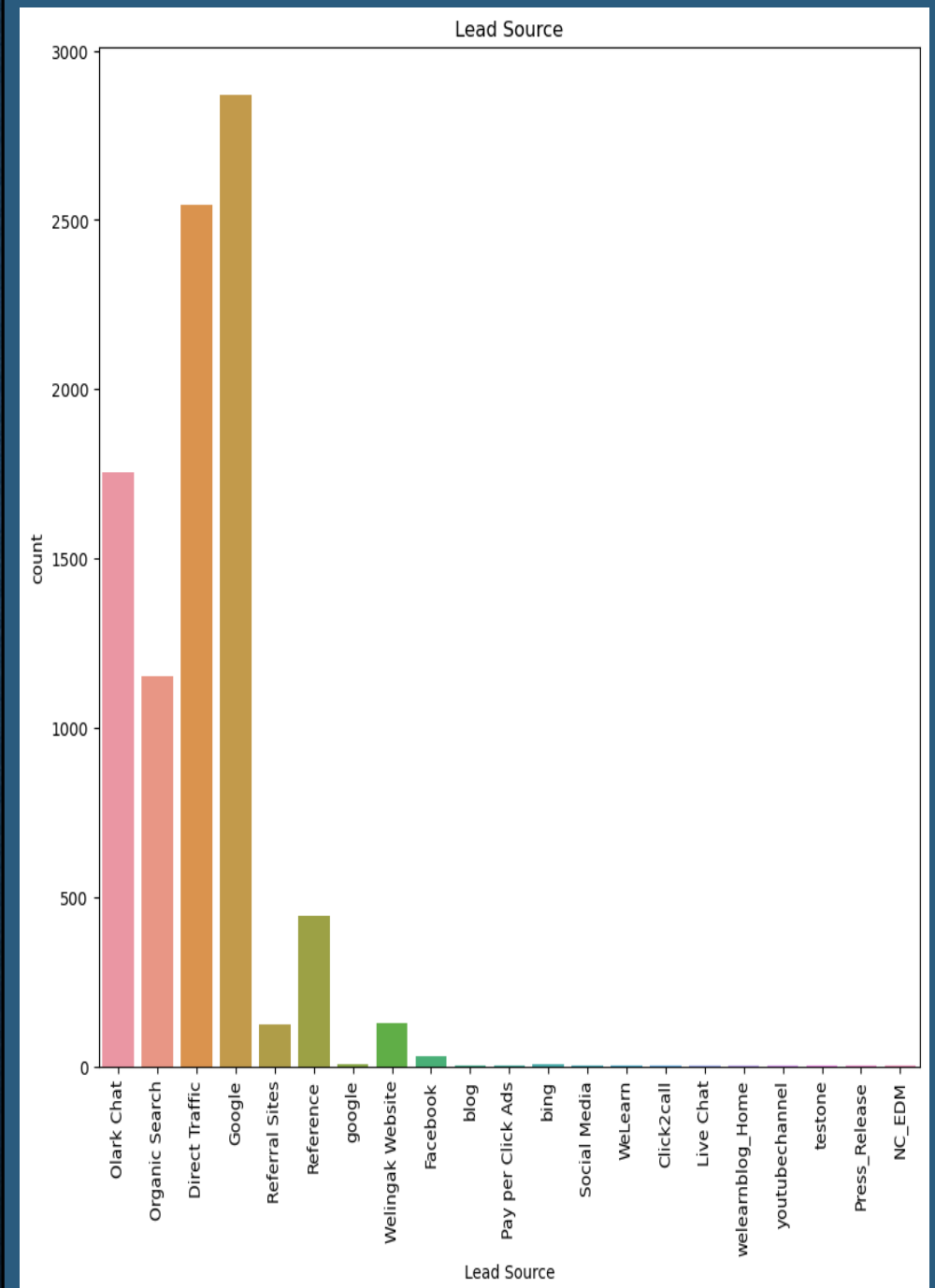
# UNDERSTANDING THE DATA

• We import all the necessary libraries for e.g. NumPy , pandas, matplotlib seaborn etc., import all the warnings.
• We read the data and check the no. of rows and columns. We also, check if there are any missing/ null values or not. Afterwards, we see the statistical summary of the data.
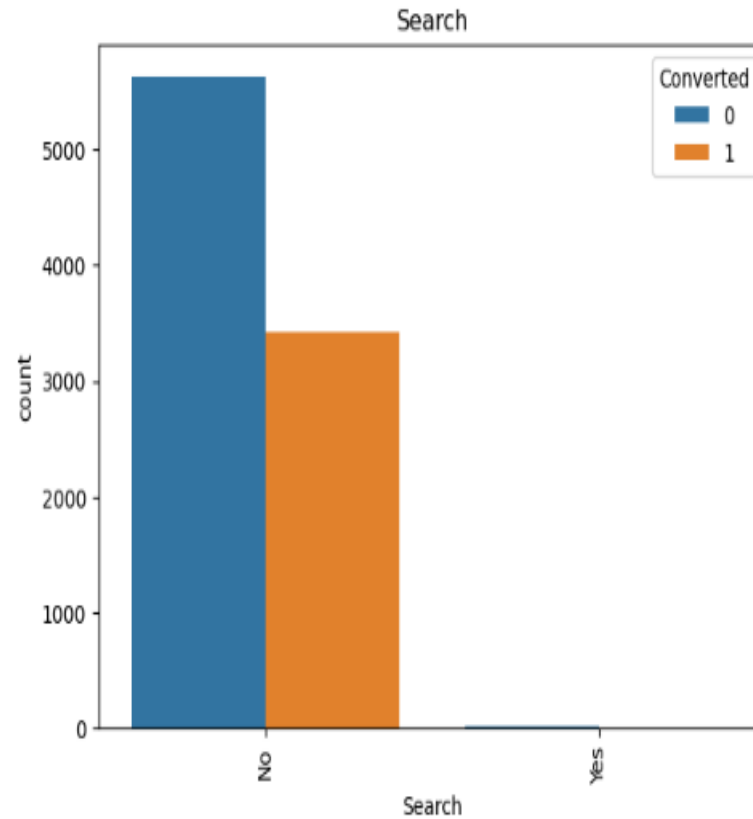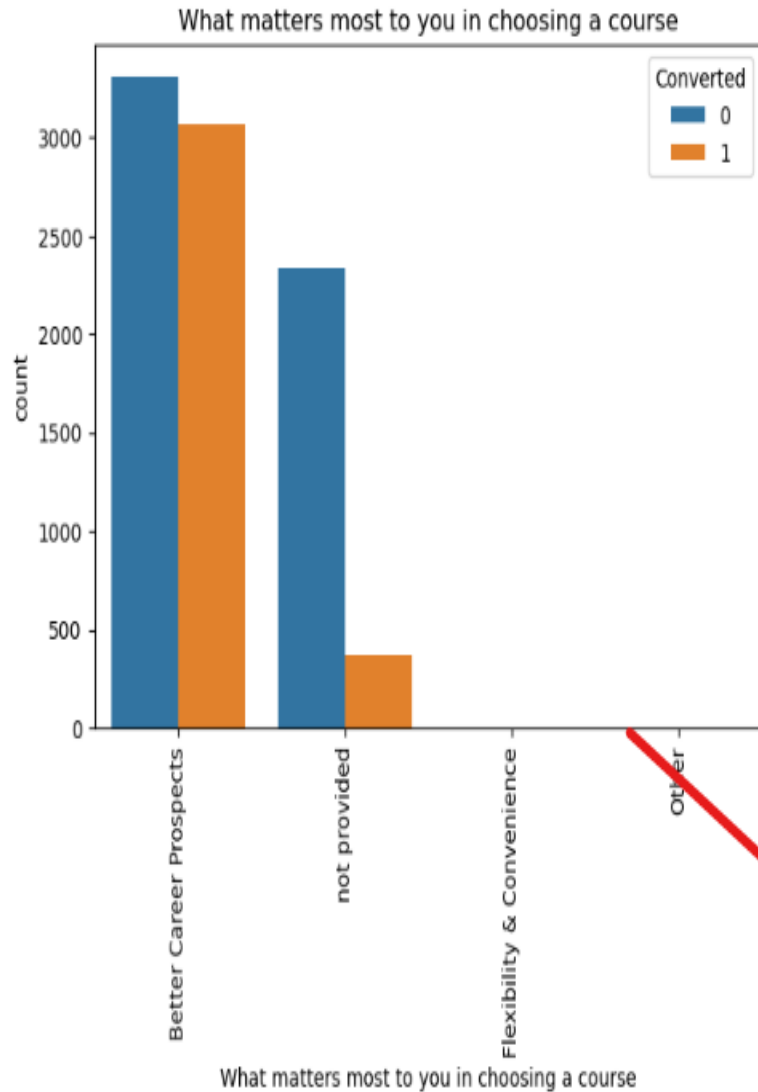
# DATA CLEANING

• We saw there were few columns with high percentage of null values, so we decided to drop those columns.
• few columns had null values, but column's were important for analysis so all null values repaced with 'not provided'.
• Few of the columns had values as 'Select' so we replaced it by 'NaN'.
• Few columns were having outliers and treatment of outlier was performed.
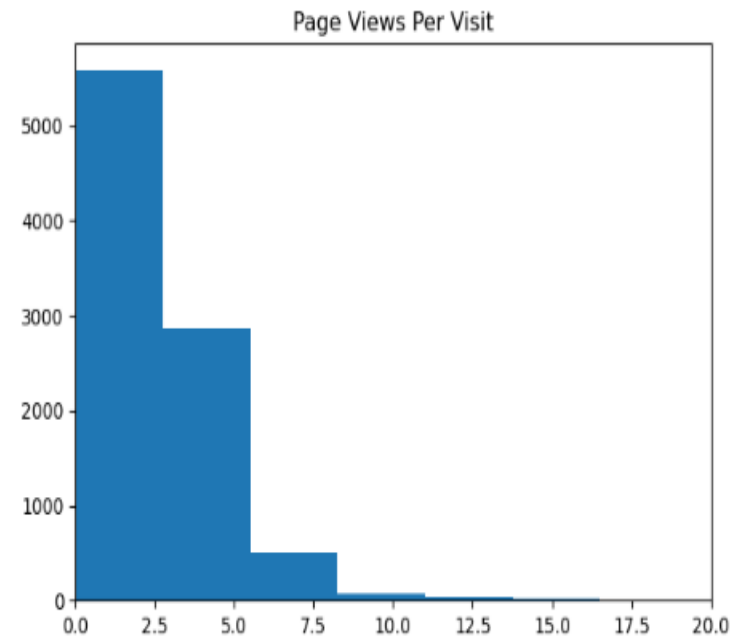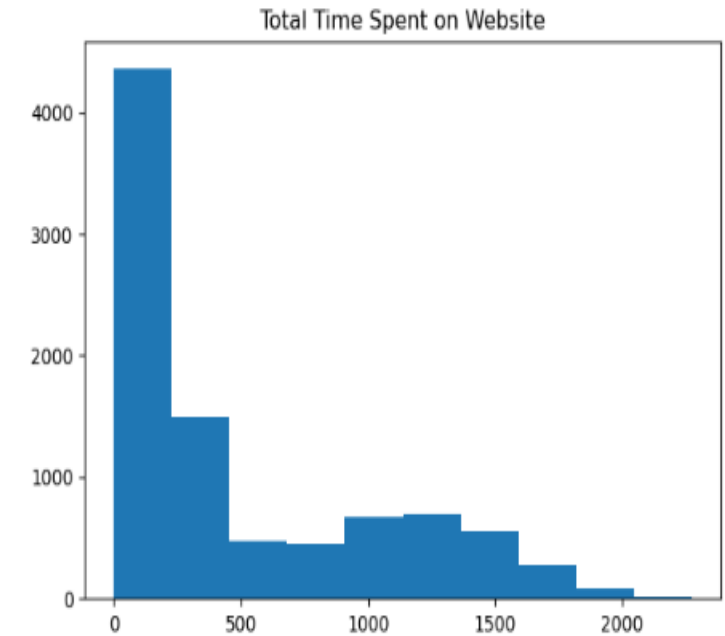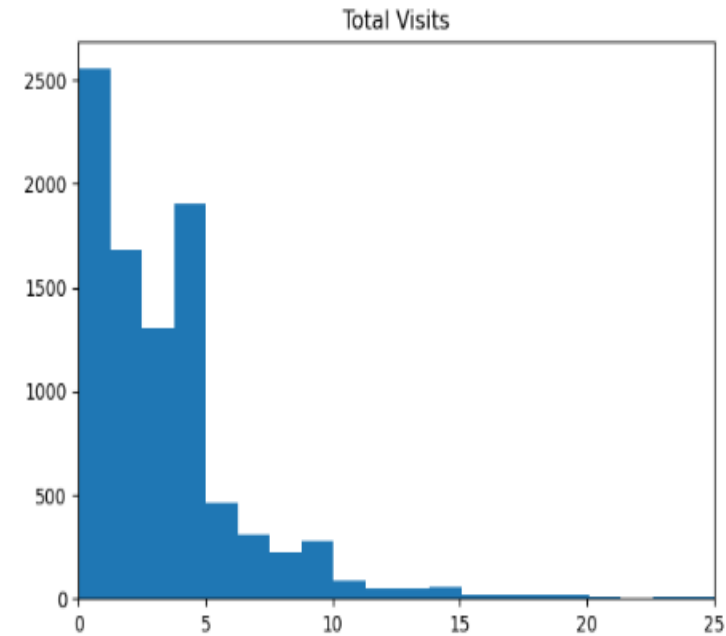
# DATA VISUALIZATION USING EDA

- For data visualization we performed exploratory data analysis (EDA). We first did the univariate analysis of categorical and continuous data.

- Quick check was done on % of null value and we dropped columns with more than 45% missing values.

- Moved ahead with the bi-variate analysis of categorical and continuous data.

- We also found in Lead Source data most of the records were from Google and few were from social media platforms and we classified as same.

- We found correlations between variables by using different plots.

❑ From this we inferred that candidates focus on better career prospectus while choosing a course.

❑ Very less number of leads saw the advertisement regarding the course which shows that less visibility of the website.

From this data visualisation graph we can infer that more than 2500 candidates visit the website but spend around two minutes and visit only two pages. Therefore, website should be improved and made more enthusiastic so that candidates spend more time on the website and see maximum pages possible.

# DATA PREPARATION FOR MODELLING

- Data preparation for multiple linear regression involves handling the categorical variables first and then performing dummy encoding.

- We then performed the train and test split using 70%-30% rule and then performed the scaling of variables. Since scaling of variables is an important step, we may have different variables of different scales. So, it's important to have everything on the same scale for the model to be easily interpretable.

- Therefore, we used **MinMaxScaler** for the same.

# MODEL BUILDING

- We follow the bottom up approach for this, i.e. we start by building the model with just one variable. Hence, the choice of variables becomes very crucial.

- RFE was done to attain the top 15 relevant variables. First, we will look at the significance of variables and based on the significance.

- We checked VIF values and p-value (The variables with VIF < 5 and p-value < 0.05 were kept). In our model for all the feature VIF < 5 and p-value < 0.05.

# FEATURE SELECTION USING RFE
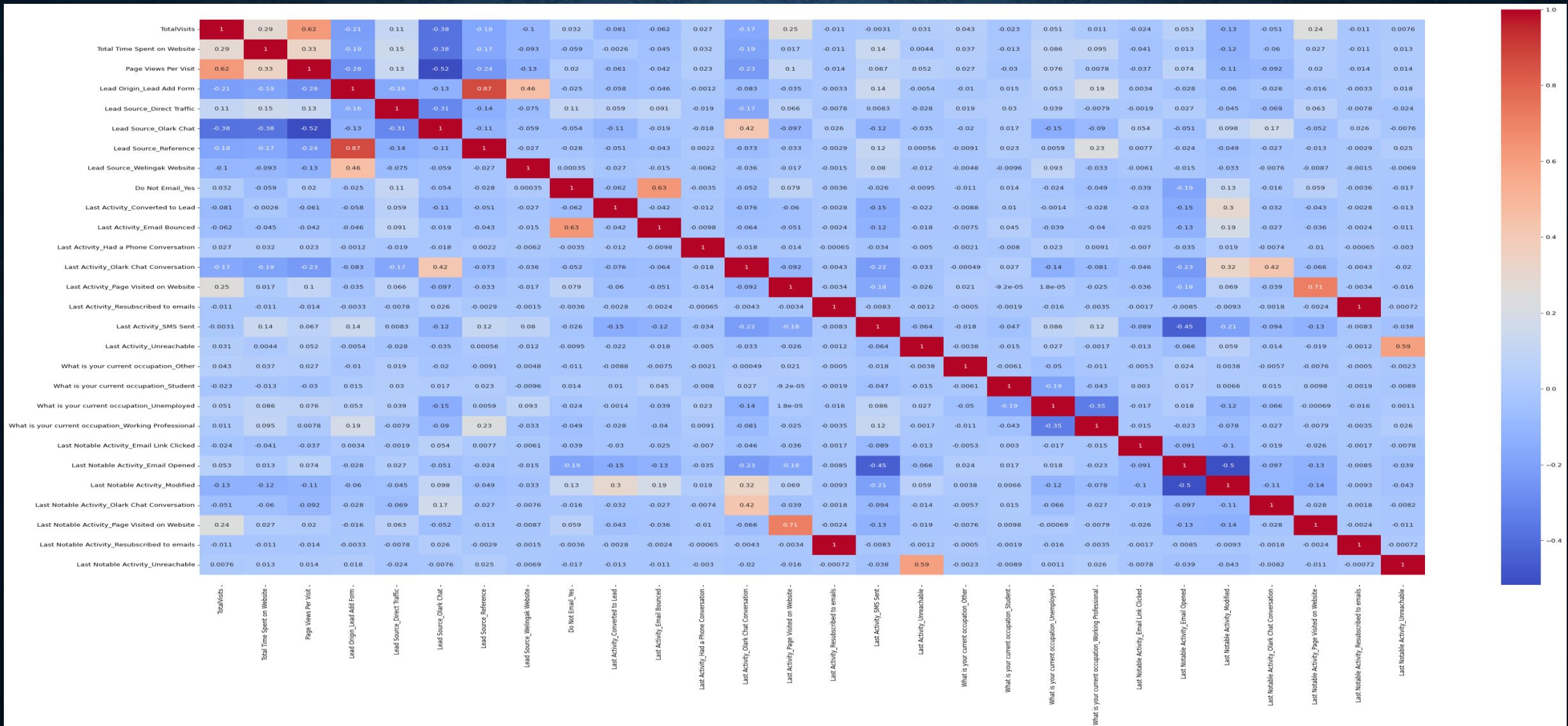
| | TotalVisits | Total Time Spent on Website | Page Views Per Visit | Lead Origin_Lead Add Form | Lead Source_Direct Traffic | Lead Source_Olark Chat | Lead Source_Reference | Lead Source_Welingak Website | Do Not Email_Yes | Activity_Converted to Lead | ... | WI occupa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1289 | 0.014184 | 0.612676 | 0.083333 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | |
| 3604 | 0.000000 | 0.000000 | 0.000000 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | ... | |
| 5584 | 0.042553 | 0.751761 | 0.250000 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | ... | |
| 7679 | 0.000000 | 0.000000 | 0.000000 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | ... | |
| 7563 | 0.014184 | 0.787852 | 0.083333 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | ... | |

5 rows × 30 columns

- Recursive feature elimination is an optimization technique for finding the best performing subset of features.

- It is based on the idea of repeatedly constructing a model and choosing either the best (based on coefficients), setting the feature aside and then repeating the process with the rest of the features.

- This process is applied until all the features in the dataset are exhausted. Features are then ranked according to their time of elimination.

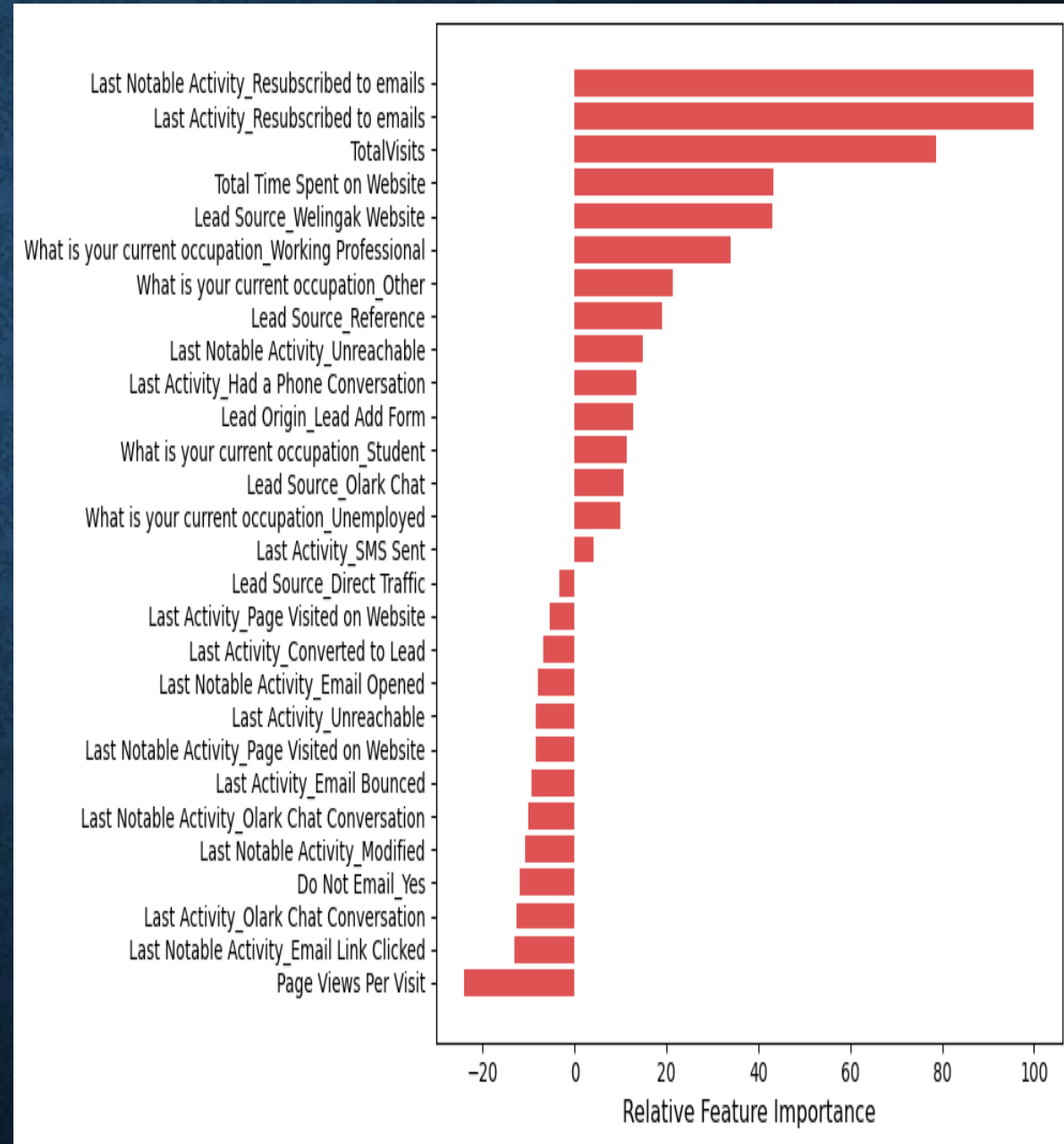# TO CHECK CORRELATION AMONG VARIABLES

# LEAD SCORE CALCULATION

- Lead score is calculated for all the leads in the original data frame.

- Formula for Lead Score calculation is Lead Score = 100* Conversion Probability.

- Higher the lead score, higher is the probability of a lead getting converted and vice versa.

- Since, we had used 0.42 as our final probability threshold for deciding if a lead will convert or not, any lead with a lead score of 34 or above will have a value of '1' in the final predicted column.

| | Lead Number | Converted | Conversion_Probability | final_predicted | Lead_Score |
|---|---|---|---|---|---|
| 0 | 660737 | 0.0 | 0.190369 | 0.0 | 19.0 |
| 1 | 660728 | 0.0 | 0.422565 | 1.0 | 42.0 |
| 2 | 660727 | 1.0 | 0.751855 | 1.0 | 75.0 |
| 3 | 660719 | 0.0 | 0.063021 | 0.0 | 6.0 |
| 4 | 660681 | 1.0 | 0.542724 | 1.0 | 54.0 |
| 5 | 660680 | 0.0 | 0.035828 | 0.0 | 4.0 |
| 6 | 660673 | 1.0 | 0.772298 | 1.0 | 77.0 |
| 7 | 660664 | 0.0 | 0.035828 | 0.0 | 4.0 |
| 8 | 660624 | 0.0 | 0.043511 | 0.0 | 4.0 |
| 9 | 660616 | 0.0 | 0.054540 | 0.0 | 5.0 |

## FEATURE DETERMINATION

◈ The Coefficient (beta) values for each of these features from the model parameters are used to determine the order of importance of these features.

◈ Features with high positive beta values are the ones that contribute most towards the probability of a lead getting converted.

◈ Similarly, features with high negative beta values contribute the least

# MODEL EVALUATION

- We first created confusion matrix to find the TP, TN, FP and FN. We calculated the sensitivity and specificity.

- We plotted the graph of sensitivity, accuracy and specificity for each level of probability. We found that 0.42 was the cut-off point.

- We then used the cut-off point 0.42 to select the person and see if he would be converted or not then ,

- We again created the confusion matrix to calculate TP, TN, FP & FN calculated the sensitivity and specificity and got 0.84984 and 0.78669 respectively.

Thank You