

# Final Report

## Problem Statement

In the financial industry, a term deposit is a type of savings account offered by banks and financial institutions. It involves depositing a fixed amount of money for a specific period, often ranging from a few months to several years, at a fixed interest rate.

A Portuguese bank is aiming to improve its customer acquisition marketing strategies for term deposits. By identifying potential subscribers, the bank will tailor its marketing efforts to increase subscription rates.

The insights gained from the model can be used to drive targeted campaigns for reaching out to clients most likely to subscribe, optimizing marketing expenses and improving overall customer satisfaction. This project will be working on predicting whether a customer will subscribe to a term deposit or not. This project also identifies the biggest predictors of term deposit subscription for the bank.

- *In summary, the classification task primarily aims to assist the financial institution in optimizing its marketing efforts by targeting potential customers who are more likely to subscribe to term deposits, thereby increasing the efficiency of their campaigns and ultimately improving the conversion rate.*

## Data

The raw dataset contains 45,211 records and 17 attributes, organized in a long format. It contains the following:

- **Bank client data:** age (numeric), Job, marital status , education level, whether the client has credit in default, average yearly balance in euros, whether the client has housing loan, or other loans.
- **Data related to the last contact of the current campaign:** contact communication type , last contact day of the month, last contact month of year and last contact duration in seconds.
- **other attributes:** number of contacts performed during this campaign, number of days that passed by after the client was last contacted from a previous campaign, number of contacts performed for the client before this campaign and outcome of the previous marketing campaign
- **Output variable (desired target):** subscription

## Key Data source:

Moro, S., Rita, P., and Cortez, P.. (2012). Bank Marketing. UCI Machine Learning Repository. <https://doi.org/10.24432/C5K306>.

Link: (<https://archive.ics.uci.edu/dataset/222/bank+marketing>)

The dataset is related with direct marketing campaigns (phone calls) of a Portuguese banking institution.

## Data Wrangling

The raw data is more or less clean and doesn't have any duplicate values, however some features like 'previous outcome' have a high number of 'unknown' values which are essentially missing values.

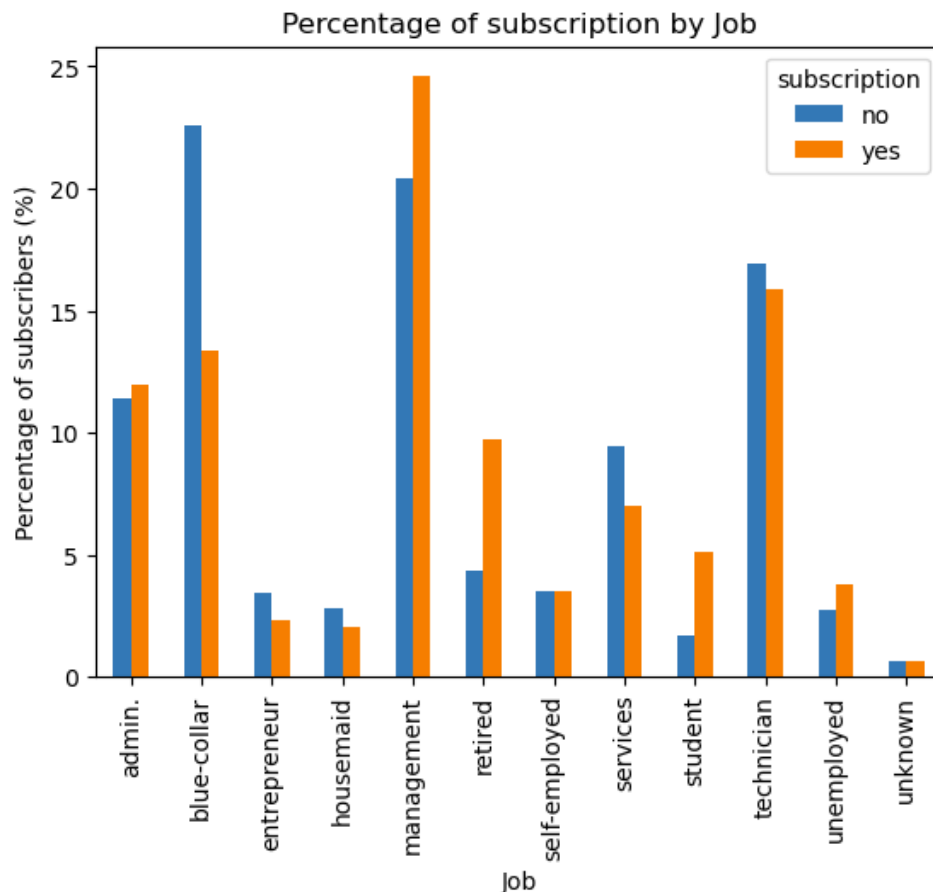
By analyzing the numerical variables' statistics and frequency distribution, I've identified a few potential issues in the dataset such as:

- **out-of-range data or outliers.** '*previous*' representing past contacts, has an outlier with the value of '275' which seems like an error.  
'*balance*' representing customers' balance in their bank account has outliers on both ends.  
'*campaign*' i.e the number of contacts performed during this campaign - shows a few instances with values over 30, possibly due to a longer campaign duration.
- **Category imbalance:** looking into the categorical variables, they appear to be mostly balanced, except '*defaults*', which is highly imbalanced.  
'*job*' has many unique categories, which may need potential grouping before modeling.  
The target feature '*subscription*' is imbalanced, possibly requiring techniques to address this issue later on before modeling.
- **Skewness and low variability** - Variables representing the number of days that passed by after the client was last contacted from a previous campaign, the number of contacts performed during this campaign for each client, and duration of campaign calls are highly skewed and have clustered values towards the lower end.  
The variable representing past contacts, has limited variability, mostly having 0 values. Around 81.7% of clients were never contacted previously making the variable one with the least variability which was eventually removed before modeling.

# Data Exploration

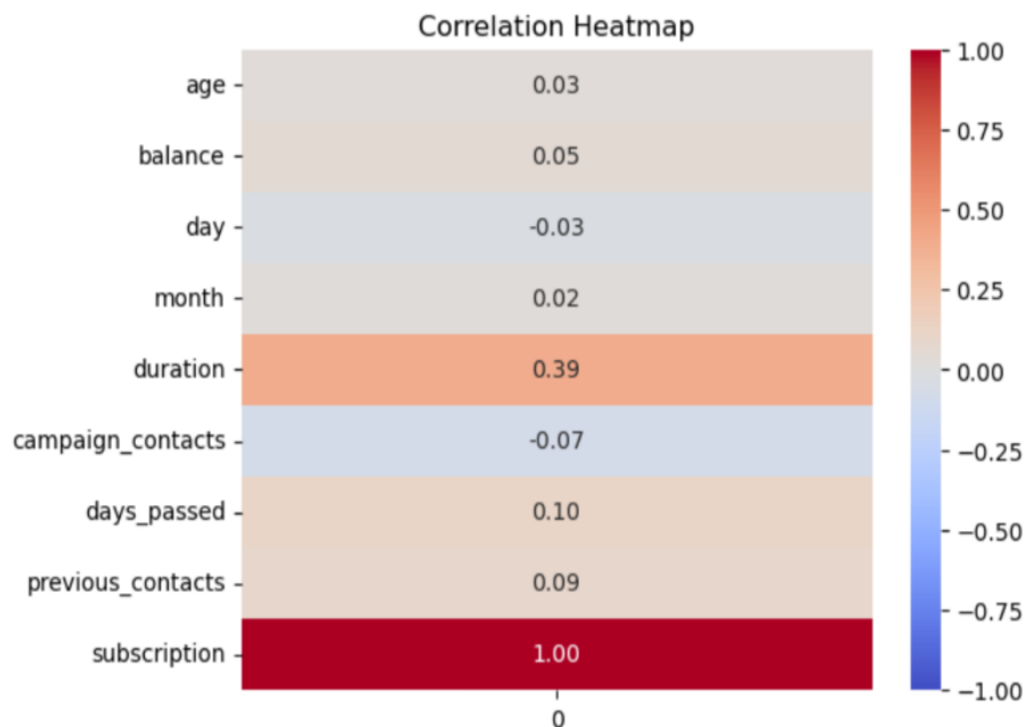
The initial data exploration revealed that there isn't a direct correlation among the independent variables. The following are the relationships that were found between the independent variables and the target.

- 'Age' of the customer has a nonlinear relationship with subscription i.e. customers above the age of 60 and those below 30, are likely to be associated with a higher number of subscriptions.
- 'Job' of the customer - Those in management as well as retirees and students are associated with a higher number of subscriptions. There appears to be a connection between age and job, where older individuals aged above 60 are likely to be retired and those under 30 students, aligning with the expected demographics.



- Customers without existing loans show a higher rate of subscription, aligning with intuitive expectations.

- We can conclude that, overall, longer call **duration** and a moderate number of contacts appear to be linked with a higher number of subscription rates. 'Duration' is the only numerical variable that shows a direct correlation with subscription. *(shown in the heatmap below)*
- Generally, moderately high **balances** (low positives) are associated with higher numbers of subscriptions. However, when plotting **age\_group** against **balance** it's clear that only the age group spanning from the 20s to 50s with a higher balance demonstrate a higher mean subscription
- It also appears that higher number of **contacts** made to clients results in adverse effects on subscription.

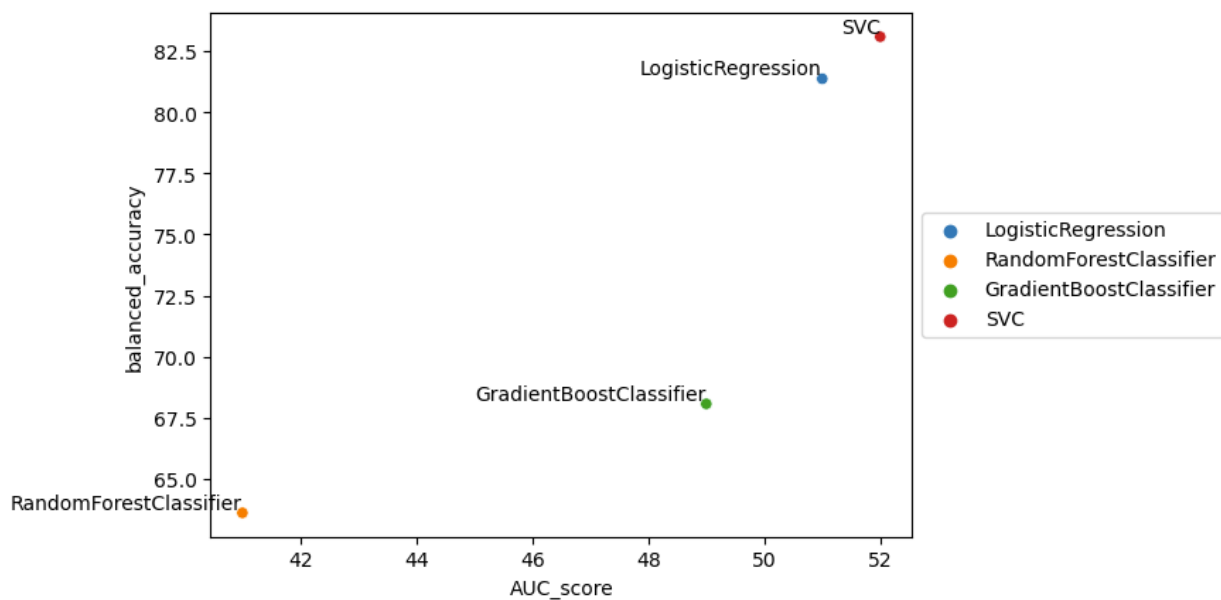


## Model Selection

After fitting 4 different models, initially, the accuracy metric indicated that Gradient Boosting and Random Forest models performed better than others. However, other evaluation metrics like balanced accuracy, confusion matrix, F1 score, and AUC and the high number of False Negatives in these models suggested otherwise.

The Logistic Regression and Support Vector Classifier (SVC) were identified as the better options for the imbalanced dataset due to their performance in handling the

imbalance. Hyperparameter tuning was conducted, which didn't significantly improve model performance.



Model Comparison between Logistic Regression and SVC found that both Logistic Regression and SVC had close balanced accuracy, but SVC showed slightly higher AUC and F1 scores. However, training and prediction with SVC took significantly more time compared to Logistic Regression. Finally, Logistic Regression was selected due to its interpretable coefficients, faster performance and reasonable performance metrics.

## Conclusion

Based on the results of the Logistic regression model - key features impacting the model the most were identified as follows:

- *'previous\_outcome\_success' (positive)*- The coefficient for previous\_outcome\_success is 2.36 and  $e^{2.36} = 10.665727$ . Therefore, those with 'previous\_outcome\_success' have **10.67 times** the odds of subscribing than those who had a different outcome.
- *'contact\_type\_unknown' (negative)* - was shown to have a significant negative impact on subscription - recalling that a very significant portion of the feature contact\_type had the value 'unknown', this association could be considered incorrect. The model might have incorrectly associated this feature's dominance with the target variable, and assume a strong relationship between them only due

to the model's bias towards the majority class within this highly imbalanced feature.

- *'duration' (positive)* - The coefficient of 'duration' is 1.285004 and  $e^{1.285004} = 3.62$ . An increase of 1 second in the phone call duration in the campaign calls multiplies the odds of subscribing to term deposits by **3.62**.
- *'housing loan' (negative)* - The coefficient of 'housing loan' is -0.910418, and  $e^{-0.910418} = 0.4036$ . Those with 'housing loans' have **0.4 times** the odds of subscribing than those who don't.
- Different *job* types had varied effects on subscription, aligning with the exploratory data analysis results - i.e. being a student or retired increased subscription tendency, while being a housemaid, blue-collar worker, or entrepreneur decreased it.
- *Month* and *Day* of contact, had minimal impact on subscription according to the model.

In conclusion, previous campaign's success is the biggest predictor of term deposit subscription followed by duration of campaign calls and housing loans.

## Future Research

Considering the imbalanced dataset, Logistic Regression appeared as the preferred model due to its interpretability, reasonable performance metrics, and quicker training/prediction times compared to SVC. However, further analysis or feature engineering may enhance the model's performance, especially in handling the imbalanced nature of the data and improving predictive capabilities.

To mitigate the issue of having incorrect coefficients of `contact_type_unknown`, we could consider removing this feature from the training set before fitting the model.

The threshold for prediction can be adjusted (lowered) to allow the model to capture more positive instances, potentially reducing missed positive cases which could enable the bank to maximize efforts in targeting more customers leading to a higher potential subscription. Similarly, the threshold can be increased to make the model more conservative in classifying instances as positive, allowing the bank to use less resources in the campaign and focus on a smaller number of customers that are highly likely to subscribe. The decision to adjust this threshold depends on the bank's Subscription goals for future campaigns.