



LEAD SCORE CASE STUDY

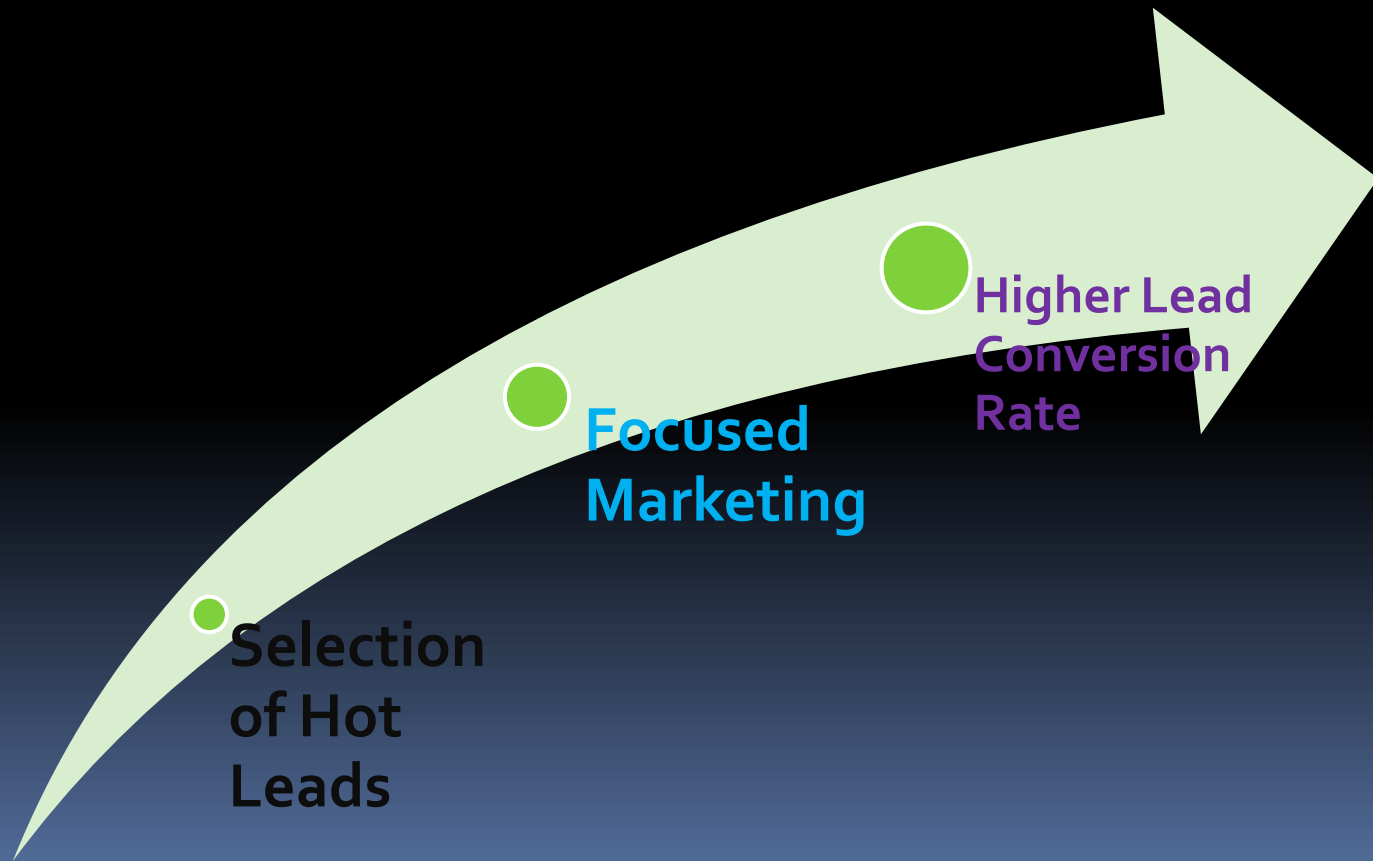
Focused Business approach using logistic regression technique

Vaishnavi Patil



Business objective

To help X Education select most promising leads (Hot Leads), i.e the leads that are most likely to convert into paying customers.





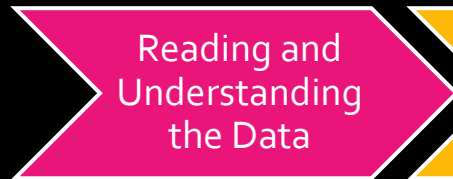
METHODOLOGY

To build a Logistic Regression model that assigns lead scores to all leads such that the customers with higher lead score have a higher conversion chance and vice versa.

Target Lead Conversion Rate = 80%



- Importing and Observing the past data provided by the Company



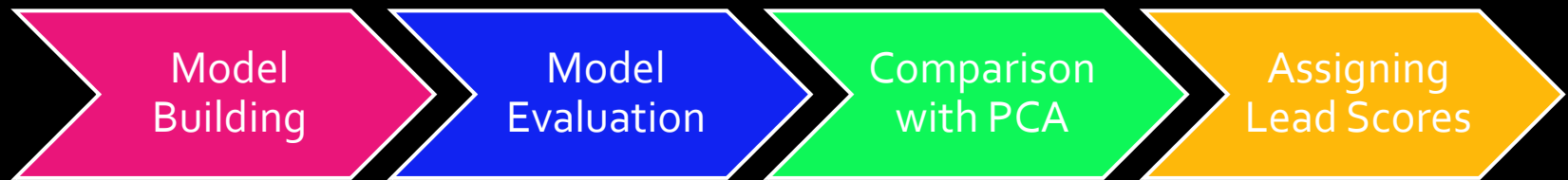
- Univariate and Bivariate analysis



- Missing Value imputation
- Removing duplicated data and other redundancies



- Outlier treatment
- Dropping unnecessary columns
- Dummy variable creation
- Feature standardization



- Feature selection using RFE
- Manual feature elimination based on p-values and VLFs

- Evaluating model based on various evaluation metrics
- Finding the optimal probability threshold

- Building another model using PCA
- Comparing the two models

- Finalizing the First model
- Using Predicted probabilities to calculate Lead Scores: $\text{Lead Score} = \text{Probability} * 100$

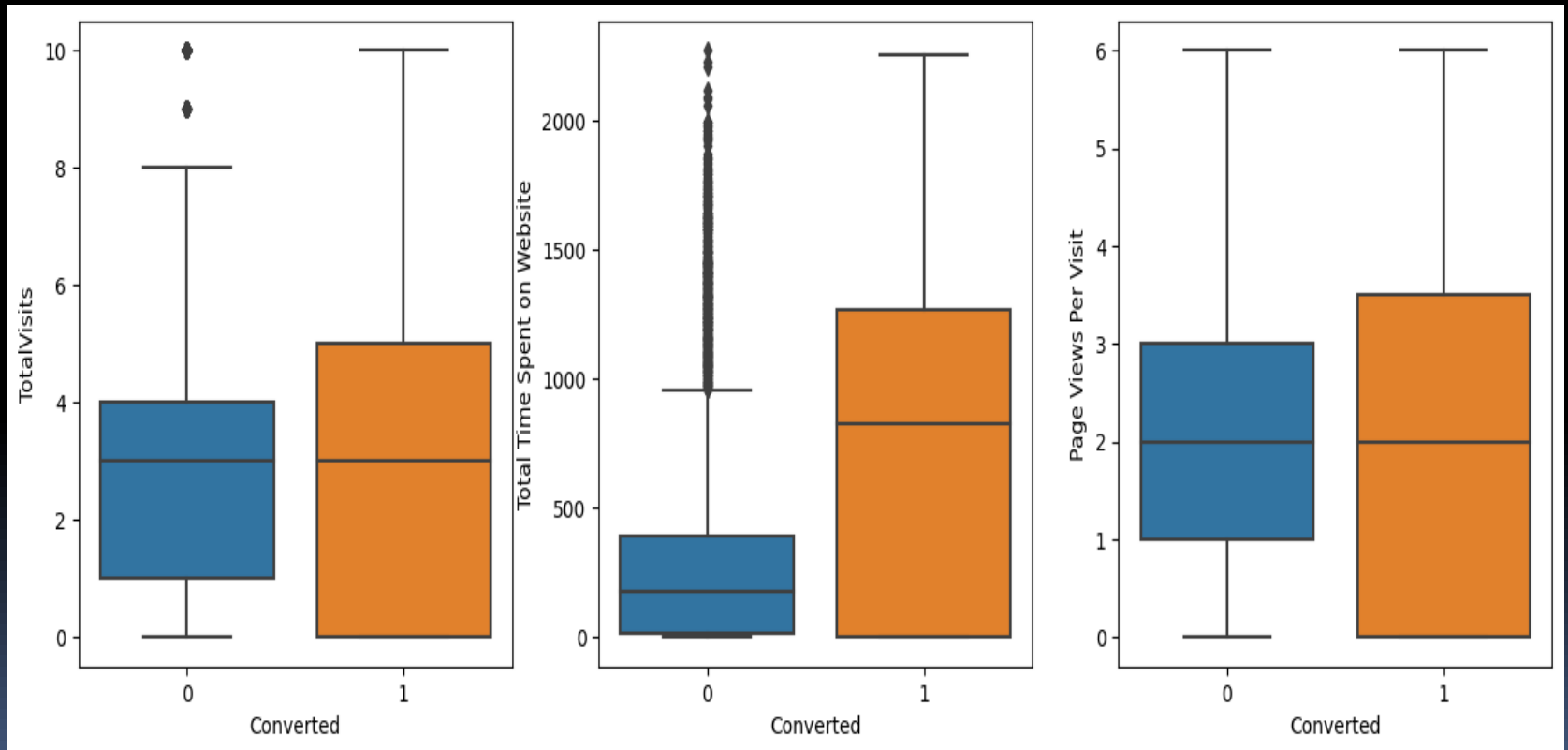


DATA VISUALIZATION

- To identify important features
 - To get Insights
- 

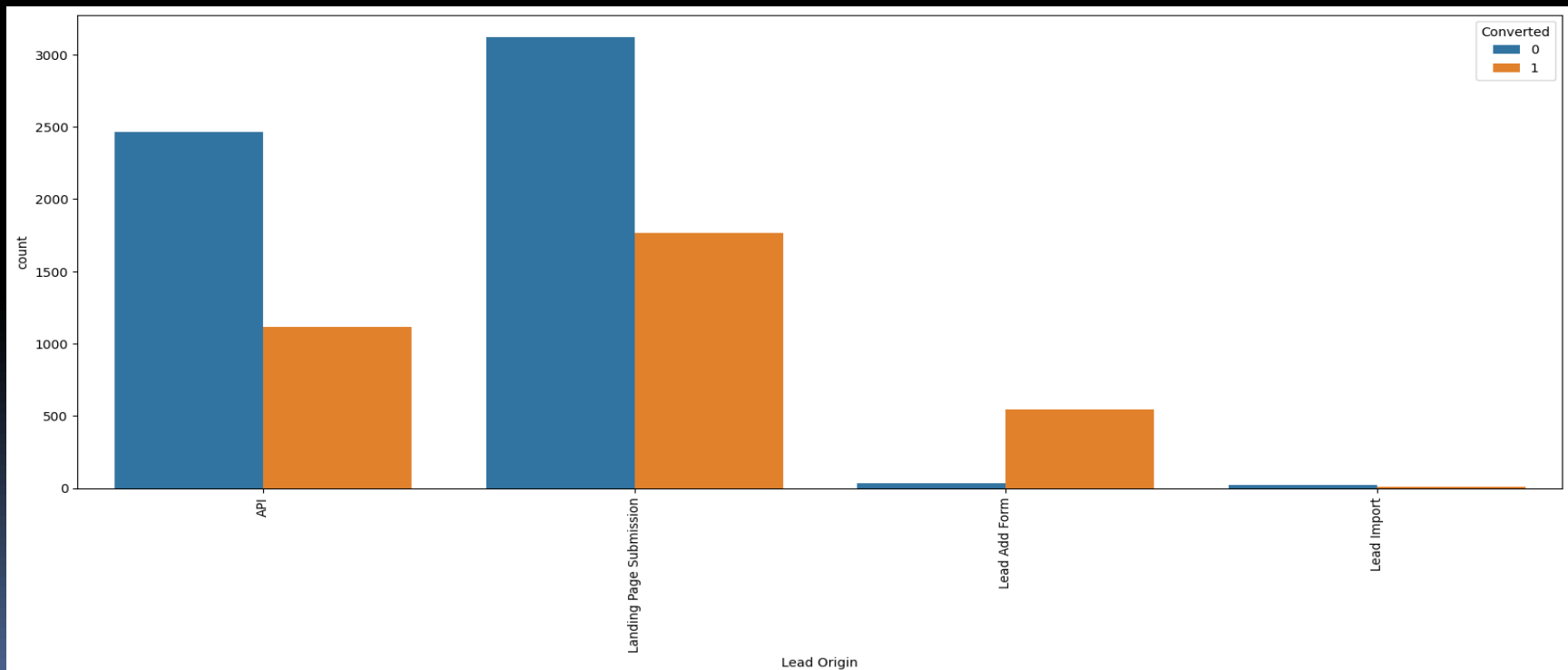
Numerical Variables

People spending more time on website are more likely to get converted.



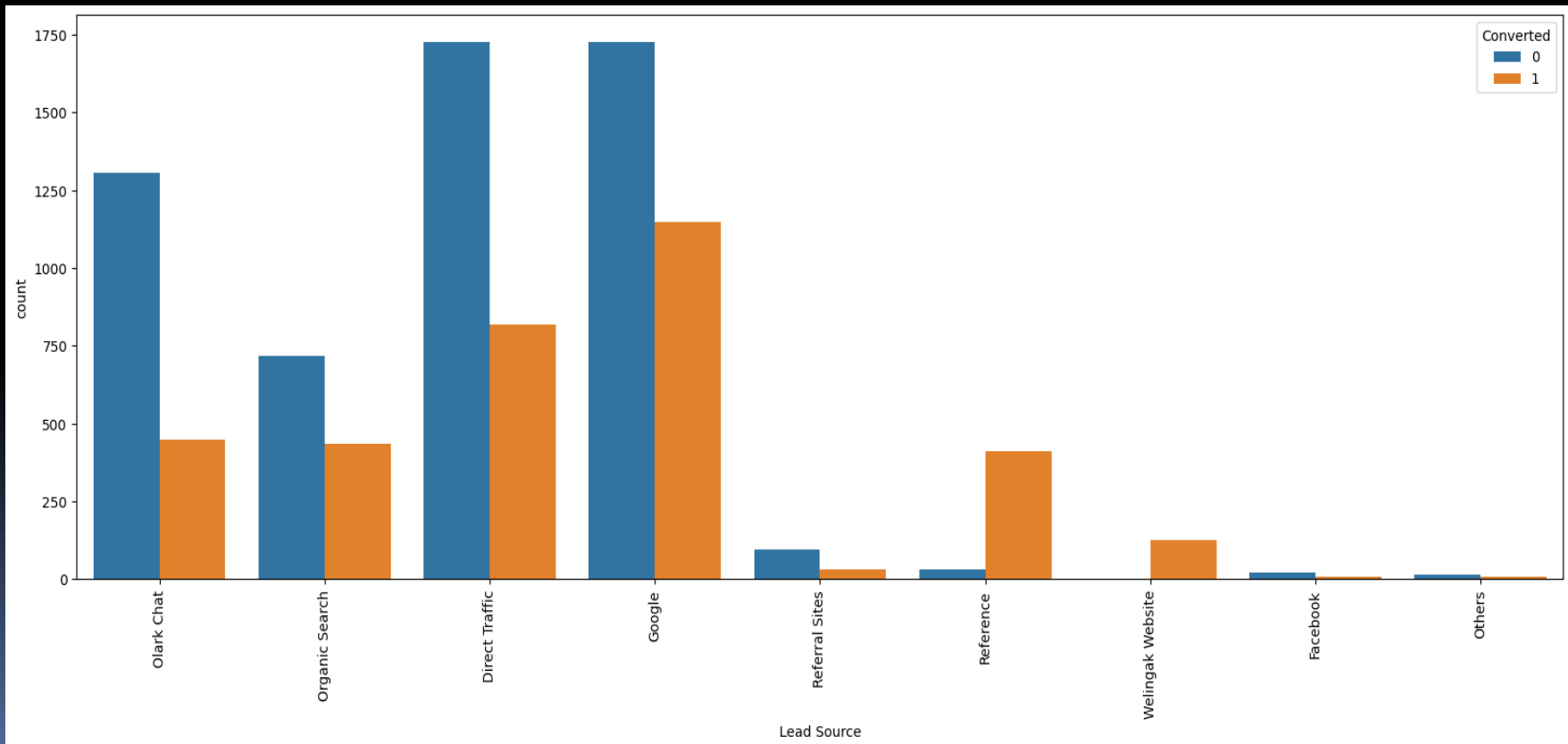
Lead Origin

- **'API' and 'Landing page Submission'** generate the most leads but have less conversion rates, whereas **'Lead Add Form'** generates less leads but conversion rate is great.
- **Try to increase conversion rate for 'API' and 'Landing Page Submission' , and increase leads generation using 'Lead Add Form'.**



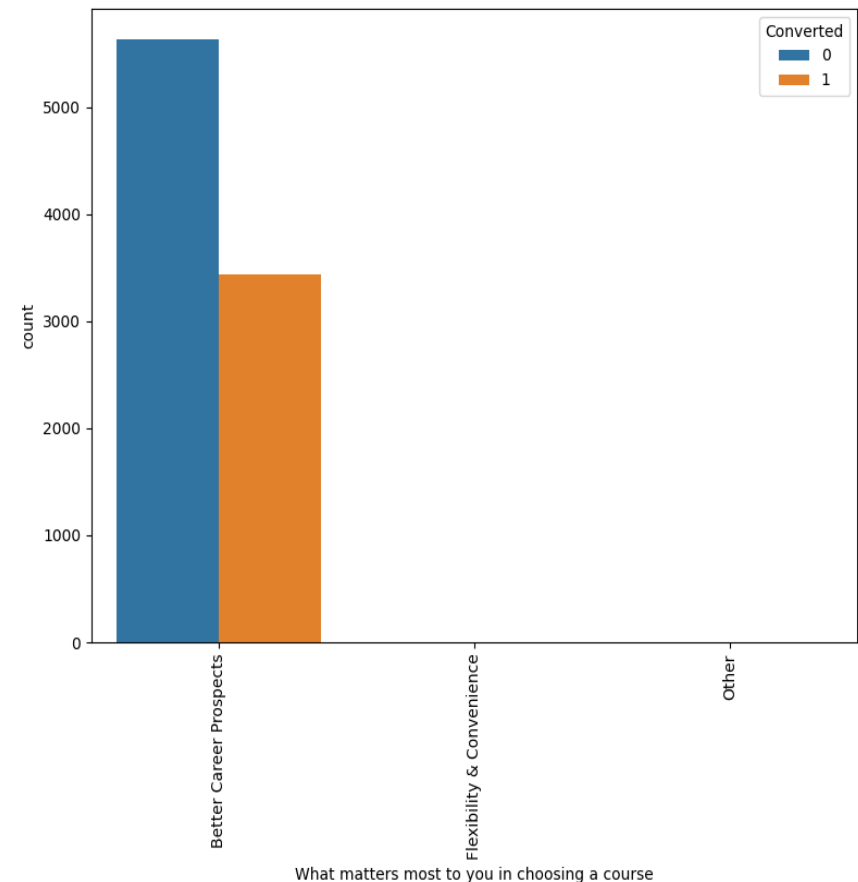
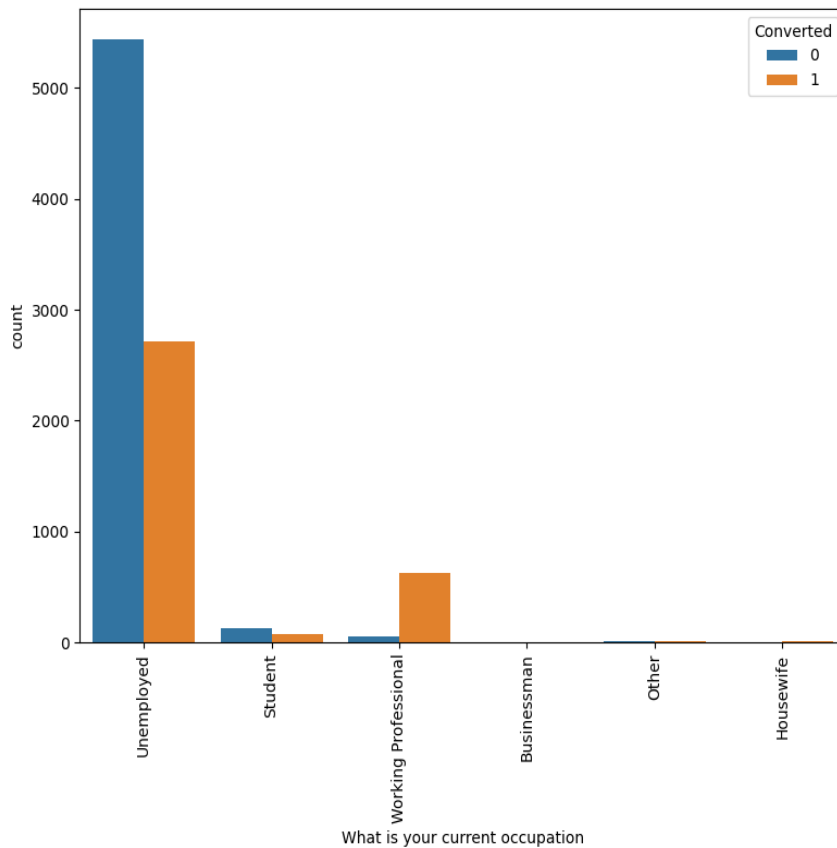
Lead Source

- Very high conversion rates for lead sources '**Reference**' and '**Welingak Website**'.
- Most leads are generated through '**Direct Traffic**' and '**Google**'.



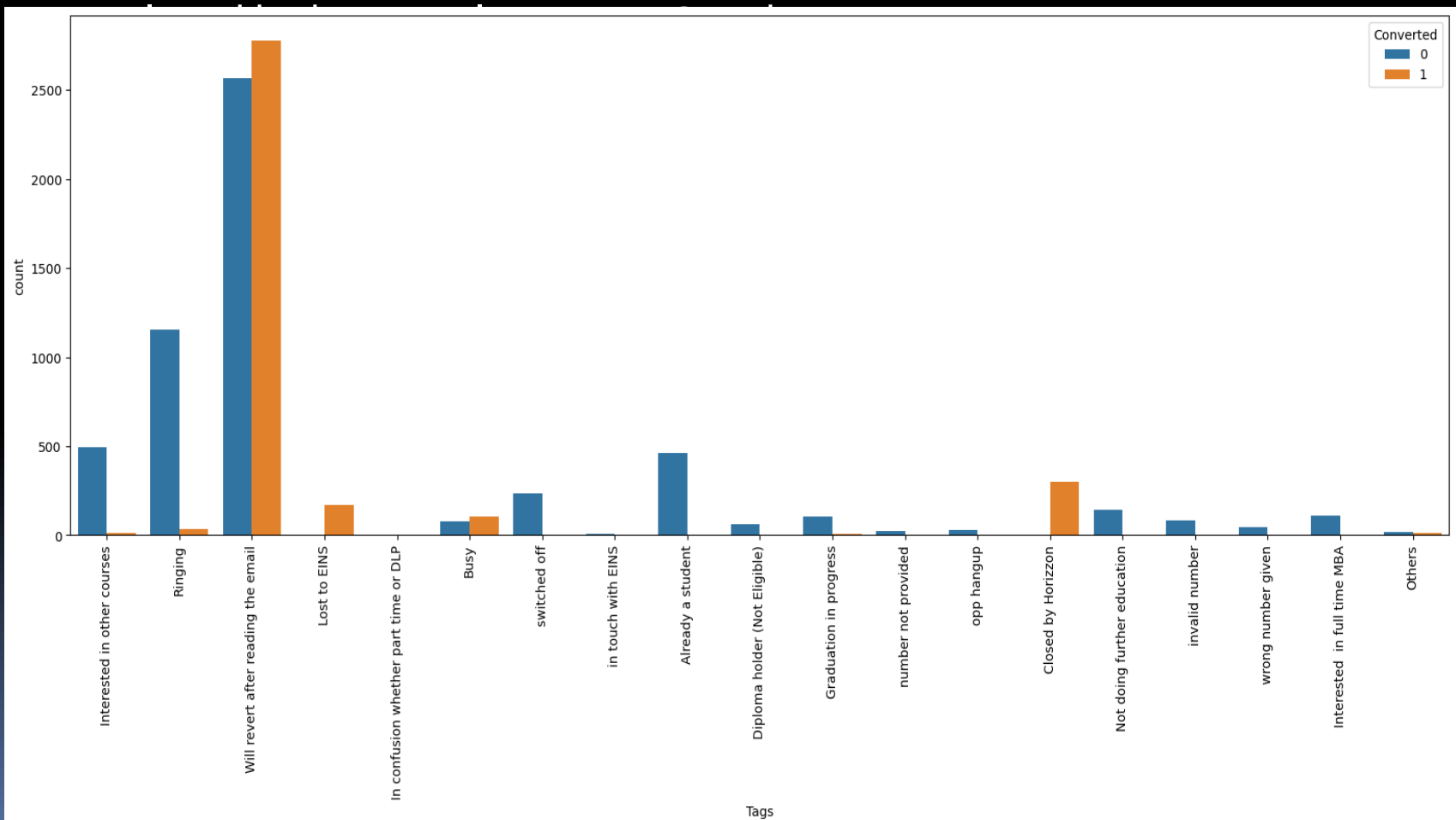
Current Occupation

- **Working professionals** are most likely to get converted



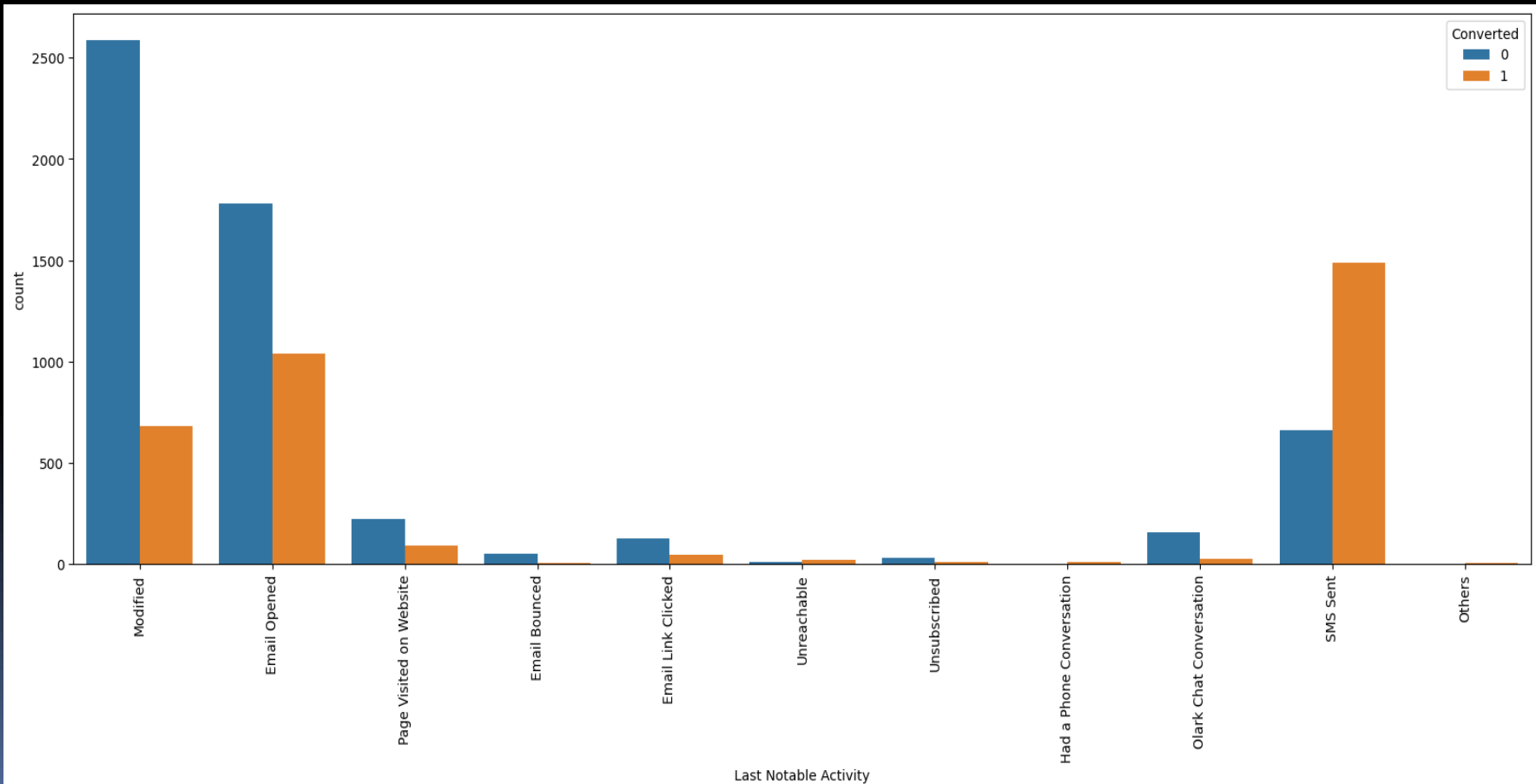
Tags

- High conversion rates for tags will revert after reading the email,



Last Notable Activity

- Highest conversion rate is for the last notable activity 'SMS Sent'.

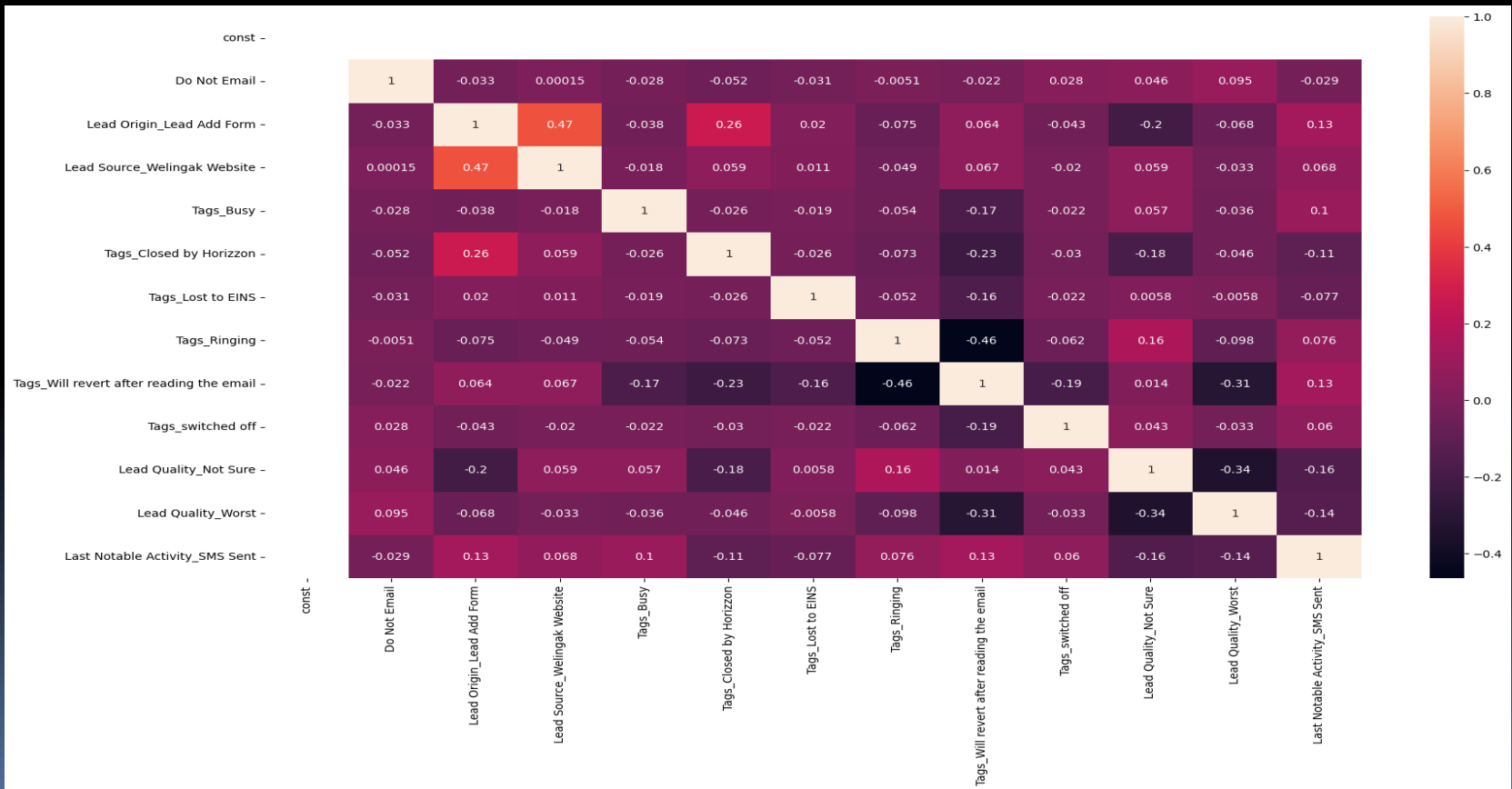




MODEL EVALUATION

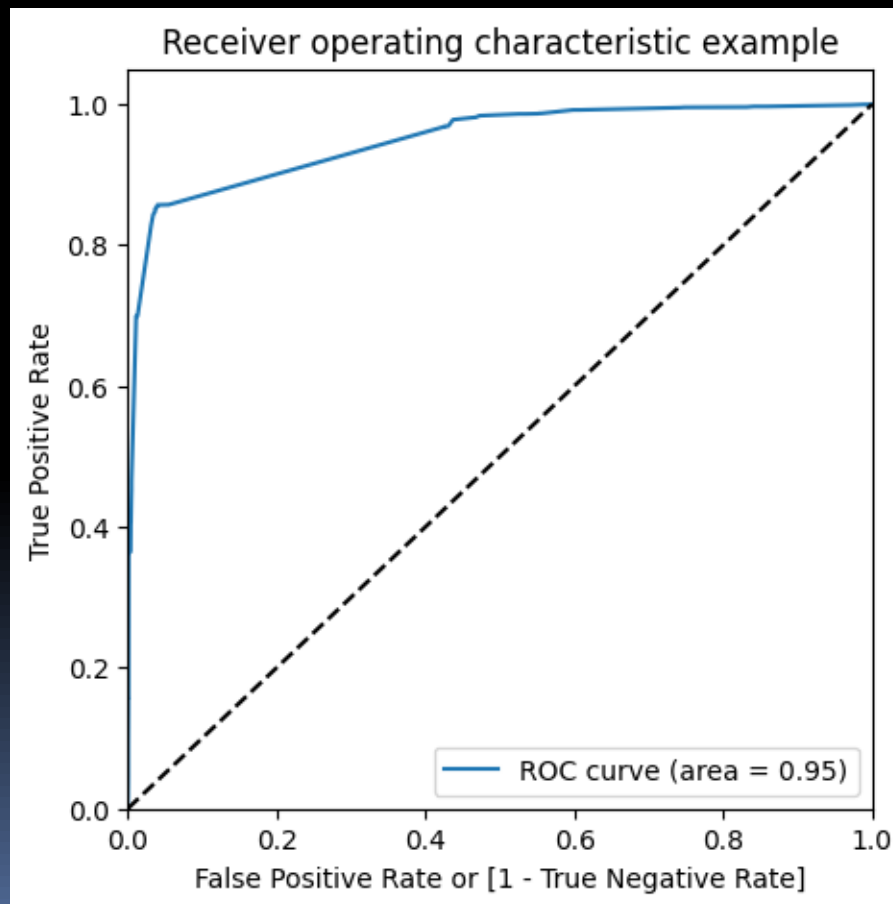
Heatmap

- Correlations between features in the final model are negligible.



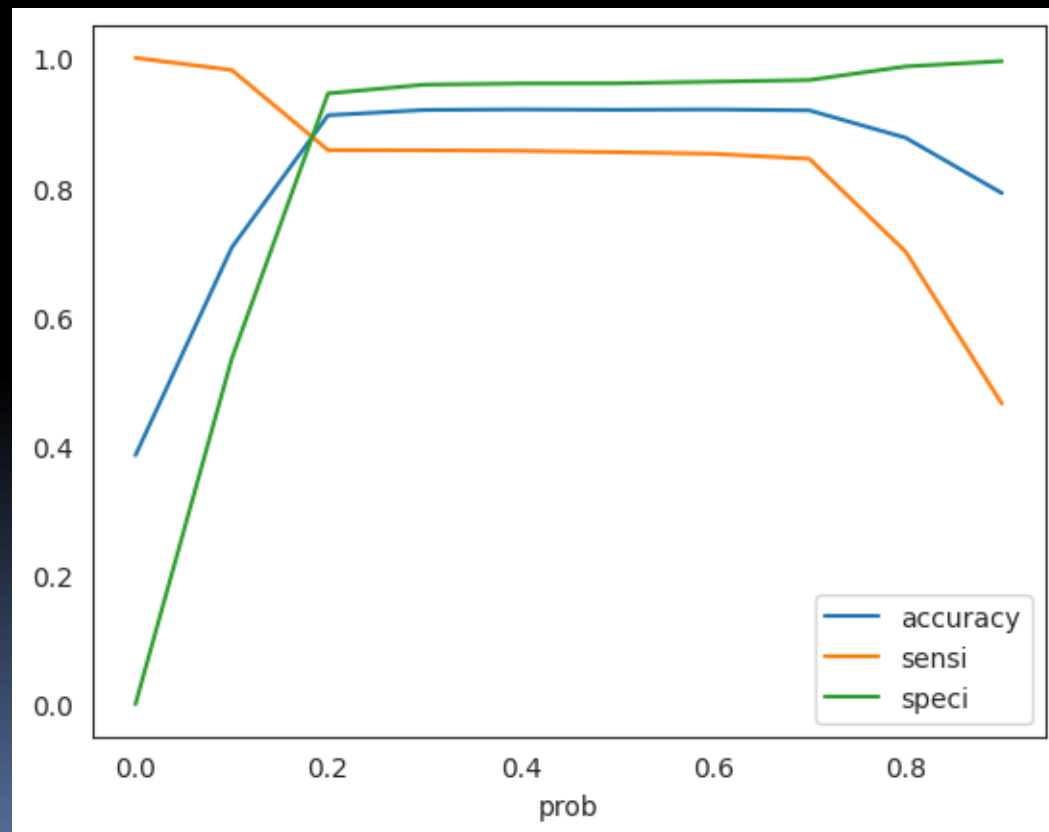
ROC curve

- Area under curve = 0.95



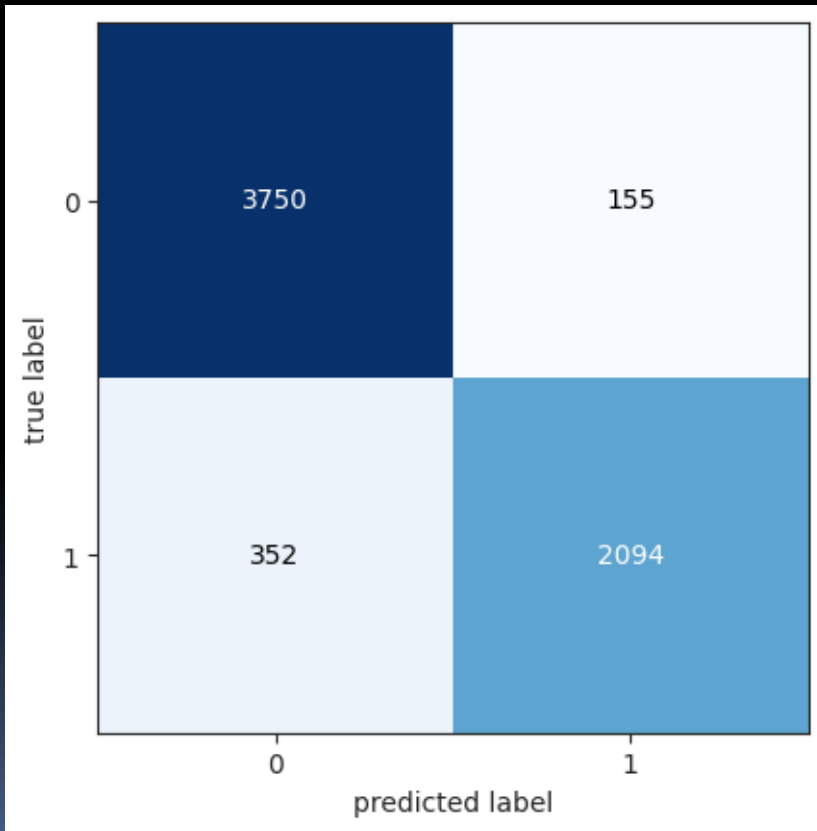
Finding Optimal Threshold

- Graph showing changes in sensitivity, Specificity and Accuracy with changes in the probability threshold values
- Optimal cutoff = 0.20

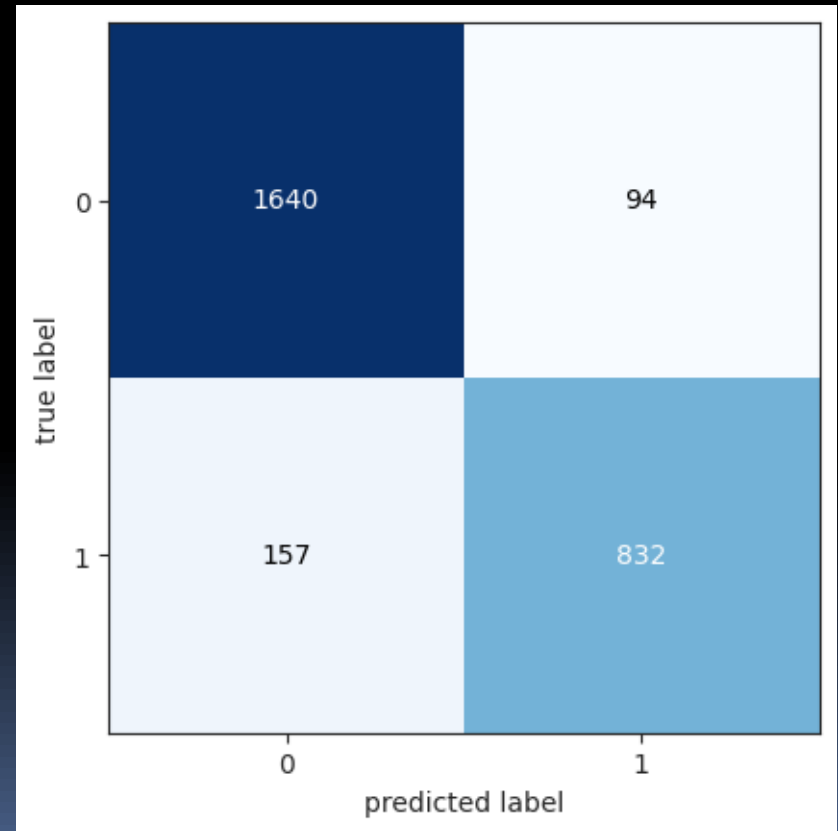


Confusion Matrix

- For train set



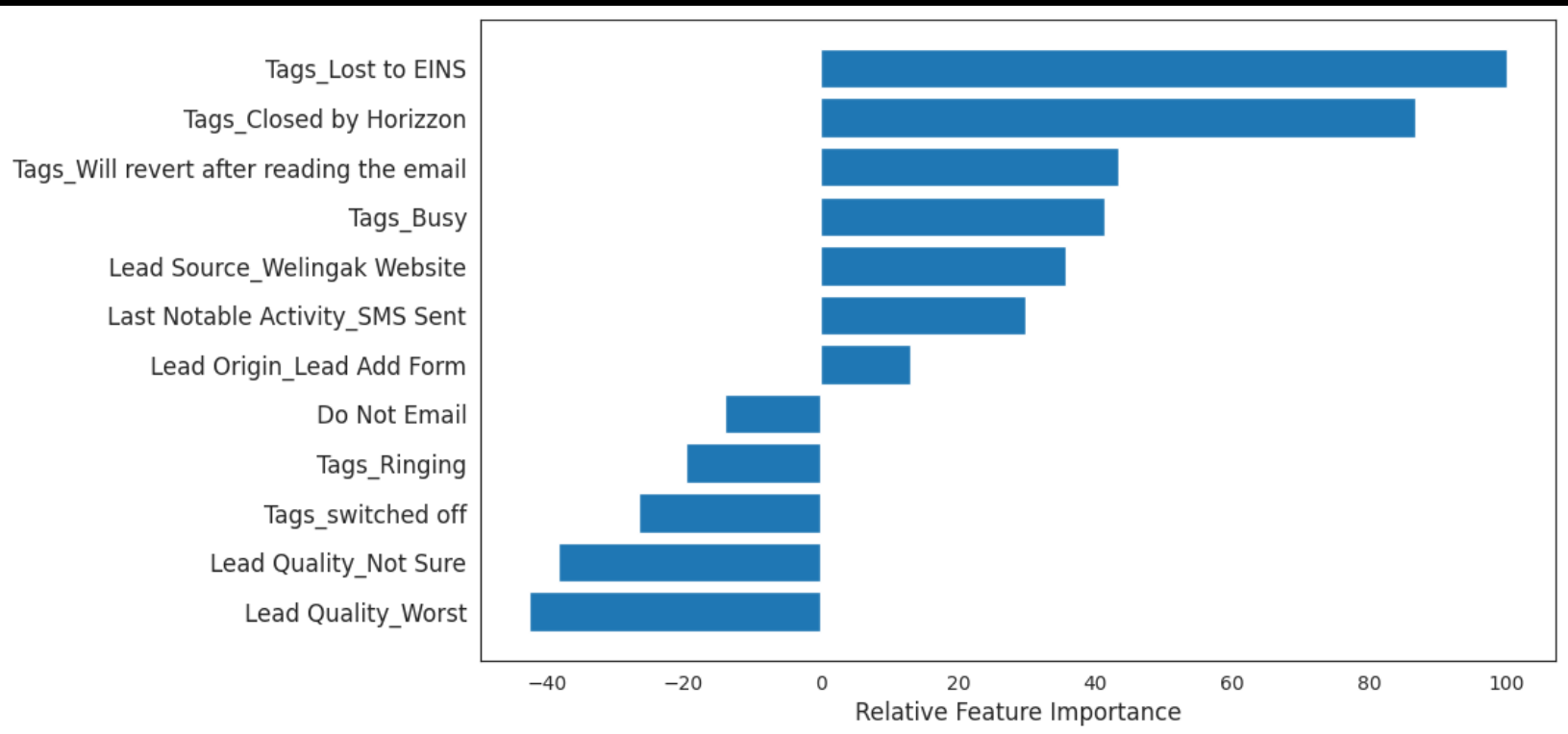
- For test set



Final Results

- Sensitivity: 0.8412537917087968
- Specificity: 0.9457900807381776
- False positive rate - predicting the lead conversion when the lead does not convert: 0.05420991926182238
- Positive predictive value: 0.8984881209503239
- Negative predictive value: 0.9126321647189761

Relative Importance Of Features





INFEEENCES



Feature Importance

- Three variables which contribute most towards the probability of a lead conversion in decreasing order of impact are
 1. Tags_Lost to EINS
 2. Tags_Closed by Horizzon
 3. Tags_will revert after reading the email
- These are dummy features created from the categorical variable Tags.
- All three contribute positively towards the probability of a lead conversion.
- These results indicate that the company should focus more on the leads with these three tags

- Situation 1: Company has interns for 2 months. They wish to make lead conversion more aggressive. They want almost all of the potential leads to be converted and hence, want to make phone calls to as much of such people as possible.
- Solution:
 - $\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$
 - Sensitivity can be defined as the number of actual conversions predicted correctly out of total number of actual conversions. As we saw earlier, sensitivity decreases as the threshold increases.
 - High sensitivity implies that our model will correctly predict almost all leads who are likely to convert. At the same time, it may overestimate and misclassify some of the non-conversions as conversions.
 - As the company has extra man-power for two months and wants to make the lead conversion more aggressive, it is a good strategy to go for high sensitivity. To achieve high sensitivity, we need to choose a low threshold value

- Situation 2: At times, the company reaches its target for a quarter before the deadline. It wants the sales team to focus on some new work. So during this time, the company's aim is to not make phone calls unless it's extremely necessary.
- Solution:
 - $\text{Specificity} = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$
 - Specificity can be defined as the number of actual non-conversions predicted correctly out of total number of actual non-conversions. It increases as the threshold increases.
 - High specificity implies that our model will correctly predict almost all leads who are not likely to convert. At the same time, it may misclassify some of the conversions as non-conversions.
 - As the company has already reached its target for a quarter and doesn't want to
 - make unnecessary phone calls, it is a good strategy to go for high specificity.
 - It will ensure that the phone calls are only made to customers who have a very high probability of conversion. To achieve high specificity, we need to choose a high threshold value.



Recommendations

- By referring to the data visualizations, focus on
 - Increasing the conversion rates for the categories generating more leads and
 - Generating more leads for categories having high conversion rates.
- Pay attention to the relative importance of the features in the model and their positive or negative impact on the probability of conversion.
- Based on varying business needs, modify the probability threshold value for identifying potential leads.