

A Data-Driven Approach to Predict the Success of Bank Telemarketing

Curated By:
Pyimoe Than
&
Zain Mirza

San Francisco State University

Table of Contents

- 1. Executive Summary**
- 2. Introduction**
- 3. Data Description**
 - 3.1. Correlation Matrix Graph
 - 3.2. Data Cleaning
- 4. Feature Engineering**
 - 4.1. Forward Selection
 - 4.2. Backward Elimination
- 5. Model Selection + Results (using K-fold cross-validation)**
 - 5.1. Logistic Regression with logit link
 - 5.2. Logistic regression with prob link
 - 5.3. Proposed Solution
- 6. Appendix**
 - 6.1. Data Repo
 - 6.2. Code - GitHub

This document outlines our approach to help a Portuguese bank overcome its challenges in effectively advertising its ‘bank term deposit’ service. Through the use of machine learning techniques, we are confident in our ability to predict which individuals are most likely to subscribe to this service based on various factors in a dataset. Our project involves identifying the most effective ML model for predicting responses and ensuring that our phone call campaigns are personalized and targeted toward the proper audience.

The dataset, acquired from the University of California, Irvine Repository, from this project contains information on direct marketing campaigns conducted by a Portuguese banking institution. The ML techniques aim to enhance the bank to understand the predicting factors that potentially determine a classified “yes” or “no” response to a ‘term deposit’ subscription.

The data utilized in this project is sourced from the UCI Repository, specifically from the Bank Marketing Data Set. This dataset comprises 45,211 observations including 17 columns.

Correlation coefficients between many pairs of predictor variables in the data appear to be less than 0.5. Thus, we do not need to worry about multicollinearity in this dataset.

Variable overview:

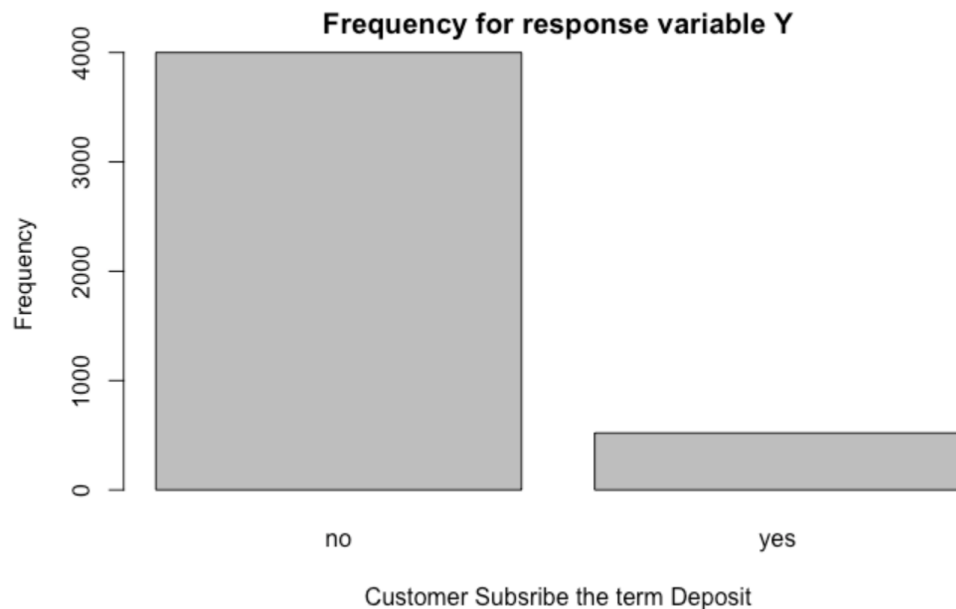
1. **age**: (numeric)
 2. **job**: type of job (categorical: 'admin.', 'blue-collar', 'entrepreneur', 'housemaid', 'management', 'retired', 'self-employed', 'services', 'student', 'technician', 'unemployed', 'unknown')
 3. **marital**: marital status (categorical: 'divorced', 'married', 'single', 'unknown'; note: 'divorced' means divorced or widowed)
 4. **education** (categorical:)
 5. **default**: has credit in default? (categorical: 'no', 'yes', 'unknown')
 6. **Balance**: (numeric)
 7. **housing**: has a housing loan? (categorical: 'no', 'yes', 'unknown')
 8. **loan**: has a personal loan? (categorical: 'no', 'yes', 'unknown')
 9. **contact**: contact communication type (categorical: 'cellular', 'telephone')
 10. **month**: last contact month of the year (categorical: 'Jan', 'Feb', 'mar', ..., 'Nov', 'Dec')
 11. **day_of_week**: last contact day of the week (categorical: 'mon', 'Tue', 'wed', 'Thu', 'Fri')
 12. **duration**: last contact duration, in seconds (numeric). Important note: this attribute highly affects the output target (e.g., if duration=0 then y='no'). Yet, the duration is not known before a call is performed. Also, after the end of the call, y is known. Thus, this input should only be included for benchmark purposes and should be discarded if the intention is to have a realistic predictive model.
 13. **campaign**: number of contacts performed during this campaign and for this client (numeric, includes the last contact)
 14. **pdays**: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means the client was not previously contacted)
 15. **previous**: number of contacts performed before this campaign and for this client (numeric)
 16. **poutcome**: outcome of the previous marketing campaign (categorical: 'failure', 'nonexistent', 'success')
- Output variable (desired target)
17. **y**: has the client subscribed to a term deposit? (binary: 'yes', 'no')

3.2 Data Cleaning

During the data cleaning process, we conducted a series of integrity checks to ensure the dataset's accuracy, completeness, and consistency. We were pleased to discover that there were no missing values or duplicated data in the dataset. However, we did notice an unbalanced response variable,

with a higher count of "No" responses than "Yes" responses. This class imbalance has the potential to introduce bias in our prediction model, which is something we must address.

Figure 1: Frequency of Responses



Percentage of 'yes' subscription of term deposit: 11.524 %
Percentage of 'no' subscription of term deposit: 88.476 %

The figure above depicts the total number of 'yes' responses and 'no' responses in the dataset. From the graph above we can see there are far more 'no' responses to the bank service than there are 'yes'.

An oversampling method was applied to address this issue to balance the data. The result of this technique exhibits a balanced distribution in our dataset and is now ready for further analysis using machine learning techniques. Below, in Figure 2, you can visualize the frequency of responses after including the oversampling method.

Figure 2: Frequency of Responses (After Oversampling Method)

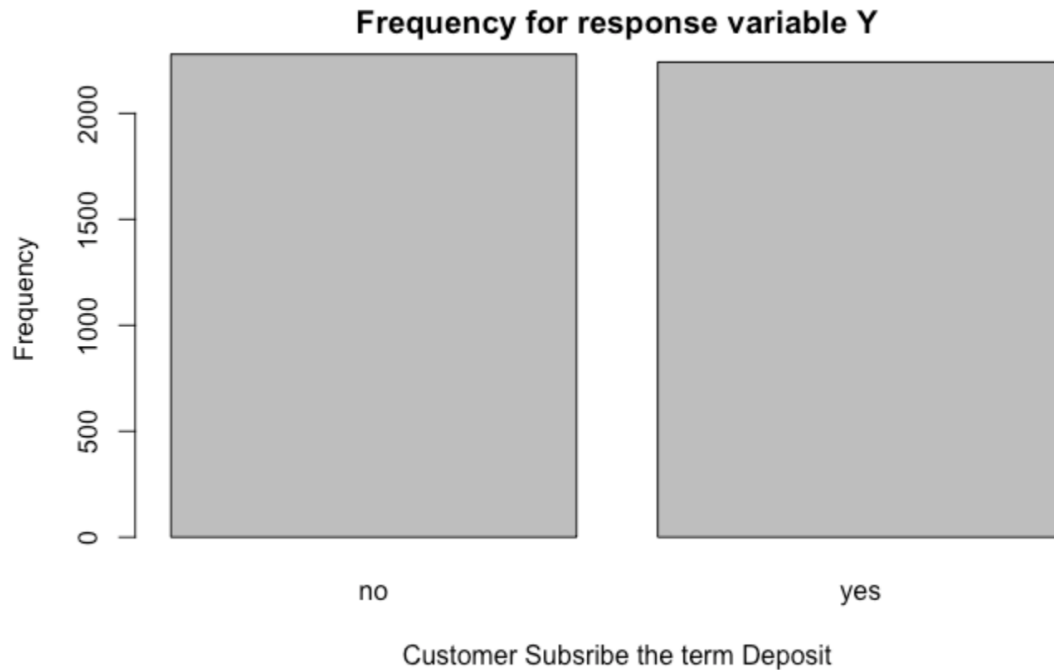
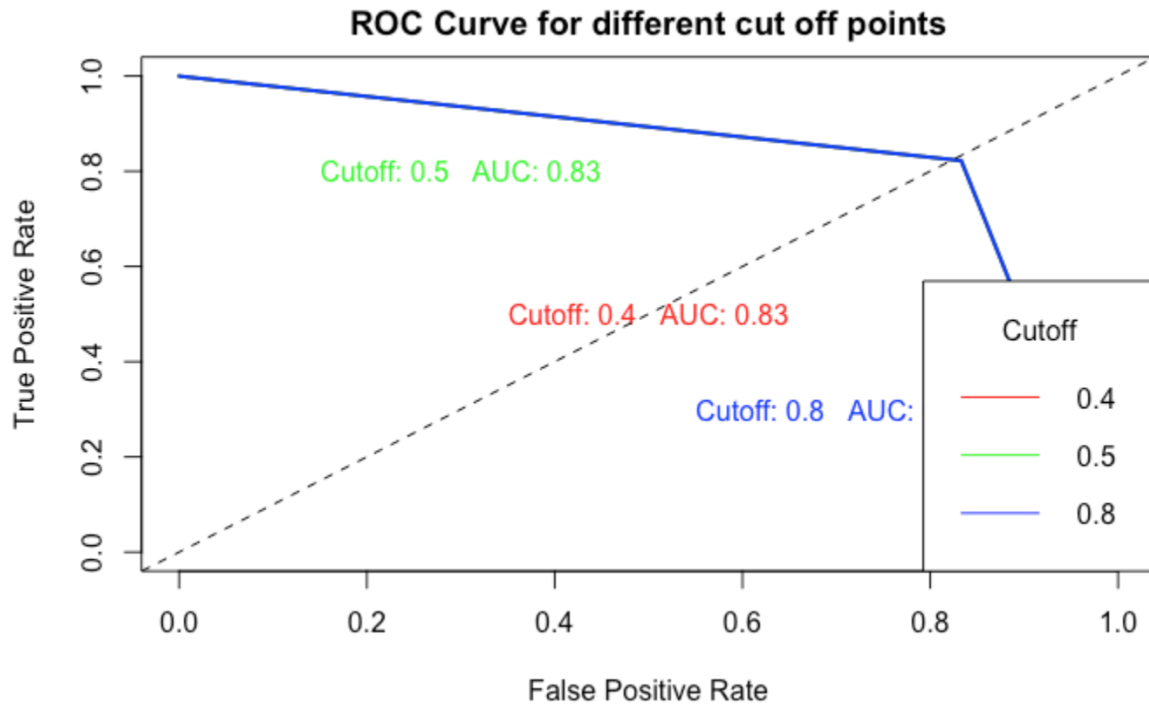


Figure 2 above displays a balanced number of ‘yes’ and ‘no’ responses, after incorporating the oversampling method.

4. Feature Engineering

Following the data cleaning process, we divided the data into 80% training and 20% testing data prior to fitting. The Logistic Regression model was trained using both k-fold cross-validation and Leave-One-Out Cross-validation (LOOCV) on the training data. This comprehensive validation approach provides robust estimates of the model's performance. We generated Confusion Matrices for different cutoff points, specifically 0.4, 0.5, and 0.8, and were surprised to find that all three cutoff points resulted in the same accuracy and Area Under the Curve (AUC) values. Therefore, we confidently recommend using 0.5 as the best cutoff point.

Cutoff points	Accuracy	AUC
0.4	0.8276	0.83
0.5	0.8276	0.83
0.8	0.8276	0.83



We also found a significant number of predictor variables from the null model that can be used to predict the response for a ‘term deposit’ bank service. However, this approach carries the risk of over-fitting since all features are used to fit the model. To mitigate this risk, we will use feature selection to identify the most important predictors responsible for our response outcome. By selecting only the most relevant predictors, we can boost the accuracy and interpretability of our predictive models. In this analysis, we will ensure that our selection of predictors is optimal to produce the best results. We will use the forward and backward elimination process to determine the proper feature selection.

4.1 Forward Selection

The Forward Selection method starts with an empty feature set and iteratively adds one feature at a time based on its performance and contribution to the model’s predictive power. From the features selection using the forward method, we found these variables to be most significant:

job	education	loan	day	duration	pdays	poutcome	marital	housing	contact	month	campaign	previous
-----	-----------	------	-----	----------	-------	----------	---------	---------	---------	-------	----------	----------

The table above displays all the significant predictor variables found using the forward selection process.

4.2 Backward Selection

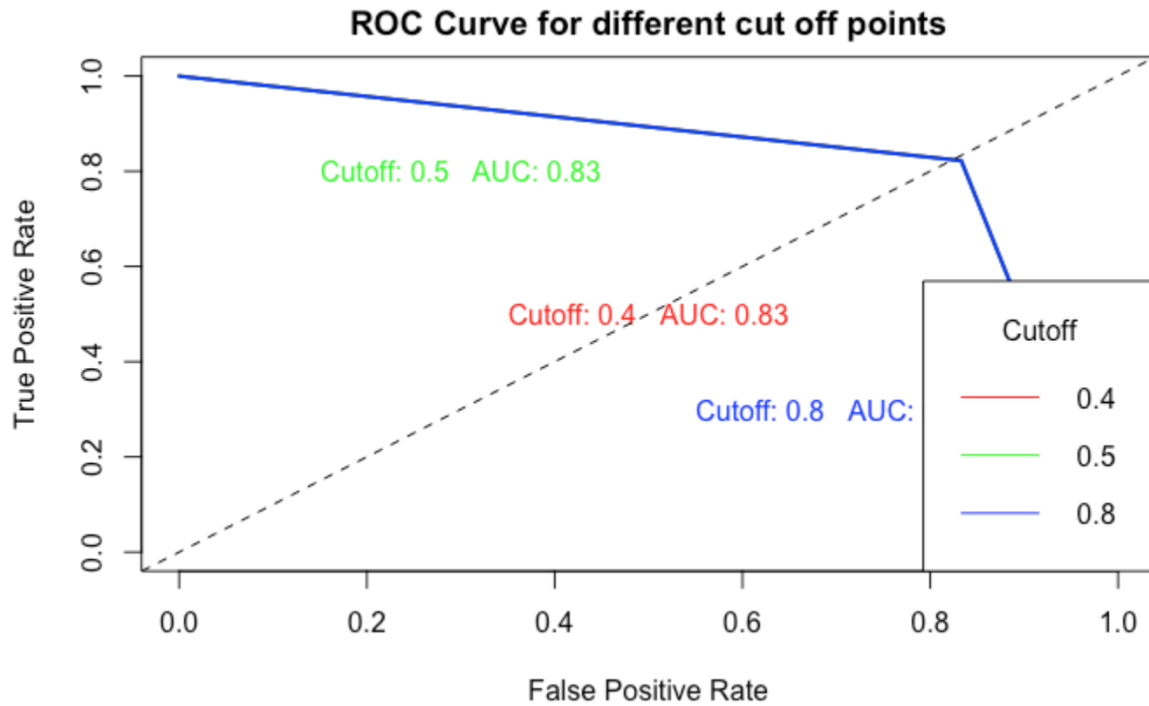
The Backward Selection method, in contrast to the Forward Selection, commences with a complete set of features and gradually eliminates one feature at a time based on its performance and impact on the model's predictive ability. Our previous findings have demonstrated that the Backward Elimination approach produces the same outcome as the Forward Selection. Both Forward and Backward methods have identified the same 13 predictor variables from the dataset's 16.

5. Model Selection

5.1 Logistic Regression with Logit Link

To ensure a reliable evaluation, we divided the data into 80% training and 20% testing data prior to fitting. The Logistic Regression model was trained using both k-fold cross-validation and Leave-One-Out Cross-validation (LOOCV) on the training data. This comprehensive validation approach provides robust estimates of the model's performance. We generated Confusion Matrices for different cutoff points, specifically 0.4, 0.5, and 0.8, and were surprised to find that all three cutoff points resulted in the same accuracy and Area Under the Curve (AUC) values. Therefore, we confidently recommend using 0.5 as the best cutoff point.

Cutoff points	Accuracy	AUC
0.4	0.8276	0.83
0.5	0.8276	0.83
0.8	0.8276	0.83

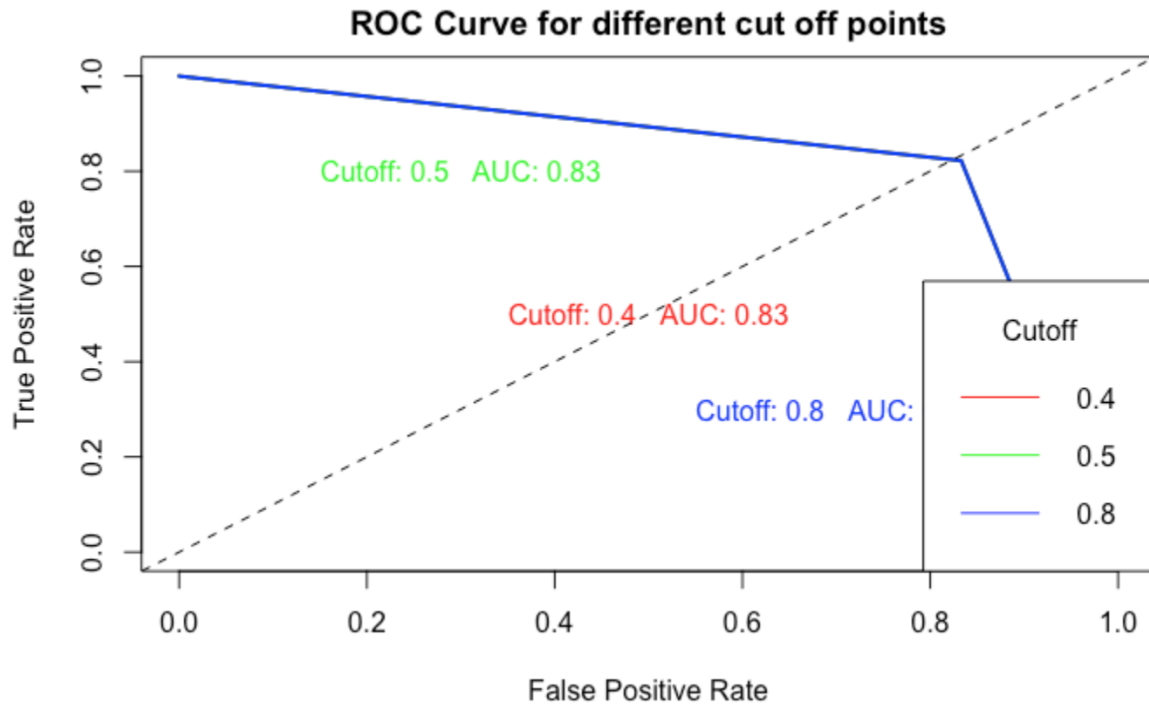


AUC score greater than 0.5. It means the model performs better than random guessing.

5.2 Logistic Regression with Prob Link

Confusion Matrices were generated for different cutoff points, specifically, 0.4, 0.5, and 0.8. Surprisingly, all three cutoff points resulted in the same accuracy and Area Under the Curve(AUC) values. For this particular reason, we should use 0.5 as the best cutoff point.

Cutoff points	Accuracy	AUC
0.4	0.821	0.82
0.5	0.821	0.82
0.8	0.821	0.82



When the AUC score is greater than 0.5, it indicates that the model's performance is superior to random guessing.

5.3 Proposed Solution

Based on the results, it can be confidently concluded that Logistic Regression with Logit link is a more suitable choice as it improves model accuracy and AUC score compared to Logistic regression with Prob link. This improvement is a clear indication of the effectiveness of the Logit link in enhancing the model's performance.

6. Appendix

6.1 Data Source - UCI Repository

<https://archive.ics.uci.edu/ml/datasets/bank+marketing>.

6.2 Code - GitHub

[Project Code](#)

