

# R- Classification

*Project Report -*

## Predicting Customer Purchase Behavior in Bank Marketing

*Wai Ying Wong, Winnie*

Introduction:.....	4
1. Discovery .....	4
1.1. Business Domain and Problem Area .....	4
1.2. Developing Initial Hypotheses.....	4
1.3. Dataset: .....	5
2. Data Planning .....	5
2.1. Data Inventory.....	5
2.2. Exploratory Data Analysis.....	6
2.3. Data Processing and Transformation .....	7
Removing missing data:.....	7
Target Variable Transformation: .....	8
Feature Selection:.....	8
Feature Removal: .....	8
Data Binning: .....	8
Categorical Data Handling: .....	9
Scaling and Special Consideration: .....	9
Balancing the Dataset:.....	9
Dataset Split for training and testing:.....	9
Data shuffling: .....	9
2.4 Feature Observation & Hypothesis .....	10
2.4.1. Preliminary analysis: Correlation Analysis .....	10
3. Data Analytic Model Planning and Implementation .....	11
3.1. Logistic Regression .....	11
3.2. Other Classification Techniques .....	12
4. Results Interpretation and Business Implications .....	13
4.1.a. Logistic regression model Result .....	13
Individual Coefficients: .....	14
▫ Continuous variables: .....	14
Significance: .....	14
Sign: .....	14
Size:.....	14

▫ Categorical variables:.....	15
Significance: .....	15
Sign: .....	15
Size:.....	15
4.1.b. Business implications of the results .....	16
4.2.a. Classification regression model Result .....	17
Confusion Matrix.....	17
4.1.b. Business implications of the results .....	18
5. Application of the Model for Predictive Analysis .....	19
5.1 Out Of Sample Prediction .....	19
6. Conclusion .....	20
7. Challenging Factors .....	21
7.1. Dataset Quality and Completeness .....	21
7.2. Imbalanced Data Distribution.....	21
7.3. Skillful Feature Selection and Engineering .....	21
7.4. Constrained Model Selection .....	21
7.5 Limited by the Model Assumptions.....	22
7.6 Data Extrapolation.....	22

## Introduction:

This project aims to revolutionize customer purchase behavior prediction in the bank marketing sector. In the retail banking industry, blind marketing techniques like cold calling have been widely employed. However, the effectiveness of such approaches is inherently limited. To overcome this limitation, our objective is to harness the power of data analytics techniques to develop highly accurate predictive models by R. These models will enable us to identify customers who are most likely to respond positively to marketing campaigns and subscribe to the bank products, term deposit in this study. Furthermore, the developed classification model can serve as a means to measure campaign effectiveness and evaluate employee performance. By leveraging these insights, we can optimize our marketing efforts, refine customer targeting strategies, and ultimately enhance the overall response rate.

## 1. Discovery

### 1.1. Business Domain and Problem Area

The problem area we wish to work on is improving the effectiveness of blind marketing strategies in the retail banking industry. Through the analysis of historical customer data, we aim to segment the customers and develop a model that can predict the likelihood of positive responses from customers towards marketing campaigns, with the aim of optimizing marketing campaigns and increasing conversion rates. The significance of this initiative becomes evident when considering limited customer interest observed in our bank marketing dataset, where only approximately 12% of customers expressed interest in the bank's offers. This indicates a substantial gap between the marketing efforts and the customers' receptiveness to those efforts. By targeting the right customers, we can significantly improve campaign effectiveness and drive business growth.

### 1.2. Developing Initial Hypotheses

This hypothesis seeks to explore the intricate relationship between several key factors, including customer demographics, portfolio characteristics, the bank's engagement approach, social and economic factors and their influence on customers' purchase behavior in bank marketing.

We propose the null hypothesis that there is no statistically significant relationship between these variables. In other words, it is assumed that customer demographics, portfolio characteristics, and the engagement approach employed by the bank, social and economic factors do not have a substantial impact on customers' response rate.

### 1.3. Dataset:

For this project, we will utilize a publicly shared bank marketing dataset from UCI Machine Learning Repository (<http://archive.ics.uci.edu/dataset/222/bank+marketing>). This dataset describes the direct marketing campaigns (phone calls) of a Portuguese banking institution. These campaigns conducted were based mostly on direct phone calls and offering bank customers the chance to purchase a term deposit. Often, more than one call needs to be made to a single customer before they either decline or agree to a term deposit subscription. After the marketing campaigns if the customer had agreed to place deposit - target variable marked 'yes', otherwise 'no'.

There are four datasets in the file:

1. bank-full.csv with all examples, ordered by date (from May 2008 to November 2010).
2. bank.csv with 10% of the examples (4521), randomly selected from bank-full.csv.
3. bank-additional-full.csv with all examples, ordered by date (from May 2008 to November 2010).
4. bank-additional.csv with 10% of the examples (4119), randomly selected from bank-additional-full.csv.

In this project, we are going to use 'bank-additional-full.csv' which has 21 columns, and 41189 rows for the analysis. The dataset contains valuable information about past marketing campaigns, customer demographics, economic indicators, and other relevant variables. We consider the variable "y" (yes or no) as the target domain, which indicates customer has subscribed a term deposit.

## 2. Data Planning

### 2.1. Data Inventory

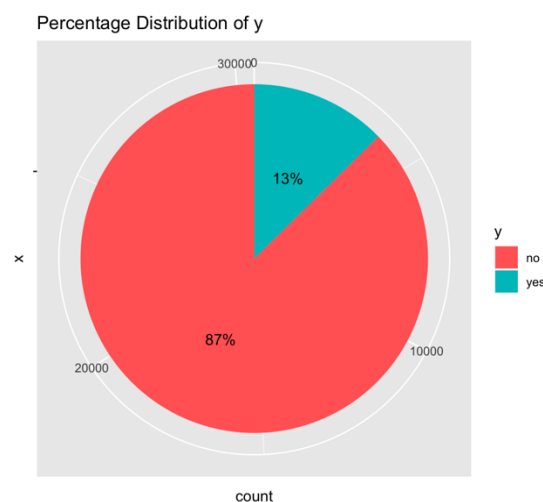
#### Summary of the variables in the dataset

Column name	Description	Type
Age	Age	Numeric
Job	Type of job	Categorical
Marital	Marital status	Categorical
Education	Level of education	Categorical
Default	Customer has credit in default?	Categorical
Housing	Customer has housing loan?	Categorical
Loan	Customer has personal loan?	Categorical
Contact	Contact communication type	Categorical
Month	Last contact month of year	Categorical
Day_of_week	Last contact day of the week	Categorical

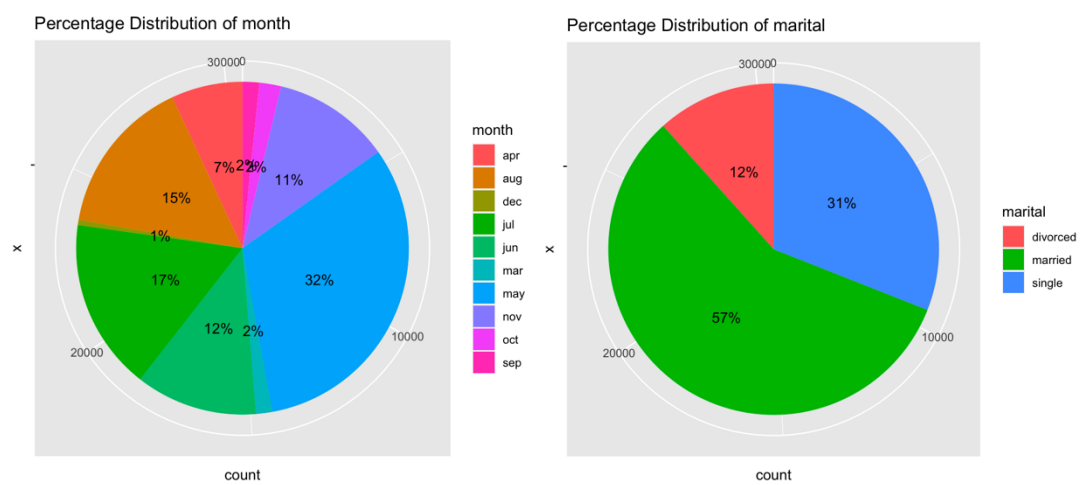
<i>Duration</i>	Last contact duration, in seconds	Numeric
<i>Campaign</i>	Number of contacts performed during this campaign and for this client	Numeric
<i>Pdays</i>	Number of days that passed by after the client was last contacted from a previous campaign (999 means client was not previously contacted)	Numeric
<i>Previous</i>	Number of contacts performed before this campaign and for this client	Numeric
<i>Poutcome</i>	Outcome of the previous marketing campaign	Categorical
<i>Emp.Var.Rate</i>	Employment variation rate - quarterly indicator	Numeric
<i>Cons.Price.Idx</i>	Consumer price index - monthly indicator	Numeric
<i>Cons.Conf.Idx</i>	Consumer confidence index - monthly indicator	Numeric
<i>Euribor3m</i>	Euribor 3-month rate - daily indicator (The euribor is the euro interbank offered rate. It generally refers to the price at which european banks lend money to each other. In the same way that people and businesses borrow money from banks, when banks need money, they borrow from other banks for which they pay interest.)	Numeric
<i>Nr.Employed</i>	Number of employees - quarterly indicator	Numeric
<i>Y</i>	Has the client subscribed to a term deposit?	Categorical-binary

## 2.2. Exploratory Data Analysis

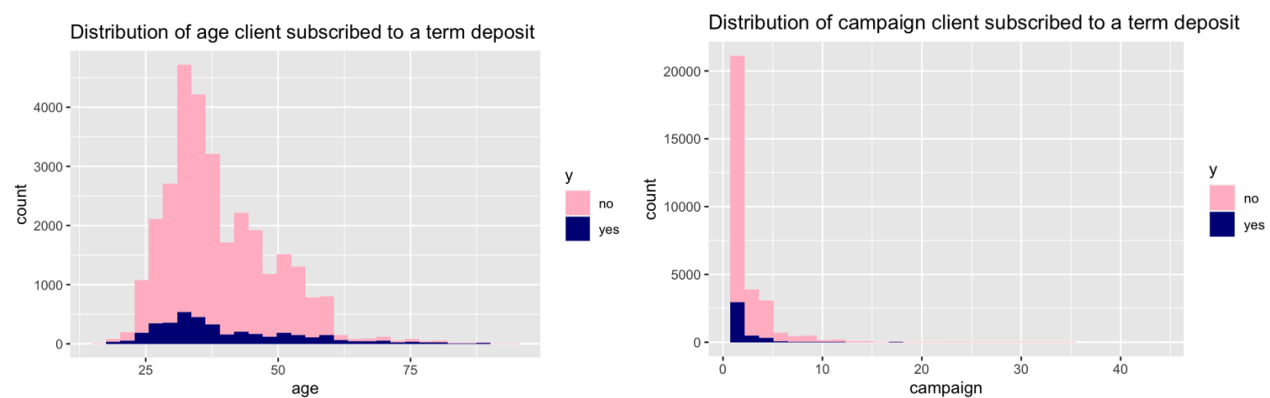
This pie chart provides insights into the distribution of our target variable, 'y' (subscribed to a term deposit). Notably, the majority, comprising 87% of the total, falls under the 'No'. In contrast, the 'Yes' category constitutes a smaller proportion, accounting for only 13%.



The pie chart on the left reveals an uneven distribution of the marketing campaign timeframe, with May accounting for 32%, while certain months like Jan and Feb saw no promotion. On the right, the marital distribution among customers indicates that over half are married, around a third are single, and only 12% are divorced.



The left histograms reveal that most customers are aged 25 to 40, with a notable portion of term deposit subscribers falling between 25 to 35. The right plot indicates that over 20,000 customers were approached for the first marketing campaign, and this group generated a significant number of positive responses.



### 2.3. Data Processing and Transformation

Our dataset underwent rigorous data cleaning and transformation to ensure the integrity of the subsequent analysis and prevent any misleading outcomes that might arise from incomplete or ambiguous information.

To prepare the dataset for analysis, we performed several key manipulations:

#### Removing missing data:

- A complete dataset without blank cells has been verified. We remove all "unknown" labeled data points in order to improve accuracy by eliminating uncertainties. The result is a remaining total of 30,488 data points after removal.

#### Target Variable Transformation:

- The binary target variable "Y" was transformed into binary values, specifically "1" and "0." This adjustment was necessary to align the target variable with the modeling requirements, facilitating accurate predictions.

#### Feature Selection:

- We use *aggregate()* function to calculate the mean age for each job category with the summary below.
- We exclude the "age" variable from our analysis. It was motivated by concerns about potential multicollinearity effects with the "job" variable. For example, the "retired" category typically involves ages above 60.

<i>Job</i>	<i>Age (mean)</i>
<i>student</i>	25.36557
<i>services</i>	36.76094
<i>admin</i>	37.54435
<i>technician</i>	37.60278
<i>blue-collar</i>	38.16423
<i>self-employed</i>	38.82143
<i>entrepreneur</i>	40.89624
<i>management</i>	41.52402
<i>housemaid</i>	44.67536
<i>retired</i>	63.19984

#### Feature Removal:

- The "default" variable, with only three positive records, was removed from the dataset. Given its limited impact on analysis, this variable's omission contributed to a more focused and effective study.

#### Data Binning:

- To enhance the interpretability of our analysis, we grouped similar categories under column- "Education" together. Notably, "basic.4year," "basic.6year," and "basic.9year" were consolidated into a single category labeled "basic". This binning strategy reduces data granularity and enables the identification of patterns that might remain obscured when dealing with individual data points.



### Categorical Data Handling:

- Categorical data were appropriately encoded using dummy variables by one-hot encoding function. It aims to transform the categorical variables into a binary format, where each unique category becomes a separate binary column (dummy variable). Each data point is assigned a '1', representing presence of the category, in the corresponding category column and '0', representing absence of the category, in the others. This approach enables statistical models to interpret categorical information, ensuring accurate and meaningful analysis.

### Scaling and Special Consideration:

- We scale the numerical data. Due to significant variations in numerical variable magnitudes, data scaling is applied to normalize their ranges so as to enhance the interpretability of regression coefficients, aids stable convergence in training, and improves overall predictive performance in both types of models.
- However, the dataset introduced a unique challenge in the form of the "pday" column. This column indicates the number of days since the client's last contact, with values ranging from 1 to 4 and a special value, 999, signifying that no such contact occurred. Given this distinctive feature and the presence of the special value 999, traditional scaling techniques are inapplicable to this context.

### Balancing the Dataset:

- Recognizing the imbalance in the distribution of positive and negative records within the targeted variable, we took steps to address this issue. With only around 15% (3859 out of 26629) of records classified as positive after some data processing, a potential bias could arise. To counteract this, we employed a random undersampling technique. Specifically, we randomly extracted an equal number of negative records (3859) to achieve a balanced representation in the dataset. This approach not only rectifies the imbalance but also ensures that the resulting dataset maintains a randomized order, preventing any unintentional sequencing effect.

### Dataset Split for training and testing:

- We partition the transformed data, allocating 70% for training and reserving 30% for testing purpose.

### Data shuffling:

- Data shuffling is employed for sampling training and testing data to mitigate any possible pre-formatting of the dataset (e.g. sorted dataset, time-series data, ...). Training and testing split allows out-of sample testing for better model evaluation that simulates real-world forecasting and mitigates overfitting.

## 2.4 Feature Observation & Hypothesis

### 2.4.1. Preliminary analysis: Correlation Analysis

Before data modeling, we initiate correlation analysis, delving into potential relationships between independent variables and the dependent variable "Y" (subscribed term deposit or not) and the direction. It serves to uncover initial insights and guide our subsequent modeling decisions.

The table provided below enumerates the dependent variables that exhibit a positive relationship with the target variable "Y" individually, with a significant p-value of less than 0.025 ( $\Pr(>|z|)$ ). Notably, the correlation coefficient between "duration" and "Y" was determined to be  $r = 0.4566$ , indicating a moderate positive linear correlation. In essence, as the duration of the last contact increases, there is an inclination for an enhanced chance of subscribing to a term deposit.

<b>Variable</b>	<b>Correlation Coefficient (r)</b>	<b>P-value</b>
<i>duration</i>	0.4566	0
<i>poutcome_success</i>	0.3027	0
<i>previous</i>	0.2547	0
<i>contact_cellular</i>	0.2361	0
<i>month_mar</i>	0.1498	0
<i>month_oct</i>	0.1486	0
<i>month_sep</i>	0.1264	0
<i>job_retired</i>	0.1095	0
<i>job_student</i>	0.1012	0
<i>cons_conf_idx</i>	0.0853	0
<i>marital_single</i>	0.0707	0
<i>month_dec</i>	0.07	0
<i>education_university_degree</i>	0.0647	0
<i>poutcome_failure</i>	0.0342	0.0026
<i>job_unemployed</i>	0.03	0.0084

The table below lists the dependent variables that individually exhibit a negative relationship with the target variable 'Y,' each accompanied by a significant p-value of less than 0.025. Particularly noteworthy, the correlation coefficient between 'nr\_employed' and 'Y' was calculated as  $r = -0.4692$ , while the correlation coefficient between 'euribor3m' and 'Y' was found to be  $r = -0.4475$ . These findings denote a moderate negative linear correlation. In essence, as the number of employees at the bank ('nr\_employed') and the 3-month term deposit rate ('euribor3m') increase, there is an associated tendency for a decrease in the likelihood of subscribing to a term deposit.

<b>Variable</b>	<b>Correlation Coefficient</b>	<b>P-value</b>
<i>nr_employed</i>	-0.4692	0
<i>euribor3m</i>	-0.4475	0
<i>emp_var_rate</i>	-0.4286	0
<i>pdays</i>	-0.3102	0
<i>cons_price_idx</i>	-0.1799	0
<i>month_may</i>	-0.1766	0
<i>campaign</i>	-0.121	0
<i>job_blue-collar</i>	-0.0984	0
<i>month_jul</i>	-0.0616	0
<i>marital_married</i>	-0.0585	0
<i>job_services</i>	-0.0518	0
<i>day_of_week_mon</i>	-0.039	0.0006
<i>month_nov</i>	-0.0358	0.0016
<i>job_technician</i>	-0.0285	0.0122

Moreover, correlation analysis assumes that the relationship between the variables being analyzed is approximately linear. It may not accurately capture nonlinear associations. In other words, If the relationship between the variable and the target in your logistic regression is non-linear, it might lead to a positive coefficient even if there is a negative correlation in a linear correlation analysis.

### 3. Data Analytic Model Planning and Implementation

Considering the mixed data types (numerical and categorical) and binary target variable described in part 2, we leverage the below analytic techniques to understand the relationship between various predictor variables and the likelihood of a specific outcome, in this case, successful subscriptions to term deposits and to develop an accurate predictive model for upcoming clients, enabling us to anticipate the success rate based on their demographic background and account portfolio.

#### 3.1. Logistic Regression

We employ logistic regression, which is suited to our binary classification objective – the identification of potential customers inclined to positively respond to marketing campaigns (variable-“Y”). Unlike correlation analysis, which primarily assesses the strength and direction of linear relationships between a dependent variable to the targeted variable individually, logistic regression delves deeper into understanding how multiple predictor variables collectively influence the binary outcome. By estimating the probability of a positive response, logistic regression enables us to model the complex interplay between various factors, providing a more comprehensive and predictive understanding of customer behavior.

For implementation, after the data preprocessing in Part 2, we construct a logistic regression model using the *glm()* function with the target variable "Y." By specifying "*family = binomial*," we indicate its binary outcome prediction. We then utilize *summary()* and *anova()* functions to provide insights. *summary()* presents coefficients, standard errors, z-scores, and p-values for predictor variables. *anova()* summarizes variable contributions to deviance reduction.

In term of optimization, we employ backward-stepwise selection, iteratively refining the model by removing less significant predictors. The "Overall Fit (Omnibus Test)" assesses model goodness-of-fit by comparing null and residual deviance. The resulting statistic follows a Chi-squared distribution, highlighting the collective significance of predictors.

### 3.2. Other Classification Techniques

Furthermore, we enrich our approach by integrating an array of diverse classifiers stated in the table. This incorporation ensures a comprehensive exploration of intricate predictive patterns inherent within the dataset. By thoughtfully blending these techniques, our model aims to not only deliver accurate predictions but also to unravel the underlying factors that influence customer responses.

Classifiers	Features
Logistic Classifier	<ul style="list-style-type: none"> <li>• Able to predict a binary response from a binary predictor.</li> <li>• Highly interpretable result</li> </ul>
Random Forest / Decision Tree	<ul style="list-style-type: none"> <li>• Good for prediction but a little bit difficult to interpret.</li> </ul>
Support Vector Machine (SVM-linear & non-linear)	<ul style="list-style-type: none"> <li>• Effective in creating linear or non-linear decision boundaries for classification.</li> </ul>
XGBoost / AdaBoost	<ul style="list-style-type: none"> <li>• A boosting algorithm that excels in predictive accuracy and handles complex relationships.</li> <li>• Tackle overfitting with simple model structure</li> </ul>
Neural Network	<ul style="list-style-type: none"> <li>• Able to compute a probabilistic output like logistic classifier and also give a more powerful non-linear decision boundary.</li> </ul>
KNN Classifier	<ul style="list-style-type: none"> <li>• Logistic regression poses assumptions of its linear relationship on the log-odds scale and stable and reliable predictions from low variance.</li> <li>• KNN applies fewer assumptions and is able to implement non-linear boundaries for testing and training.</li> </ul>

We use a loop to fit piped processed data in all the selected models and the results are appended to a table. After that, we choose the best for modelling based on accuracy and the result table is listed below ranked by accuracy.

<i>Classifier</i>	<i>Accuracy</i>
<i>Random Forest</i>	0.8778066
<i>XGBoost</i>	0.8743523
<i>AdaBoost</i>	0.8674439
<i>Decision Tree</i>	0.8644214
<i>Non-linear SVM</i>	0.8605354
<i>Neural Network Classifier</i>	0.8588083
<i>Logistic Regression</i>	0.8583765
<i>SVM</i>	0.8419689
<i>KNN</i>	0.8419689
<i>Linear SVM</i>	0.8372193

We train classifiers on the training data, predict outcomes for the testing dataset using *predict()*, and calculate accuracy for each through a loop. The resulting insights are presented using confusion matrices, highlighting performance patterns and disparities among classifiers. This visual analysis enhances our comprehension of classifier effectiveness.

## 4. Results Interpretation and Business Implications

### 4.1.a. Logistic regression model Result

We conducted both a logistic regression and a backward-stepwise logistic regression model. After evaluating the results, we decided to proceed with the backward-stepwise approach for interpretation. The computed Chi-squared statistic is 5709.254 with a p-value of 0, indicating strong statistical significance.

The combination of independent variables we selected plays a valuable role in predicting positive response rates. Collectively, these variables led to a reduction in the deviance statistic from the initial null model value of 10699.4 to a final value of 4990.2. The deviance statistic gauges the "Lack of fit," and in this case, the notable reduction signifies an improved fit of the model. This indicates that the chosen variables enhance the model's ability to capture and explain patterns in the data.

The small p-values ( $\Pr(>|z|)$ ) associated with these variables emphasize their meaningful contribution to the model. In essence, the results affirm that our selected variables are important for predicting positive responses and for improving the overall fit of the model.

## Individual Coefficients:

- Continuous variables:

### *Significance:*

Last contact duration, Consumer Price Index, the three-month Euro Interbank Offered Rate and Employment Variation Rate are statistically significant.

### *Sign:*

All attributes, except for Employment Variation Rate, are positively linked to customers' favorable responses to the marketing campaign and their interest in term deposits. However, employment variation rate demonstrates a negative association with the response, indicating that as this rate decreases, customers' positive response likelihood diminishes.

### *Size:*

Based on our analysis, the most important factor in successful bank marketing is the duration of the last contact. Our data shows that longer call durations have a strong correlation with increased likelihood of positive customer responses. In fact, our analysis revealed a substantial positive coefficient of 1.97131 for this variable. So, it's essential for bank marketers to focus on extending the duration of their last contact with customers to increase the chances of a successful outcome.

Following closely in impact is the Consumer Price Index, boasting a coefficient of 1.56842. This index gauges inflation, representing the average alteration over time in prices paid by urban consumers for a selection of goods and services. Consequently, when inflation surges or consumer prices rise, customers appear more predisposed to positively engage with marketing campaigns promoting term deposits.

Furthermore, the three-month Euro Interbank Offered Rate has a positive coefficient has 1.10320 which is the interest rate at which European banks lend funds to one another for a three-month period. Hence, it indicates that the likelihood of a positive response to the marketing campaign is higher when the three-month Euribor rate increases. However, employment variation rate, which is an economic indicator that reflects changes in the employment situation, such as the job market and employment conditions.

A negative coefficient -3.90430 for employment variation rate suggests that a decrease in the employment variation rate is associated with a positive response to the marketing campaign. This means that when the economy is stable or the job market is doing well (with high rates of employment variation), customers may not be as interested in responding positively to marketing efforts for term deposits.

- Categorical variables:

#### *Significance:*

Contact communication type, Last contact month of year, outcome of the previous marketing campaign and different contact months are all statistically significant  $\Pr(>|z|)$ .

#### *Sign:*

New customers or those who successfully converted to complete a term deposit in the previous marketing campaign exhibit a positive association with placing a term deposit, particularly if customer were contacted in August or March. Conversely, contact via telephone or being contacted in June, May, or November shows a negative association with customer response.

#### *Size:*

Customers who positively responded to the previous campaign hold a substantial coefficient of 2.07076, coupled with an exceptionally low p-value. Similarly, customers who have not previously been contacted also wield significant influence, as evidenced by the coefficient of 0.53257.

The month of August emerges as a notable contributor, boasting a significant coefficient of 2.28519. Furthermore, March closely follows with a coefficient of 1.10942. Both months are accompanied by remarkably low p-values, underscoring their potent impact. Specifically, customers contacted during March and August exhibit a heightened likelihood of delivering positive responses. In contrast, customer engagement during May, June, and November is associated with a diminished propensity to subscribe to a deposit.

Conversely, customers contacted via telephone, as opposed to cellular communication, exhibit a negative coefficient of -0.56481. This highlights their reduced probability of responding positively, emphasizing the relatively more favorable outcomes associated with alternative contact methods.

Additional attributes, including education with a degree or professional course, as well as job statuses like retired or student, exhibit p-values that approach 0. These values suggest a subtle yet discernible positive influence on the customer response. In other words, customers with educational backgrounds involving degrees or professional courses, along with those who are retired or students, are more likely to exhibit a favorable response to the marketing campaign.

#### 4.1.b. Business implications of the results

Our analysis has highlighted key factors that strongly influence positive responses in our marketing campaigns. These factors encompass how we approach customers, economic indicators, customer segments, and their backgrounds.

In terms of customer approach, the duration of interactions emerges as a critical aspect. It's important to train our staff to engage effectively with customers, understand their needs, and communicate the benefits of the term deposit product. This can lead to extended conversations and improved success rates. Additionally, using cellphones for communication proves more effective than traditional home phones, possibly due to customers' mobility. Collecting mobile numbers during account setup is essential for targeted promotions. Exploring alternatives like text and email communications can also be beneficial. Moreover, certain months, like March and Aug, exhibit stronger impacts, likely due to company bonuses or customer savings preferences. On the other hand, caution is advised in June and May, which coincide with summer vacation periods and lead to higher expenditures. During these times, promoting liability products such as credit cards or line of credit may yield better results.

From an economic perspective, close attention to key indicators is crucial. The Consumer Price Index, Euro Interbank Offered Rate, and Employment Variation Rate play significant roles. Our findings suggest that customers are more likely to consider term deposits during times of higher inflation and rising interest rates. This insight encourages increased term deposit promotions during such periods. Conversely, stable job markets may lead to fewer term deposit purchases, necessitating a balanced approach.

The choice of customers and their backgrounds also impacts campaign outcomes. Repeat customers and new entrants are more likely to respond positively. Regular communication with existing customers as their term deposits mature, accompanied by attractive promotions, can boost retention rates. Customers with higher education levels also tend to engage more. Specialized strategies for students and retirees, who show interest in term deposits, can further enhance campaign effectiveness.

In conclusion, our analysis offers valuable insights into driving positive responses in marketing campaigns. By refining our customer approach, adapting to economic indicators, and tailoring strategies to customer preferences, we can optimize campaign success and increase term deposit subscriptions.



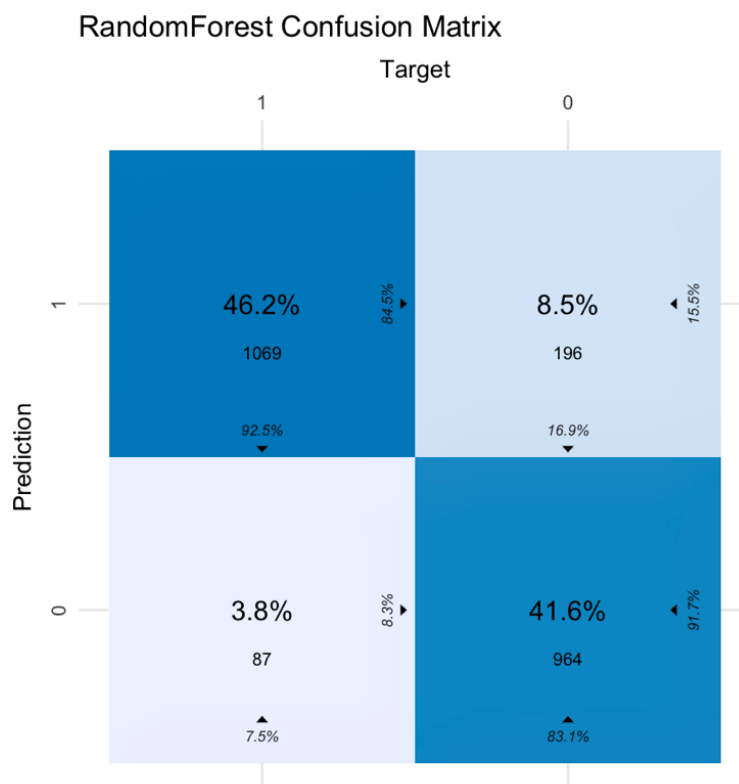
#### 4.2.a. Classification regression model Result

After optimizing for accuracy and conducting a thorough assessment, we determined that the Random Forest model with the highest accuracy at 0.8778. An accuracy of 0.8778, indicates a relatively good classifier performance. It suggests that the model is correctly predicting the outcome for approximately 87.7% of the instances in the test dataset.

##### Confusion Matrix

An examination of the confusion matrix sheds light on the model's predictive prowess, as illustrated in the figure. An analysis of the matrix reveals that the model excels in predicting instances of 1 (indicating a positive response to term deposit subscription), with 1069 accurate predictions. Similarly, when forecasting instances of 0 (indicating a negative response), 964 predictions align accurately with actual outcomes.

However, it's worth noting some instances of misclassification. Specifically, there are 196 cases where the actual outcome should be 1, but the model predicted 0 (No). Conversely, 87 instances were predicted as 0, whereas the true outcome is 1 (Yes).



The classification approach demonstrates strong overall accuracy. However, the discrepancies we observed indicate areas where the model can further excel. To enhance performance, I recommend refining the model's hyperparameters through systematic experimentation. Fine tuning parameters

like the number of trees in a Random Forest or the learning rate in gradient boosting algorithms could yield better outcomes. Moreover, employing cross-validation methods would help gauge the model's reliability and general applicability, reducing the impact of random occurrences or overfitting.

#### 4.1.b. Business implications of the results

In real-world scenarios, the implementation involves effectively using customer profiles, tailoring marketing strategies to individual preferences, and considering factors like contact duration and the prevailing economic environment. By including these various elements as inputs to the predictive model, the system gains the ability to automatically predict whether a customer's response is likely to be positive (coded as 1) or negative (coded as 0).

For customers predicted to respond positively, an automated reminder email could be sent and sales may call in the next few days for follow up, encouraging them to consider and potentially act on the term deposit offer. This approach aims to enhance the overall response rate and conversion.

Furthermore, the predictive model is a useful tool for future evaluations. By comparing its predictions with the actual outcomes of campaigns in regular reviews, the company gets a clear measure of how well the campaigns are working. If the results are not as expected, the company should investigate and make changes to improve the methods, strategies, or target audience of the campaigns. This helps ensure that the company's marketing efforts are effective and on track.

The predictive model is also important for evaluating how well employees are doing. By using the model's predictions as a guide, the company can fairly measure employee achievements. By comparing what the model predicts with how things actually turn out with customers, we can figure out how well each employee is doing. This helpful information makes evaluations fairer and helps the company decide how to reward and help employees grow.

## 5. Application of the Model for Predictive Analysis

### 5.1 Out Of Sample Prediction

The provided figures in this context pertain to predictions made for data outside the model's training set. From this out-of-sample data, customer is retrieved who holds a university degree, was contacted during the promotion in August, has previously subscribed to a term deposit, and experienced a lengthy promotion contact duration. During this period, the economy had high consumer prices and interest rates, but low employment changes. These factors made it seem likely that the customer would sign up for a term deposit.

```
## job_blue-collar job_entrepreneur job_housemaid job_management job_retired
## 1 0 0 0 0 1
## job_self-employed job_services job_student job_technician job_unemployed
## 1 0 0 0 0 0
## education_high_school education_illiterate education_professional_course
## 1 0 0 0
## education_university_degree contact_telephone month_aug month_dec month_jul
## 1 1 0 1 0 0
## month_jun month_mar month_may month_nov month_oct month_sep
## 1 0 0 0 0 0
## poutcome_nonexistent poutcome_success duration campaign emp_var_rate
## 1 0 1 5 1 -0.8
## cons_price_idx euribor3m nr_employed
## 1 5 3 0.2
```

```
# Make predictions on the out-of-sample data
outOfSample_pred <- predict(rf_model, as.matrix(outOfSample))
outOfSample_pred <- ifelse(outOfSample_pred >= 0.5, 1, 0)
print(outOfSample_pred)
```

```
## 1
## 1
```

By plugging these variables into the model, the system can independently forecast whether the outcome is inclined towards positivity (labeled as 1) or negativity (labeled as 0). This predictive ability empowers us to make well-informed choices and enhances our marketing efforts for better results.

Based on a result of 1, indicating a positive response, the system can trigger a follow-up email, and within the next two days, staff can reach out to the customer for a reminder, effectively boosting the conversion rate.

## 6. Conclusion

The aim of this project is to build models to optimize bank marketing strategy by utilizing data analytics to anticipate customer purchasing patterns. We followed a rigorous process of data cleaning, transformation, and analysis, leading us to valuable insights.

Employing logistic regression and ten classification models are trained, we accurately predicted positive customer responses to marketing campaigns, achieving an 87.78% accuracy rate with Random Forest. Our findings highlighted the importance of tailored communication, economic indicators, and customer demographics.

Our findings recommend banks can optimize marketing by focusing on personalized communication, economic trends, and targeted segmentation. These insights offer actionable strategies to enhance campaign success and drive business growth.

Ultimately, our project demonstrates the power of data analytics to reshape bank marketing, bridging the gap between customer response and marketing efforts. Informed by data-driven insights, our approach promises enhanced campaign outcomes, more accurate employee evaluations, and sustained business growth. This project showcases the transformative potential of predictive modeling in modern banking practices.

## 7. Challenging Factors

Making our ambitious project to change how we predict customer buying habits in the banking industry a success involves overcoming several tough challenges.

### 7.1. Dataset Quality and Completeness

One of the foremost challenges pertains to ensuring the quality and completeness of the dataset. Any inaccuracies, incompleteness, or inconsistencies within the dataset can compromise the reliability of predictions. The presence of erroneous or missing data points may lead to skewed results, ultimately impacting the effectiveness of the predictive models. A meticulous data preprocessing and cleaning strategy is paramount to address these issues, guaranteeing that the data accurately reflects the real-world scenarios it aims to capture.

### 7.2. Imbalanced Data Distribution

The inherent imbalance in the distribution of positive and negative outcomes within the dataset introduces a significant challenge. We have a lot less examples of positive outcomes (customers buying) compared to negative outcomes (customers not buying). The risk of overlooking the minority class and focusing excessively on the dominant class could undermine the predictive power of the models. Hence, we employed technique under-sampling, can help mitigate this challenge and restore balance to the dataset.

### 7.3. Skillful Feature Selection and Engineering

Crafting predictive models that can discern meaningful patterns from a complex dataset is a skillful task. Effective feature selection and engineering are pivotal to identify influential variables while filtering out noise. This process demands a deep understanding of the domain and an insightful exploration of the data. The challenge lies in striking the right balance between inclusivity and exclusivity, ensuring that the chosen features align with the objectives of predicting customer purchase behavior accurately.

### 7.4. Constrained Model Selection

The selection of suitable models for this project is constrained due to the characteristics of the available data. An exploration was conducted clustering model, which aims to group similar data points together. However, regrettably, this approach yielded inconclusive outcomes with limited utility.

## 7.5 Limited by the Model Assumptions

The effectiveness of certain analytical techniques, such as logistic regression and specific classifiers, is contingent upon adherence to underlying model assumptions. For example, in random forest, individual decision trees should be built independently, without being influenced by other trees. Otherwise, the correlation between trees can lead to overfitting. In logistic regression, logistic regression assumes a linear relationship between the independent variables and the log-odds of the dependent variable. Deviations from these assumptions can compromise the accuracy and reliability of the results, potentially leading to misinterpretation and flawed conclusions.

## 7.6 Data Extrapolation

While data extrapolation allows us to extend predictions beyond observed ranges, it introduces inherent limitations. When the model is tasked with making predictions or inferences for scenarios lying outside the scope of the original data, risks arise. For instance, if a new variable emerges that is beyond the observed data range, like the introduction of an unaccounted age group or a novel value of the 3-month term deposit rate ('euribor3m'), the assumed continuous patterns might not accurately apply. This underscores the need for cautious consideration and rigorous validation of extrapolated results. Ensuring the reliability and relevance of such projections is crucial to prevent drawing erroneous or misleading conclusions.

In conclusion, the journey to revolutionize customer purchase behavior prediction is not without its share of challenges. By proactively addressing these multifaceted obstacles, we can pave the way for the successful implementation of our transformative project. The pursuit of improved customer purchase behavior prediction and optimization in the bank marketing sector demands a strategic and collaborative approach, guided by a deep understanding of these challenges and their potential solutions.