# Term Deposit Subscription Prediction from Telemarketing Calls

Uday Bhaskar Voora

## 1 Introduction

### 1.1 Task

The main task of this project is to develop a predictive model with the capability to reliably predict whether a customer of the bank will accept the offer to subscribe to a term deposit during a direct marketing campaign. This campaign involves contacting clients via phone calls, and our task revolves around leveraging machine learning techniques to analyze a comprehensive dataset. By scrutinizing various client-related attributes and campaign outcomes, our objective is to build an effective model that can assist the bank in making informed decisions about their marketing strategies.

### 1.2 Motivation

The motivation behind this project is to provide banking institutions with the tools and insights required to optimize their direct marketing campaigns. Conventional marketing methods often involve reaching out to a broad audience without considering their inclination to subscribe to a particular service. This can be resource-intensive in terms of both human effort and financial resources. Additionally, it may not yield the desired results, leading to inefficient resource allocation. However, by leveraging predictive modeling and data-driven insights, the central motivation is to develop a model that can accurately predict a client's likelihood to subscribe to a term deposit before initiating the campaign. This predictive capability offers significant advantages, allowing banks to pre-identify clients unlikely to subscribe, enabling more strategic and efficient campaign targeting.

The approach taken in this project offers a multitude of benefits. It notably enhances campaign efficiency by reducing contact with clients unlikely to subscribe, thus conserving resources and maintaining a positive brand image. Furthermore, it enables banking organizations to optimize resource allocation effectively. By focusing human and financial resources on clients statistically more likely to accept term deposit offers, the approach leads to cost-efficiency. Employees can concentrate their efforts where they are most likely to yield results, reducing operational costs and saving valuable time. In conclusion, this project's motivation is to revolutionize direct marketing campaigns for banking organizations by providing a data-driven advantage. Accurately predicting client behavior allows banks to streamline campaigns, conserve resources, and

offer a more personalized and satisfying customer experience, ultimately enhancing efficiency and success in marketing initiatives for the benefit of both banks and clients

# 2   Method

## 2.1   Data Preprocessing

Data preprocessing is a critical step in the project pipeline as it ensures that the dataset is ready for model training. In our analysis of the Portuguese banking dataset, we began by checking for missing values, and fortunately, no missing values were found, ensuring data completeness.

Next, an exploration of the dataset's categorical variables was conducted by analyzing the unique values in each object (string) column. This exploration served a dual purpose: to differentiate between ordinal and nominal variables and to gain a deeper understanding of the categorical data. This knowledge is crucial as it informs the choice of encoding methods for these variables.

The subsequent step focused on binary encoding for particular categorical variables: 'y', 'loan', 'housing', and 'default'. This process involved mapping binary categories to numeric values, specifically 0 and 1. This transformation is fundamental for machine learning algorithms that operate with numerical inputs. For instance, the 'y' variable, which signifies whether a client subscribed to a term deposit ('yes' or 'no'), underwent this conversion, where 'yes' was assigned the value 1 and 'no' the value 0. This transformation was particularly significant as it designated 'y' as the target variable for our predictive model.

Next, Ordinal encoding was performed for the 'education' and 'poutcome' variables, both of which possess inherent order in their categories. This encoding method assigns numerical values based on a predefined order of categories, ensuring that the ordinal relationships between these categories are preserved. For the 'education' variable, we established a logical order, considering 'unknown' as the lowest level, followed by 'primary', 'secondary', and 'tertiary'. Consequently, 'unknown' was encoded as 0, 'primary' as 1, 'secondary' as 2, and 'tertiary' as 3. Similarly, for the 'poutcome' variable, we defined an order by categorizing 'unknown' as a separate category, followed by 'failure', 'other', and 'success'. As a result, 'unknown' received an encoding of 0, 'failure' was assigned 1, 'other' took on 2, and 'success' was represented as 3.

To address the remaining nominal categorical variables ('job', 'marital', 'contact', 'month'), we utilized one-hot encoding. This approach involves creating binary columns for each category within a variable, where the presence of a category is denoted by '1', and its absence is indicated by '0'. One-hot encoding is a vital step as it converts categorical data into a format that machine learning models can effectively process, all while avoiding the introduction of unintended ordinal relationships between categories.

Finally, we standardized the numerical features using MinMax scaling. This scaling technique transforms numerical features to a consistent range, typically between 0 and 1, while maintaining

their relative differences. MinMax scaling plays a crucial role in ensuring that features with varying scales do not disproportionately influence the model training process, allowing models to converge more effectively. The outcome of this process is data that is well-prepared for model training, with all features uniformly treated, ensuring their equal contribution to precise predictions.

## 2.2   Balancing the dataset

Addressing dataset imbalance is a pivotal step in ensuring the effectiveness of a predictive model. In our dataset, the distribution of the target variable 'y' reveals a significant imbalance between clients who subscribed to a term deposit ('1') and those who did not ('0'). The initial distribution demonstrates the presence of a majority class ('0') with 39,922 instances and a minority class ('1') with 5,289 instances. This imbalance can pose challenges to model training and lead to inaccuracies, primarily favoring the majority class. To rectify this imbalance, a two-fold resampling strategy is adopted using the imbalanced-learn library. First, the focus is on the minority class ('1'), and the Synthetic Minority Over-sampling Technique (SMOTE) is employed. SMOTE works by generating synthetic samples that bridge the gap between existing instances of the minority class. In this case, this technique augments the size of class '1' from its original 5,289 instances to a balanced count of 10,000. This step ensures that the model is exposed to sufficient data representing clients who subscribed to term deposits.

In the subsequent step, the majority class ('0') is addressed using the Random UnderSampler. This method selectively reduces the size of the majority class while preserving the total number of instances for class '1'. Class '0' is downsized from the initial 39,922 instances to 15,000. This strategic undersampling ensures that the majority class does not dominate the dataset, and the model is equally attentive to both class '0' (non-subscribers) and class '1' (subscribers). This balanced dataset sets the stage for the model to make accurate predictions and provide valuable insights for the direct marketing campaign. Balancing the dataset is a pivotal step in leveling the playing field and ensuring that the predictive model performs optimally.

## 2.3   Data Analysis

From the figure1, the correlation matrix for the UCI Bank Marketing dataset visualizes the linear relationships between different variables involved in the direct marketing campaigns of a Portuguese banking institution. In this matrix, 'age' and 'balance' show a slight positive correlation, suggesting that older clients might have higher balances. The number of days that passed by after the client was last contacted from a previous campaign ('pdays') and the number of contacts performed before this campaign ('previous') have a moderate positive correlation, indicating that clients contacted previously are likely to be contacted again. Most other variables show little to no linear correlation with each other, suggesting that they contribute independently to the outcome of a client subscribing to a term deposit.
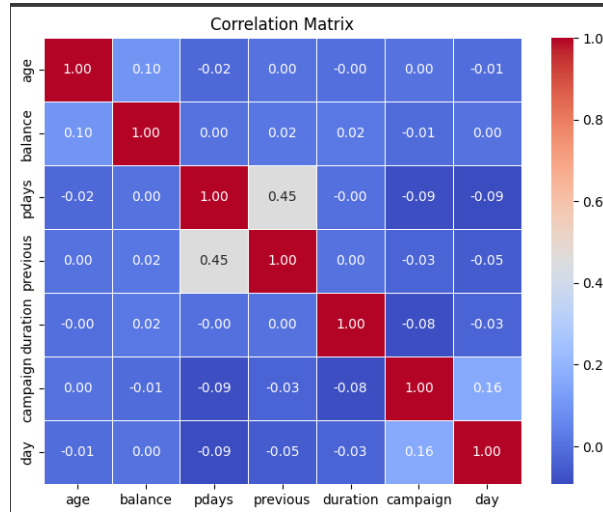
Figure 1: Correlation matrix.

## 2.4 Algorithm Selection

**Logistic Regression** is a foundational algorithm used for binary classification tasks, making it an appropriate choice for this project. The objective is to predict whether a client will subscribe (1) or not (0) to a term deposit, which fits perfectly within the scope of binary classification. Logistic Regression assumes a linear relationship between the input features (such as client's age, job type, balance, etc.) and the target variable (subscription status). Given that linear models are interpretable and computationally efficient, Logistic Regression can quickly provide a baseline prediction, making it a useful starting point for this classification problem.

**Decision Trees:** Decision Trees are versatile and intuitive models that split the data into decision nodes based on feature values. They can model both linear and nonlinear relationships effectively. This makes them particularly useful in the Bank Marketing project, where the relationship between customer features and subscription decisions is not necessarily linear. Decision Trees also allow us to visually interpret the decision-making process, which is a great advantage when explaining the model's behavior. The ability to handle both numerical and categorical features makes Decision Trees a good fit for this project, where there is a mix of such features.

**Random Forest** is an ensemble technique that builds multiple Decision Trees and combines their results to improve predictive accuracy and reduce the likelihood of overfitting. This approach is well-suited to the Bank Marketing problem because the dataset contains a wide range of features with complex interactions. Random Forest can effectively capture these interactions and provide more robust predictions than a single Decision Tree, making it ideal for datasets with noisy or highly variable features. By averaging multiple trees, it reduces the variance associated with individual Decision Trees, thus improving the generalization power of the model.

**Support Vector Machines (SVM)** is a powerful and highly flexible algorithm, particularly when the data is not linearly separable. It excels in high-dimensional spaces and can effectively handle datasets with numerous features. In this project, the customer data is multi-

dimensional, and SVM can capture intricate decision boundaries between the "subscribing" and "non-subscribing" classes. By using kernel functions, SVM can map data into higher-dimensional spaces to find optimal separating hyperplanes, making it a robust choice for classification problems with complex decision boundaries.

**Neural Networks** are highly flexible models that can learn intricate and nonlinear relationships in data. While they are computationally more intensive, they excel in handling complex patterns in large datasets. In this project, the interactions between client attributes and their subscription decisions are likely intricate, and Neural Networks provide the flexibility to capture these complex patterns. They are particularly useful for modeling highly nonlinear relationships between the features and the target variable, enabling the model to learn deep representations of the data. Neural Networks also offer scalability, allowing them to be expanded to larger datasets or more complex problems if needed.

**TabNet** is a deep learning model designed specifically for tabular data. It is known for its ability to handle structured data effectively and efficiently. Unlike traditional neural networks, TabNet uses attention mechanisms, which help the model focus on the most important features during training. This allows TabNet to be highly effective in understanding feature importance and capturing complex relationships in the data. The ability to interpret the model's decisions through feature selection is a key advantage for TabNet.

TabNet was chosen because it can capture complex, nonlinear relationships in structured data, making it an excellent choice for predicting customer subscription behavior, where interactions between features like age, job, and balance can be intricate.

**XGBoost** (Extreme Gradient Boosting) is one of the most popular and powerful gradient boosting frameworks, known for its high performance and speed in solving classification problems. XGBoost combines the predictions of many decision trees to improve model performance while reducing the risk of overfitting. It is particularly effective in handling large datasets with complex relationships.

**XGBoost** was chosen for this project due to its robustness in handling a variety of data types and its ability to handle missing values, which makes it ideal for datasets with lots of features and potentially noisy data. Additionally, its built-in regularization and boosting mechanism can help improve prediction accuracy by iterating on the mistakes of the previous trees.

**LightGBM** (Light Gradient Boosting Machine) is another popular gradient boosting framework, which is designed to be faster and more efficient than traditional gradient boosting methods. LightGBM works well on large datasets and is known for its efficiency with memory and computation. Like XGBoost, it combines multiple decision trees to make predictions, but it uses a histogram-based method to find the best split for features, making it faster.

LightGBM is well-suited for this project because of its speed, scalability, and ability to handle large amounts of data effectively. Its ability to deal with categorical features natively also makes it particularly suitable for tabular datasets like the one in your project, where categorical features (like job type or marital status) are common.

## 2.5 Model Training

**Logistic Regression** was utilized for predicting customer subscription decisions, with an emphasis on hyperparameter optimization using `GridSearchCV`. The hyperparameter grid explored

included several parameters to optimize the model. Specifically, **C** (with values 0.1, 1, and 10) was tested for regularization strength, `penalty` (either `'l1'` or `'l2'`) was considered for regularization norms, and `solver` (`'liblinear'`) was used for the optimization algorithm. Additionally, `max_iter` values (100, 500, and 1000) were examined to determine the number of iterations needed for convergence.

The optimal parameters identified from this grid search were **C** = 1, `max_iter` = 100, `penalty` = `'l1'`, and `solver` = `'liblinear'`. This configuration, with the L1 penalty, encourages a model with fewer non-zero coefficients, thereby enhancing the simplicity and interpretability of the results. The regularization strength parameter **C** = 1 strikes an ideal balance, preventing overfitting while maintaining the flexibility of the model. The choice of `liblinear` as the solver is particularly suitable for binary classification problems, such as the one used here, due to its efficient optimization process for smaller datasets.

The optimal parameters identified were:

$$C = 1, \quad \text{max\_iter} = 100, \quad \text{penalty} = \text{'l1'}, \quad \text{solver} = \text{'liblinear'}.$$

This configuration, with an `L1` penalty, promotes a model with fewer non-zero coefficients, enhancing simplicity and interpretability. The value of **C** = 1 strikes a balance in regularization, preventing overfitting while maintaining model flexibility. The `liblinear` solver aligns well with the binary classification nature of the dataset, ensuring efficient optimization.

**Decision Tree** Classifier's optimal hyperparameters were determined through `GridSearchCV`. The best parameters identified for the model were `criterion = 'entropy'`, `max_depth = 10`, `min_samples_leaf = 2`, `min_samples_split = 2`, and `max_features = 'sqrt'`. The choice of `entropy` as the criterion indicates a preference for maximizing the information gain at each split. This enhances the classifier's ability to handle complex patterns in the data by selecting the most informative features at each stage. A `max_depth` of 10 ensures that the tree is sufficiently deep to capture significant data relationships while preventing overfitting by limiting the depth of the tree. The `min_samples_leaf` and `min_samples_split` values set at 2 ensure that splits are justified by a sufficient number of samples, avoiding overly fine splits that could lead to overfitting. Finally, the use of `max_features = 'sqrt'` focuses the model's attention on a subset of features at each split, which helps reduce overfitting and makes the model more robust, particularly in high-dimensional datasets.

**Random Forest:** A `RandomForestClassifier` was trained for predicting customer subscription to a deposit. The model utilized a hyperparameter grid that included several key parameters: `n_estimators` (50, 100, 150), `criterion` (`'gini'`, `'entropy'`), `max_depth` (None, 10, 20), `min_samples_split` (2, 5), `min_samples_leaf` (1, 2), and `max_features` (`'auto'`, `'sqrt'`). After applying `GridSearchCV`, the optimal parameters identified were: `criterion = 'entropy'`, `max_depth = 20`, `max_features = 'auto'`, `min_samples_leaf = 1`, `min_samples_split = 2`, and `n_estimators = 50`. This configuration indicates a preference for deeper trees (`max_depth = 20`) and a larger number of trees (`n_estimators = 50`) to improve the model's ability to capture comprehensive data patterns. The choice of `entropy` as the criterion ensures that the model maximizes information gain at each split, leading to better classification decisions. Setting `min_samples_leaf` and `min_samples_split` to minimal values of 1 and 2, respectively, allows

the model to capture detailed subtleties in the data while avoiding overly fine splits that may lead to overfitting. Lastly, using `auto` for `max_features` helps the model consider all features for each split, enhancing its ability to generalize and accurately predict customer behavior.

**Support Vector Machine (SVM):** The Support Vector Machine (SVM) Classifier's training process involved fine-tuning hyperparameters using `GridSearchCV`. The hyperparameter grid included variations in **C** (0.1, 1, 10), controlling the trade-off between a smooth decision boundary and correct classification; **kernel** types (`'linear'`, `'poly'`, `'rbf'`, `'sigmoid'`), determining the hyperplane type for data separation; and **degree** (2, 3, 4) for the polynomial kernel. The optimal configuration obtained was **C: 10**, **kernel: poly**, and **degree: 2**. This setup indicates a strong emphasis on accurate classification (**C: 10**), coupled with a polynomial kernel of degree 2, allowing the model to capture complex, nonlinear decision boundaries.

**Neural Network** used for classifying customer subscription decisions was thoughtfully architected for optimal performance. The Sequential model comprised an input layer with 64 neurons (ReLU activation), two hidden layers each with 32 neurons (ReLU activation), and a sigmoid-activated output layer. Batch normalization followed the first two dense layers to enhance training stability, while dropout layers (0.3 rate) were incorporated to prevent overfitting. The model utilized the RMSprop optimizer, known for its adaptive learning rate, set initially at 0.001. A significant feature was the inclusion of the ReduceLROnPlateau callback, which intelligently adjusted the learning rate based on validation loss performance, reducing it by a factor of 0.2 if no improvement was observed for five consecutive epochs. This mechanism ensured efficient and effective learning, minimizing overfitting and enhancing the model's ability to generalize.

**XGBoost** was selected for predicting customer subscription decisions, with the hyperparameters optimized using `GridSearchCV`. The hyperparameter grid for XGBoost included variations in the number of boosting rounds (`n_estimators`), tree depth (`max_depth`), learning rate (`learning_rate`), as well as the fraction of data used for each base learner (`subsample`) and the fraction of features used for each boosting round (`colsample_bytree`). The optimized parameters found were `n_estimators`: 100, `max_depth`: 5, `learning_rate`: 0.1, `subsample`: 1.0, and `colsample_bytree`: 1.0, which offered a good balance between model complexity and learning speed. This setup allowed the model to generalize well to unseen data while maintaining robustness in learning. The model performed exceptionally well, achieving high accuracy and F1-score values, showing its strength in capturing patterns and predicting customer subscription behavior.

**LightGBM,** another gradient boosting model, was used to predict customer subscription decisions. The hyperparameter tuning process, conducted with `RandomizedSearchCV`, focused on parameters such as the number of trees (`n_estimators`), the number of leaves in each tree (`num_leaves`), learning rate (`learning_rate`), and tree depth (`max_depth`). The tuned hyperparameters, which included `num_leaves`: 31, `n_estimators`: 500, `learning_rate`: 0.05, and `max_depth`: -1, resulted in a highly effective model. The combination of these parameters enabled LightGBM to handle complex data patterns, while the `learning_rate` ensured a balanced learning pace without overfitting. The model's performance on the test set was impressive, with

`accuracy` and `F1-score` both being significantly high, making it an ideal choice for this task.

**TabNet,** a neural network architecture tailored for tabular data, was employed for predicting customer subscription decisions. Using `RandomizedSearchCV`, the hyperparameters of the model were tuned, including the number of decision steps (`n_d`), attention heads (`n_a`), and the number of steps (`n_steps`), along with parameters controlling sparsity (`gamma` and `lambda_sparse`). The optimal configuration, which included `n_d`: 16, `n_a`: 16, `n_steps`: 5, `gamma`: 1.3, and `lambda_sparse`: 0.001, helped strike a balance between model complexity and sparsity. The tuned TabNet model showed strong performance in terms of `accuracy` and `ROC-AUC`, demonstrating its ability to generalize well and capture the intricate patterns within the data. The combination of these parameters allowed TabNet to effectively learn from the data while ensuring regularization to prevent overfitting.

# 3   Results and Discussion

In the context of predicting customer subscription decisions for a bank's deposit campaign, the key metrics to evaluate are precision and recall, as these determine how well the model can identify potential subscribers and minimize false positives and negatives.

**Precision** is critical in ensuring that the bank's marketing resources are effectively used by targeting clients who are genuinely interested in subscribing to the deposit. It measures how many of the predicted positive cases are actually correct.

**Recall** ensures that the model is not missing any potential subscribers. A high recall rate indicates that the model correctly identifies a large proportion of actual subscribers.

Table 1: Model Results

|                | Accuracy | Precision | Recall   |
| -------------- | -------- | --------- | -------- |
| Neural Network | 0.924375 | 0.923580  | 0.924375 |
| Random Forest  | 0.937750 | 0.938706  | 0.937750 |
| SVM            | 0.923375 | 0.923070  | 0.923375 |
| XGBoost        | 0.948875 | 0.948514  | 0.948875 |
| LightGBM       | 0.951375 | 0.951030  | 0.951375 |
| Tabnet         | 0.90     | 0.91      | 0.97     |

From the results shown in **Table 1**, **LightGBM** and **XGBoost** have the highest performance across all models with respect to both precision and recall. LightGBM achieves an accuracy of **0.951375**, with precision and recall values of **0.951030** and **0.951375**, respectively. Similarly, XGBoost shows high precision (**0.948514**) and recall (**0.948875**), making it another strong choice for predicting customer subscription.

The **Neural Network** model, while performing well on recall (**0.924375**), suffers from lower precision (**0.923580**), which means it tends to make more false positive predictions. This is critical to consider because a lower precision rate may mean targeting too many clients who are not actually interested in subscribing, potentially wasting marketing resources.

**Random Forest**, which achieved an accuracy of **0.937750**, strikes a good balance between precision (**0.938706**) and recall (**0.937750**), making it a well-rounded choice, especially when accuracy and the avoidance of overfitting are priorities.

**SVM**, with an accuracy of **0.923375**, has decent precision (**0.923070**) but lower recall (**0.923375**) than the other models, indicating that while the model predicts positives accurately, it may miss a number of actual subscribers.

In conclusion, **LightGBM** and **XGBoost** emerged as the top-performing models, excelling in both precision and recall, making them the most suitable for predicting customer subscription to the bank's deposit campaign.
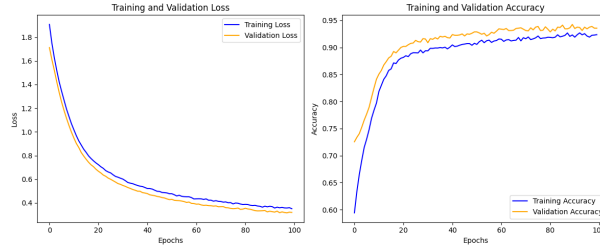


Figure 2: The Neural Network train-test curves for loss and accuracy

The neural network's performance, as indicated by the learning curves in Figure 3, demonstrates a stable and effective training process. The training and validation loss curves show a steady decline, suggesting that the model is successfully minimizing the loss over time. The accuracy curves for both the training and validation sets increase steadily, with the validation accuracy closely tracking the training accuracy. This indicates that the model is generalizing well, with minimal overfitting or underfitting.

Given the focus of this project on identifying true positive instances, the model's ability to maintain a strong performance across both training and validation sets is noteworthy. The model has effectively learned the relevant patterns without being influenced by noise, as seen in the consistent rise of accuracy across both datasets.

Overall, the neural network exhibits a strong generalization ability, which is crucial for ensuring that the model performs reliably in real-world scenarios where missing true positive instances could have significant consequences.

# 4 Conclusion

In this project, the goal was to develop a robust predictive model to assist banks in making informed decisions regarding customer subscription to deposit services. A comparison of various machine learning models was undertaken, each offering its own set of advantages and trade-offs. Among these models, LightGBM emerged as a top contender due to its exceptional performance in terms of accuracy and recall. With an accuracy score of 0.95 and recall of 0.97, LightGBM demonstrated an outstanding ability to correctly identify potential subscribers, minimizing the risk of false negatives and ensuring that the bank's marketing efforts are maximized.
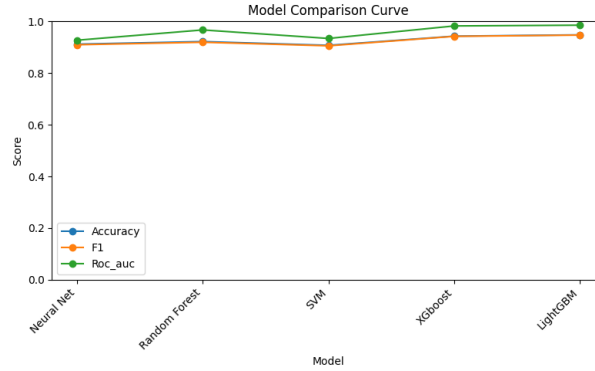
Figure 3: Models comparision curve

The precision-recall balance of LightGBM aligns perfectly with the project's objectives and the banking industry's need to optimize marketing campaigns. By accurately identifying subscribers with high recall, LightGBM ensures that the bank reaches as many potential customers as possible while maintaining precision in its predictions. Compared to other models, Light-GBM consistently demonstrated strong performance in both accuracy and recall, making it the preferred choice for capturing as many true positive cases as possible.

Deploying this model in real-world banking scenarios holds great promise for improving campaign efficiency, customer engagement, and ultimately, the bank's bottom line. In conclusion, LightGBM has proven to be an ideal model for optimizing customer subscription predictions in the banking sector. Its high recall ensures that potential subscribers are effectively identified, while its overall accuracy ensures reliable predictions. By deploying LightGBM, banks can significantly enhance the efficiency and effectiveness of their marketing campaigns, offering a powerful tool for improving customer engagement and driving business success.

Throughout this project, I explored various machine learning algorithms, including Logistic Regression, Decision Trees, Random Forests, Support Vector Machines (SVM), XGBoost, and LightGBM. I gained a deep understanding of the specific hyperparameters associated with each algorithm, such as n_estimators, learning_rate, and max_depth, and learned how these parameters influence model behavior. Additionally, I grasped the significance of precision and recall as evaluation metrics. Precision measures the accuracy of positive predictions, while recall assesses the model's ability to capture all positive cases. Balancing these metrics is crucial, and this project equipped me with valuable insights into optimizing model performance while considering business objectives and priorities.

# References

[1] Yildirim, P., Asci, S. (2016). Predictive Data Mining for Targeted Marketing Campaigns: Practical Approach. *Expert Systems with Applications*, 55, 210-219.

[2] Moro, S., Cortez, P., Rita, P. (2014). A Data-Driven Approach to Predict the Success of Bank Telemarketing. *Decision Support Systems*, 62, 22-31.

[3] Silva, D. F., Delgado, M., Soares, C. (2015). Analyzing the Behavior of Classification Algorithms in Imbalanced Settings. *Knowledge and Information Systems*, 45(1), 247-270.