# Lead Scoring Case Study using logistic regression

**Submitted by**

**1.Karishma Sahay**

**2.Keshav Raj**

**3. Swati Mehta**

# Contents

- Problem Statement
- Problem Approach
- EDA
- Correlations
- Model Evaluation
- Observations
- Conclusion

# Problem Statement

▶ An education company named X Education sells online courses to industry professionals. Many professionals who are interested in the courses land on their website and browse for courses. They have process of form filling on their website after which the company that individual as a lead.

▶ Once these leads are acquired, employees from the sales team start making calls, writing emails, etc. Through this process, some of the leads get converted while most do not.

▶ The typical lead conversion rate at X education is around 30%. Now, this means 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, that are 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.
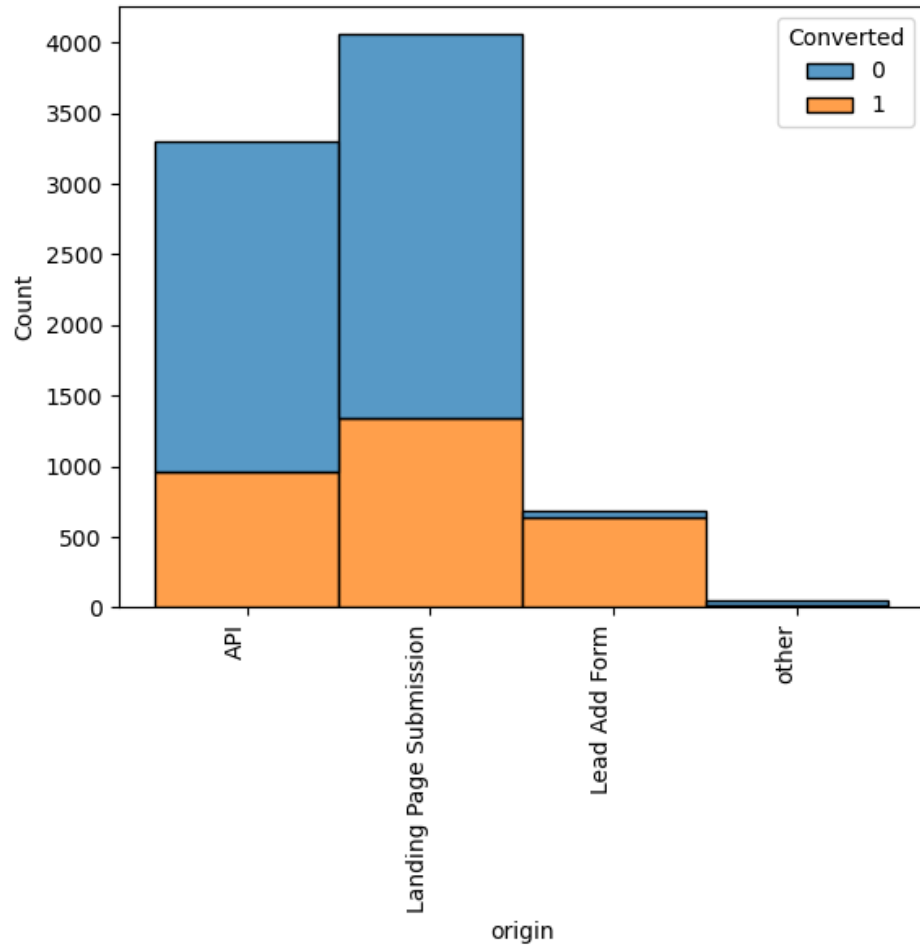
# Objective

- Lead X wants us to build a model to give every lead a lead score between 0 -100.

- So that they can identify the Hot leads and increase their conversion rate as well.

- The CEO want to achieve a lead conversion rate of 80%.

- They want the model to be able to handle future constraints as well like Peak time actions required, how to utilize full man power and after achieving target what should be the approaches.
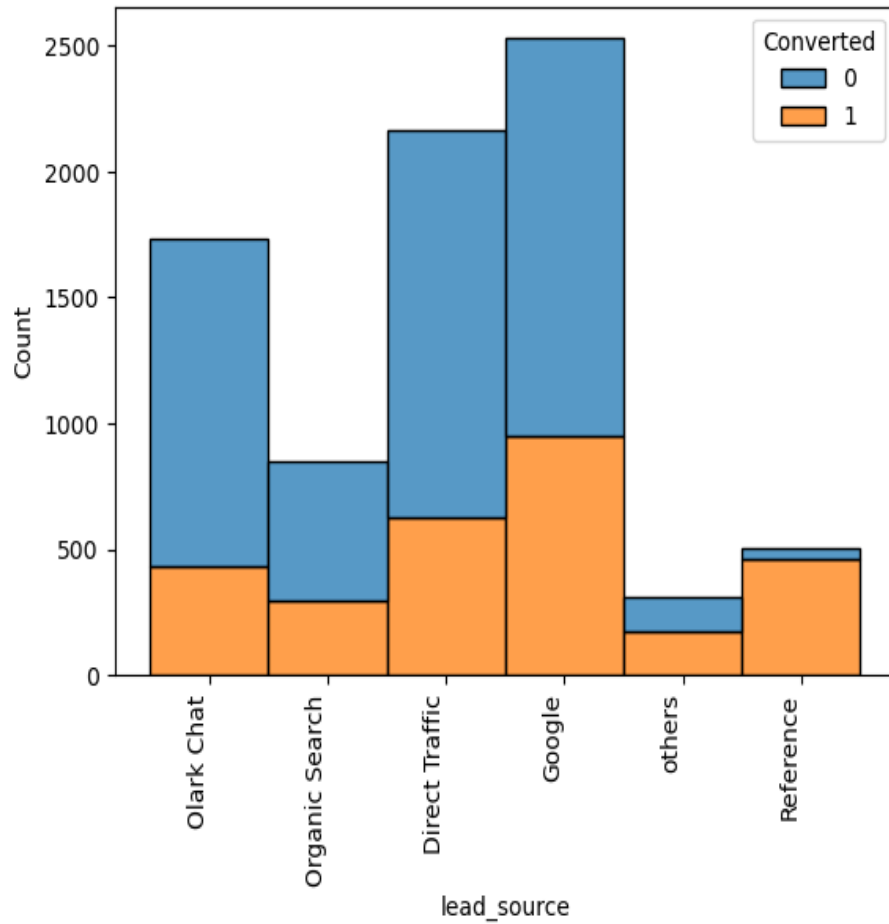
# Problem Approach

- Importing the data and inspecting the
- data frame
- Data preparation
- EDA
- Dummy variable creation
- Test-Train split
- Feature scaling
- Correlations
- Model Building (RFE,VIF and pvalues)
- Model Evaluation (Specificity, Sensitivity,Accuracy)
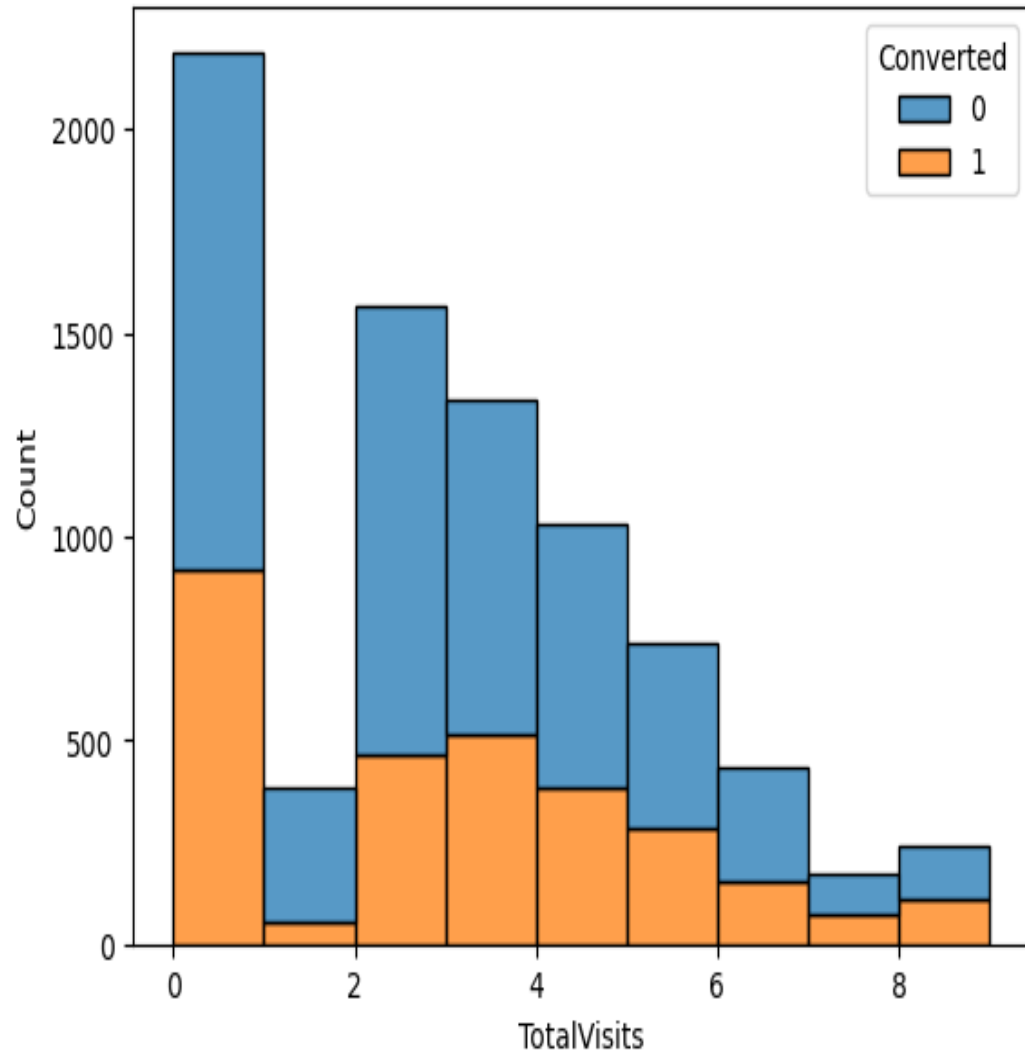- Making predictions on test set

# Exploratory Data Analysis



- In these graph we can see the higher conversation rate is in 'Lead Add Form'.
- And there is no conversation in 'other'.

# Lead Source



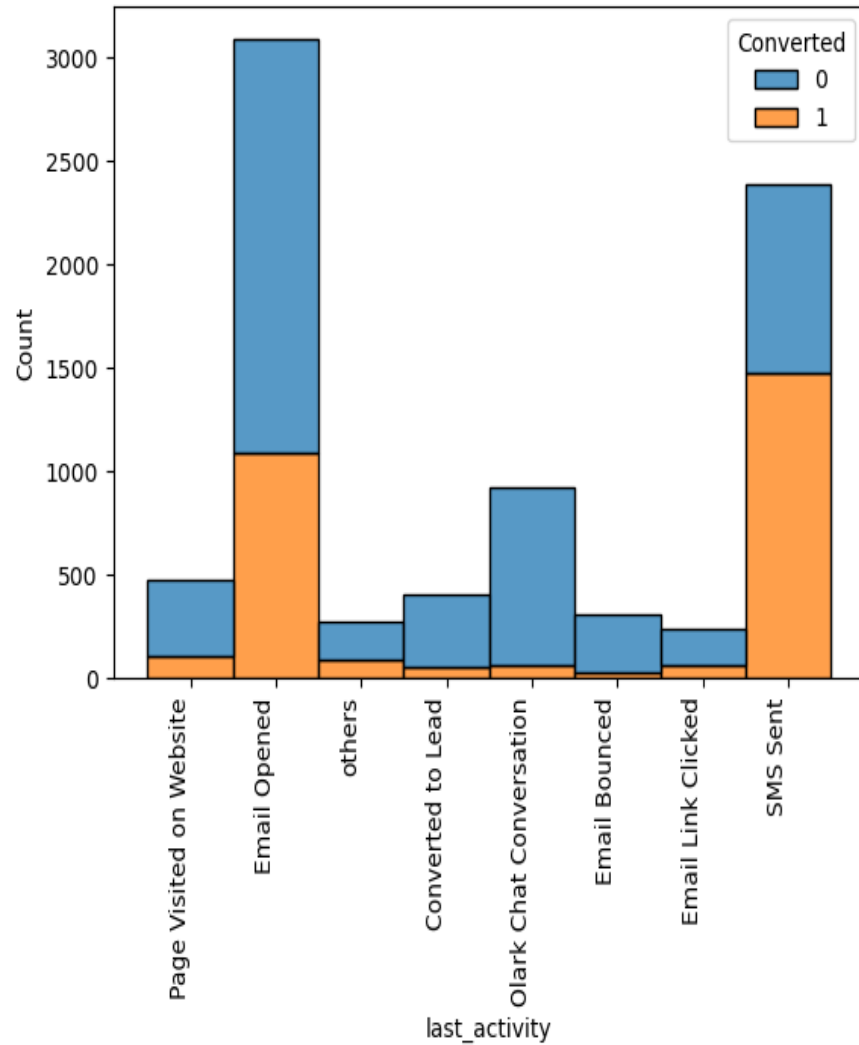- In these graph lead_source, we can observe that the higher rate of conversation is in 'Reference'.

# Total Visit



▶ In these graph, Total visit does not show a great impact on whether a lead was converted or not.
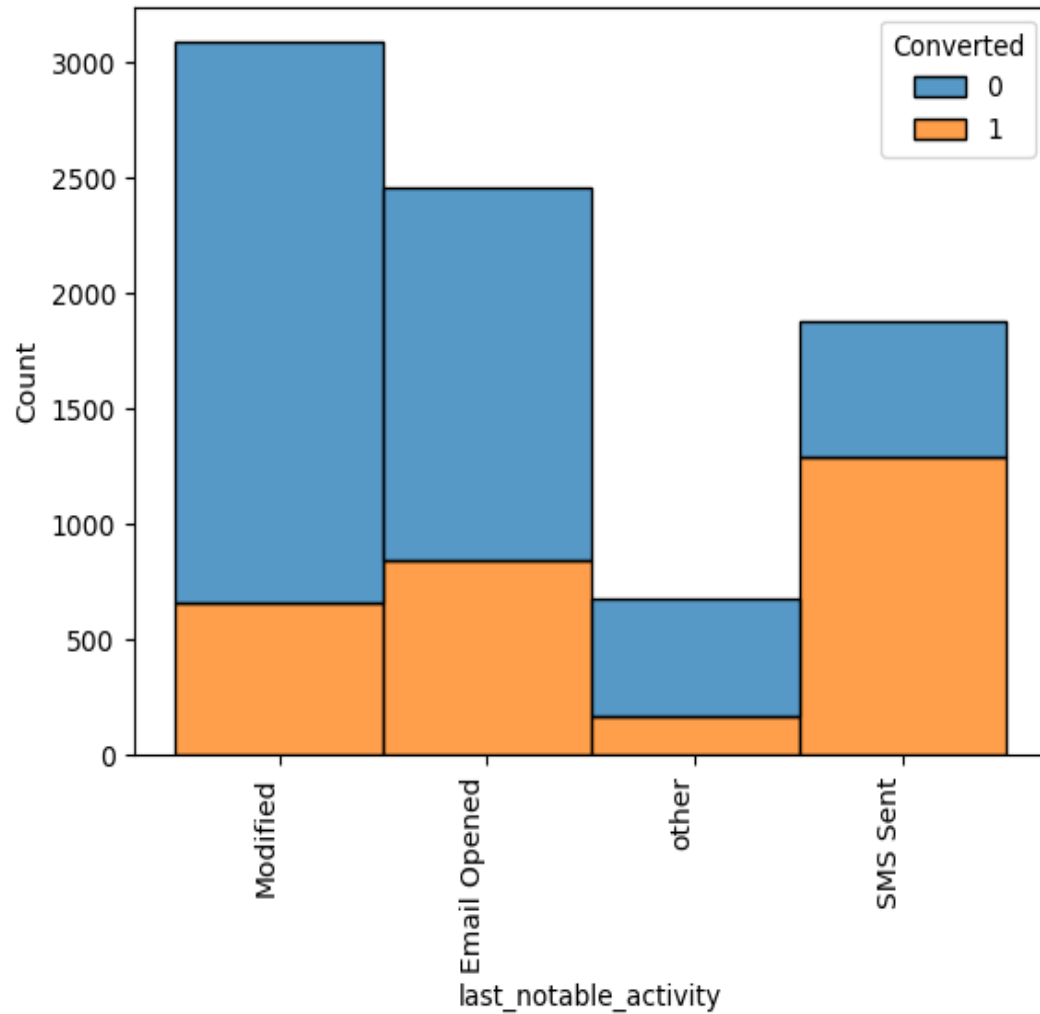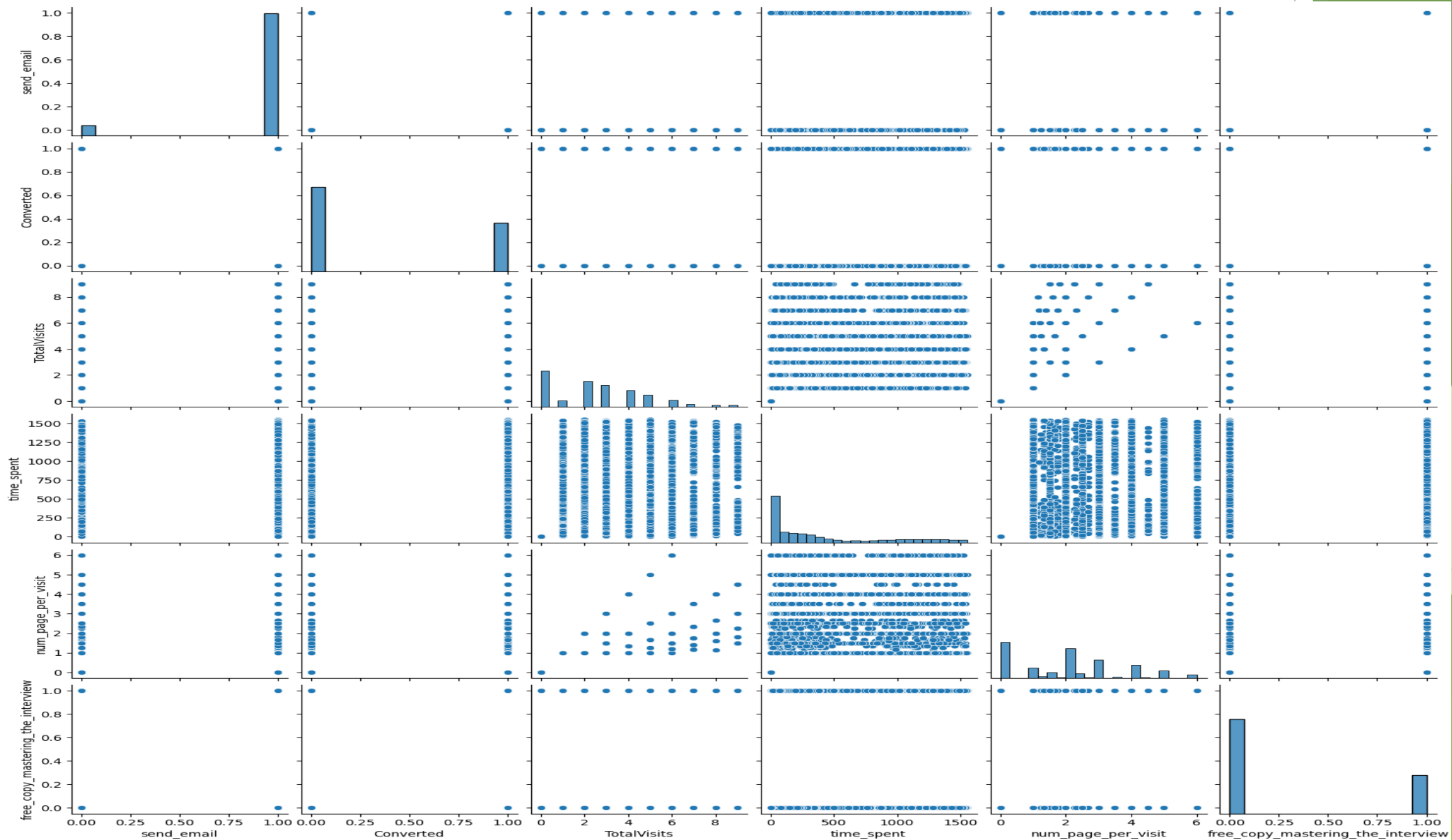
# Last Activity



- By these graph we understand that ,one whose last activity was SMS sent are converted more.

- Also Email opened has also almost 33% conversion rate.

# Last Notable Activity



▶ In these graph of last_notable_activity 'SMS Sent' seems to have higher impact on conversion rate.
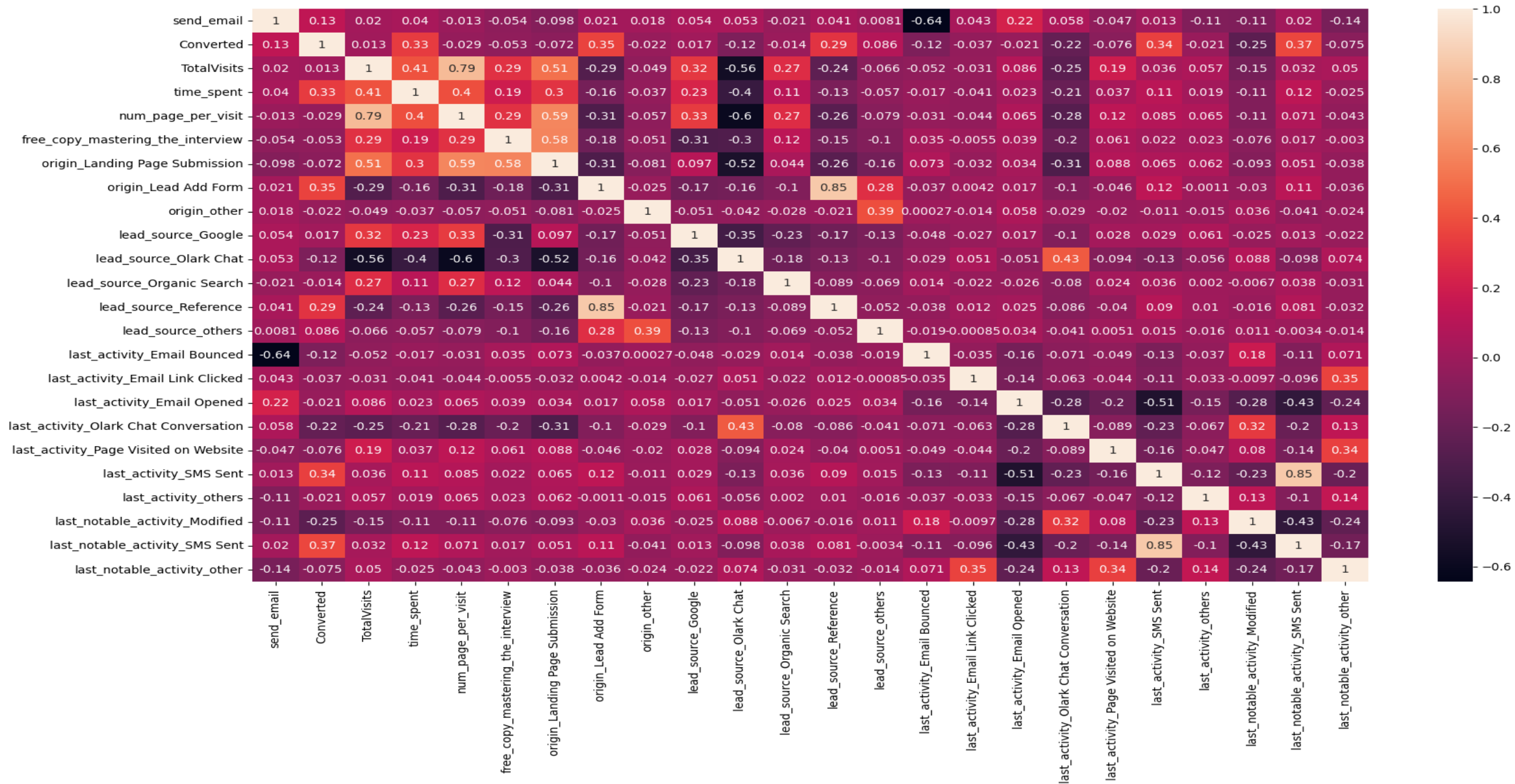
# Observation After EDA

# Observation After EDA

➤ Conversion rate is very less. i.e. almost 38%

➤ We had to drop a lot of features because either null_percentage > 10% or more than 95% values were same.

➤ Most of the lead origin is API or Landing Page Submission. But their conversion rate is very low. However Lead Add Form is having very high conversion rate

➤ Google, OLARK Chat and Direct traffic are major lead sources but having low conversion rates. Referral is having very high conversion rate.

➤ If user has asked not to send the email, then he is not going to convert for sure.

➤ Total visits, Time spent and num of pages per visit was having outliers. We have taken values up to the 98th percentile.

➤ Users who have spent more time on the website are more likely to be converted.

➤ Users who have opted for free copy of mastering the interview does not seem.
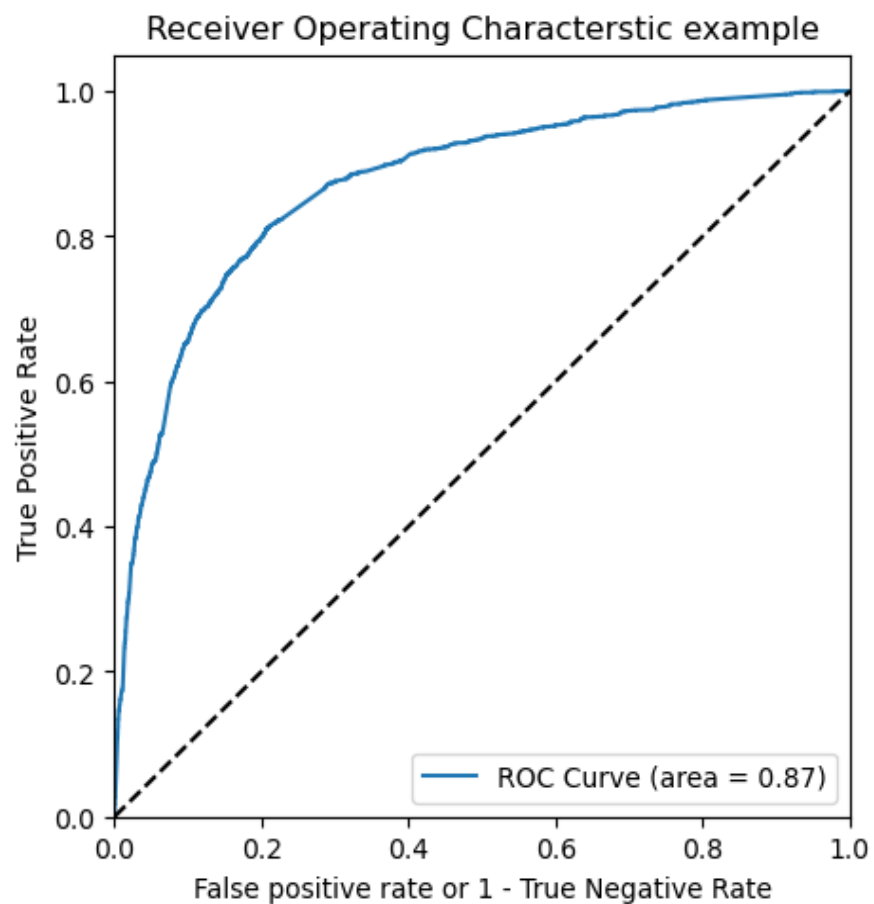
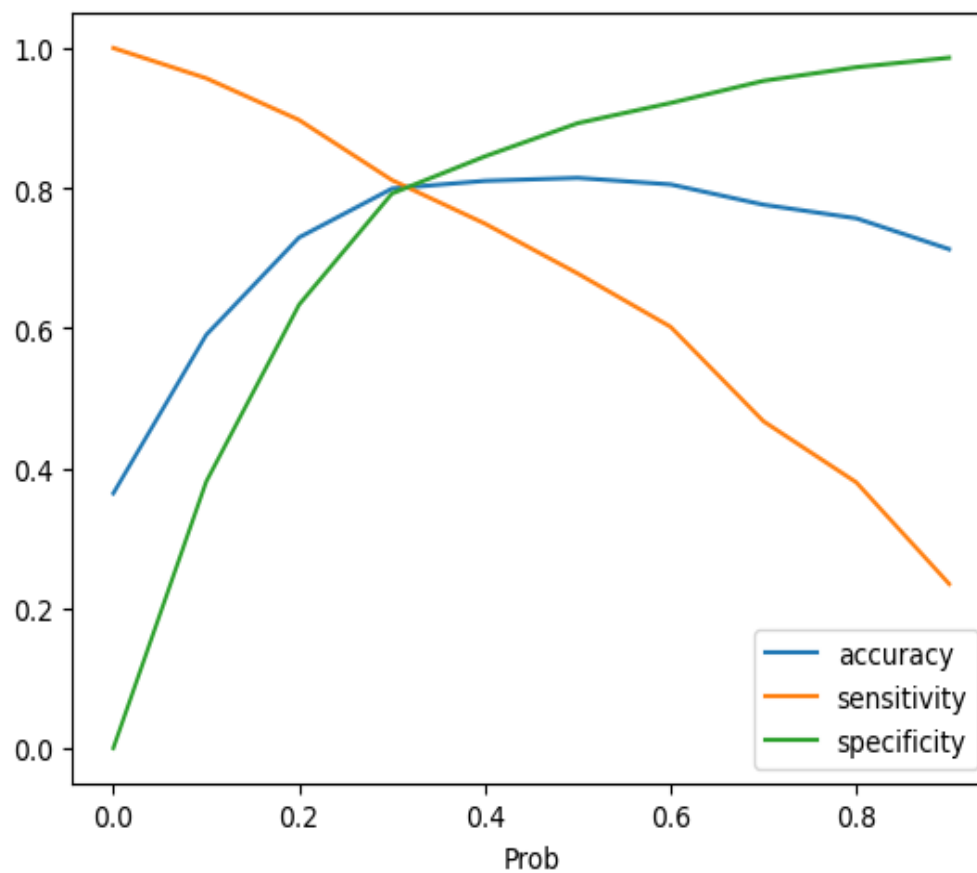# Correlations

# Observation from correlation

➢ Send_email and last_activity_email_bounced high correlation

➢ lead_source_olark_chat has high correlation with origin_landing_page_submission, num_pages_per_visit and total_visits

➢ num_page per visit has high correlation with num_pages_per_visit

➢ origin_landing_page_submission has high correlation with free_copy_mastering_the_interview and num_pages_per_visit

➢ Hope that these collinearity will be dealt while RFE and VIF Step. Just to keep eye.

# Model Evaluation

▶ Here model looks good as AUC is 0.87 for ROC

▶ As looking at the below plot we choose 0.3 as cutoff point

# Conclusion

➢ We created a model with finally 9 Features.

➢ We choose 0.3 as cutoff value.

▶ Model Evaluation on Train data :

1. accuracy = 0.79

2. sensitivity = 0.81

3. specificity = 0.79

▶ Model Evaluation on Test Data :

1. accuracy = 0.79

2. sensitivity = 0.80

3. specificity = 0.79

▶ We see max number of leads are generated by google/direct traffic.

▶ Most common last activity is email opened, highest rate = SMS Sent.

▶ Finally we merged predicted data frame with test data on index and assigned lead score by multiplying with 100

# Business Recommendation

- In order to increase probability of more lead conversion. we should focus on what actions on organization end result in more lead conversion. As per our model –

1) Lead source Olark chat has a positive coefficient of 1.4738. So, we should try to put more people on Olark chat to generate more lead source from there.

2) Last activity Olark chat conversation has negative coefficient of -1.3043. It means a lot of our leads drop off after the Olark chat conversation. So, we need to find out if our employees are not trained enough who are there on Olark chat and decrease the drop off from there itself.

3) Email bounced has a negative coefficient of -1.5312. So, it means we are not validating the email and entering the wrong email on the lead form. Let it be website or past referrals. On website we should check if email validation is there before form submission.

# Business Recommendation

▶ When business wants to go more aggressive on lead conversion , they should increase the sensitivity and recall of the model. In simpler words, No of actual leads predicted to total no of actual leads and our model does not predict a hot lead in to cold lead should be more.

▶ Also when company wants to minimize the rate of useless phone calls,  we should focus more on specificity and precision(probability that a predicted lead is actual lead). We can compromise on sensitivity of the model. In our model at cutoff = 0.6 sensitivity = 0.60 specificity = 0.92 precision = 0.81.