

Evaluating Email Outreach Performance with Advanced Machine Learning Techniques

Riya, Sanika,
Vaishnavi, Aswad

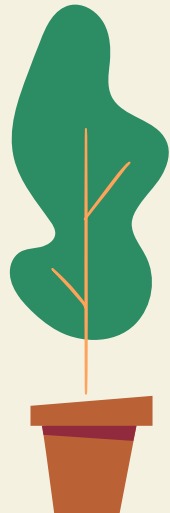


Table of Contents

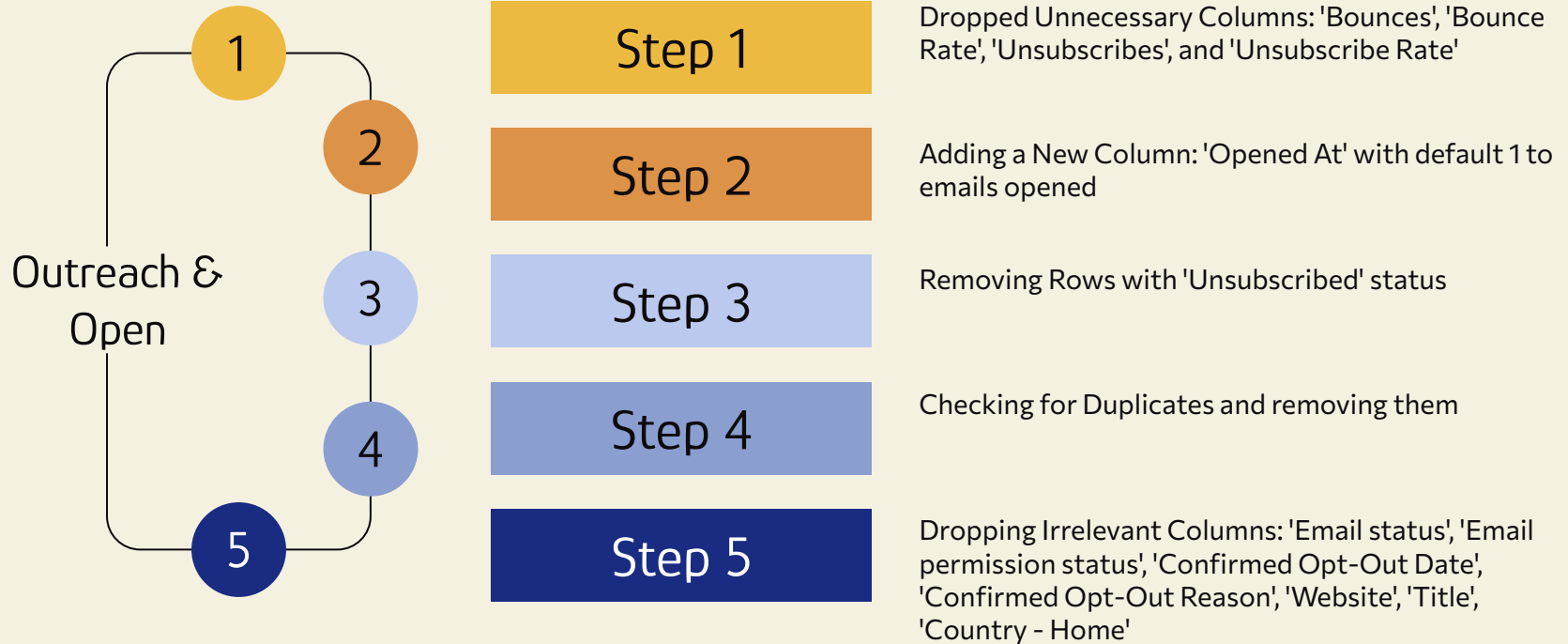
<u>Project Description</u>	Background
<u>Data Preprocessing</u>	Data cleaning and merging
<u>Data Analysis</u>	Model Building
<u>Results</u>	Insights from the analysis
<u>Suggestions / Recommendations</u>	Suggestions for improvement to get better results
<u>Summary</u>	Project reflection
<u>Reference</u>	Resources used

Project Description

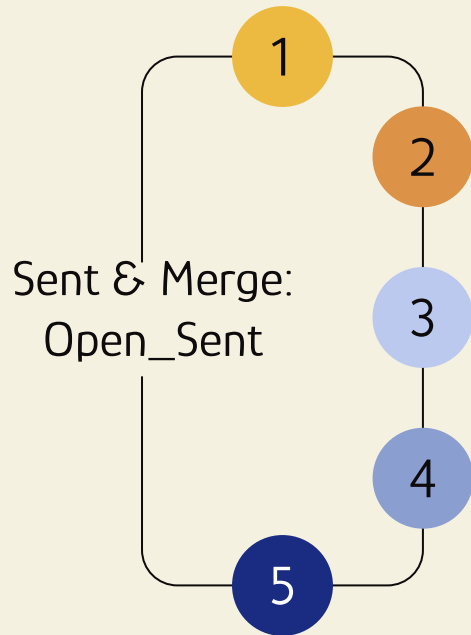
- Background:
 - This project analyzes email outreach data to improve campaign effectiveness by understanding effect of the time email sent, revenue of the company, the date, company type and number of employees, and location of the company.
- Problem Statement:
 - Identify factors influencing the likelihood of recipients engaging with the email (opens) and increase the outreach efficiency.
 - Dependent Variable (DV): Email open status (opened or not).
 - Independent Variables (IVs): City - Home, Date Sent, Time Sent, State, Company type, Revenue, Headcount.



Data Preprocessing



Data Preprocessing



Step 1

Removed Unsubscribed Records, Remove Duplicates: 'Campaign Name', 'Company', 'City - Home', and 'Created At'

Step 2

Dropping Irrelevant Columns just like Open dataset

Step 3

Left join Sent and Open datasets using 'Campaign Name' and 'Company' columns, Dropping Duplicate Columns with suffixes _x for 'sent' and _y for 'open'

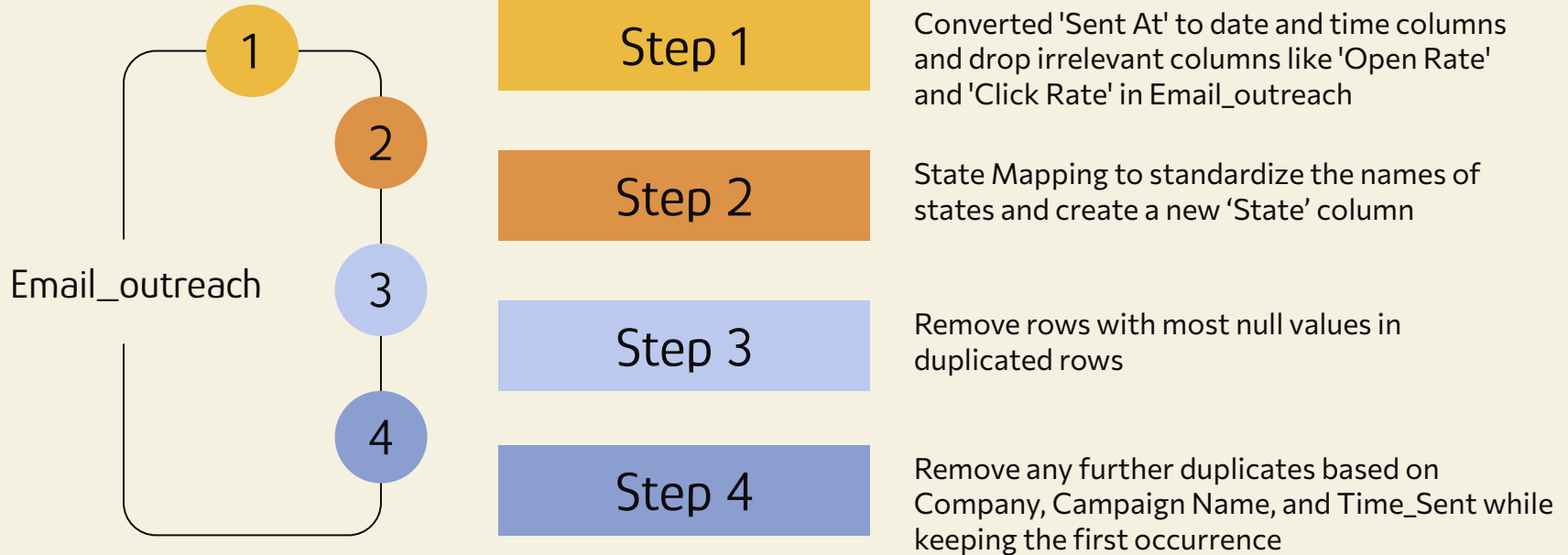
Step 4

Renaming columns like 'City - Home_x': 'City - Home', 'State/Province - Home_x': 'State/Province - Home', 'Created At_x': 'Created At', 'Updated At_x': 'Updated At'

Step 5

Filled missing values in 'Opened At' with 0 and ensure there are no duplicates in Open_Sent which is merged with Outreach summary dataset using 'Campaign Name'

Data Preprocessing



Data Preprocessing

1

Company names from various sources were normalized by converting them to lowercase and removing special characters to ensure consistency.

2

Missing city, state, revenue, headcount, and other key columns were populated using the mode into new columns.

3

Duplicate rows based on the Company column were removed and unnecessary columns were removed.

Master
dataset

1

2

3

1

2

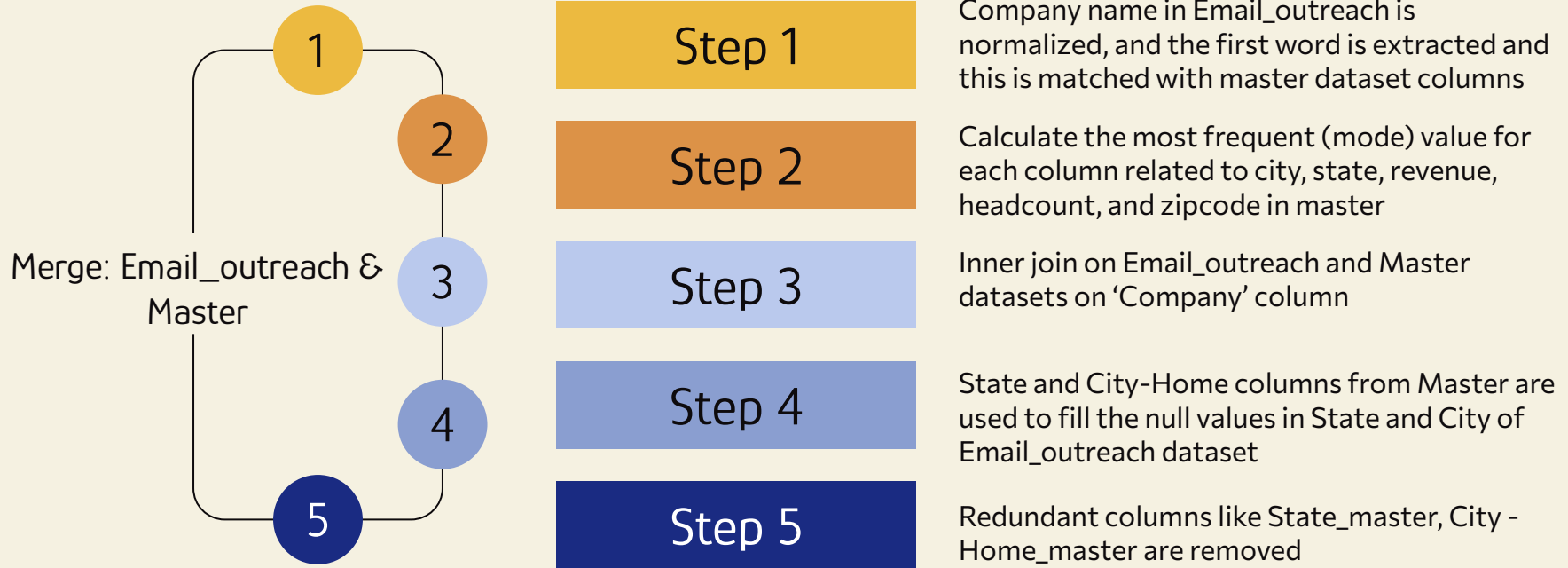
1

2

After normalization, stemming was applied to company names using the Porter Stemmer to reduce words to their root forms (e.g., "advanced" to "advanc").

The TF-IDF vectorizer was used to transform the stemmed company names into numerical vectors. Cosine similarity scores were calculated between Email_outreach and master datasets to identify the most likely matches based on the similarity of their company names.

Data Preprocessing



Final dataset: Merged1

01

Campaign
Name

02

Company

03

City - Home

04

Date Sent

05

State

06

Revenue

07

Opened

Binary Independent
Variable

08

Mobile
open rate

09

Desktop
open rate

10

Time_Sent

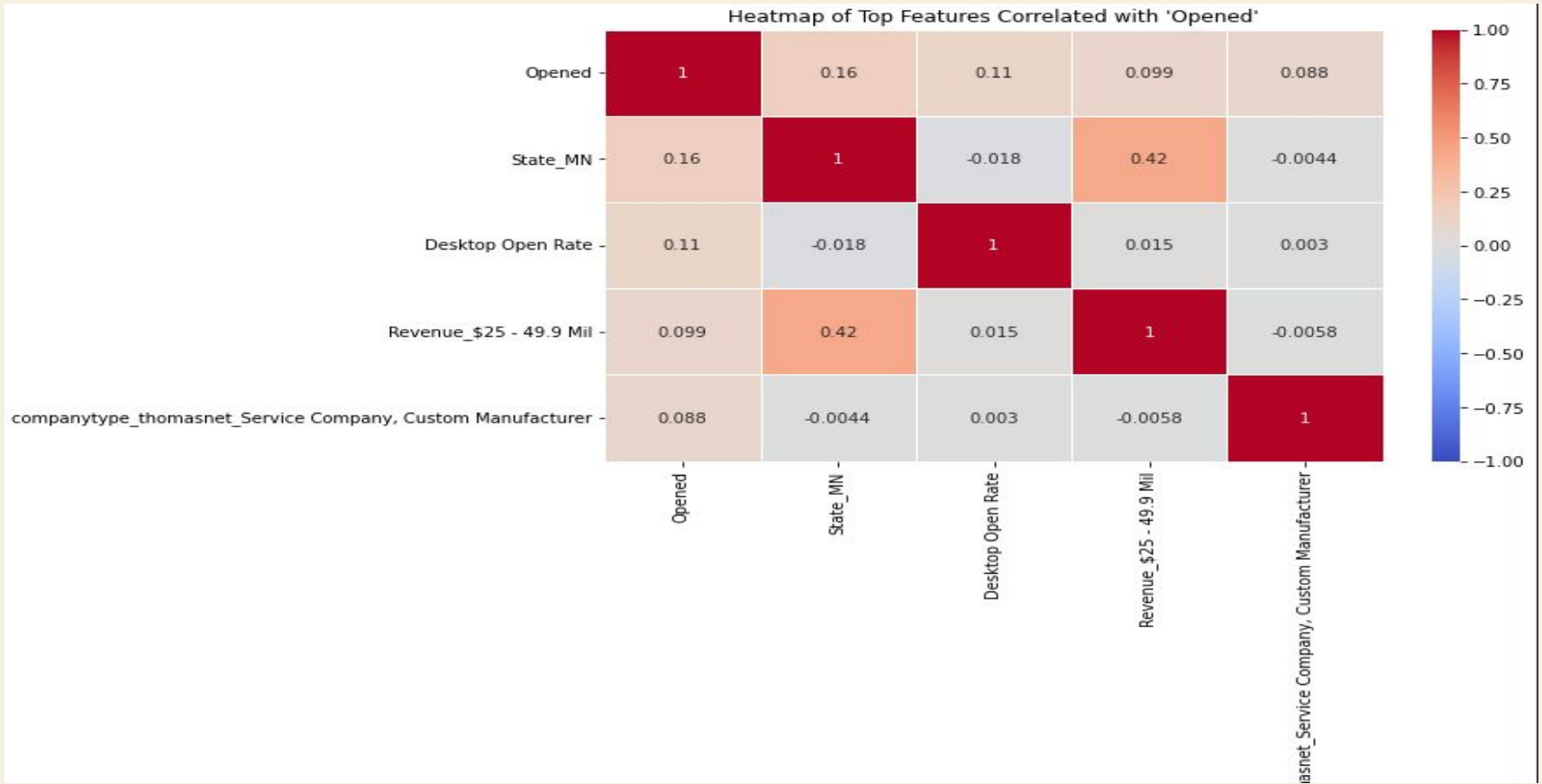
11

Company_type

12

Headcount

Correlation Matrix



Data Analysis : Training the model

Random Forest Model: It is a learning method used for both classification and regression tasks. It operates by constructing a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees.

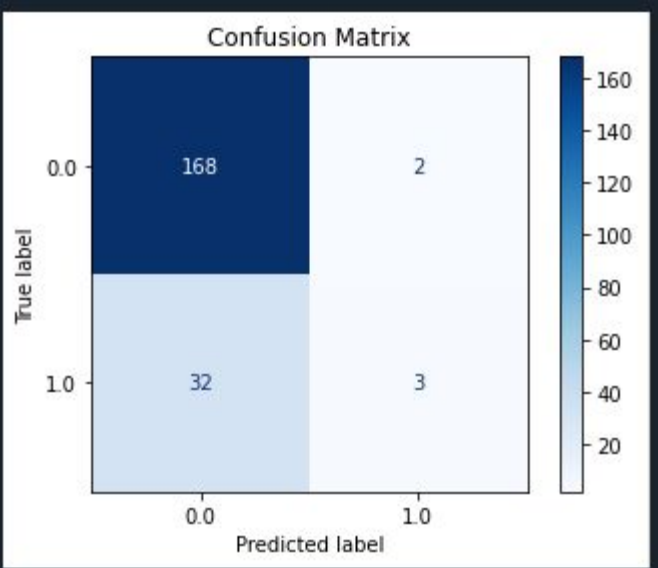
Our Findings:

Accuracy: 0.8341463414634146

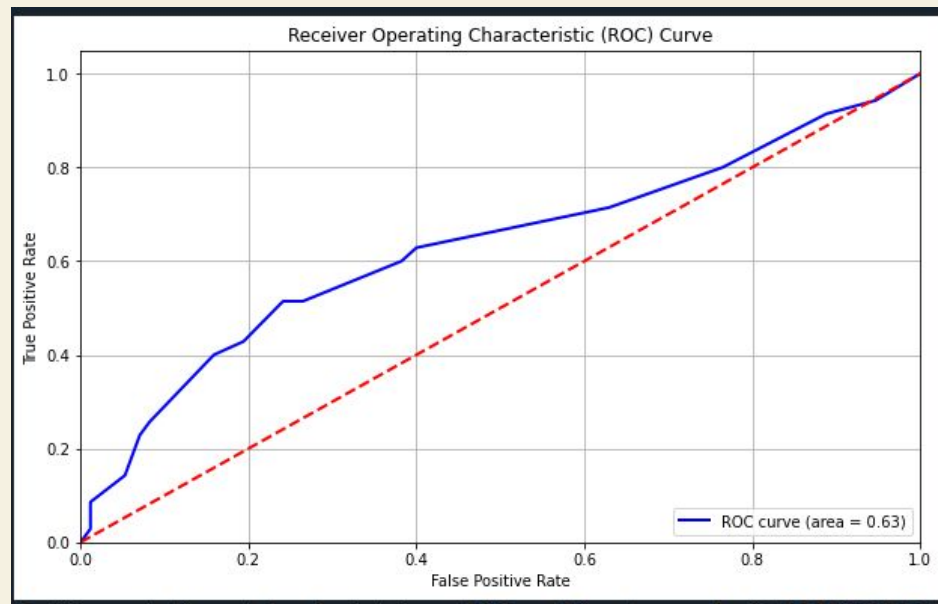
ROC-AUC Score: 0.6298319327731092

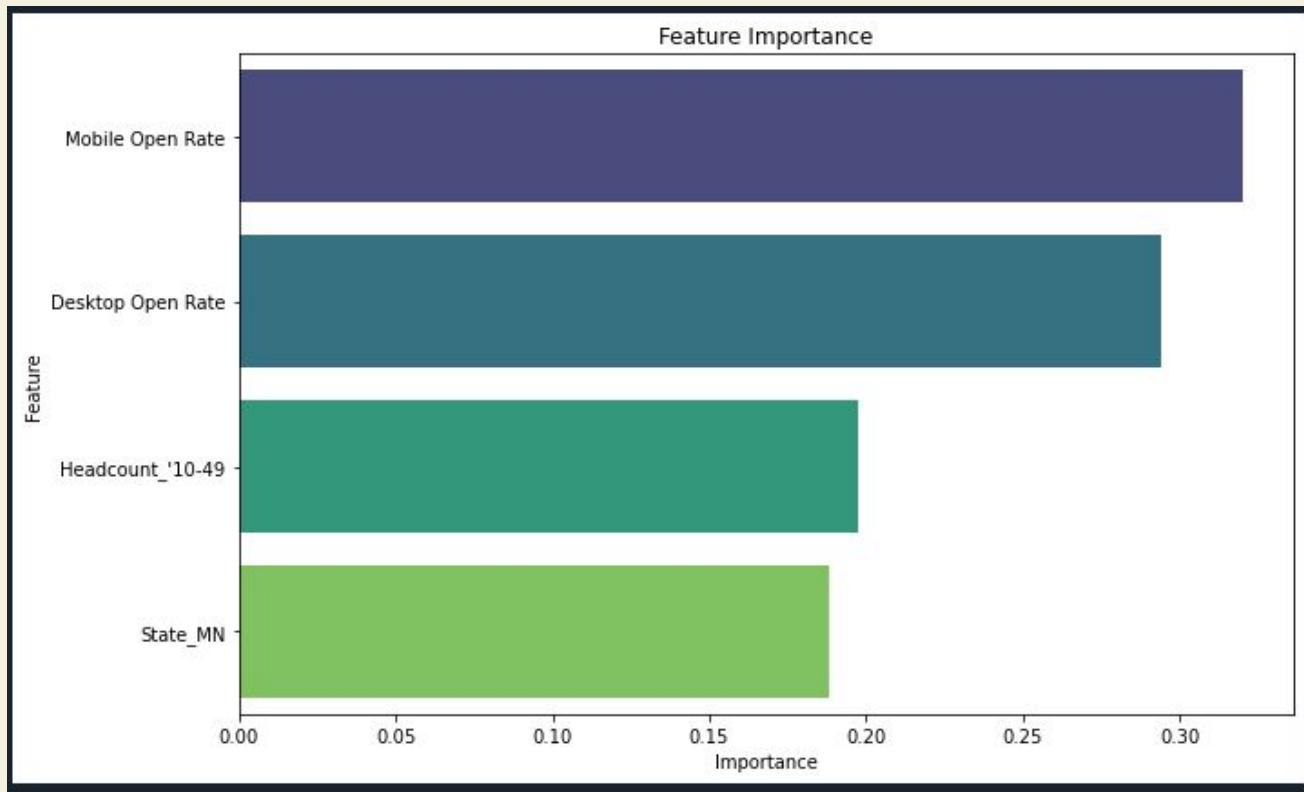
	precision	recall	f1-score	support
0.0	0.84	0.99	0.91	170
1.0	0.60	0.09	0.15	35
accuracy			0.83	205
macro avg	0.72	0.54	0.53	205
weighted avg	0.80	0.83	0.78	205

Confusion Matrix



Receiver Operating Characteristic (ROC) Curve





Feature Importance

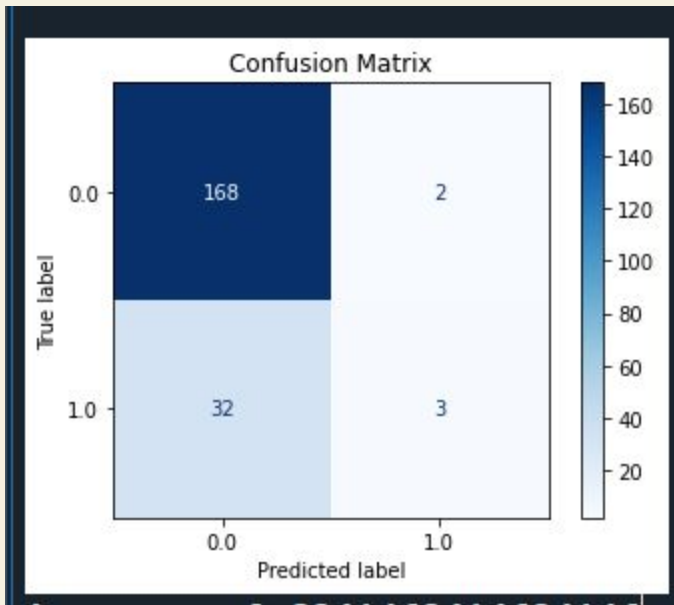
Gradient Boosting Machine(GBM): It is an ensemble learning algorithm used for classification tasks. It builds a strong predictive model by combining the predictions of multiple weaker models, typically decision trees, in a sequential manner.

Our Findings:

Accuracy: 0.8341463414634146

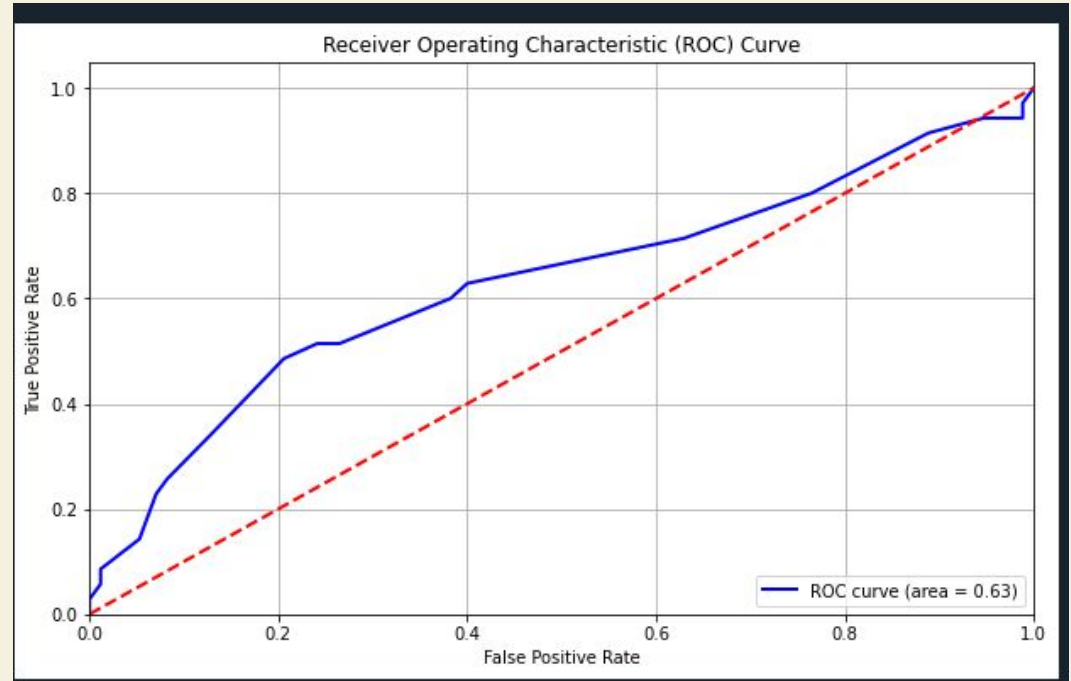
ROC-AUC Score: 0.6306722689075631

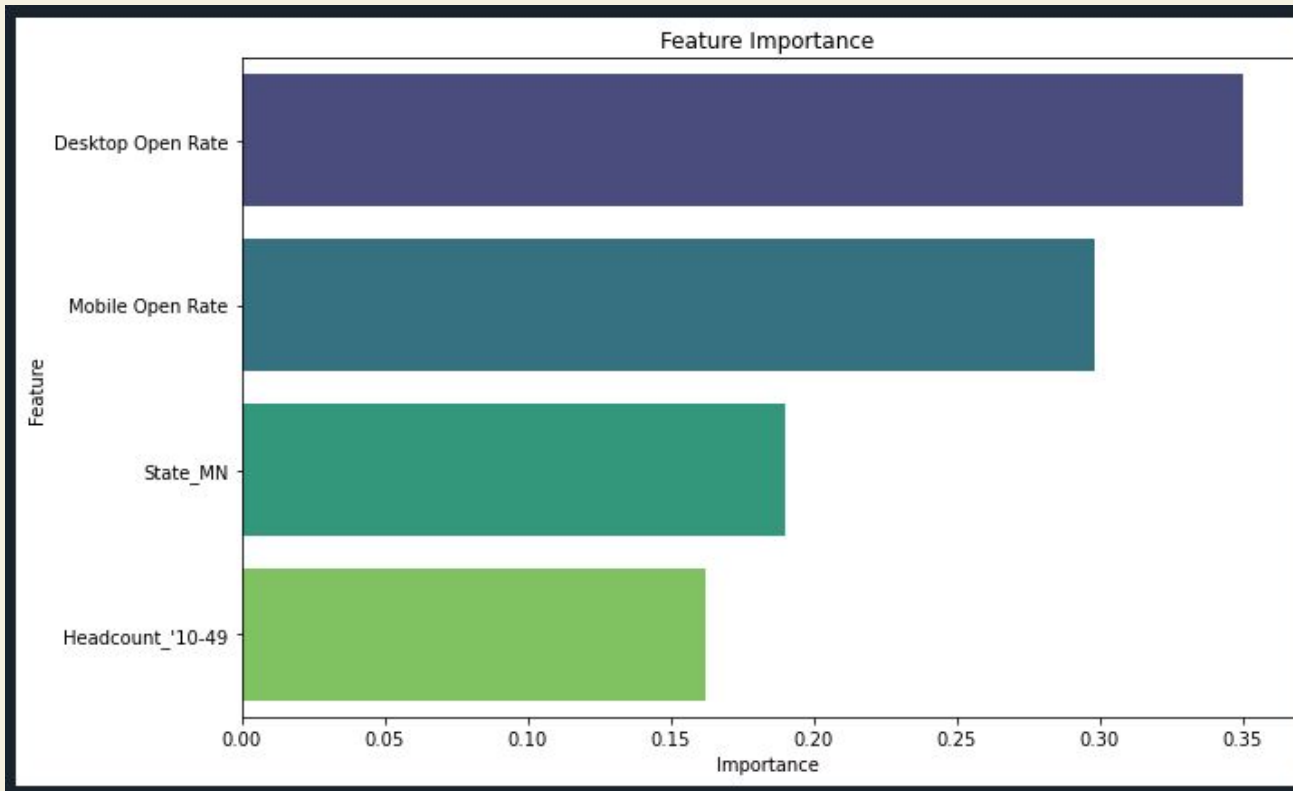
	precision	recall	f1-score	support
0.0	0.84	0.99	0.91	170
1.0	0.60	0.09	0.15	35
accuracy			0.83	205
macro avg	0.72	0.54	0.53	205
weighted avg	0.80	0.83	0.78	205



Confusion Matrix

Receiver Operating Characteristic (ROC) curve





Feature Importance

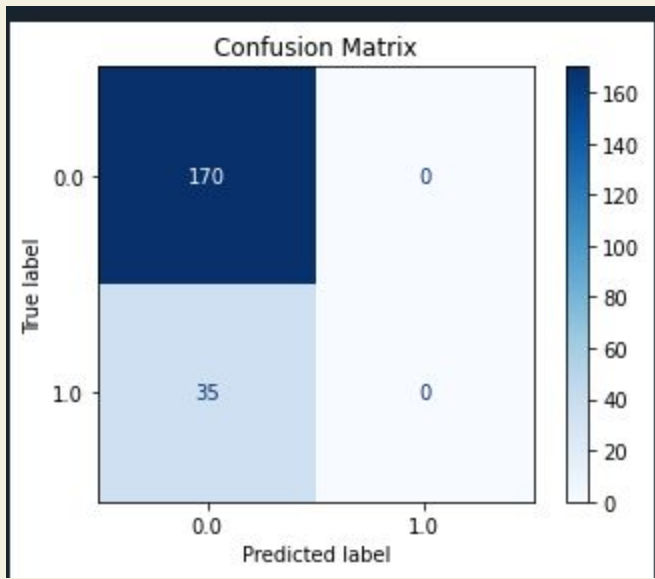
Support Vector Machine (SVM) a machine learning algorithm that classifies data and solves regression tasks. SVMs are particularly good at binary classification problems, where data is separated into two groups.

Our Findings:

Accuracy: 0.8292682926829268

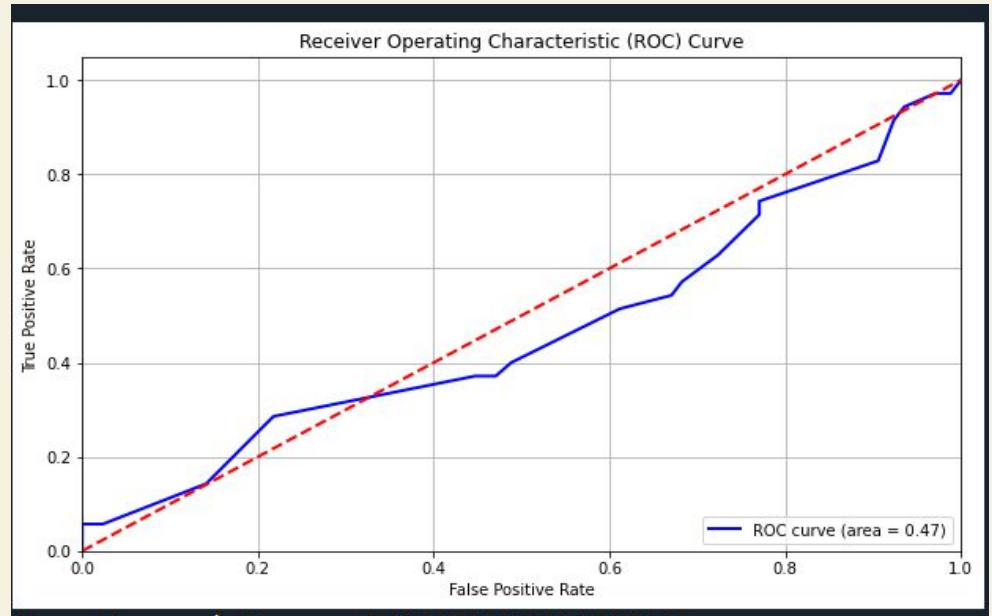
ROC-AUC Score: 0.4659663865546218

	precision	recall	f1-score	support
0.0	0.83	1	0.91	170
1.0	0.00	0.00	0.00	35
accuracy			0.83	205
macro avg	0.41	0.50	0.45	205
weighted avg	0.69	0.83	0.75	205



Confusion Matrix

Receiver Operating Characteristic (ROC) Curve



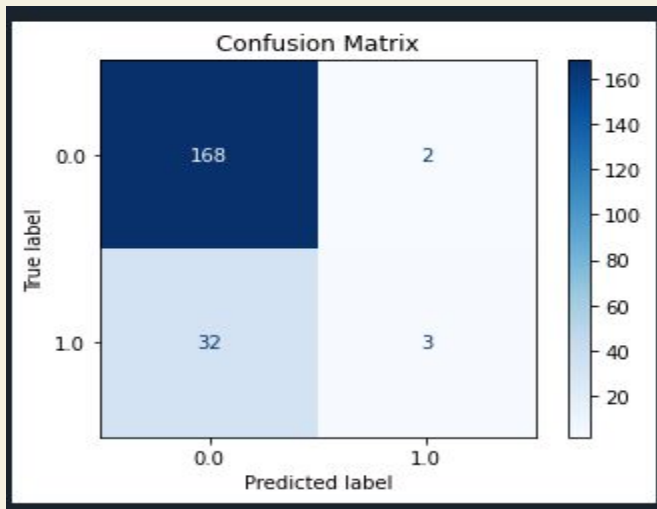
Decision Tree is a non-parametric supervised learning algorithm, which is utilized for both classification and regression tasks.

Our Findings:

Accuracy: 0.8341463414634146

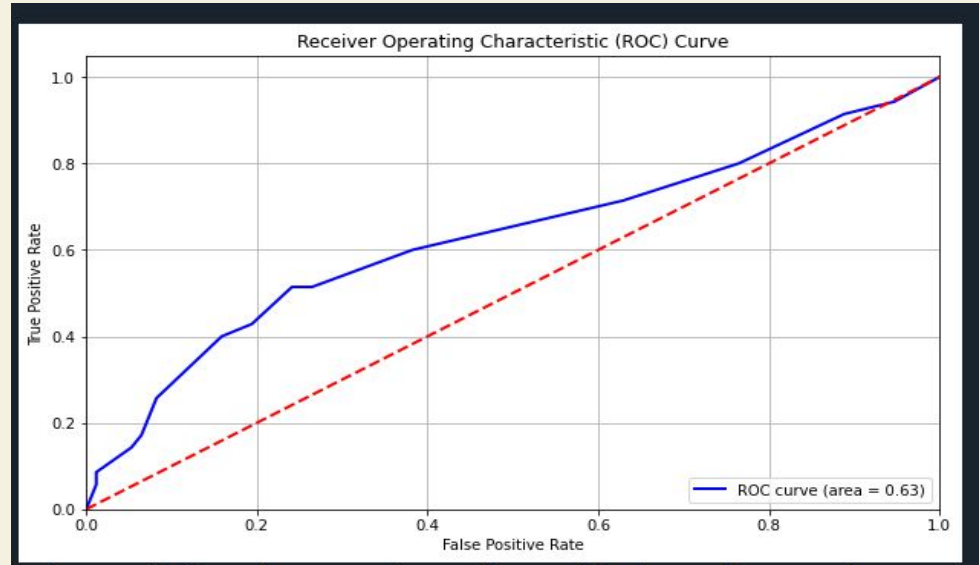
ROC-AUC Score: 0.6269747899159664

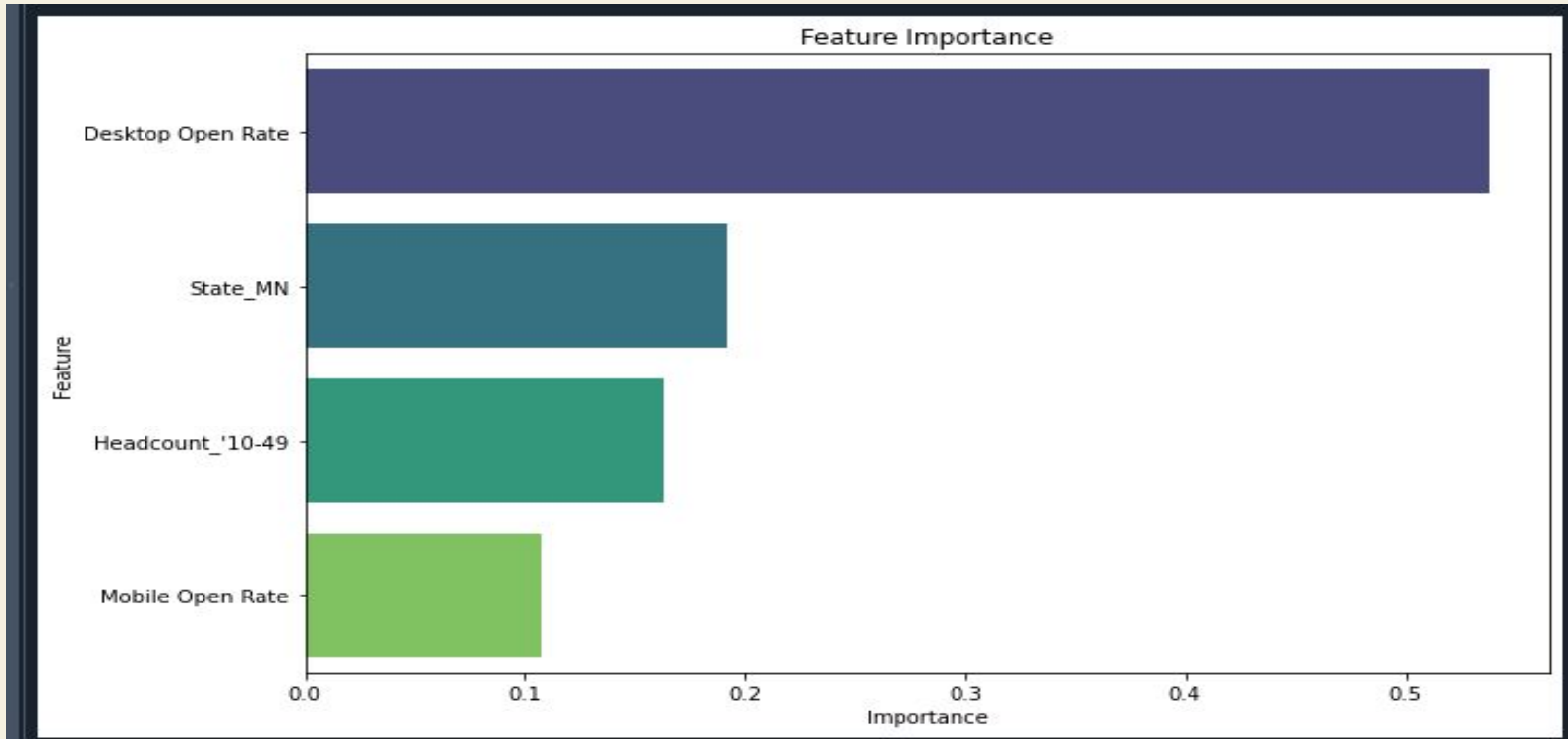
	precision	recall	f1-score	support
0.0	0.84	0.99	0.91	170
1.0	0.60	0.09	0.15	35
accuracy			0.83	205
macro avg	0.72	0.54	0.53	205
weighted avg	0.80	0.93	0.78	205



Confusion Matrix

Receiver Operating Characteristic (ROC) Curve





Feature Importance

Results

Random Forest	Gradient Boosting Machine	Support Vector Machine(SVM)	Decision Tree
Accuracy: 0.8341 ROC-AUC Score: 0.6298	Accuracy: 0.8341 ROC-AUC Score: 0.6307	Accuracy: 0.8293 ROC-AUC Score: 0.4660	Accuracy: 0.8341 ROC-AUC Score: 0.6270

- ❑ Since, correlation was strong between Desktop Open Rate, State_MN, Headcount '10-49 and Mobile Open Rate and Opened, we chose them for our independent variable.
 - ❑ We found out that after comparing all the models, Training the company data with Gradient Boosting Machine (GBM) is the best fit as it has the more ROU - AUC score in compare to other models.
-

Suggestions / Recommendations

- ❑ Final dataset can be increased by matching companies using other methods like Stemming, or other algorithms.
 - ❑ Stemming was tried which reduced number of rows to 706.
- ❑ Consistent way of recording data can make this analysis more focussed.
- ❑ Follow-up emails could be targeted more effectively to non-openers to increase interaction.
- ❑ Content of the email: Subject and going to spam can also impact the results of this analysis.



Summary of the Group Project



Biggest Challenge

Merging datasets especially master



Most Well-Done Aspect

Merging and modeling



Aspect Needing Improvement

The merging process could be further refined



Reference

- GeeksforGeeks: Provided guidance on data preprocessing, merging datasets, and applying specific Python functions.
 - Website: <https://www.geeksforgeeks.org>
- W3Schools: Helped with syntax, function usage, and understanding pandas, regex, and other Python-related concepts.
 - Website: <https://www.w3schools.com>
- For training the models like SVM, Decision tree and others:
 - Website: <https://scikit-learn.org/stable/modules/svm.html>
 - <https://scikit-learn.org/stable/modules/tree.html>
 - <https://scikit-learn.org/stable/modules/ensemble.html>
 -



Questions



Task Division and Communication

<u>Data Preprocessing: Outreach, Sent, and Open datasets</u>	Handled by Riya and Vaishnavi.
<u>Data Preprocessing: Master dataset</u>	Handled by Sanika, Riya, Aswad and Vaishnavi.
<u>Merging</u>	Collaborative effort by Sanika, Riya, and Vaishnavi.
<u>Modeling</u>	Managed by Riya and Sanika.
<u>Presentation</u>	Created by Vaishnavi and Aswad.
<u>Communication</u>	Whatsapp, Google Colab, Google docs, Gmail

