



# Assessment #3

## Telecom Marketing Campaign

Camilo Andres Uribe Guerra  
ID: 25416518  
MSc. Data Science for Innovation

# Table of Contents

<b>Table of Contents .....</b>	<b>2</b>
<b>1. Introduction .....</b>	<b>3</b>
Problem Statement .....	3
Rationale.....	3
Project aims and objective .....	3
<b>2. Methodology .....</b>	<b>4</b>
Methods overview.....	4
Methods detail .....	4
<b>3. Results.....</b>	<b>5</b>
Key Findings.....	5
In depth results.....	5
<b>4. Conclusions .....</b>	<b>8</b>
<b>5. References.....</b>	<b>9</b>
<b>6. Appendix .....</b>	<b>10</b>
A. Data Dictionary.....	10
B. Confusion Matrix Results for RF Best Model. ....	11
C. Logistic Stats Model Analysis (No tuning) .....	12
D. Feature Engineering Binned Creation .....	13

# 1. Introduction

## Problem Statement

In the competitive Australian telecommunications sector, marketing campaigns are vital for customer retention and profitability. TelecomSyd company aims to identify customer segments likely to respond positively to campaigns to optimize future strategies.

## Rationale

Effective marketing boosts profitability and engagement. By leveraging past and current campaign data, TelecomSyd seeks to allocate resources efficiently and enhance outcomes through predictive insights.

## Project aims and objective

The goal of this research is to build predictive models to estimate customer responses to TelecomSyd's marketing campaigns. The main research questions include:

1. Which predictive model—parametric or non-parametric—best estimates customer response?
  - H0: Both model types perform equally.
2. Do customer demographics or behaviors influence model performance?
  - H0: Demographic or behavioral features do not significantly affect performance.
3. Do past campaign outcomes impact future predictions?
  - H0: Previous outcomes do not significantly influence future predictions.

The objectives are:

- Develop and compare parametric and non-parametric models for customer response prediction.
- Identify key demographic and behavioral features affecting campaign success.
- Analyze the influence of previous campaign outcomes on future predictions.

## 2. Methodology

### Methods overview

The project followed a structured approach involving EDA, data preprocessing, feature selection, and model development. Both parametric (Logistic Regression) and non-parametric (Random Forest) models were used, with hyperparameter tuning and cross-validation to identify the best-performing model for further analysis.

### Methods detail

#### Feature Selection and Engineering

Feature selection involved evaluating numerical features using biserial correlation and categorical features with Cramér's V. The duration variable was excluded, as it only becomes known after campaign interactions and cannot be used for future predictions. Variables such as `cons.price.idx` and `emp.var.rate` were also excluded due to multicollinearity, being redundant with `euribor3m` and `nr.employed`. Additional feature engineering included binning previous contact counts and campaign interactions to enhance interpretability.

#### Data Transformation for Modeling

Data transformation included standard scaling for numerical features to standardize their ranges and one-hot encoding for categorical variables to ensure they were suitable for modeling. The processed data was then split into training, validation, and test sets, maintaining class distribution using stratified sampling.

#### Model Development

An initial logistic regression model was created using a basic stats model setup to assess its statistical properties, convergence, and characteristics without hyperparameter tuning, establishing a foundation for the parametric model. Both logistic regression and random forest classifiers underwent hyperparameter tuning through grid search with cross-validation, utilizing maximum likelihood estimation (MLE) for fitting the parametric model (Scikit-learn Developers, 2024) and employing ensemble methods for the non-parametric model. Precision and F1 score were the main evaluation metrics, with an emphasis on precision to accurately identify responsive customers, crucial for targeted marketing.

#### Model Comparison and Feature Importance

Models were compared based on cross-validation scores. The evaluation highlighted the best model, which was further analyzed for feature importance to guide business insights.

**Note:** The final set of features and parameters used for modeling are provided in the appendix for reference.

## 3. Results

### Key Findings

1. The initial logistic regression model did not converge, indicating challenges in model fitting without hyperparameter tuning (Appendix C).
2. The random forest model outperformed logistic regression in precision (36.6% validation precision), supported by the Mann-Whitney U test which showed a significant difference in precision ( $p < 0.05$ ) using 10 folds for best model.
3. Economic indicators (nr.employed, euribor3m, and cons.conf.idx) were the most influential features in predicting campaign success, aligning with the importance of financial metrics. In contrast, customer demographic and behavioral features had less impact.
4. The final random forest model demonstrated consistency across training, validation, and testing datasets (Figure 1), though the overall precision indicates further data refinement and acquisition is needed for better performance.
5. Feature importance analysis confirmed the relevance of previous campaign outcomes (pdays and poutcome\_success), supporting the third hypothesis regarding their influence on future campaign responses.

### In depth results

#### Initial Logistic Regression Analysis

The initial logistic regression model was assessed to understand the statistical behavior and overall convergence. The model did not converge (Appendix, Table 1), indicating potential limitations in fitting the data accurately without hyperparameter tuning. Some coefficients, such as pdays ( $p < 0.001$ ) and nr.employed ( $p < 0.001$ ), were found to be statistically significant, while others like job\_blue-collar and education\_high.school were not significant. The full list of coefficients and p-values can be found in the regression results table in the appendix.

#### Model Results with Hyperparameter Tuning

**Logistic Regression** The best logistic regression model was obtained with the following parameters: {'C': 1, 'fit\_intercept': False, 'penalty': 'l2', 'solver': 'lbfgs'}. This model achieved:

- F1 Score (Training): 45.5%
- Precision (Training): 35.6%
- F1 Score (Validation): 43.1%
- Precision (Validation): 33.4%

**Random Forest** The optimal random forest model was tuned with parameters: {'bootstrap': True, 'max\_depth': 5, 'min\_samples\_leaf': 2, 'min\_samples\_split': 4, 'n\_estimators': 100}. This model produced:

- F1 Score (Training): 47.6%
- Precision (Training): 38.3%
- F1 Score (Validation): 45.5%
- Precision (Validation): 36.6%

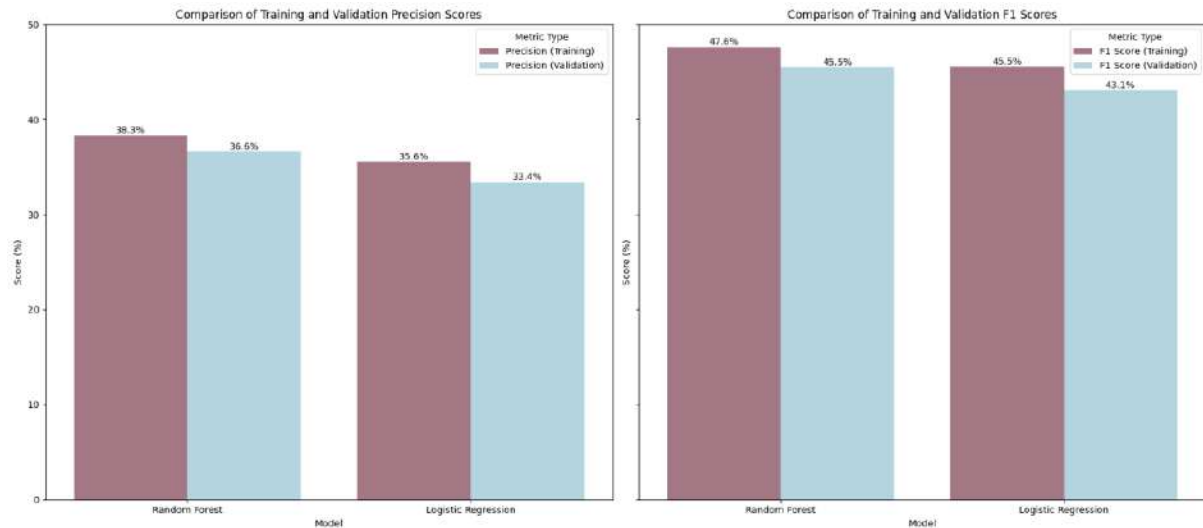


Figure 1. Best Models Performance Indicators using GridSearchCV for Logistic Regression and Random Forrest across Training and Validation Datasets.

### Mann-Whitney U Test for Model Comparison

The Mann-Whitney U test was conducted to compare the precision and F1 scores of the two models:

- Precision - U statistic: 80.0, p-value: 0.025 (significant difference)
- F1 Score - U statistic: 69.0, p-value: 0.162 (not significant)

Based on these results, the random forest model was chosen for its higher precision, with significant evidence supporting its superior performance over logistic regression.

### Model Consistency

Figure 2 illustrates the training, validation, and testing results for precision and F1 scores across the Random Forrest model, showing that it maintained consistent performance across datasets. The final confusion matrix for the training set (Appendix B) further verifies the performance.

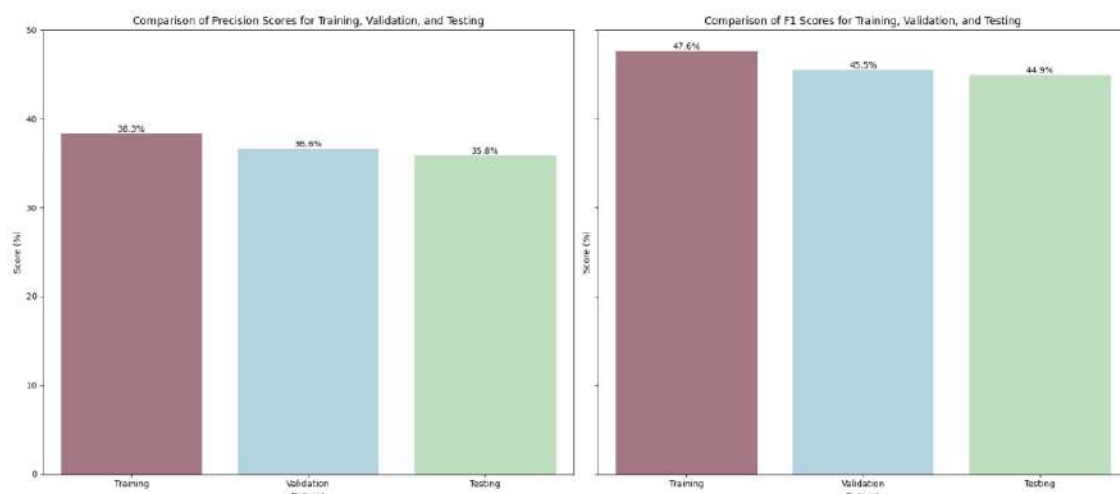


Figure 2> Final Random Forrest Model consistency performance indicators across all splitting datasets.

## Feature Importance Analysis

The feature importance analysis (Figure X) revealed that `nr.employed` was the most influential variable, highlighting the role of employment metrics. Financial indicators such as `euribor3m` and `cons.conf.idx` were also significant, confirming the impact of economic conditions. Variables related to campaign interactions, such as `pdays` and `poutcome_success`, showed moderate relevance, partially supporting the third hypothesis regarding the influence of previous campaign outcomes. Additionally, months like March, May, and October emerged as potential periods to target for campaigns. Most customer demographic and behavioral features demonstrated minimal impact, supporting the second hypothesis that these characteristics do not significantly affect model performance.

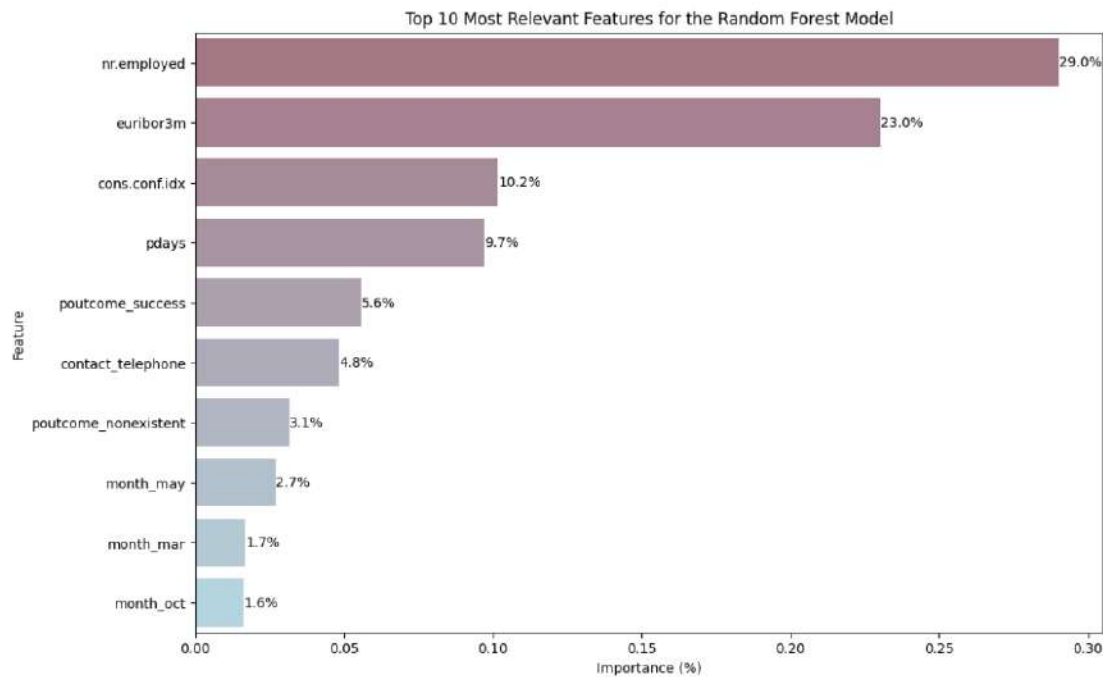


Figure 3. Feature Importance for Random Forrest best model (% by feature)



## 4. Conclusions

- The Random Forest model outperformed Logistic Regression, achieving precision (35.8%) and F1 score (44.9%) on testing dataset. However, the overall precision was still low, indicating a need for further experiments, data refinement, and data collection, showing this model is not suitable for business cases.
- Key features impacting predictions included economic indicators such as the number of employed individuals, the 3-month Euribor rate, and the consumer confidence index, underscoring the significant influence of broader economic conditions on campaign responses.
- Previous campaign interactions, such as the number of days since the client was last contacted and the success outcome of prior campaigns, held moderate importance but can influence results in future campaigns. These factors contribute to defining successful marketing targets based on past campaign experiences.



## 5. References

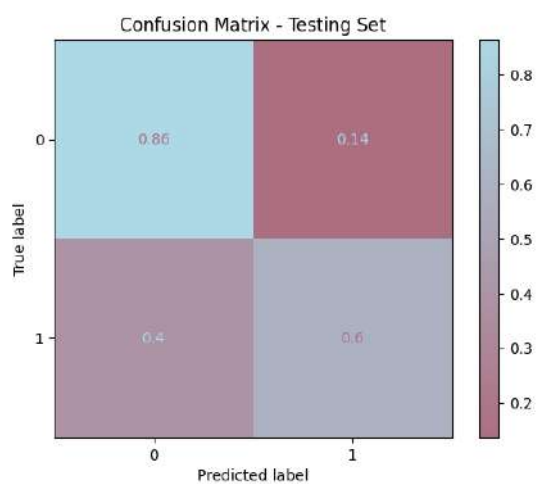
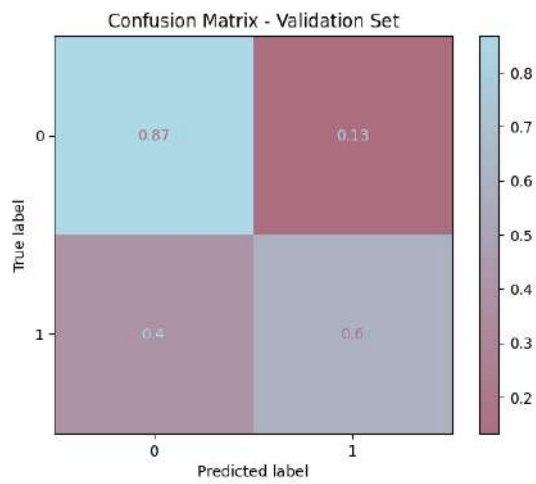
- Gao, W., Ding, Z., Cui, T., & Tao Cui. (2022). Construction of Digital Marketing Recommendation Model Based on Random Forest Algorithm. *Security and Communication Networks*, 2022, 1–9. <https://doi.org/10.1155/2022/1871060>
- Ullah, I., Raza, B., Malik, A. K., Imran, M., Islam, S. U., & Kim, S. W. (2019). A Churn Prediction Model Using Random Forest: Analysis of Machine Learning Techniques for Churn Prediction and Factor Identification in Telecom Sector. *IEEE Access*, 7, 60134–60149. <https://doi.org/10.1109/ACCESS.2019.2914999>
- Scikit-learn Developers. (2024). *LogisticRegression*. Scikit-learn: Machine Learning in Python. [https://scikit-learn.org/1.5/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/1.5/modules/generated/sklearn.linear_model.LogisticRegression.html)

## 6. Appendix

### A. Data Dictionary

Variable Name	Description
<b>age</b>	Age
<b>job</b>	Type of job
<b>marital</b>	Marital status
<b>education</b>	Level of education
<b>default</b>	Has credit in default
<b>balance</b>	Average yearly balance
<b>housing</b>	Has a housing loan
<b>loan</b>	Has a personal loan
<b>contact</b>	Contact communication type
<b>day</b>	Day of contact
<b>month</b>	Month of contact
<b>duration</b>	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 obviously known.
<b>campaign</b>	Number of contacts performed during this campaign and for this client
<b>pdays</b>	Number of days that passed by after the client was last contacted from a previous campaign (numeric, -1 means client was not previously contacted)
<b>previous</b>	Number of contacts performed before this campaign and for this client
<b>poutcome</b>	Outcome of the previous marketing campaign
<b>emp.var.rate</b>	employment variation rate - quarterly indicator (numeric)
<b>cons.price.idx</b>	consumer price index - monthly indicator (numeric)
<b>cons.conf.idx</b>	consumer confidence index - monthly indicator (numeric)
<b>euribor3m</b>	euribor 3 month rate - daily indicator (numeric)
<b>nr.employed</b>	number employed - quarterly indicator (numeric)
<b>y</b>	Did the client subscribe to a Telecom plan?

## B. Confusion Matrix Results for RF Best Model.



## C. Logistic Stats Model Analysis (No tuning)

Parameter	Value
No. Observations	26356
Model	Logit
Df Residuals	26310
Method	MLE
Df Model	45
Pseudo R-squ.	0.2169
Log-Likelihood	-7264
converged	FALSE
LL-Null	-9275.8
Covariance Type	nonrobust
LLR p-value	0

Coefficient	Coef	Std Err	z	P> z	[0.025	0.975]
const	-1.9468	nan	nan	nan	nan	nan
pdays	-0.224	<b>0.047</b>	-4.799	<b>0</b>	-0.315	<b>-0.13</b>
cons.conf.idx	0.0909	<b>0.025</b>	3.69	<b>0</b>	0.043	<b>0.139</b>
euribor3m	-0.0518	<b>0.081</b>	-0.64	<b>0.522</b>	-0.21	<b>0.107</b>
nr.employed	-0.7669	<b>0.072</b>	-10.656	<b>0</b>	-0.908	<b>-0.63</b>
job_blue-collar	-0.1164	<b>0.086</b>	-1.349	<b>0.177</b>	-0.286	<b>0.053</b>
job_entrepreneur	0.018	<b>0.131</b>	0.137	<b>0.891</b>	-0.238	<b>0.274</b>
job_housemaid	-0.129	<b>0.162</b>	-0.794	<b>0.427</b>	-0.447	<b>0.189</b>
job_management	0.0069	<b>0.092</b>	0.075	<b>0.94</b>	-0.174	<b>0.188</b>
job_retired	0.2885	<b>0.103</b>	2.792	<b>0.005</b>	0.086	<b>0.491</b>
job_self-employed	-0.0452	<b>0.128</b>	-0.354	<b>0.724</b>	-0.295	<b>0.205</b>
job_services	-0.0923	<b>0.093</b>	-0.993	<b>0.321</b>	-0.274	<b>0.09</b>
job_student	0.0696	<b>0.124</b>	0.563	<b>0.573</b>	-0.173	<b>0.312</b>
job_technician	0.0052	<b>0.078</b>	0.067	<b>0.947</b>	-0.148	<b>0.158</b>
job_unemployed	-0.0816	<b>0.143</b>	-0.569	<b>0.569</b>	-0.362	<b>0.199</b>
job_unknown	-0.0025	<b>0.253</b>	-0.01	<b>0.992</b>	-0.498	<b>0.493</b>
marital_married	-0.039	<b>0.073</b>	-0.531	<b>0.596</b>	-0.183	<b>0.105</b>
marital_single	0.0456	<b>0.08</b>	0.571	<b>0.568</b>	-0.111	<b>0.202</b>
marital_unknown	0.3784	<b>0.459</b>	0.825	<b>0.409</b>	-0.521	<b>1.277</b>

education_basic.6y	0.0712	<b>0.131</b>	0.546	<b>0.585</b>	-0.185	<b>0.327</b>
education_basic.9y	-0.0557	<b>0.103</b>	-0.54	<b>0.589</b>	-0.258	<b>0.147</b>
education_high.school	0.0014	<b>0.1</b>	0.014	<b>0.989</b>	-0.194	<b>0.197</b>
education_illiterate	1.0354	<b>0.791</b>	1.309	<b>0.191</b>	-0.515	<b>2.586</b>
education_professional.course	-0.0242	<b>0.111</b>	-0.219	<b>0.827</b>	-0.241	<b>0.193</b>
education_university.degree	0.0598	<b>0.1</b>	0.599	<b>0.549</b>	-0.136	<b>0.256</b>
education_unknown	0.1542	<b>0.131</b>	1.175	<b>0.24</b>	-0.103	<b>0.411</b>
default_unknown	-0.3165	<b>0.073</b>	-4.335	<b>0</b>	-0.46	<b>-0.17</b>
default_yes	-15.6819	<b>8019.341</b>	-0.002	<b>0.998</b>	-15700	<b>15700</b>
contact_telephone	-0.4245	<b>0.073</b>	-5.827	<b>0</b>	-0.567	<b>-0.28</b>
month_aug	-0.1013	<b>0.112</b>	-0.901	<b>0.368</b>	-0.322	<b>0.119</b>
month_dec	0.0233	<b>0.229</b>	0.102	<b>0.919</b>	-0.425	<b>0.471</b>
month_jul	0.2382	<b>0.103</b>	2.319	<b>0.02</b>	0.037	<b>0.44</b>
month_jun	0.223	<b>0.098</b>	2.268	<b>0.023</b>	0.03	<b>0.416</b>
month_mar	0.9031	<b>0.136</b>	6.627	<b>0</b>	0.636	<b>1.17</b>
month_may	-0.6775	<b>0.081</b>	-8.331	<b>0</b>	-0.837	<b>-0.52</b>
month_nov	-0.3364	<b>0.112</b>	-3.013	<b>0.003</b>	-0.555	<b>-0.12</b>
month_oct	-0.2053	<b>0.147</b>	-1.393	<b>0.163</b>	-0.494	<b>0.083</b>
month_sep	-0.4841	<b>0.157</b>	-3.078	<b>0.002</b>	-0.792	<b>-0.18</b>
previous_1 (Contacted Once)	-0.4686	nan	nan	nan	nan	nan
previous_2-3 (Few Contacts)	-0.6066	nan	nan	nan	nan	nan
previous_4+ (Many Contacts)	-0.8453	nan	nan	nan	nan	nan
campaign_2-3 Contacts	0.0047	<b>0.048</b>	0.098	<b>0.922</b>	-0.089	<b>0.098</b>
campaign_4-5 Contacts	-0.0766	<b>0.08</b>	-0.953	<b>0.341</b>	-0.234	<b>0.081</b>
campaign_6-10 Contacts	-0.3529	<b>0.115</b>	-3.061	<b>0.002</b>	-0.579	<b>-0.13</b>
campaign_11+ Contacts	-1.4928	<b>0.361</b>	-4.131	<b>0</b>	-2.201	<b>-0.78</b>
poutcome_nonexistent	-0.0264	nan	nan	nan	nan	nan
poutcome_success	0.6118	<b>0.244</b>	2.509	<b>0.012</b>	0.134	<b>1.09</b>

## D. Feature Engineering Binned Creation

