# CASE STUDY: LEAD SCORING
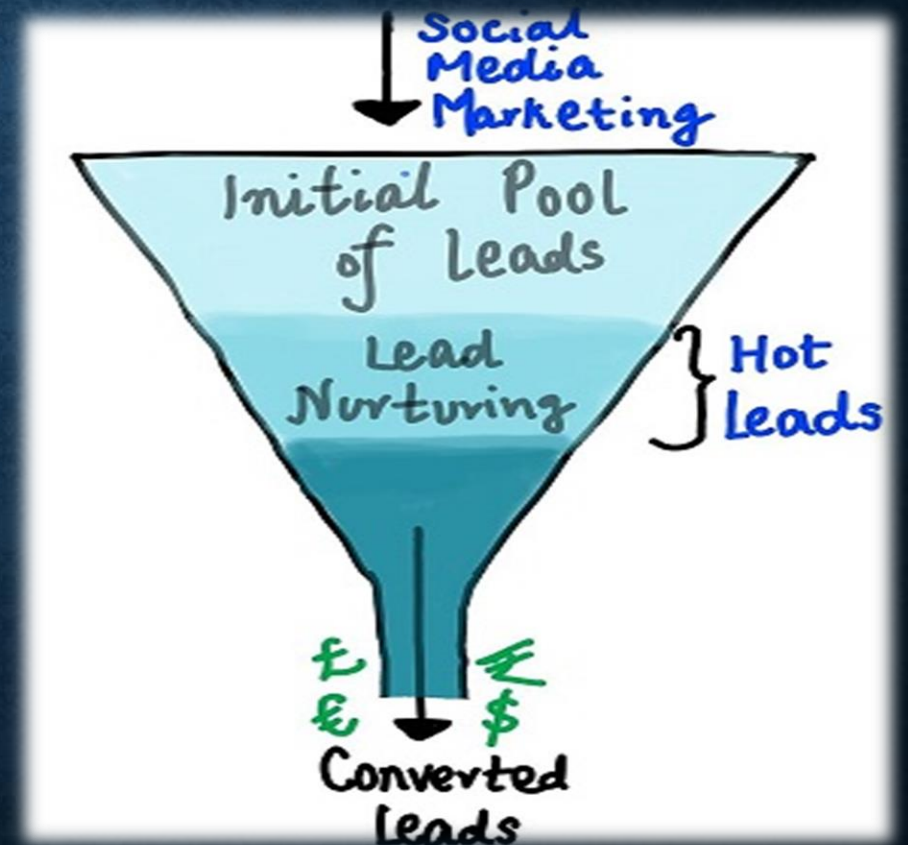
Group Members :

1.Satyam Singh

2. Somnath Nemani

3.Pushyamitra Sharma

# PROBLEM STATEMENT

- X Education Sells Online Courses To Industry Professionals And Many Professionals Browse Through Their Website For Courses. When These People Fill Up A Form And Provide Their Email Address And Their Phone Number They Are Classified As A Lead.

- The Lead Conversion Rate At X Education Is Around 30% ( If They Acquire 100 Leads In A Day Only 30 Would Get Converted) Which Clearly Is A Very Poor Rate.

- The Company Wishes To Identify The Most Potential Leads Which Are " Hot Leads".

- If They Identify This Set Of Leads, The Lead Conversion Rate Would Go Up As The Sales Team Would Focus On Communicating With The Potential Leads Rather Than Making Calls To Everyone

A Typical Lead Conversion Process Can Be Represented Using The Following Funnel:

# BUSINESS OBJECTIVE

X Education is keen on pinpointing the most promising leads, emphasizing those with a high probability of transitioning into paying customers. Our assignment involves crafting a model to allocate lead scores, strategically attributing higher scores to customers with an elevated conversion likelihood, aligning with the CEO's target conversion rate of around 80%. It's imperative that the model is designed for long-term applicability, ensuring its effectiveness in future scenarios.

# METHODOLOGY USED

**Library Import:**

Initiate the process by importing crucial libraries like Pandas, NumPy, and others essential for data analysis. Additionally, load the provided dataset to facilitate subsequent steps in the analysis.

**Data Cleaning and Manipulation:**

**Missing Values Assessment**: Conduct a thorough examination of the dataset to identify the presence and extent of missing values in various columns.

**Null Percentage Calculation**: Calculate the percentage of null values for each column, providing insights into the overall data completeness.

**Dropping Irrelevant Columns**: Evaluate columns with a substantial number of missing values, opting to eliminate those that do not contribute significantly to the analytical objectives.

# METHODOLOGY USED

**Imputation of Missing Values**: Where applicable, perform imputation of missing values to enhance the completeness of the dataset. This ensures a more robust foundation for subsequent analyses.

**Outlier Detection and Handling:** Scrutinize the dataset for outliers—data points significantly deviating from the norm. Implement appropriate strategies to handle these outliers, maintaining the integrity and reliability of the dataset for subsequent analytical procedures.

# METHODOLOGY USED

**Exploratory Data Analysis (EDA)**

**Univariate Analysis**

**Value Count**: Explore the frequency distribution of individual variables, gaining insights into the occurrence of unique values within each attribute.

**Distribution of Variables**: Analyze the distribution patterns of individual variables, providing a comprehensive understanding of the range and spread of each variable.

**Bivariate Analysis**

**Correlation Coefficients**: Investigate the relationships between pairs of variables by calculating correlation coefficients. This quantifies the strength and direction of linear associations.

**Patterns Between Variables**: Explore patterns, dependencies, and interactions between variables, unraveling insights into how changes in one variable relate to changes in another.

*This two-fold approach to Exploratory Data Analysis aims to uncover intrinsic patterns, relationships, and characteristics within the dataset, setting the stage for informed decision-making and further analysis.*

# METHODOLOGY USED

**Data Preparation**

**Creating Dummy Variables for Categorical Variables:** Generate dummy variables for categorical attributes, converting them into a format suitable for machine learning models. Exclude the first dummy variable to avoid multicollinearity issues.

**Test-Train Split**: Divide the dataset into training and testing sets to evaluate model performance. This ensures that the model is trained on one subset and tested on another, providing a reliable assessment of its generalizability.

**Feature Scaling**: Normalize or standardize the features to ensure that variables with different scales do not unduly influence the machine learning model. This step enhances the model's performance and convergence during training.

# METHODOLOGY USED

**Model Building**

**Logistic Regression:** Employ Logistic Regression for model construction and prediction, particularly suited for binary classification tasks.

**Validation of the Model**: Assess the model's performance using appropriate validation techniques to ensure its accuracy and reliability in making predictions.

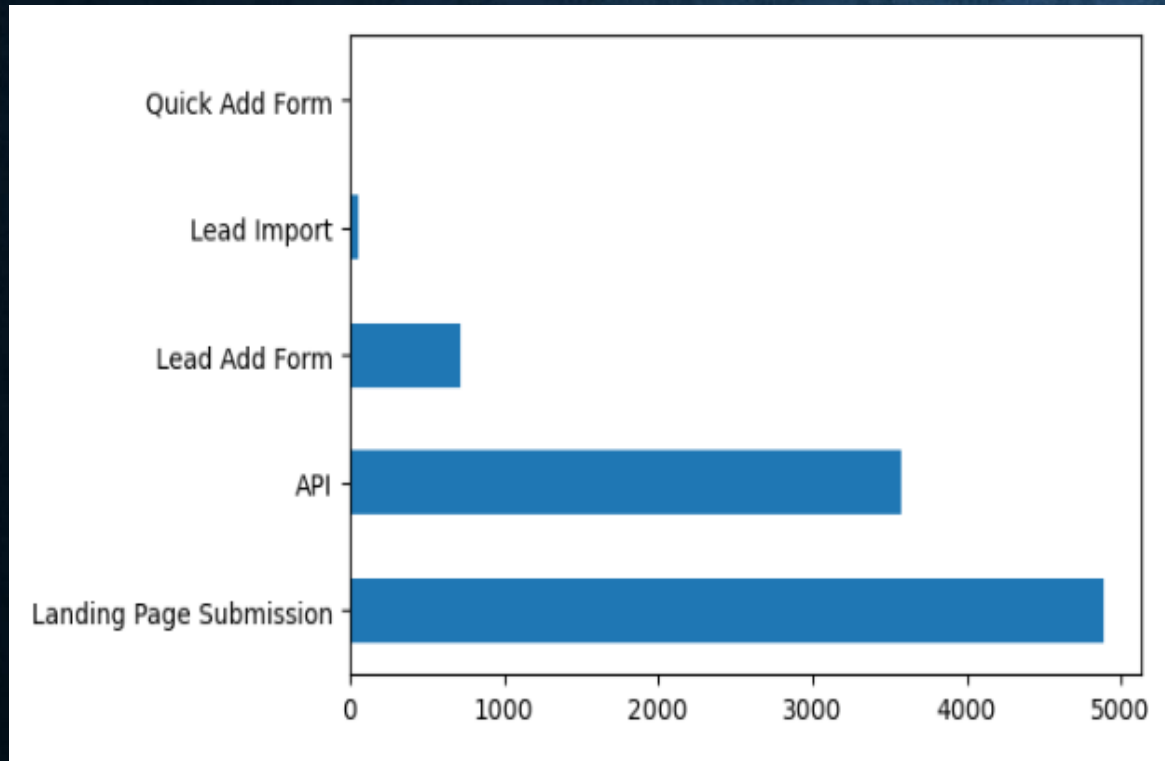**Model Presentation**: Present the finalized model, highlighting key insights, coefficients, and any pertinent information that aids in understanding its predictive capabilities.

**Conclusions and Recommendations**: Summarize findings, draw conclusions from the analysis, and provide actionable recommendations based on the model's insights to inform decision-making.

# EXPLORATORY DATA ANALYSIS

# EXPLORATORY DATA ANALYSIS
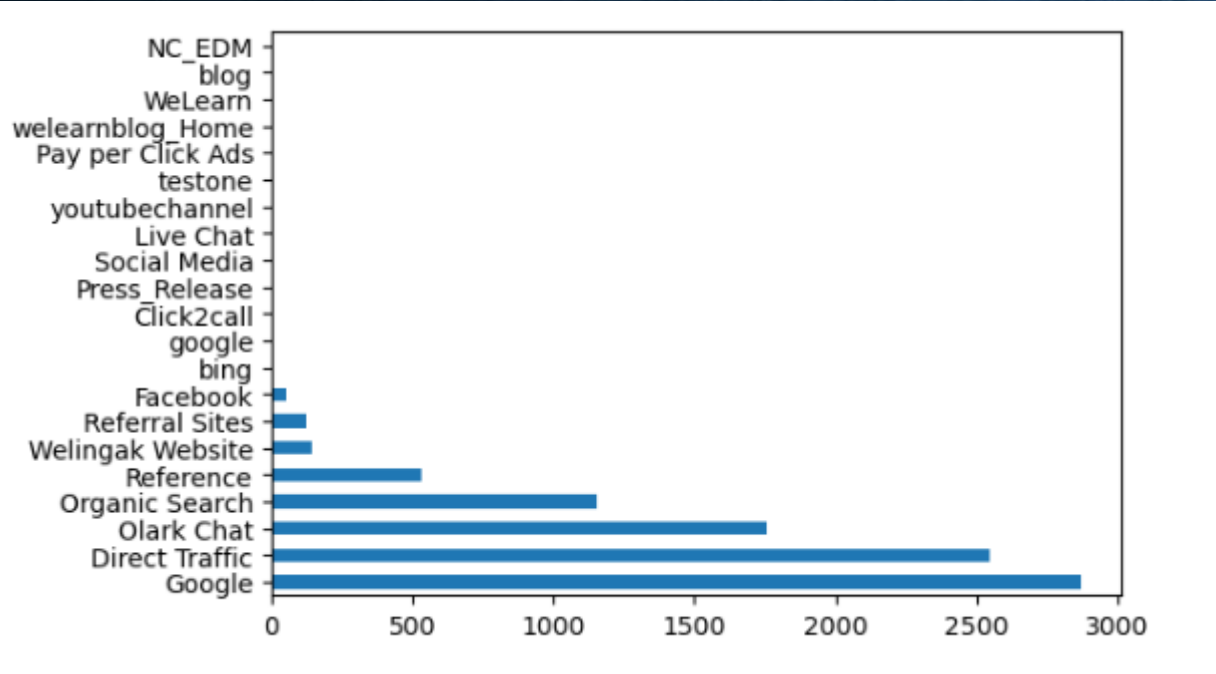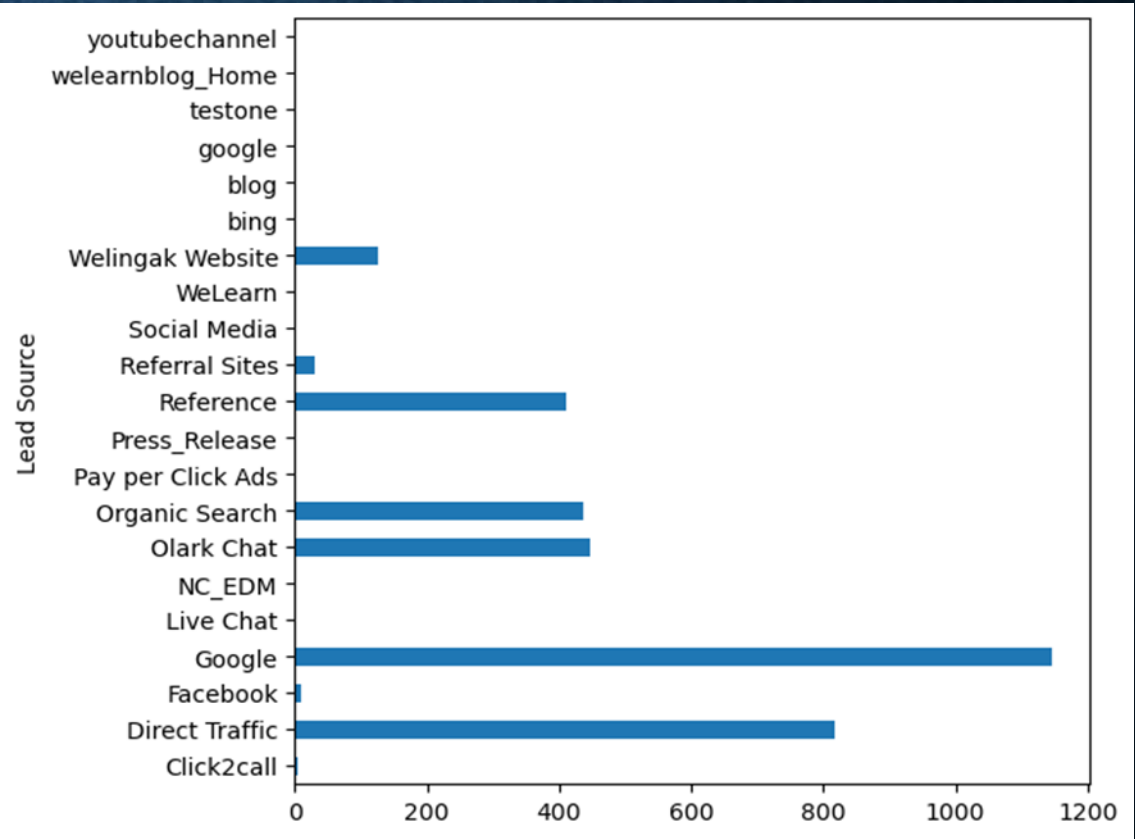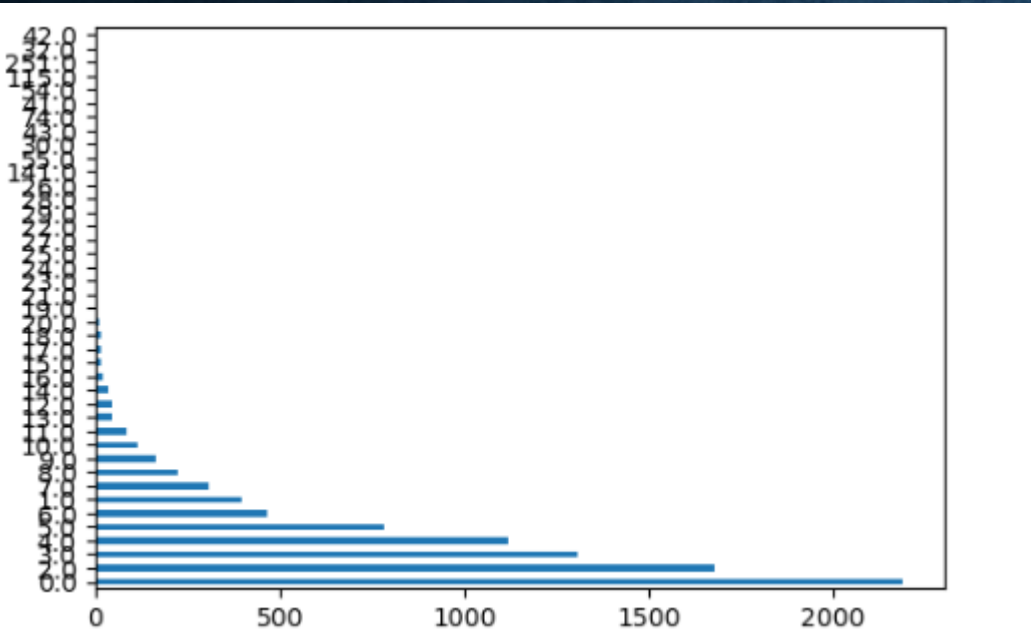
Number of people getting from the source

Number of people getting converted from source

# EXPLORATORY DATA ANALYSIS

Number of people getting from the source

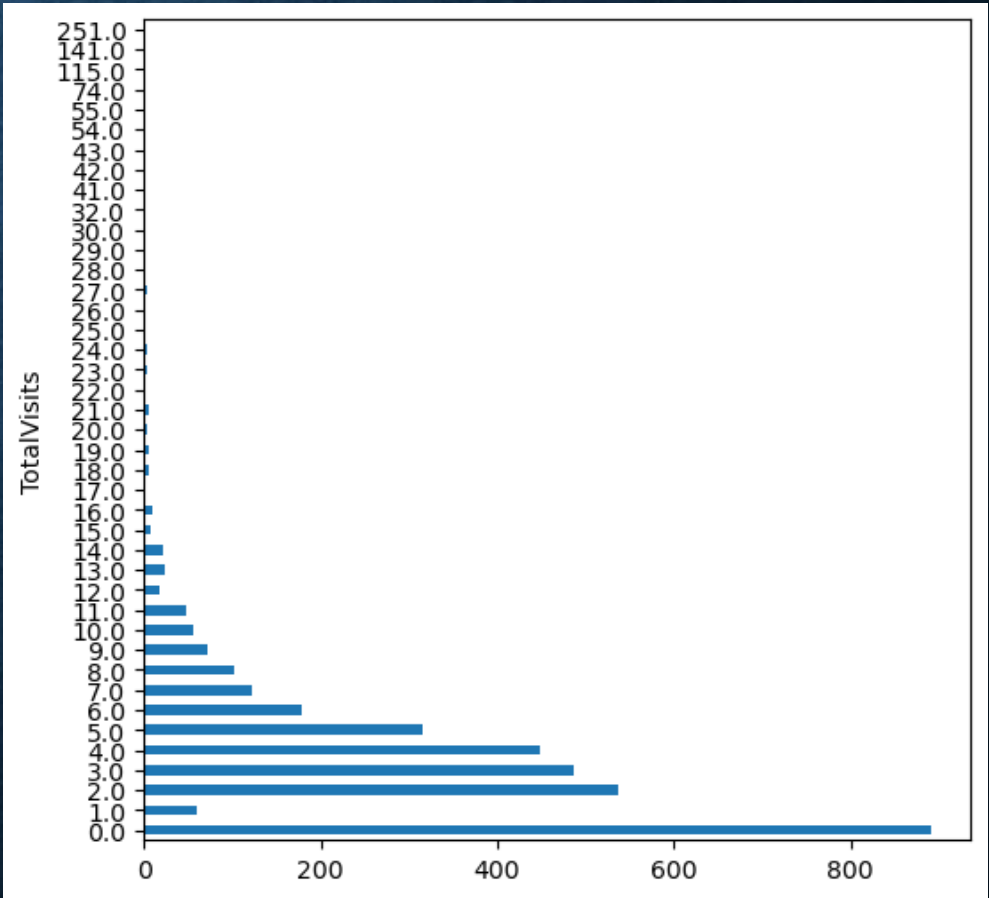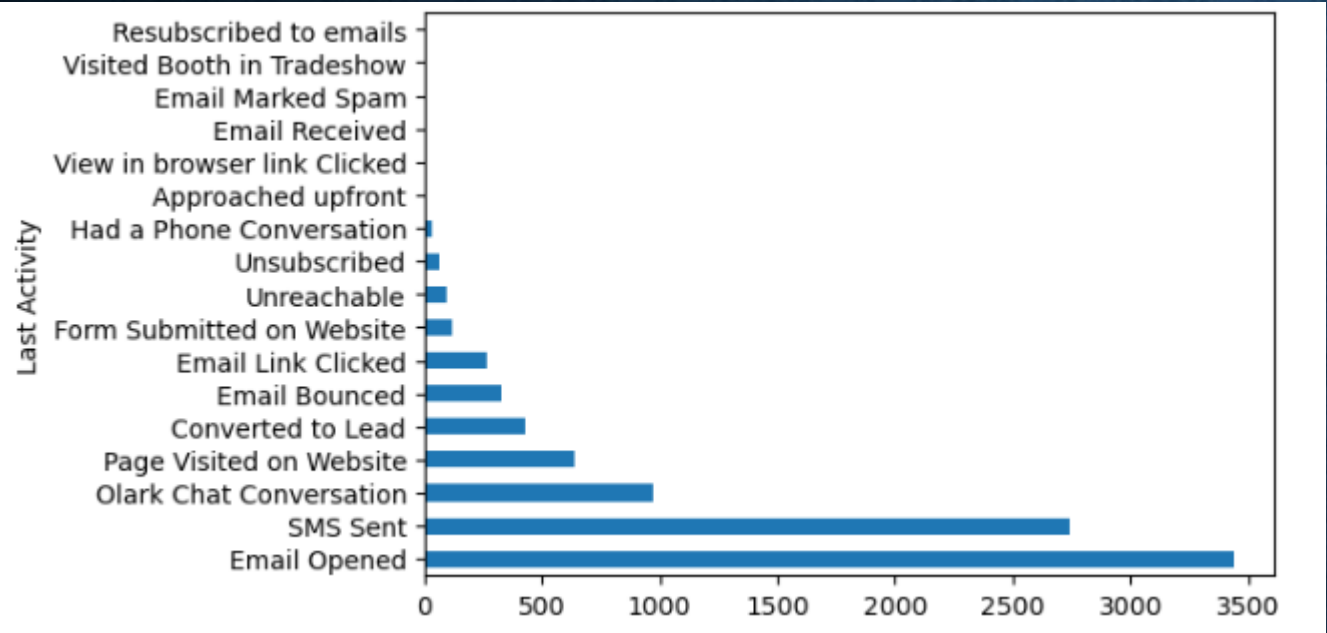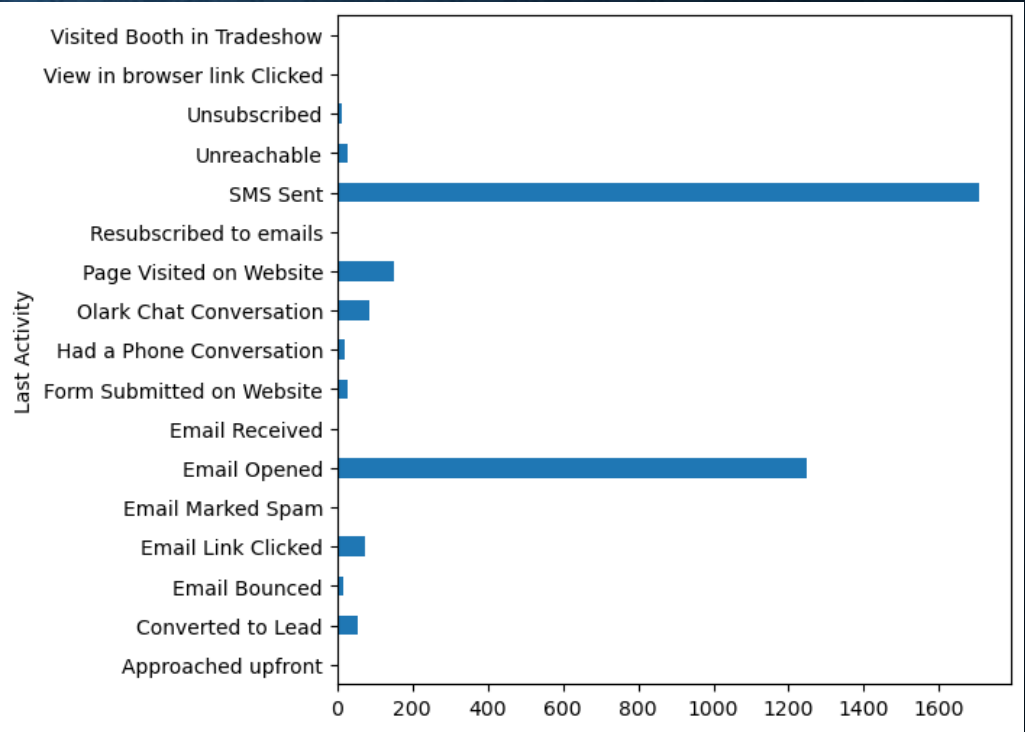Number of people getting converted from source

# EXPLORATORY DATA ANALYSIS
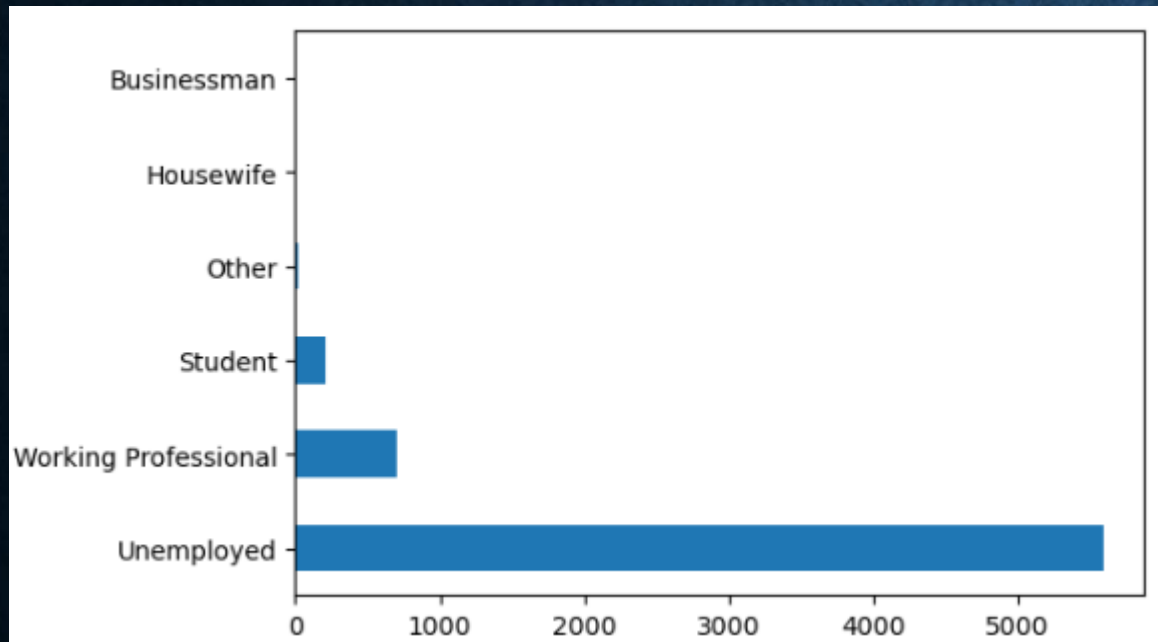
Number of people getting from the source

Number of people getting converted from source

# EXPLORATORY DATA ANALYSIS

Number of people getting from the source

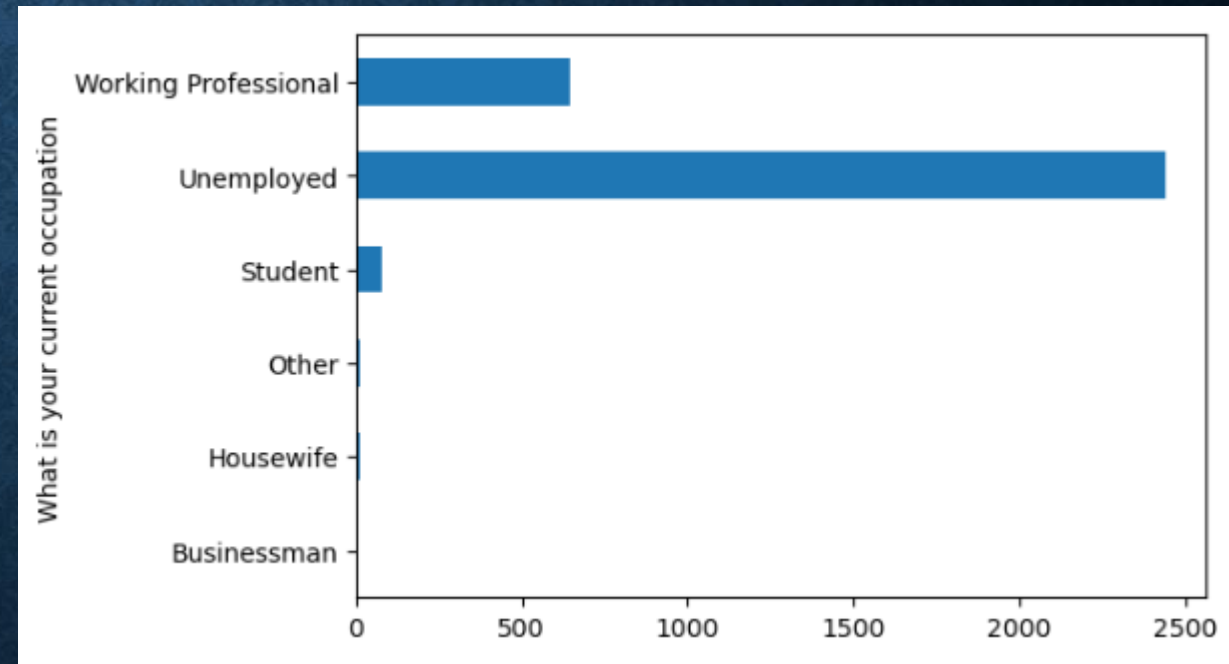Number of people getting converted from source

# EXPLORATORY DATA ANALYSIS

Number of people getting from the source

Number of people getting converted from source

# EXPLORATORY DATA ANALYSIS
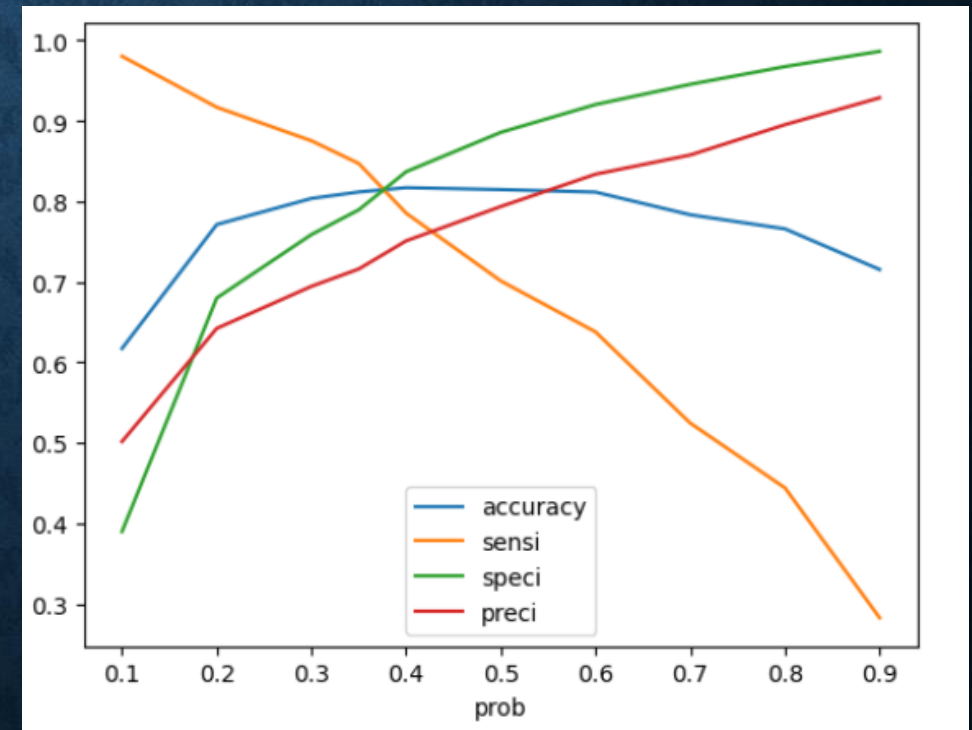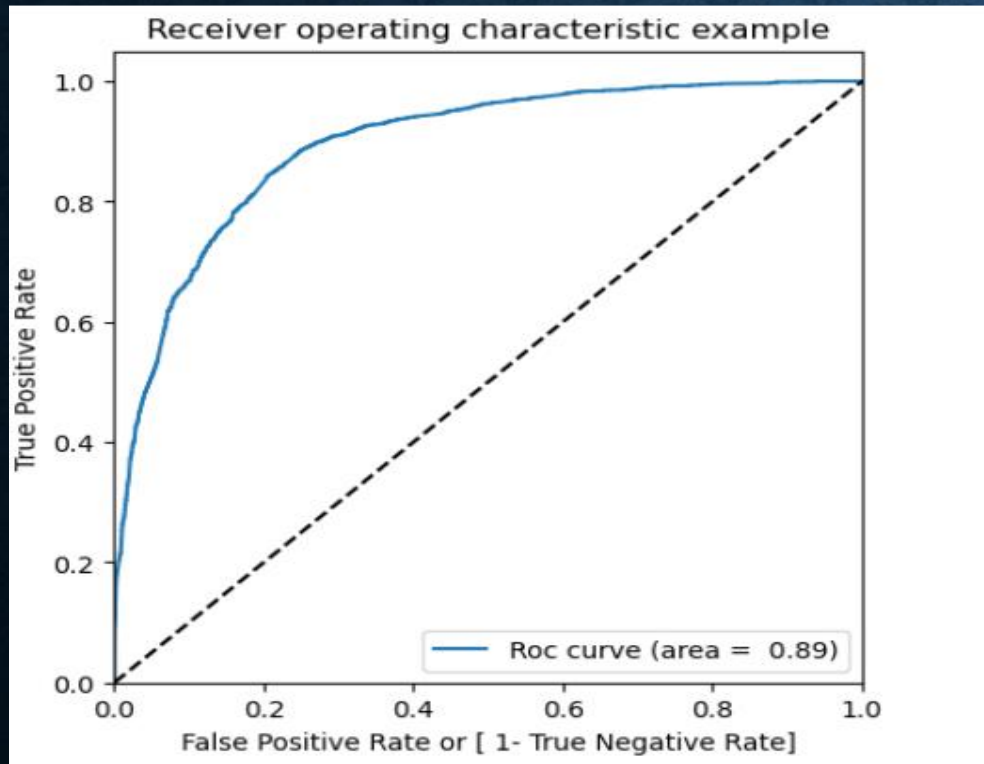
Contribution of the final variables

# RESULTS

1. Split the data into training and testing sets using a 70:30 ratio.

2. Employed RFE for feature selection, yielding 25 variables as output.

3. Built the model by eliminating variables with high VIF values.

4. Finalized the model with 12 variables based on a p-value threshold of 0.

5. Conducted predictions on the test dataset.

6. Achieved an overall accuracy of 81%.

# METHODOLOGY USED / RESULTS

1. Constructed an ROC curve (Receiver Operating Characteristic curve) for model evaluation.

2. Determined the optimal cutoff point using Accuracy, Sensitivity and Specificity.

3. The area under the ROC curve (AUC) is 0.89, indicating a strong discriminative ability of the model in distinguishing between classes.

# CONCLUSION

1. The model demonstrates adaptability to future company requirements.

2. Key variables influencing potential buyers include:

    i. Total time spent on the website.

    ii. Lead source from olark chat ,Reference and Wellingak website.

    iii. Current occupation being a working professional.

    iv. Lead origin through landing page submission.