**Data Glacier**

Your Deep Learning Partner

# Final Presentation

*Project name:* **Bank Marketing Campaign**
*Team:* **Data Science Enthusiasts**
*Date:* **August 18th, 2021**

# Agenda

Executive Summary

Data Understanding

Data Transformation

Data Dependency

Model Building

Model results

Recommendations

Data Glacier
*Your Deep Learning Partner*

# Team member's details

| Group Name: *Data Science Enthusiasts* | | | | | |
|---|---|---|---|---|---|
| | **Name** | **Email** | **Country** | **College/Company** | **Specialization** |
| **1** | Amira Asta | amira.asta02@gmail.com | Tunisia | Afrikanda | Data Science |
| **2** | Vatsal Vinesh Mandalia | vatsalvm10@outlook.com | Oman | Graduated | Data Science |

# Github Repo link:

https://github.com/AsAmira02/Bank-Marketing-Campaign-DSEnthusiasts2021

This repository includes the four datasets, model code and necessary files used in this project.

# Executive Summary

- ## **The Client:**

  ABC Bank wants to sell its term deposit product to customers and before launching the product they want to develop a model which helps them understand whether a particular customer will buy their product or not.

- ## **Problem statement:**

  Build a Classification ML model to shortlist customers who are most likely to buy the term deposit product. This would allow the marketing team to target those customers through various channels.

- ## **Analysis:**

  The Analysis of this data is divided into the following parts:

  - Data Understanding
  - Univariate analysis
  - Bivariate analysis
  - Model recommendations

# Data Understanding

- **Datasets description:**

  Four datasets provided:

  - bank-additional-full: 20 inputs (+1 target variable) and 41119 observations

  - bank-additional: 20 inputs (+1 target variable) and 4119 observations

  - bank-full: 17 inputs (+1 target variable) and 45211 observations

  - bank: 17 inputs (+1 target variable) and 4521 observations

➔ The first 7 features are Bank client data.

➔ Features 8, 9, 10 and 11 are related to the last contact of the current campaign.

➔ Features 12, 13, 14 and 15 are other attributes.

➔ Features 16, 17, 18, 19 and 20 are social and economic context attributes.

➔ The last feature is the output feature (desired target).
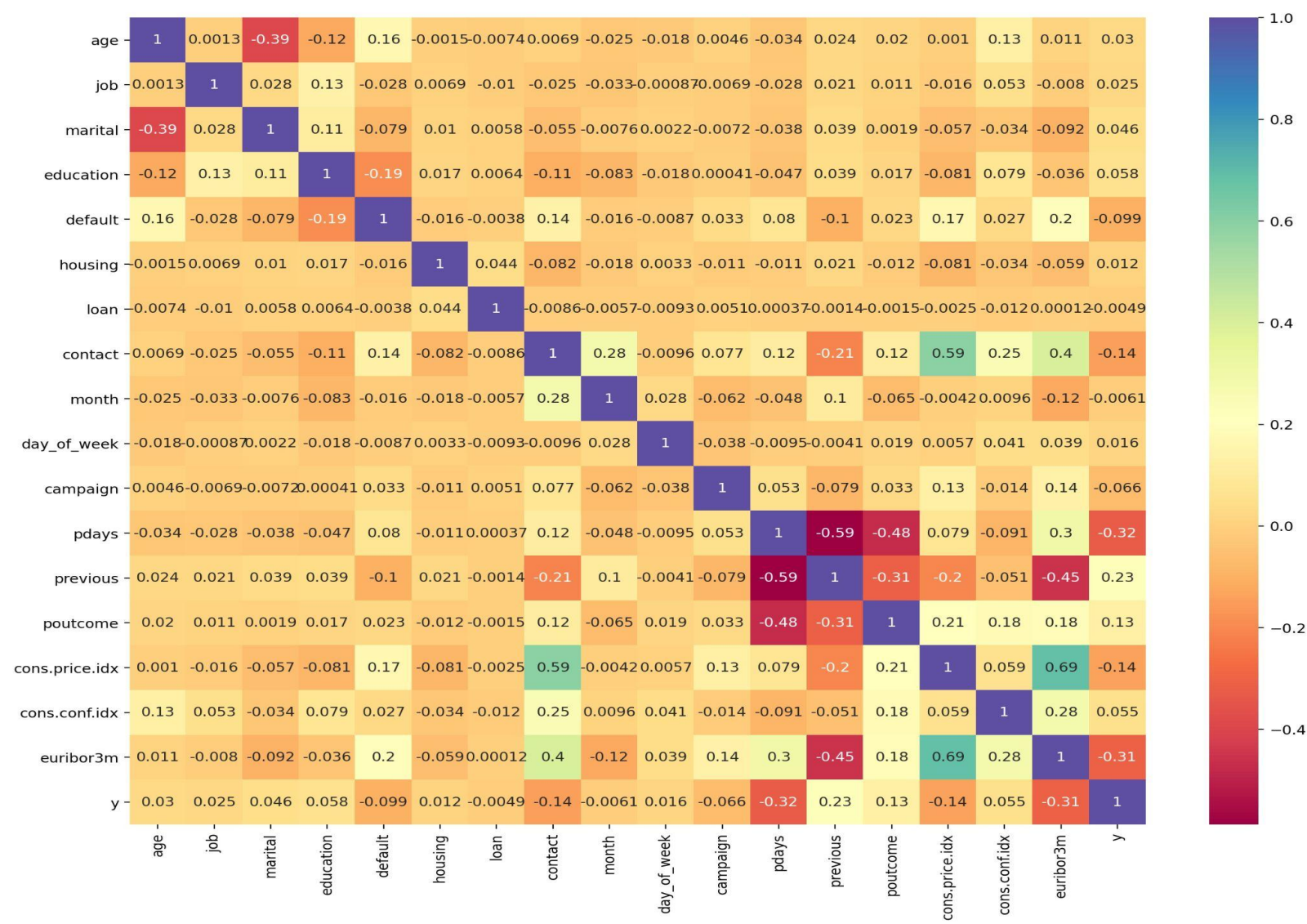
# Data Transformation

- **<u>Transformations:</u>**

  - Timeline of observations - May 2008 to November 2010.

  - 'Duration' feature is dropped to give realistic predictions from the classification model since it highly effects the output.

  - A frequently occurring missing value 'unknown' is considered as another category for the categorical features.

  - Duplicated rows were deleted from the dataset.

  - Outliers were not removed since they will help with the model generalization.

  - Handle highly correlated variables: Observations showed positive high correlation between **'emp.var.rate'**, **'nr.employed'**, and **'euribor3m'**. So we drop **'emp.var.rate'** and **'nr.employed'** as **'euribor'** can give us the price of money in the current market.

  - Encoding Categorical features using LabelEncoder since Machine learning algorithms can only read numerical values. It is therefore essential to do this step.

# Data Dependency

Now that our dataset contains all numeric variables, we could check correlation between all the features. As shown in the following figure, there are no features that are highly correlated and inversely correlated.

# Model Building

In order to predict the client subscription for a deposit term, we will use a predictive ML model to help us identify potential customers. As a start we begin by splitting the dataset, into training and testing sets in **80%** and **20%** respectively.

We choose to test out the following set of models since we don't know yet what algorithms will do well on this dataset.

The following algorithms selected for this classification problem include:

- **Linear Algorithms:**

  Logistic Regression (LR) *(Base Model)*

  Linear Discriminant Analysis (LDA).

- **Nonlinear Algorithms:**

  Classification and Regression Trees (CART),

  Support Vector Machines (SVM),

  Gaussian Naive Bayes (NB)

  k-Nearest Neighbors (KNN).

- **Ensemble Methods:**

  *Boosting Methods:* AdaBoost (AB) and Gradient Boosting (GBM).

  *Bagging Methods:* Random Forests (RF) and Extra Trees (ET).
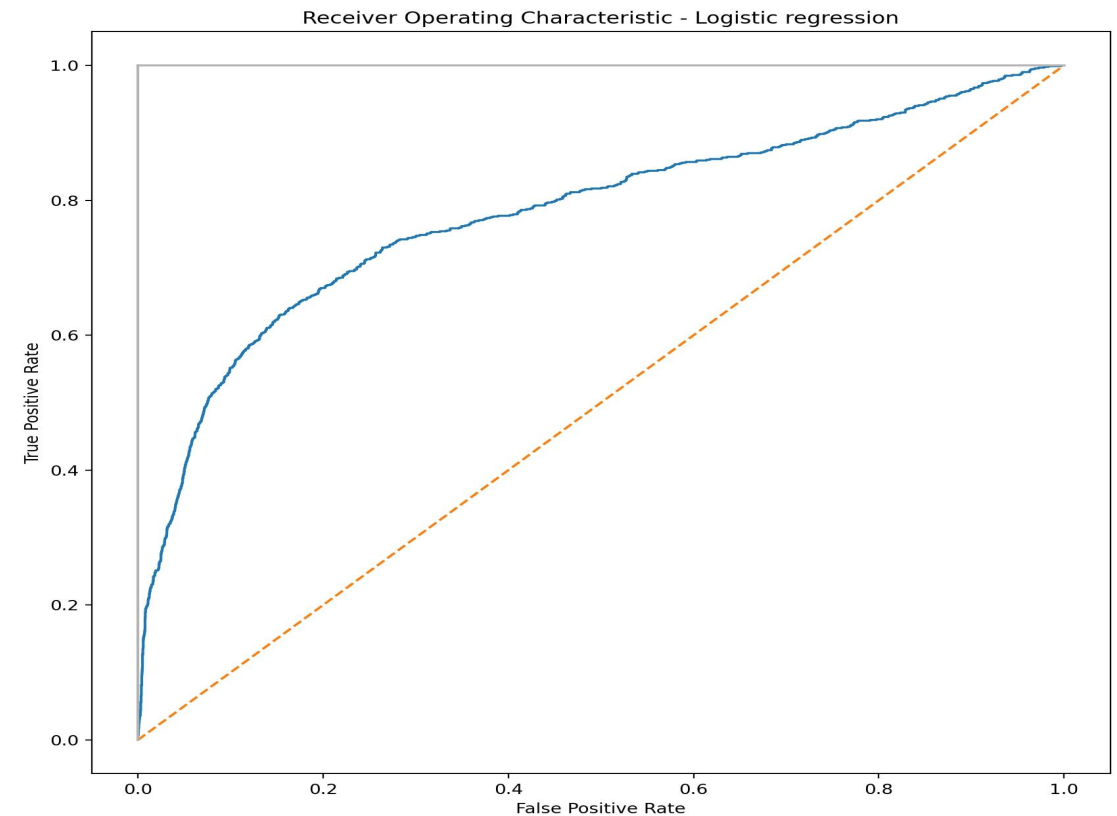
# Model Results

# Model results for Linear and nonlinear models

Logistic Regression is used as a base model in our case. As shown in the table, LDA gave better results than LR. Here are the ROC AUC Curves plotted separately for each model.
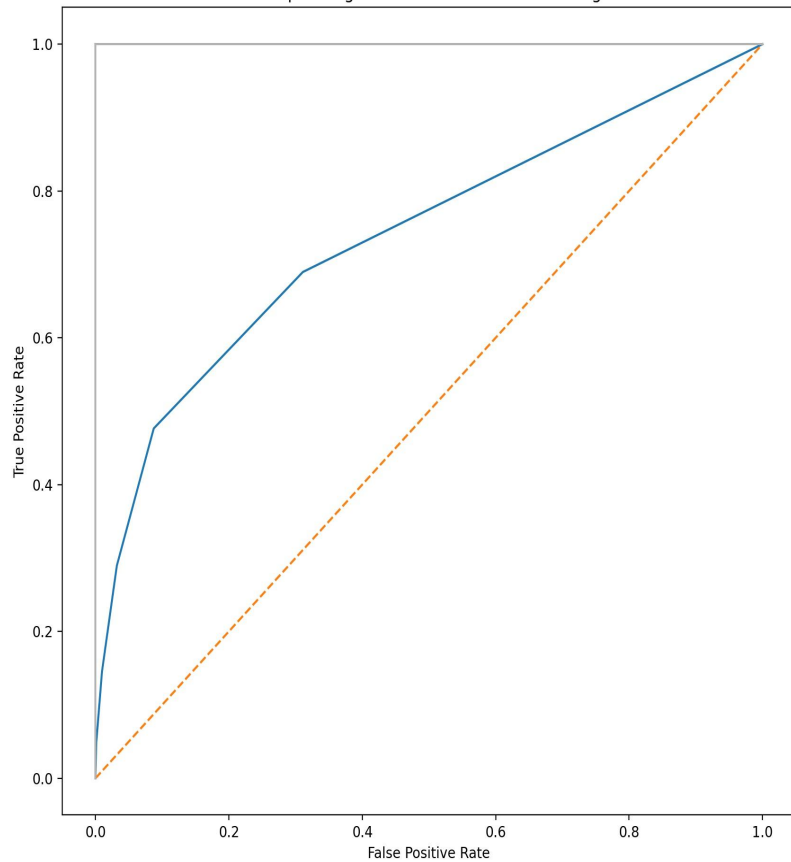
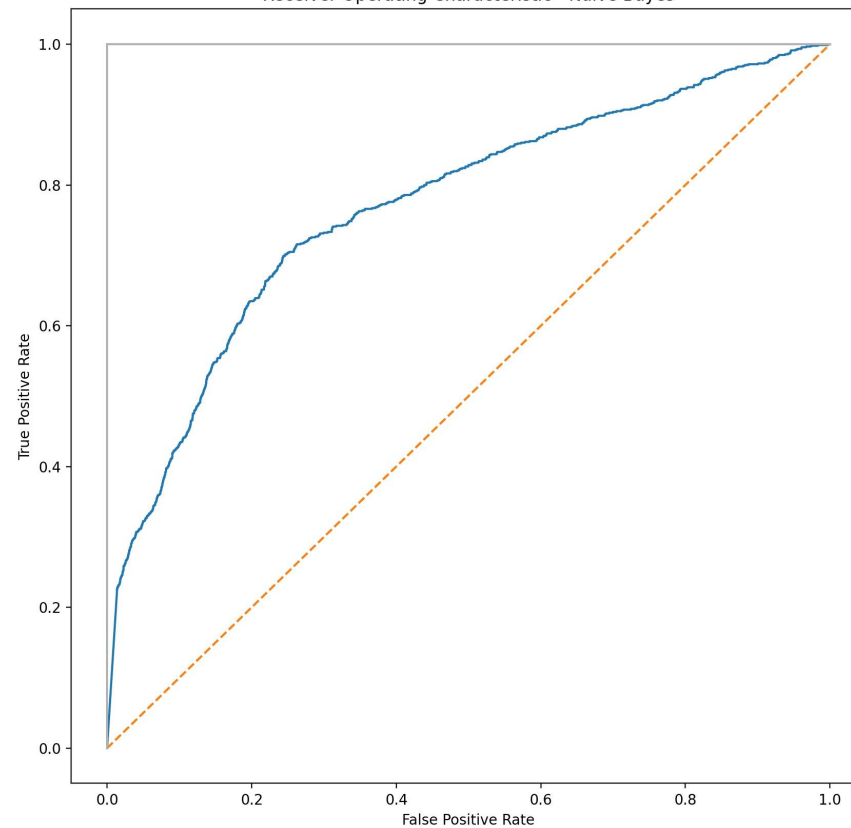## Linear Discriminant Analysis                    ## Logistic Regression

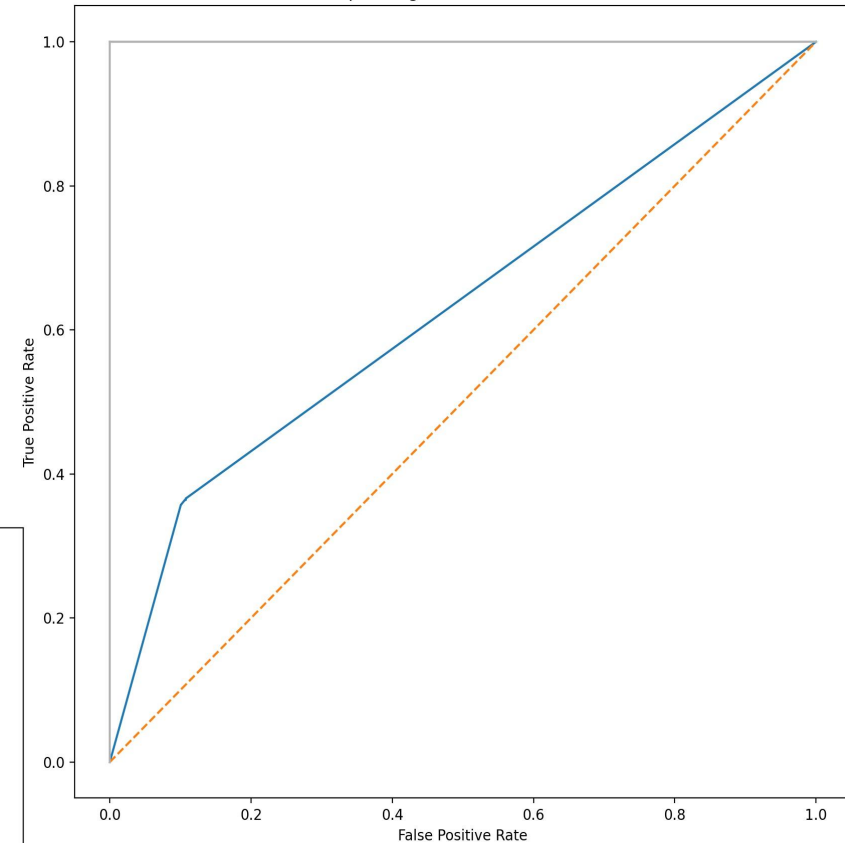Receiver Operating Characteristic - K-Nearset Neighbours

**KNN**

**Naive Bayes**
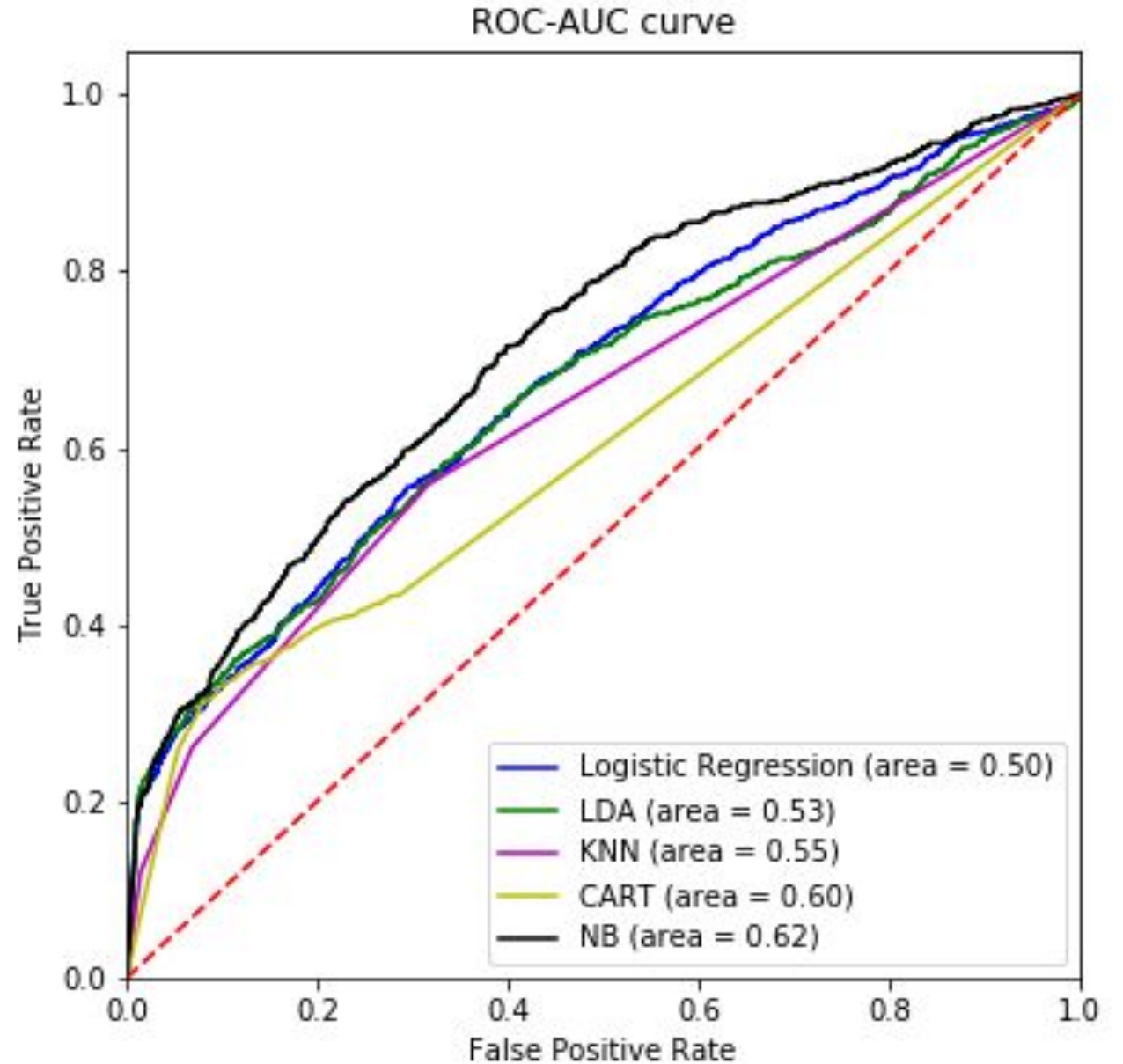
Receiver Operating Characteristic - Naive Bayes

Receiver Operating Characteristic - Decision Tree

**Decision Tree**

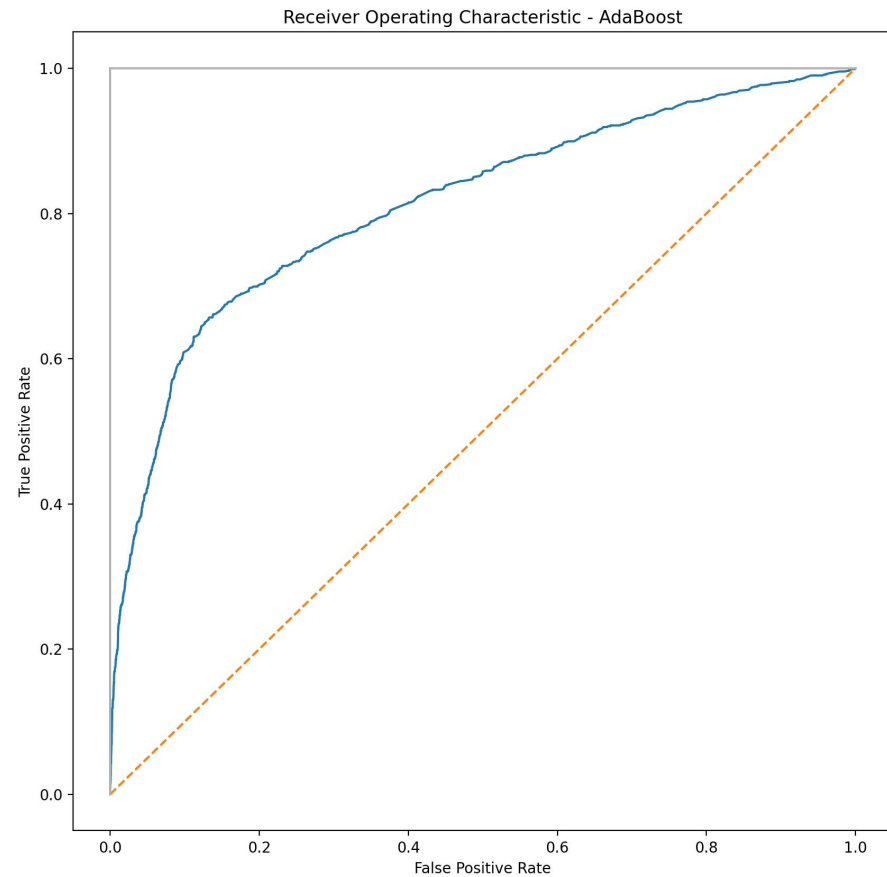The above results were gathered in one plot to compare the models.

The ROC curves are produced without Standard Scaler and Cross Validation techniques. This gives a basic result.
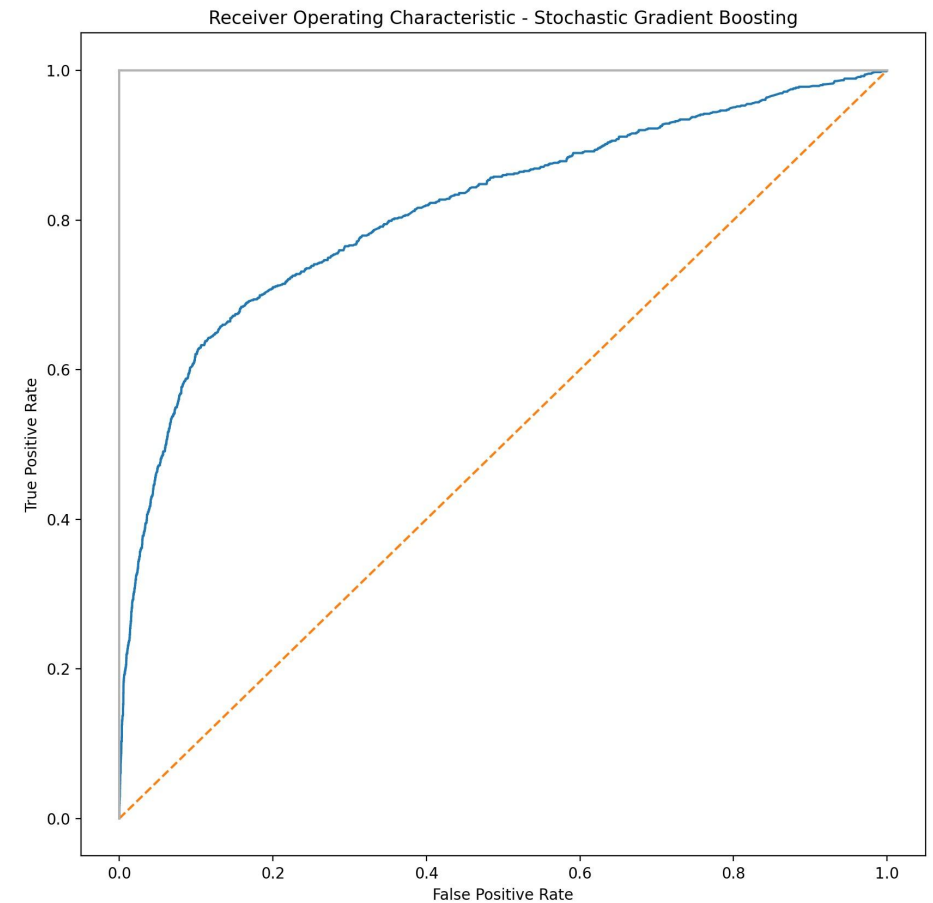


ROC-AUC curve

Legend:
- Logistic Regression (area = 0.50)
- LDA (area = 0.53)
- KNN (area = 0.55)
- CART (area = 0.60)
- NB (area = 0.62)

| | Model Name | ROC AUC Score |
|---|---|---|
| 0 | AB | 0.812690 |
| 1 | GBM | 0.814578 |
| 2 | RF | 0.791371 |
| 3 | ET | 0.771145 |

**Model results for Ensemble methods**

**AdaBoost**

**Stochastic Gradient Boosting**



Receiver Operating Characteristic - AdaBoost



Receiver Operating Characteristic - Stochastic Gradient Boosting

**Random Forest**

Receiver Operating Characteristic - Random Forest
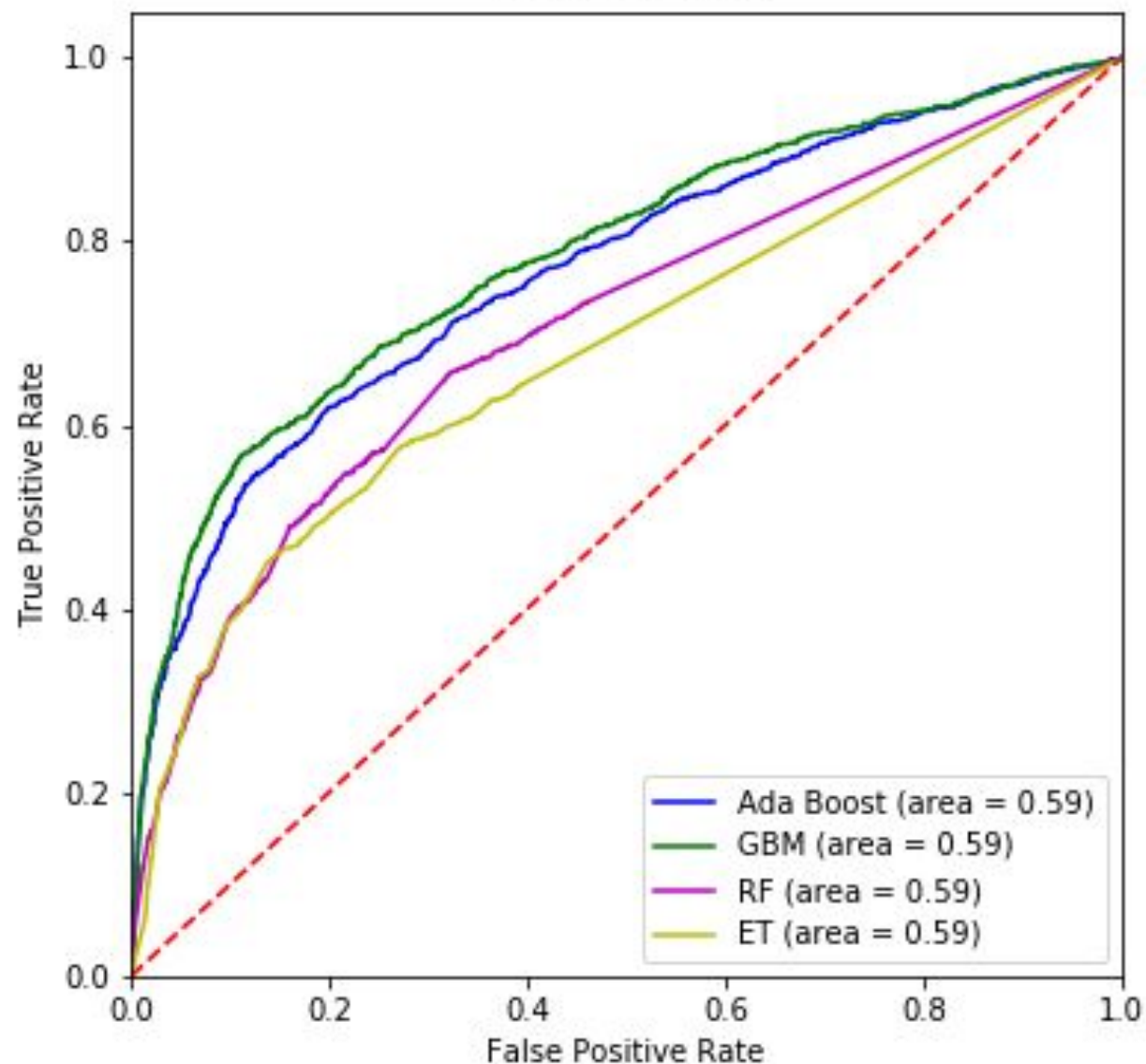
**Extra Trees**

Receiver Operating Characteristic - Extra Trees

ROC-AUC curve

We can see that both boosting techniques provide strong accuracy scores in the high 70% and even reached 80% with default configurations. The GBM model is the best model compared to the other ones. Therefore we will consider that model for production.

# Model Recommendation

# What type of ML model to use ?

Stochastic Gradient Boosting is the best model that fits with our classification task. Therefore, we proceed to save this model for later use.

# Thank You