

Capstone Project

Bank Marketing Prediction

(Supervised ML - Classification)

Bank Marketing Effective Prediction



Problem Statement

Make predictions

01

The classification Goal is to predict if the client will subscribe to a term deposit

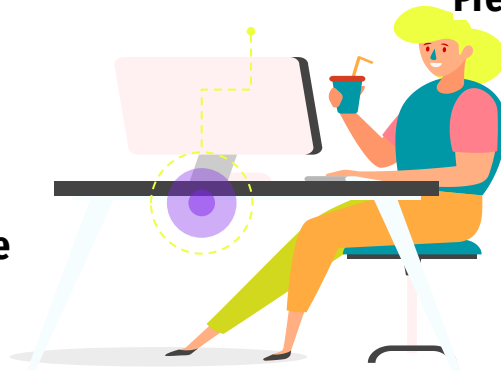
Model Development

02

Develop a **Supervised Machine Learning Model** using **Classification**.

Features

- 01 Age
- 02 Job -type of a job
- 03 Marital -Marital Status
- 04 Education
- 05 Default- Has credit in default?
- 06 Housing - Has housing Loan?
- 07 Loan - Has a personal loan
- 08 Contact - Communication type
- 09 Month- Last contact month of the year
- 10 Day_of_week -Last day of week
- 11 Duration -Last contact duration
- 12 RH5- humidity in bathroom



Campaign-No. Of contacts
performed during this campaign 13

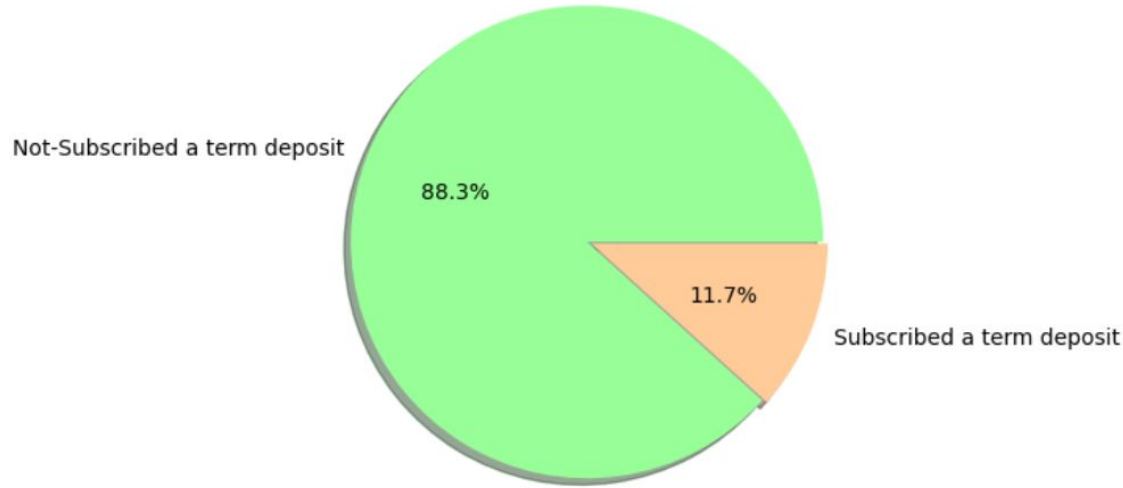
Pdays- No. of days 14

Previous-No. Of contacts performed
before this campaign 15

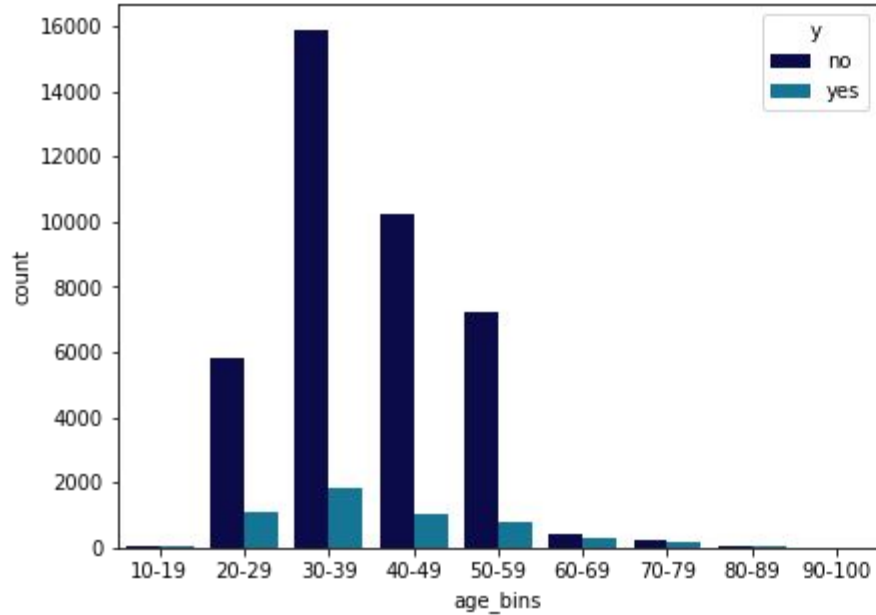
P-outcome -Outcome of
previous marketing campaign 16

Target variable(y) -Has client
subscribed to term deposit 17

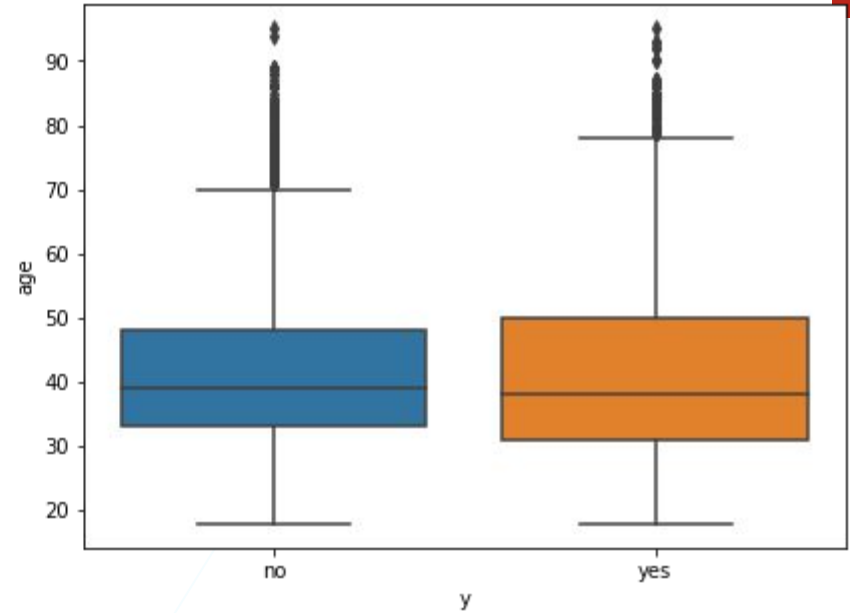
Proportion of Subscribed & Not Subscribed term Deposit



We can see from the above plot that the dataset is imbalanced, where the number of the subscribed class is close to 8 times the number of Not-subscribed class



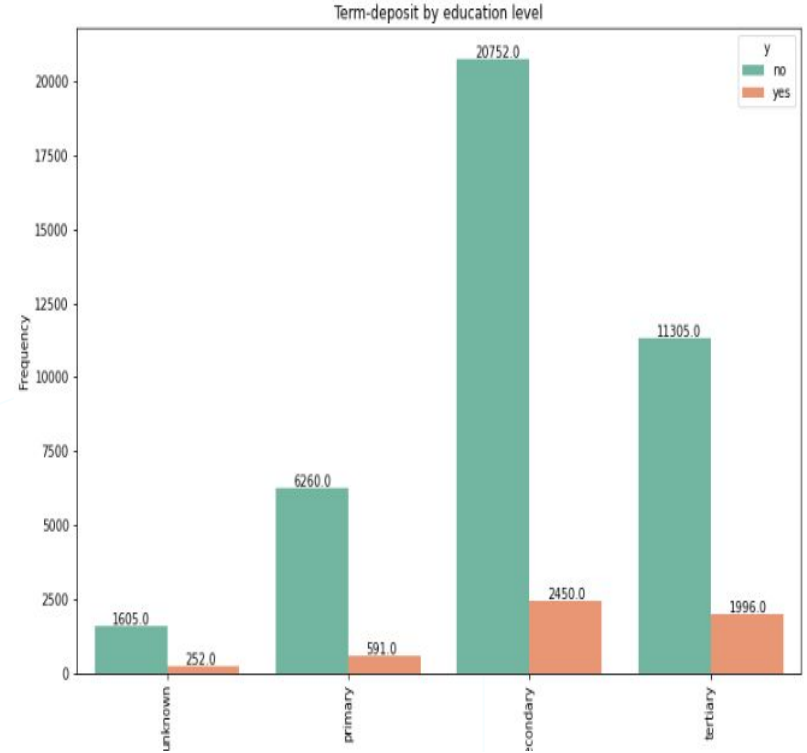
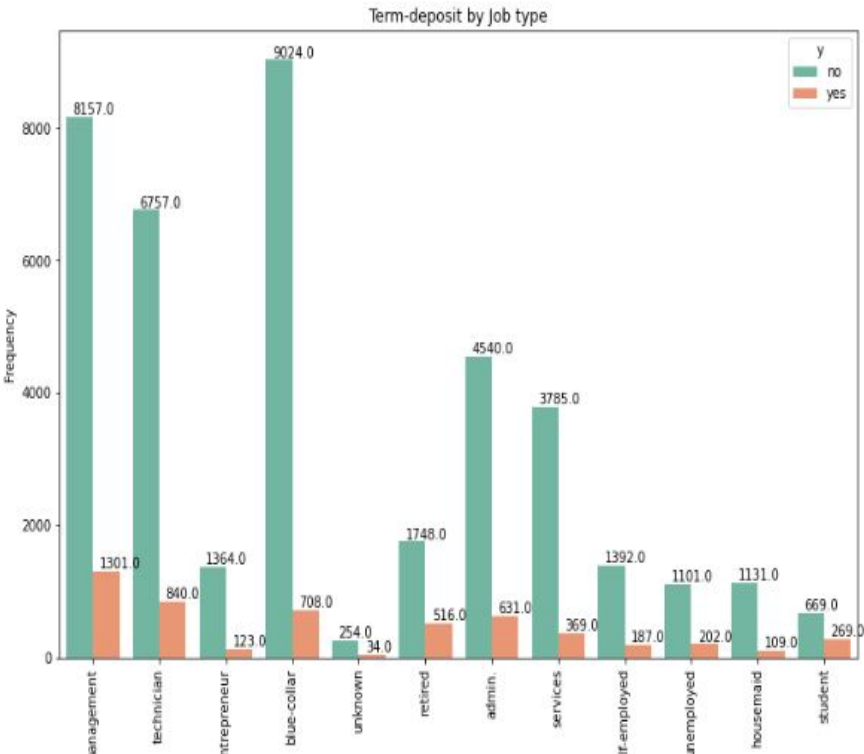
Majority of the customers are of the age group 30-39. Followed by 40-49 and 50-59



The Box Plot for the both subscribed and Not-subscribed customers looks the same

In No class, outliers are present above age 70 and For Yes class, Outliers are present above age 75

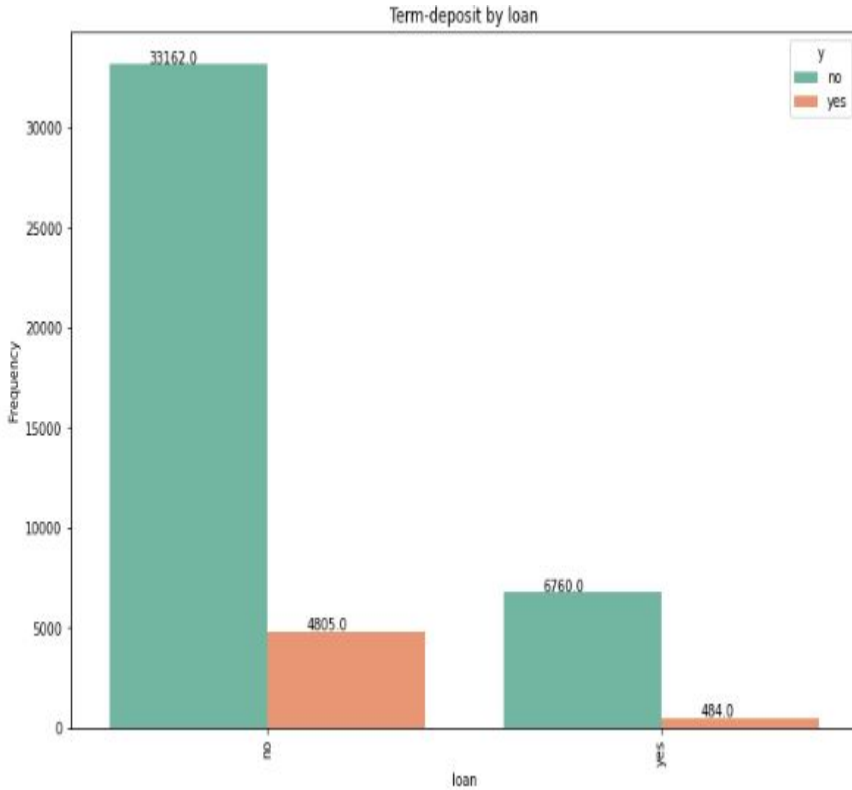
EDA



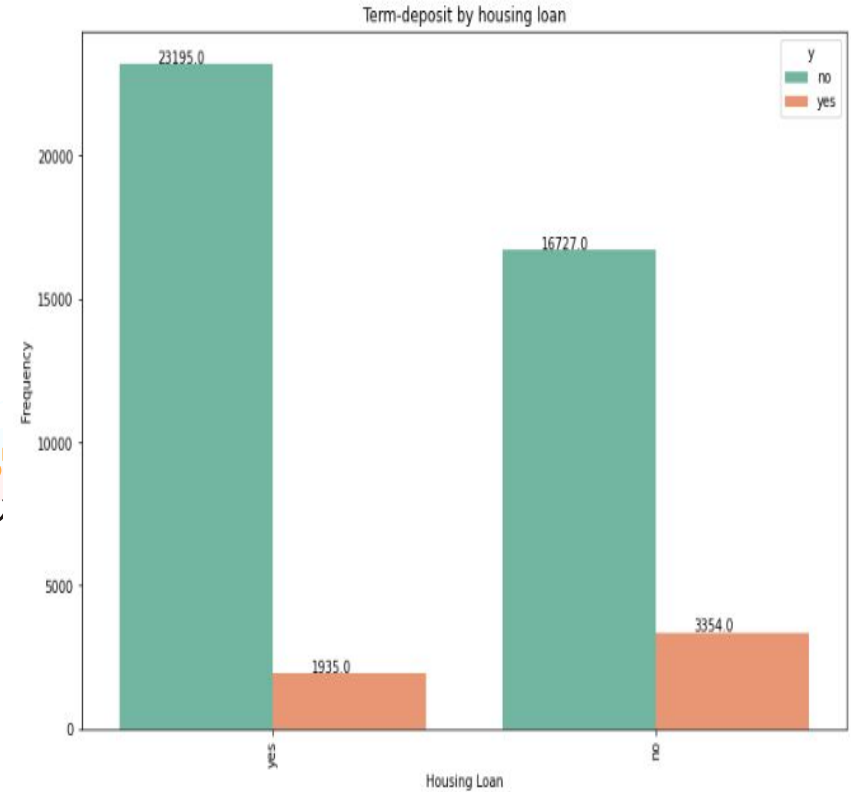
The majority of the customers who were contacted by the bank have blue-collar jobs but most of the term deposits have been taken by Management professionals

The majority of the customers have secondary education as the highest level of education

EDA

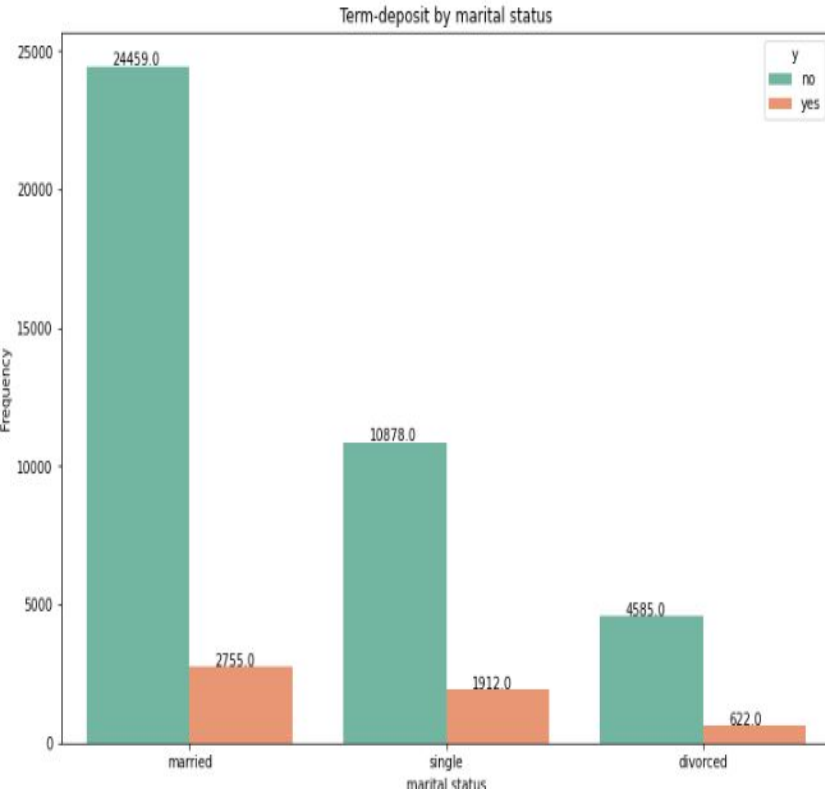


Very few people have taken a loan and these are more likely to take a term deposit.

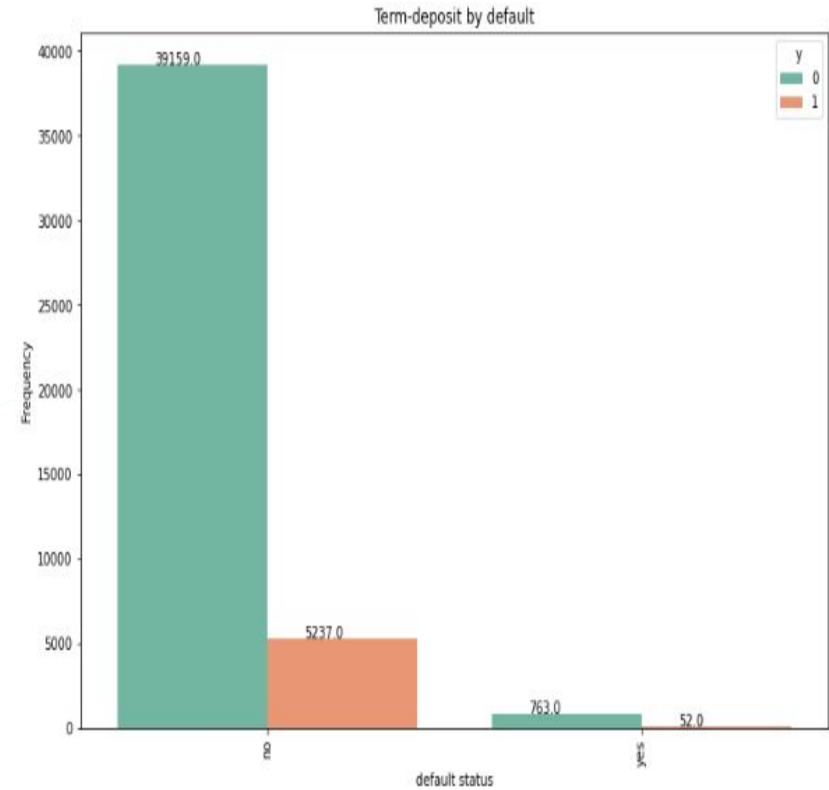


The majority of the customers have a housing loan. But those who do not have a housing loan are more likely to subscribe to a term deposit.

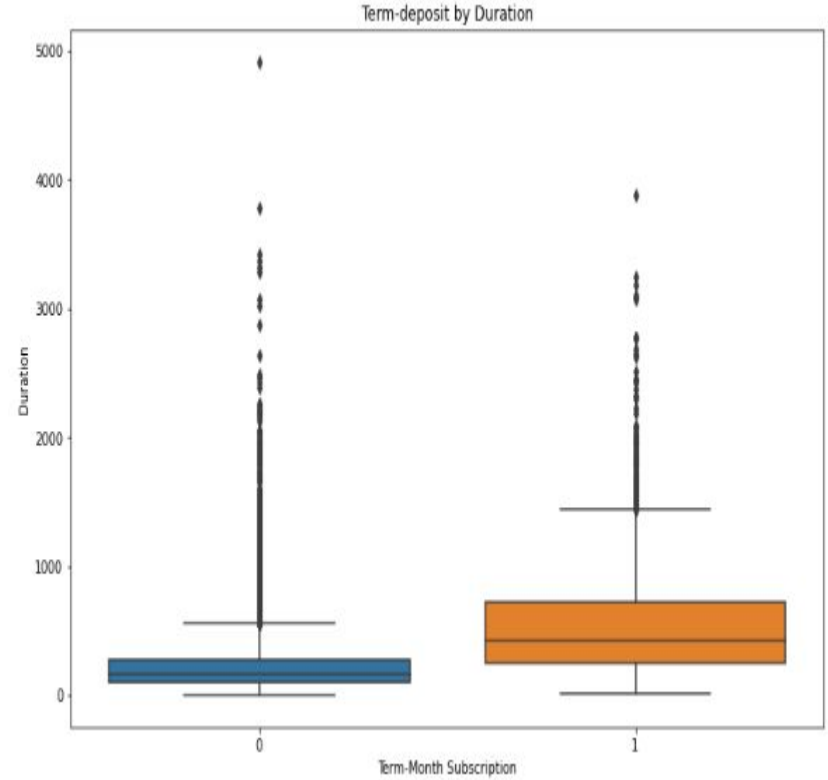
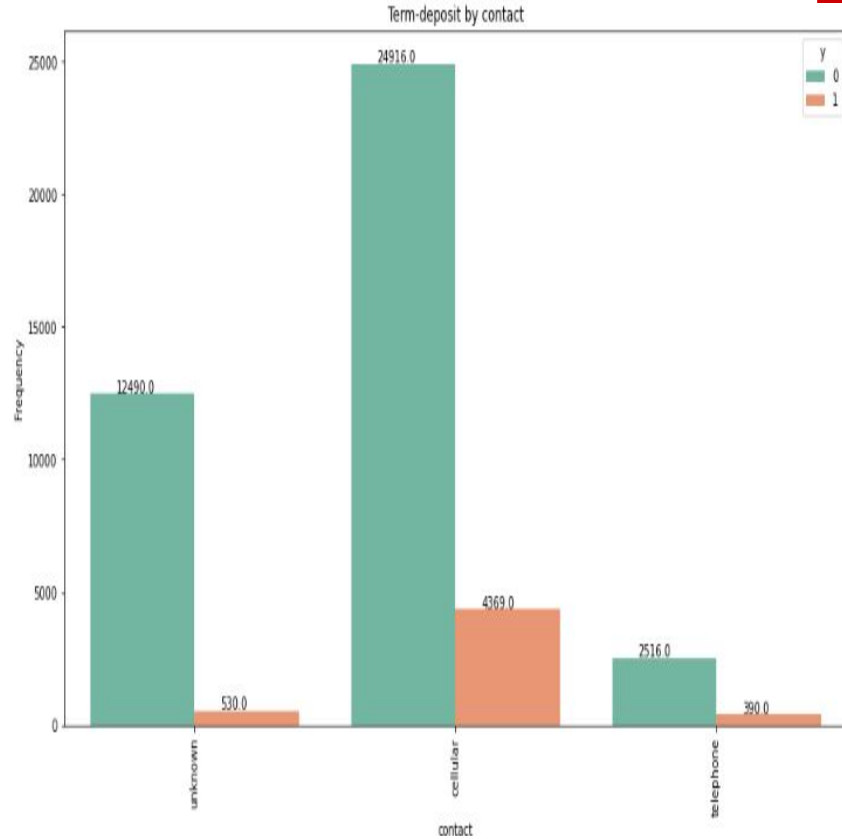
EDA



The Majority of the customers are married, followed by Single and Divorced.



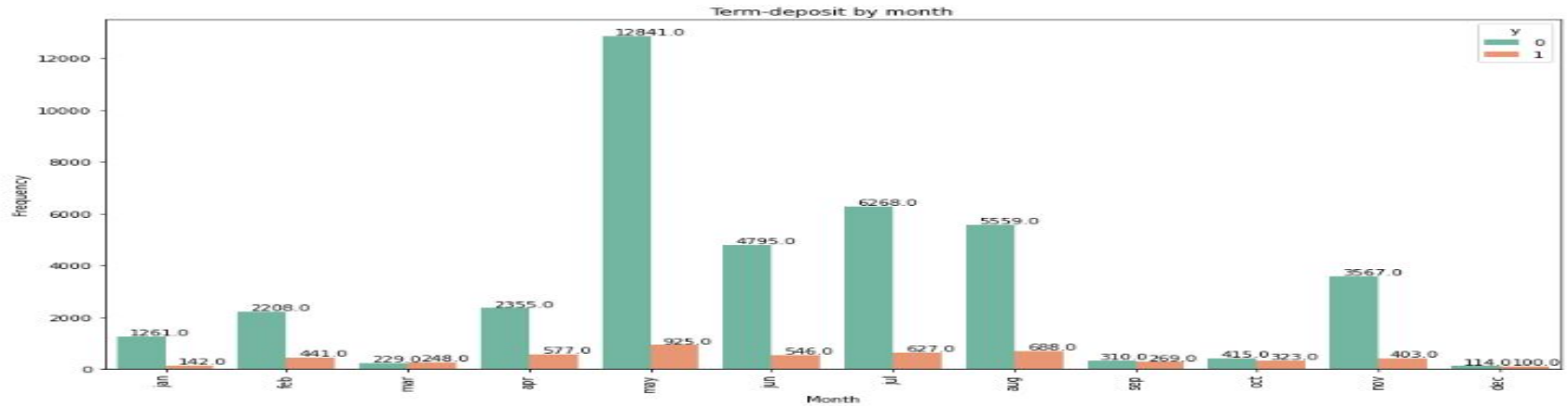
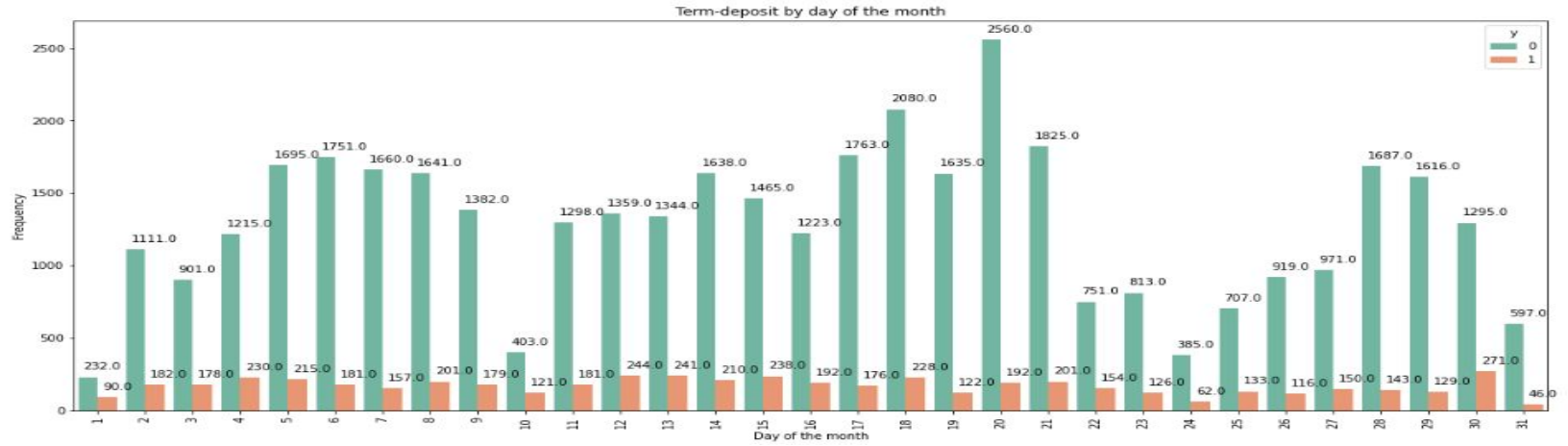
People with default status as 'no' are the most who have been contacted by the bank for the deposits.



Customers are contacted more by cellular rather than telephone.

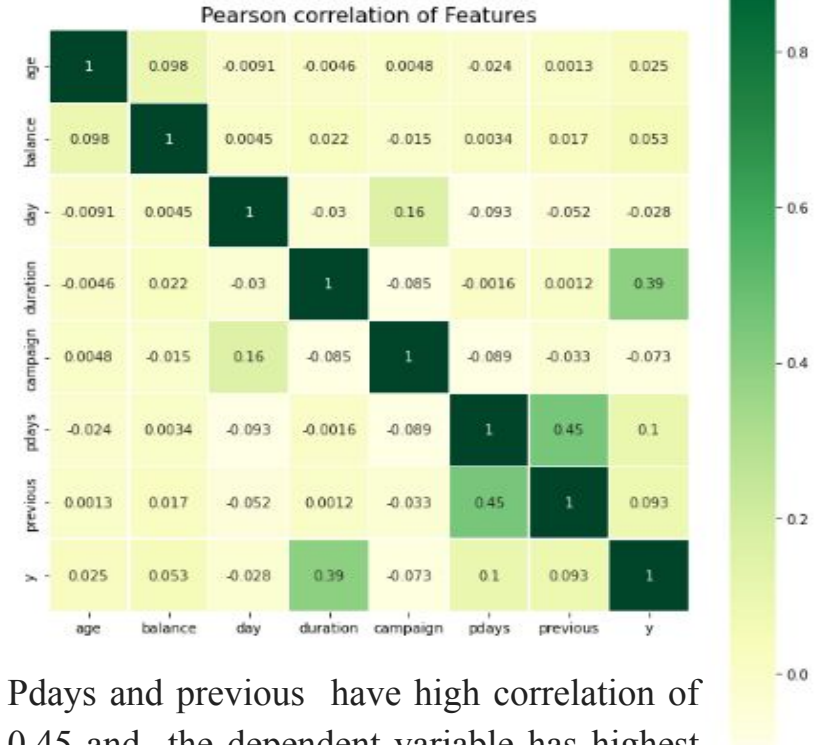
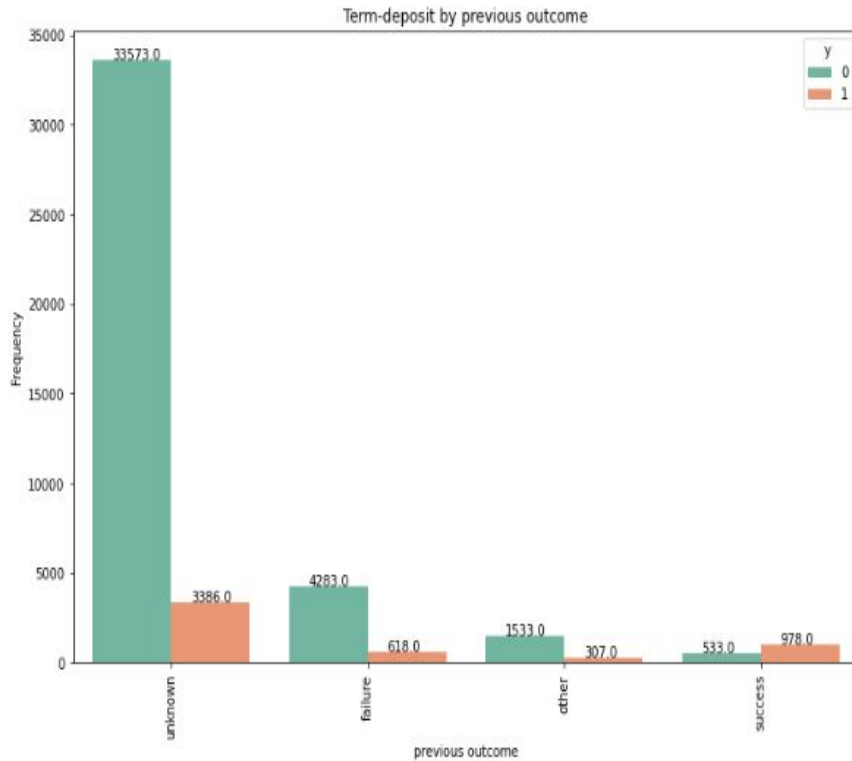
Customers who subscribed to term deposits have a relatively higher call duration.

EDA



EDA

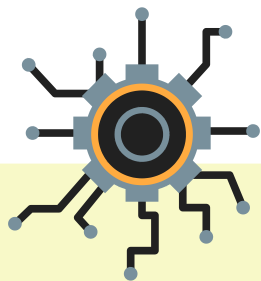
AI



Pdays and previous have high correlation of 0.45 and the dependent variable has highest correlation with duration .

The outcome of the previous campaign is unknown for most of the customers but 64% of customers who had a successful outcome in the previous campaign did subscribe for a term deposit.

Feature Engineering



Feature Engineering

01

Dropping unknowns

Education, Job, contact, outcome and balance.

02

Label Encoding

Default, Housing, Loan and y.

03

Getting Dummies

Job, Education, Marital Status, Contact, Month and Previous outcome

Model Selection

01

**Logistic
Regression**

02

Decision Trees

03

**Random Forest
Classifier**

04

**Gradient
Boosting
Classifier**

05

**XGBoost
Classifier**

06

**Neural Networks
Classifier**

07

**Naïve Bayes
Classifier**

08

Linear SVM

Cross Validation

Definition:

Cross-validation is a resampling method that uses different portions of the data to test and train a model on different iterations. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice.

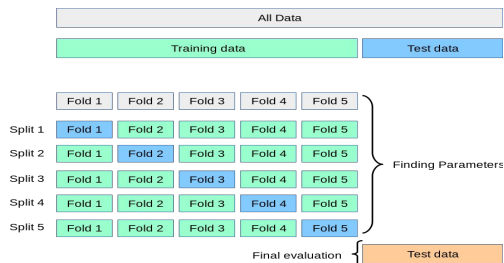


Fig A: K-Fold Cross-Validation

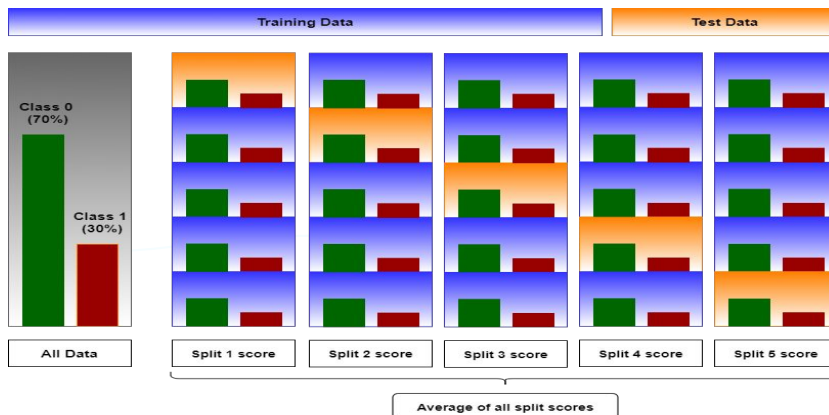
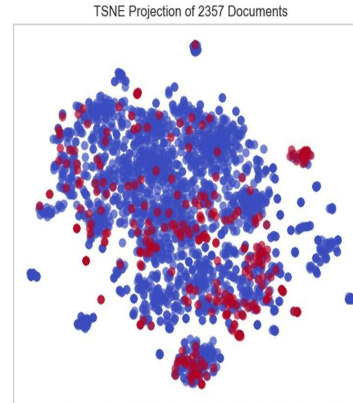
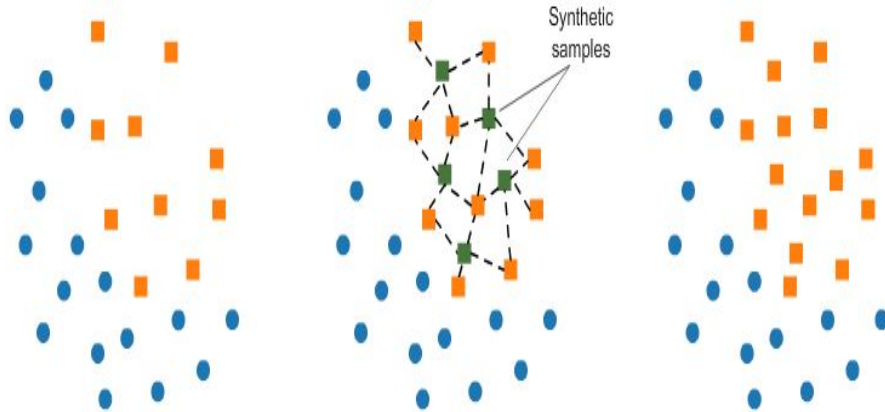


Fig B: Stratified K-Fold Cross-Validation

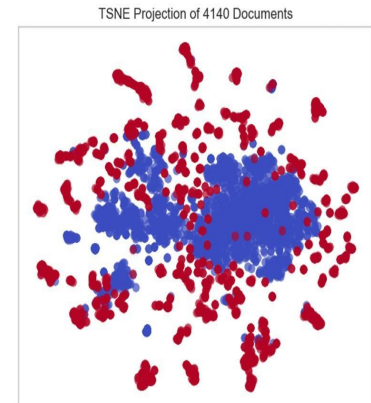
Stratified K-Fold Cross-Validation

Stratified K-Fold is an enhanced version of K-Fold cross-validation which is mainly used for imbalanced datasets. Just like K-fold, the whole dataset is divided into K-folds of equal size. But in this technique, each fold will have the same ratio of instances of target variable as in the whole datasets.

SMOTE



(a) Before SMOTE



(b) After SMOTE

SMOTE: a powerful solution for imbalanced data.

SMOTE is an algorithm that performs data augmentation by creating **synthetic data points** based on the original data points. SMOTE can be seen as an advanced version of oversampling, or as a specific algorithm for data augmentation. The advantage of SMOTE is that you are **not generating duplicates**, but rather creating synthetic data points that are **slightly different** from the original data points.

The SMOTE algorithm works as follows:

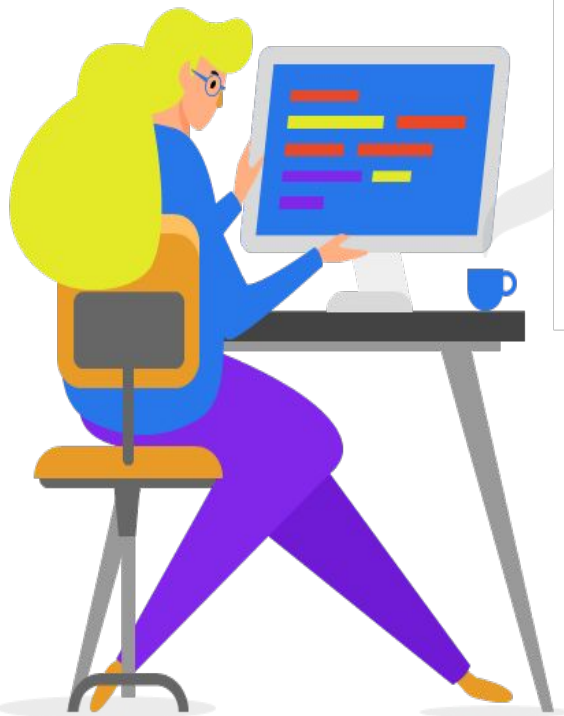
Steps Involved

1. You draw a random sample from minority class.
2. For the observations in this sample, you will identify the k nearest neighbors.
3. Take one of those neighbors and identify the vector between the current data point and the selected neighbor.
4. You multiply the vector by a random number between 0 and 1.
5. To obtain the synthetic data point, you add this to the current data point.

Tree Visualization



Confusion matrix



		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Confusion matrix terminologies

TP

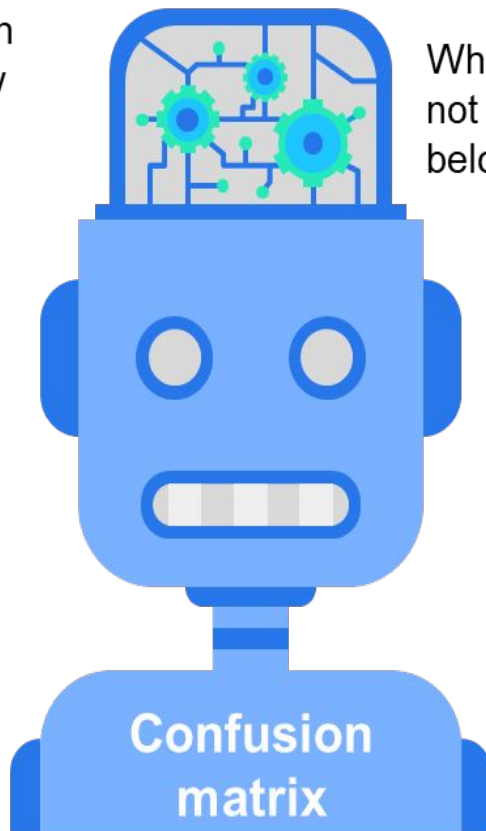
When you predict an observation belongs to a class and it actually does belong to that class.

TN

When you predict an observation does not belong to a class and it actually does not belong to that class.

FP

When you predict an observation belongs to a class and it actually does not belong to that class.



When you predict an observation does not belong to a class and it actually does belong to that class

FN**FPR**

Type 1 error
 $FP/(FP+FN)$

FNR

Type 2 error

0
6

Evaluation matrix for classification

Accuracy

0.91

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TP}}$$

F1 Score

92

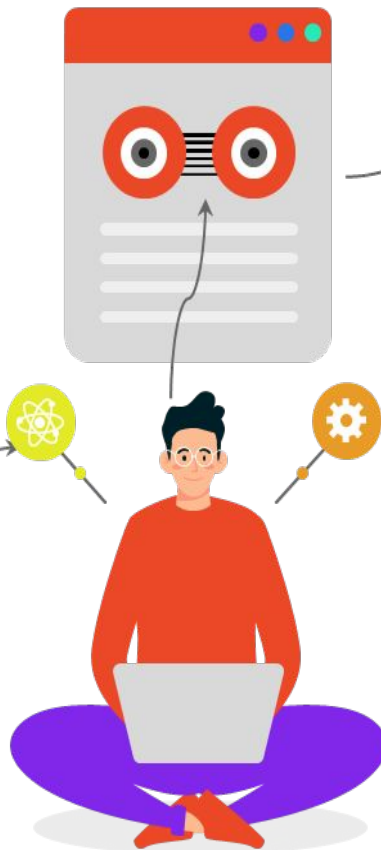
$$2 * ((\text{precision} * \text{recall}) / (\text{precision} + \text{recall}))$$

Precision

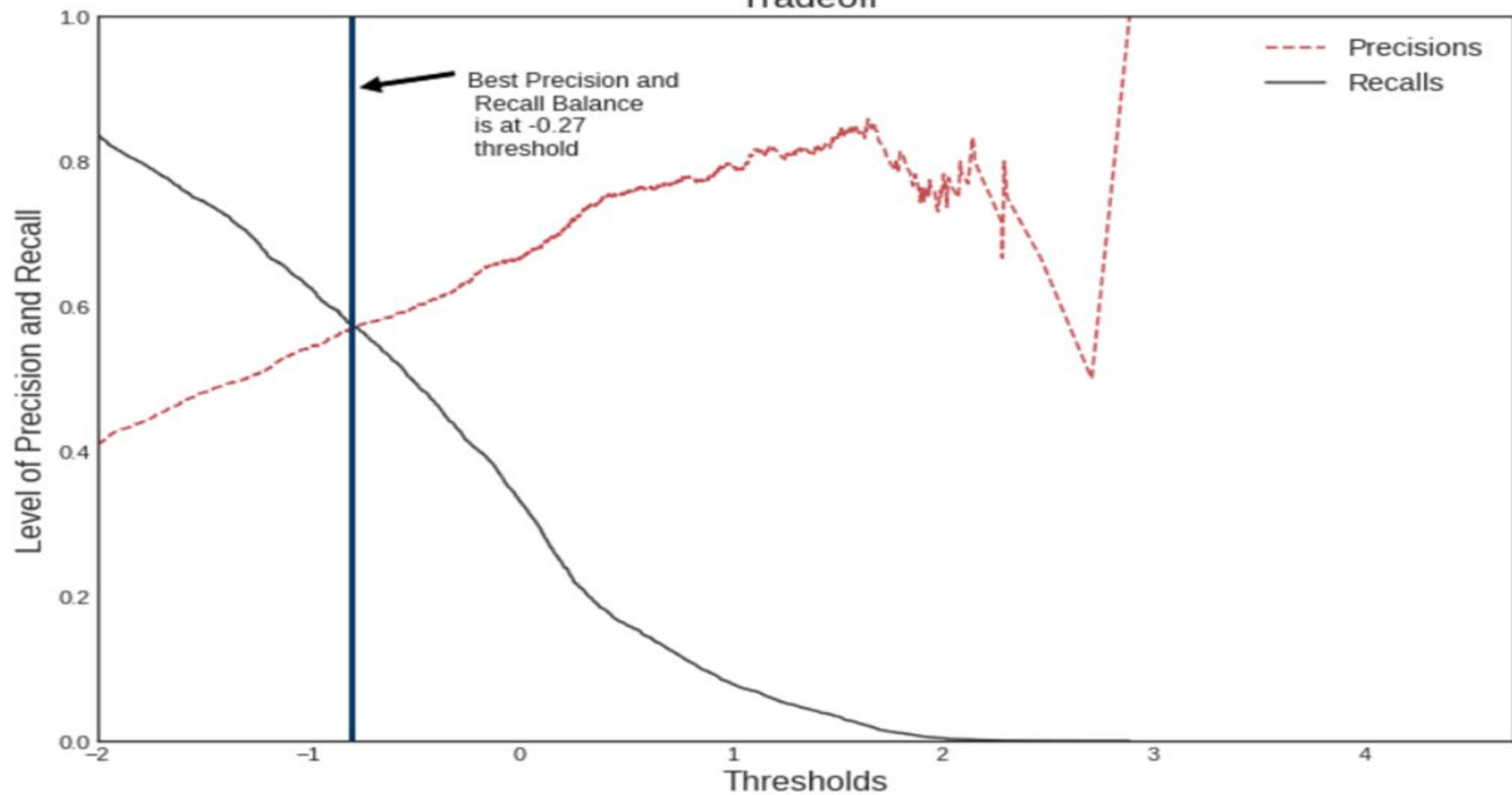
$$\text{TP} / (\text{TP} + \text{FP}) = 91$$

Recall

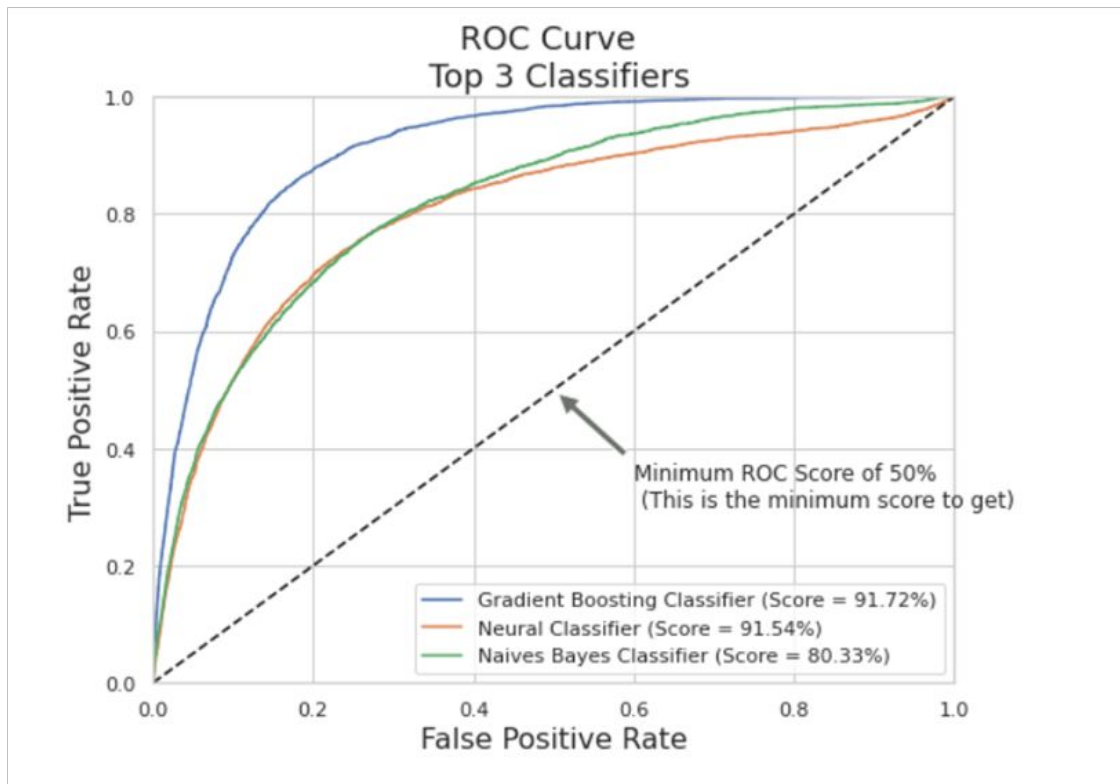
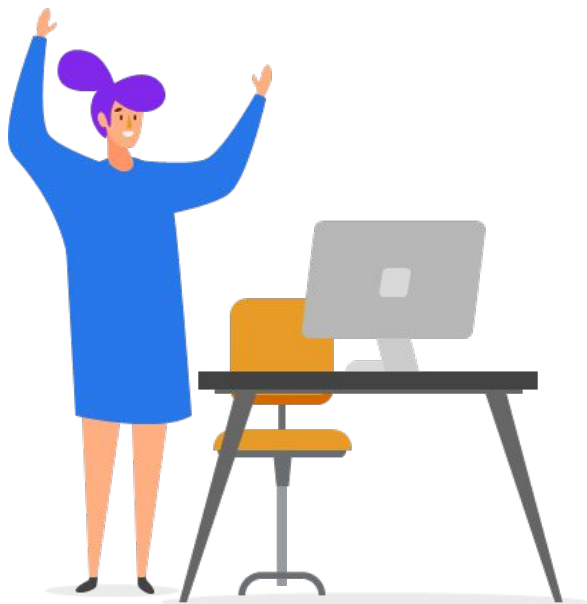
$$\text{TP} / (\text{TP} + \text{FN}) = 92$$



Precision and Recall Tradeoff



ROC AUC curve and Model performance comparison



Conclusions

Months of Marketing Activity

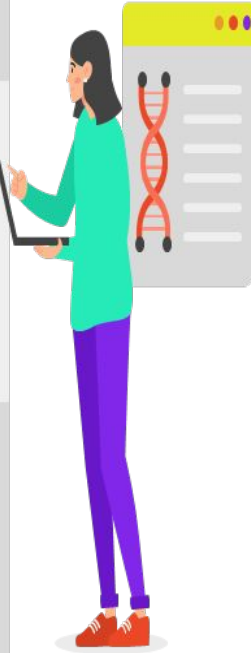
it will be wise for the bank to focus the marketing campaign during the months of **March, September, October and December.**

Campaign Calls

A policy should be implemented that states that no more than 3 calls should be applied to the same potential client in order to save time and effort in getting new potential clients. Remember, the more we call the same potential client, the more likely he or she will decline to open a term deposit.

Age Category

The next marketing campaign of the bank should target potential clients in their 20s or younger and 60s or older. The youngest category had a 60% chance of subscribing to a term deposit while the eldest category had a 76% chance of subscribing to a term deposit. It will be great if for the next campaign the bank addressed these two categories and therefore, increased the likelihood of more term deposits subscriptions.



Thank You