# \<Banking Analysis with Predicted Modeling>

## Summary

Our team of data scientists has been tasked with developing an effective telemarketing strategy to sell term deposit accounts for a Portuguese bank. This bank has been conducting marketing campaigns, but it has not been effective. Our goal then was to develop an effective marketing strategy by using machine learning and to collect correlated attributes to increase the possibility to get more subscriptions.

The dataset consisted of 4521 rows and 17 attributes – 7 of which were quantitative and 10 which were qualitative including the target attribute. The dataset was prepared by looking at the attribute types and analyzing the mean, max, min and std as well as extreme values. Furthermore, a correlation matrix was used to determine which attributes were most strongly related to the class attribute. This was then taken into account when choosing which attributes to use for the machine learning techniques. The dataset had no missing data, but it contained a lot of outliers, and it was imbalanced.

The two machine learning methods were used in this analysis – Decision Tree and Naïve Bayes. These methods were used for all attributes and the selected attributes (Age, Duration, and Poutcome). The Decision Tree was created using Python weka and Sklearn, while the Naïve Bayes method was applied using Python. The performance metrics – accuracy, recall, and precision – were used to evaluate the accuracy of the models.
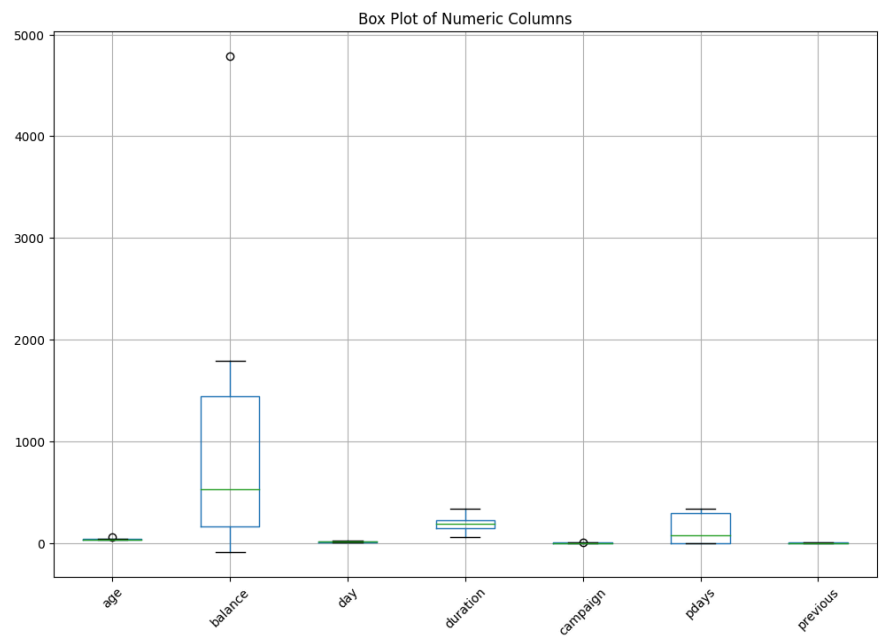
### Data Preparation

| Column Name | Attribute type | Description |
| --- | --- | --- |
| **Age** | quantitative | Age of the customer (numeric). |
| **Job** | nominal | Type of job (qualitative). |
| **Marital** | nominal | Marital status (qualitative). |
| **Education** | ordinal | Education of the customer (qualitative). |
| **Default** | nominal | Shows whether the customer has credit in default or not (qualitative). |
| **Balance** | quantitative | Average yearly balance in Euros (numeric). |
| **Housing** | nominal | Shows whether the customer has a housing loan or not (qualitative). |

| | | |
|---|---|---|
| **Loan** | qualitative | Shows whether the customer has a personal loan or not (qualitative/categorical). |
| **Contact** | nominal | Shows how the last contact for marketing campaign has been made (qualitative) |
| **Day** | quantitative | Day: Shows on which day of the month the last time a customer was contacted (numeric). |
| **Month** | ordinal | Month: Shows on which month of the year the last time a customer was contacted (qualitative). |
| **Duration** | quantitative | Shows the last contact duration in seconds (numeric). |
| **Campaign** | quantitative | Number of contacts performed during the marketing campaign and for this customer (numeric). |
| **Pdays** | quantitative | Pdays: Number of days that passed by after the client was last contacted from a previous campaign (numeric, -1 means client was not previously contacted). |
| **Previous** | quantitative | Previous: Number of contacts performed before this campaign and for this client (numeric). |
| **Poutcome** | ordinal | Outcome of the previous marketing campaign (qualitative). |
| **Y** | nominal | Class attribute showing whether the client has subscribed a term deposit or not (binary: "yes","no"). |

## Table of max, min, mean and standard deviation of attributes:

| | Age | Balance | Day | Duration | Campaign | Pdays | Previous |
|---|---|---|---|---|---|---|---|
| Max | 87.000000 | 71188.000000 | 31.000000 | 3025.000000 | 50.000000 | 871.000000 | 25.000000 |
| Min | 19.000000 | -3313.000000 | 1.000000 | 4.000000 | 1.000000 | -1.000000 | 0.000000 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Mean | 41.170095 | 1422.657819 | 15.915284 | 263.961292 | 2.793630 | 39.766645 | 0.542579 |
| Std | 10.576211 | 3009.638142 | 8.247667 | 259.856633 | 3.109807 | 100.121124 | 1.693562 |

Box Plot of Numeric Columns

Correlation Matrix:



## Correlations:

- Duration is the most correlated with y. poutcome, month. contact and age are also slightly correlated.
- Previous is highly positively correlated with pdays.
- poutcome and previous are highly negatively correlated with each-other

Based on the BestFIRst algorithm in Weka, we decided to filter for the selected attributes of age, duration, poutcome and the class attribute.

## Bar Charts/Histograms of Selected Attributes:

poutcome
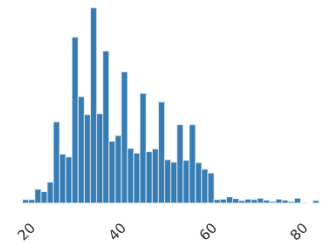Categorical

HIGH CORRELATION   IMBALANCE

| | | | |
|---|---|---|---|
| Distinct | 4 | unknown | 3705 |
| Distinct (%) | 0.1% | failure | 490 |
| Missing | 0 | other | 197 |
| Missing (%) | 0.0% | success | 129 |
| Memory size | 35.4 KiB | | |

## age
Real number (ℝ)

| | | | | |
|---|---|---|---|---|
| Distinct | 67 | Minimum | 19 | |
| Distinct (%) | 1.5% | Maximum | 87 | |
| Missing | 0 | Zeros | 0 | |
| Missing (%) | 0.0% | Zeros (%) | 0.0% | |
| Infinite | 0 | Negative | 0 | |
| Infinite (%) | 0.0% | Negative (%) | 0.0% | |
| Mean | 41.170095 | Memory size | 35.4 KiB | |

## duration
Real number (ℝ)

| | | | | |
|---|---|---|---|---|
| Distinct | 875 | Minimum | 4 | |
| Distinct (%) | 19.4% | Maximum | 3025 | |
| Missing | 0 | Zeros | 0 | |
| Missing (%) | 0.0% | Zeros (%) | 0.0% | |
| Infinite | 0 | Negative | 0 | |
| Infinite (%) | 0.0% | Negative (%) | 0.0% | |
| Mean | 263.96129 | Memory size | 35.4 KiB | |

## y
Boolean

| | |
|---|---|
| Distinct | 2 |
| Distinct (%) | < 0.1% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 4.5 KiB |

| | |
|---|---|
| False | 4000 |
| True | 521 |

## Predictive Modeling/Classification

### Predictive Modeling/Classification
Classification is a supervised machine learning technique used to train models which predict classes based on the training from the training dataset. In this project, the Decision Tree and Naïve Bayes classification techniques were used to train, test, and evaluate the performance of the dataset. The train-test split method was used for the classification.

### Classification using Decision Tree
A Decision Tree is a classification technique that uses attributes in a dataset to predict a class based on decisions. Decision trees are referred to as classification for qualitative attributes or regression trees for continuous variables (Gupta, 2017). For this project Python Weka Package and Sklearn were used to create the decision tree.

The current sklearn decision tree graph has depth of 5 and is a decision tree based on all attributes.

duration <= 628.5
gini = 0.198
samples = 3164
value = [2812, 352]
class = no

True

poutcome_success <= 0.5
gini = 0.139
samples = 2890
value = [2673, 217]
class = no

duration <= 222.5
gini = 0.114
samples = 2821
value = [2650, 171]
class = no

housing_yes <= 0.5
gini = 0.444
samples = 69
value = [23, 46]
class = yes

month_oct <= 0.5
gini = 0.046
samples = 1841
value = [1798, 43]
class = no

pdays <= 373.0
gini = 0.227
samples = 980
value = [852, 128]
class = no

duration <= 83.0
gini = 0.315
samples = 46
value = [9, 37]
class = yes

balance <= 114.5
gini = 0.476
samples = 23
value = [14, 9]
class = no

month_mar <= 0.5
gini = 0.041
samples = 1821
value = [1783, 38]
class = no

duration <= 146.0
gini = 0.375
samples = 20
value = [15, 5]
class = no

age <= 61.5
gini = 0.215
samples = 970
value = [851, 119]
class = no

month_nov <= 0.5
gini = 0.18
samples = 10
value = [1, 9]
class = yes

gini = 0.0
samples = 2
value = [2, 0]
class = no

pdays <= 96.5
gini = 0.268
samples = 44
value = [7, 37]
class = yes

gini = 0.0
samples = 5
value = [0, 5]
class = yes

duration <= 176.0
gini = 0.346
samples = 18
value = [14, 4]
class = no

day <= 9.5
gini = 0.457
samples = 51
value = [18, 33]
class = yes

gini = 0.036
samples = 1800
value = [1767, 33]
class = no

gini = 0.363
samples = 21
value = [16, 5]
class = no

gini = 0.0
samples = 13
value = [13, 0]
class = no

gini = 0.408
samples = 7
value = [2, 5]
class = yes

gini = 0.2
samples = 943
value = [837, 106]
class = no

gini = 0.499
samples = 27
value = [14, 13]
class = no

gini = 0.0
samples = 9
value = [0, 9]
class = yes

gini = 0.0
samples = 1
value = [1, 0]
class = no

gini = 0.0
samples = 18
value = [0, 18]
class = yes

gini = 0.393
samples = 26
value = [7, 19]
class = yes

gini = 0.0
samples = 8
value = [8, 0]
class = no

gini = 0.48
samples = 10
value = [6, 4]
class = no

gini = 0.48
samples = 20
value = [12, 8]
class = no

gini = 0.312
samples = 31
value = [6, 25]
class = yes

False

marital_married <= 0.5
gini = 0.5
samples = 274
value = [139, 135]
class = no

duration <= 1026.5
gini = 0.473
samples = 120
value = [46, 74]
class = yes

day <= 3.5
gini = 0.478
samples = 154
value = [93, 61]
class = no

day <= 19.5
gini = 0.493
samples = 86
value = [38, 48]
class = yes

balance <= 766.5
gini = 0.36
samples = 34
value = [8, 26]
class = yes

gini = 0.0
samples = 4
value = [0, 4]
class = yes

pdays <= 203.5
gini = 0.471
samples = 150
value = [93, 57]
class = no

job_blue-collar <= 0.5
gini = 0.49
samples = 35
value = [20, 15]
class = no

age <= 65.0
gini = 0.172
samples = 21
value = [2, 19]
class = yes

age <= 42.0
gini = 0.497
samples = 13
value = [6, 7]
class = yes

campaign <= 11.0
gini = 0.458
samples = 141
value = [91, 50]
class = no

previous <= 3.5
gini = 0.346
samples = 9
value = [2, 7]
class = yes

gini = 0.493
samples = 25
value = [11, 14]
class = yes

gini = 0.18
samples = 10
value = [9, 1]
class = no

gini = 0.095
samples = 20
value = [1, 19]
class = yes

gini = 0.0
samples = 1
value = [1, 0]
class = no

gini = 0.375
samples = 8
value = [6, 2]
class = no

gini = 0.0
samples = 5
value = [0, 5]
class = yes

gini = 0.449
samples = 138
value = [91, 47]
class = no

gini = 0.0
samples = 3
value = [0, 3]
class = yes

gini = 0.0
samples = 6
value = [0, 6]
class = yes

gini = 0.444
samples = 3
value = [2, 1]
class = no

## Classification using Naive Bayes

Naive Bayes: Naive Bayes is a machine learning technique. For this project we used the Weka Python Package to complete a Naive Bayes from the Banking data.

## Summary Table:

\<Baseline\>

|  | Position 0 Precision | Position 0 Recall | Position 1 Precision | Position 1 Recall |
|---|---|---|---|---|
| Decision Tree | 0.91805887032 61734 | 0.96327212020 03339 | 0.55555555555 55556 | 0.34810126582 278483 |
| Naive Bayes | 0.92497938994 22918 | 0.93656093489 14858 | 0.46853146853 146854 | 0.424050632911 3924 |

\<Selected Features\>

|  | Position 0 Precision | Position 0 Recall | Position 1 Precision | Position 1 Recall |
|---|---|---|---|---|
| Decision Tree | 0.92469635627 53037 | 0.95325542570 95158 | 0.53719008264 46281 | 0.411392405063 29117 |
| Naive Bayes | 0.911302982731 5542 | 0.969115191986 6444 | 0.54878048780 48781 | 0.28481012658 22785 |

## Machine Learning Comparison: Decision Tree vs Naive Bayes

Evaluation from the perspective of class at position 1 for recall was decided to be the most valuable. This is because we are interested in customers who will say YES to the subscription and we want to recall all potential customers, whether or not it is precise everytime.

Naive Bayes Recall at Position 1 is greater than Decision Tree Recall at Position 1

0.4240506329113924 > 0.34810126582278483

Machine Learning Baseline "all attributes" and "selected features" Comparison

Selected Features: Age, Duration, Poutcome, Yes or No

| Machine Learning Recall | Baseline | Selected Features |
|---|---|---|
| Naive Bayes | 0.4240506329113924 | 0.2848101265822785 |
| Decision Tree | 0.34810126582278483 | 0.41139240506329117 |

Highest Recall

Highest Recall with Selected Features

## Conclusions and Recommendations

### Major findings from different sections

It was found that duration is the most correlated with y. poutcome, month. contact and age are also correlated. Taking this into account and using the BestFIRst algorithm in Weka, we decided to filter for the selected attributes of "age", "duration", "poutcome" and the class attribute.

Recall for the decision tree improved by selecting for the attributes "Age", "Duration", "Poutcome", "Yes or No", but decreased for Naive Bayes. Naive Bayes is more accurate for recall with all attributes. However, the Decision Tree Recall for the Selected Features was increased to a similar recall percentage.

### Recommendations:

Because the recall for the decision tree was increased to a similar percentage when selecting for the specific attributes as the Naive Bayes recall for all attributes, it is recommended that when speaking to the company we would recommend using the Decision Tree Machine Learning model using the features: "Age", "Duration", "Poutcome", "Yes or No" to recruit potential subscribers.

As well, since banks would typically have millions of clients, instead of having to go through 16 different attributes and process the dataset using Naive Bayes, it would be faster and have a similar percentage outcome to use Decision Tree on the selected features mentioned above.