

Udacity Machine Learning Nanodegree

Bank Marketing Campaign Predictive Analysis

Tarun Parmar
August 15th, 2017

I. Definition

Project Overview

In banks, huge data records information about their customers. This data can be used to create and keep clear relationship and connection with the customers in order to target them individually for definite products or banking offers. Usually, the selected customers are contacted directly through: personal contact, telephone cellular, mail, and email or any other contacts to advertise the new product/service or give an offer, this kind of marketing is called direct marketing. In fact, direct marketing is in the main a strategy of many of the banks and insurance companies for interacting with their customers [1].

Historically, the name and identification of the term direct marketing suggested first time in 1967 by Lester Wunderman, which he is considered to be the father of direct marketing [2]. In addition, some of the banks and financial-services companies may depend only on strategy of mass marketing for promoting a new service or product to their customers. In this strategy, a single communication message is broadcasted to all customers through media such as television, radio or advertising firm, etc. [3]. In this approach, companies do not set up a direct relationship to their customers for new-product offers. In fact, many of the customers are not interesting or respond to this kind of sales promotion [4].

Accordingly, banks, financial-services companies and other companies are shifting away from mass marketing strategy because its ineffectiveness, and they are now targeting most of their customers by direct marketing for specific product and service offers [1, 4]. Due to the positive results clearly measured; many marketers attractive to the direct marketing. For example, if a marketer sends out 1,000 offers by mail and 100 respond to the promotion, the marketer can say with confidence that the campaign led immediately to 10% direct responses. This metric is known as the 'Response Rate', and it is one of many clear quantifiable success metrics employed by direct marketers.

From the literature, the direct marketing is becoming a very important application in data mining these days. The data mining has been used widely in direct marketing to identify prospective customers for new products, by using purchasing data, a predictive model to measure that a customer is going to respond to the promotion or an offer [5]. Data mining has gained popularity for illustrative and predictive applications in banking processes.

Problem Statement

All bank marketing campaigns are dependent on customers' huge electronic data. The size of these data sources is impossible for a human analyst to come up with interesting information that will help in the decision-making process. Data mining models are completely helping in the performance of these campaigns.

The purpose is increasing the campaign effectiveness by identifying the main characteristics that affect a success (the deposit subscribed by the client) based on a handful of algorithms that we will test (e.g. Logistic Regression, Random Forests, Decision Trees and others). With the experimental results we will demonstrate the performance of the models by statistical metrics like accuracy, sensitivity, precision, recall, etc. We the higher scoring of these metrics, we will be able to judge the success of these models in predicting the best campaign contact with the clients for subscribing deposit. The aim of the marketing campaign was to get customers to subscribe to a bank term deposit product. Whether they did this or not is variable 'y' in the data set. The bank in question is considering how to optimize this campaign in future. What would your recommendations to the marketing manager be?

Metrics

The evaluation metrics proposed are appropriate given the context of the data, the problem statement, and the intended solution. The performance of each classification model is evaluated using three statistical measures; classification accuracy, sensitivity and specificity. It is using true positive (TP), true negative (TN), false positive (FP) and false negative (FN). The percentage of Correct/Incorrect classification is the difference between the actual and predicted values of variables. True Positive (TP) is the number of correct predictions that an instance is true, or in other words; it is occurring when the positive prediction of the classifier coincided with a positive prediction of target attribute. True Negative (TN) is presenting a number of correct predictions that an instance is false, (i.e.) it occurs when both the classifier, and the target attribute suggests the absence of a positive prediction. The False Positive (FP) is the number of incorrect predictions that an instance is true. Finally, False Negative (FN) is the number of incorrect predictions that an instance is false. Table below shows the confusion matrix for a two-class classifier.

	Predicted No	Predicted Yes
Actual No	TN	FN
Actual Yes	FP	TP

Classification accuracy is defined as the ratio of the number of correctly classified cases and is equal to the sum of TP and TN divided by the total number of cases (TN + FN + TP + FP).

$$Accuracy = \frac{TP + TN}{TN + FN + TP + FP}$$

Precision is defined as the number of true positives (TP) over the number of true positives plus the number of false positives (FP).

$$Precision = \frac{TP}{TP + FP}$$

Recall is defined as the number of true positives (TP) over the number of true positives plus the number of false negatives (FN).

$$Recall = \frac{TP}{TP + FN}$$

Sensitivity refers to the rate of correctly classified positive and is equal to TP divided by the sum of TP and FN. Sensitivity may be referred as a True Positive Rate.

$$Sensitivity = \frac{TP}{FN + TP}$$

Specificity refers to the rate of correctly classified negative and is equal to the ratio of TN to the sum of TN and FP

$$Specificity = \frac{TN}{TN + FP}$$

II. Analysis

Data Exploration

The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product (bank term deposit) would be ('yes') or not ('no') subscribed. The classification goal is to predict if the client will subscribe (yes/no) a term deposit (variable y).

Input variables:

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: 'basic.4y', 'basic.6y', 'basic.9y', 'high.school', 'illiterate', 'professional.course', 'university.degree', 'unknown')
5. **default**: has credit in default? (categorical: 'no', 'yes', 'unknown')
6. **housing**: has housing loan? (categorical: 'no', 'yes', 'unknown')
7. **loan**: has personal loan? (categorical: 'no', 'yes', 'unknown')
8. **contact**: contact communication type (categorical: 'cellular', 'telephone')
9. **month**: last contact month of year (categorical: 'jan', 'feb', 'mar', ..., 'nov', 'dec')
10. **day_of_week**: last contact day of the week (categorical: 'mon', 'tue', 'wed', 'thu', 'fri')
11. **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 obviously 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.
12. **campaign**: number of contacts performed during this campaign and for this client (numeric, includes last contact)

13. **pdays**: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means client was not previously contacted)
14. **previous**: number of contacts performed before this campaign and for this client (numeric)
15. **poutcome**: outcome of the previous marketing campaign (categorical: 'failure','nonexistent','success')
16. **emp.var.rate**: employment variation rate - quarterly indicator (numeric)
17. **cons.price.idx**: consumer price index - monthly indicator (numeric)
18. **cons.conf.idx**: consumer confidence index - monthly indicator (numeric)
19. **euribor3m**: euribor 3 month rate - daily indicator (numeric)
20. **nr.employed**: number of employees - quarterly indicator (numeric)

Output variable (desired target):

1. **y** - has the client subscribed a term deposit? (binary: 'yes','no')

The dataset has 21 columns and 41188 rows, with 20 features, and one response variable. The number of customers who subscribed is 4640 and the number of customers who did not subscribe is 36548. Based on this the response rate of customers is about 11.27%, this makes the dataset very imbalanced. Fig.1 shows age distribution in the data and is a normal distribution with slight left skewness. This hints that majority of the responses are in 25-40 age group.

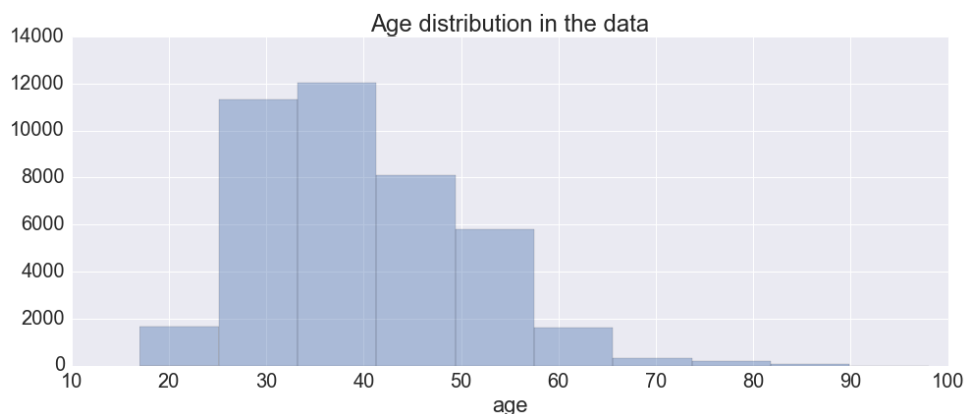


Fig. 1 Age distribution

Below table shows descriptive statistics on the data. `duration` seems to have a lot more variance, it could be a good predictor.

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx
count	41188.00000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000
mean	40.02406	258.285010	2.567593	962.475454	0.172963	0.081886	93.575664
std	10.42125	259.279249	2.770014	186.910907	0.494901	1.570960	0.578840
min	17.00000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000
25%	32.00000	102.000000	1.000000	999.000000	0.000000	-1.800000	93.075000
50%	38.00000	180.000000	2.000000	999.000000	0.000000	1.100000	93.749000
75%	47.00000	319.000000	3.000000	999.000000	0.000000	1.400000	93.994000
max	98.00000	4918.000000	56.000000	999.000000	7.000000	1.400000	94.767000

Exploratory Visualization

Fig. 2 shows data distribution of `job` factors with `duration` and split by `marital` status. It is hard to spot relevances between factors from this plot. We can pick a few numeric variables and plot them to see if we can find any different patterns.

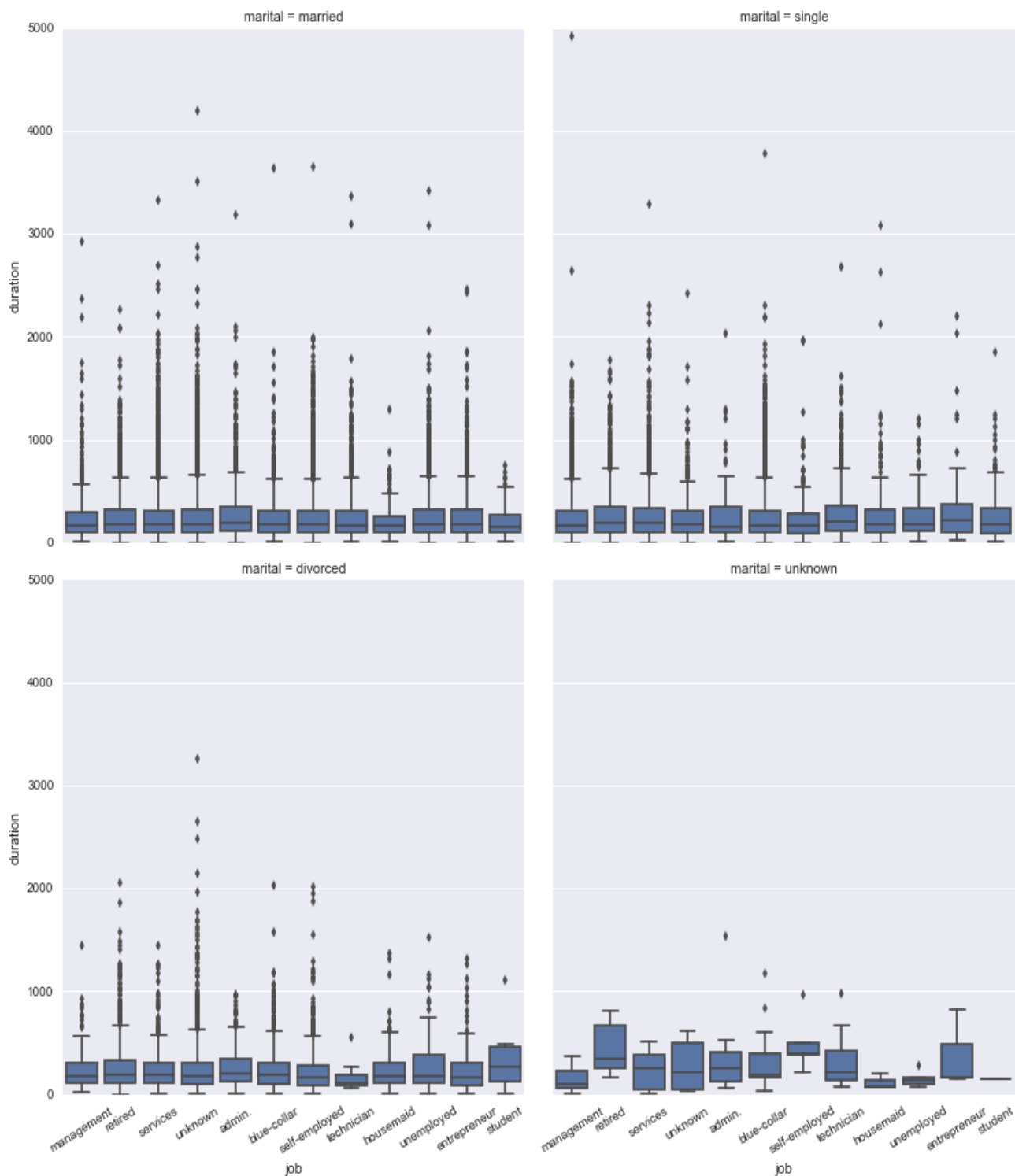


Fig. 2 Job vs duration split by marital status



Fig. 3 Scatterplot matrix of `age`, `duration`, `campaign` and `euribor3m`

We can see some interesting patterns in Fig. 3, there are splits in data for `euribor3m` and also note that the durations are smaller for age > 60 and other such relevant information. We could also check data distributions of each variable but I will leave it out for now and move ahead. We also applied scatterplot matrix on whole dataset to visualize inherent relationships between variables in the data shown in Fig. 4. But concluded that there was not much of a strong relationship.

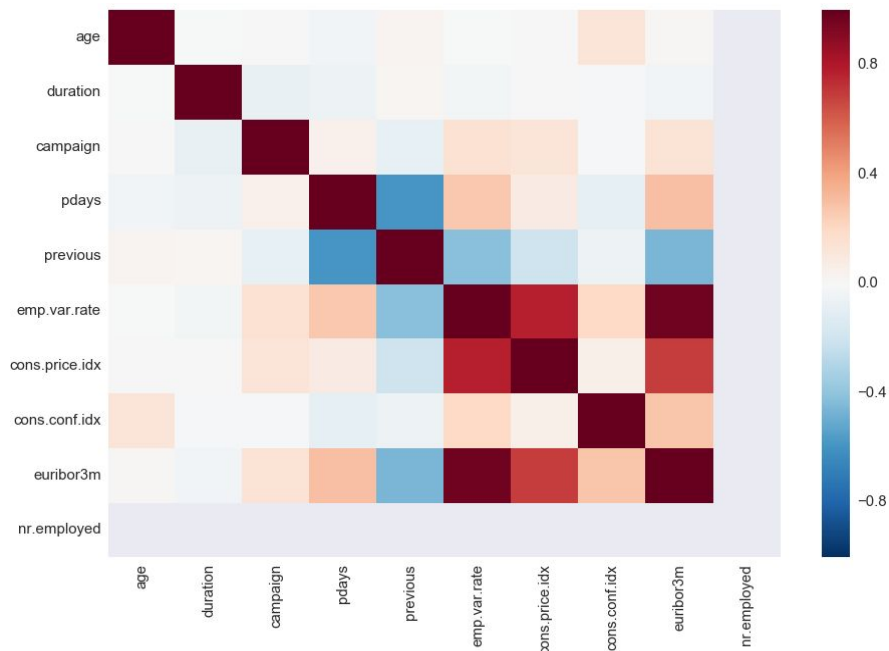


Fig. 4 Scatterplot matrix

Algorithms and Techniques

A Decision Tree is a robust and transparent Machine Learning model. The tree starts with a single node and then branches out, with a decision being made at every branch point. It can be used to predict whether a particular variable would have mattered in the customer's decision to subscribe or not to the bank's term deposit. The given dataset is a typical supervised learning problem for which tree type models perform a lot better than the rest. We don't know which algorithms would be fit for this problem or what configurations to use. So let's pick a few algorithms to evaluate.

- **Logistic Regression (LR):** is the appropriate regression analysis to conduct when the dependent variable is dichotomous (binary).
- **Classification and Regression Trees (CART):** The CART algorithm is structured as a sequence of questions, the answers to which determine what the next question, if any should be. The result of these questions is a tree like structure where the ends are terminal nodes at which point there are no more questions.
 - **Classification Trees:** where the target variable is categorical and the tree is used to identify the "class" within which a target variable would likely fall into.

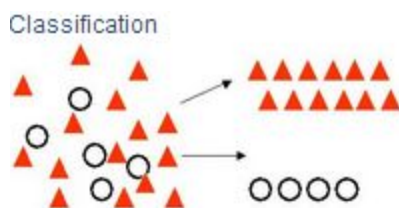


Fig. 5 Classification Tree

- **Regression Trees:** where the target variable is continuous and tree is used to predict it's value.



Fig. 6 Regression Tree

- **Random Forests (RF):** The idea behind RF, is to build many small Trees that use a random subset of the features and then combine them into a "Forest" of Trees. Each Tree by itself is a weak predictor, but when combining the many weak predictors, you can often (but not always!) get a stronger model.
- **Adaptive Boosting (AB):** AdaBoost is a type of "Ensemble Learning" where multiple learners are employed to build a stronger learning algorithm. AdaBoost works by choosing a base algorithm (e.g.

decision trees) and iteratively improving it by accounting for the incorrectly classified examples in the training set.

- **Extreme Gradient Boost (XGB):** The idea behind GBM is a more sophisticated. Roughly, the idea is to again, combine weak predictors. The trick is to find areas of misclassification and then “boost” the importance of those incorrectly predicted data points. And repeat. The result is a single Tree unlike RF.

XGBoost is a model based on tree ensemble which is a set of classification and regression trees (CART). XGBoost classifies the members of a family into different leaves and assigns scores on corresponding leaf. Usually, a single tree is a weak learner which is not strong enough to use in practice. Boosted trees are typically shallow which are high in bias and low in variance. A tree ensemble model sums prediction of multiple tree. Below is an example of tree ensemble of two trees that classifies whether someone likes play computer game (Fig.7).

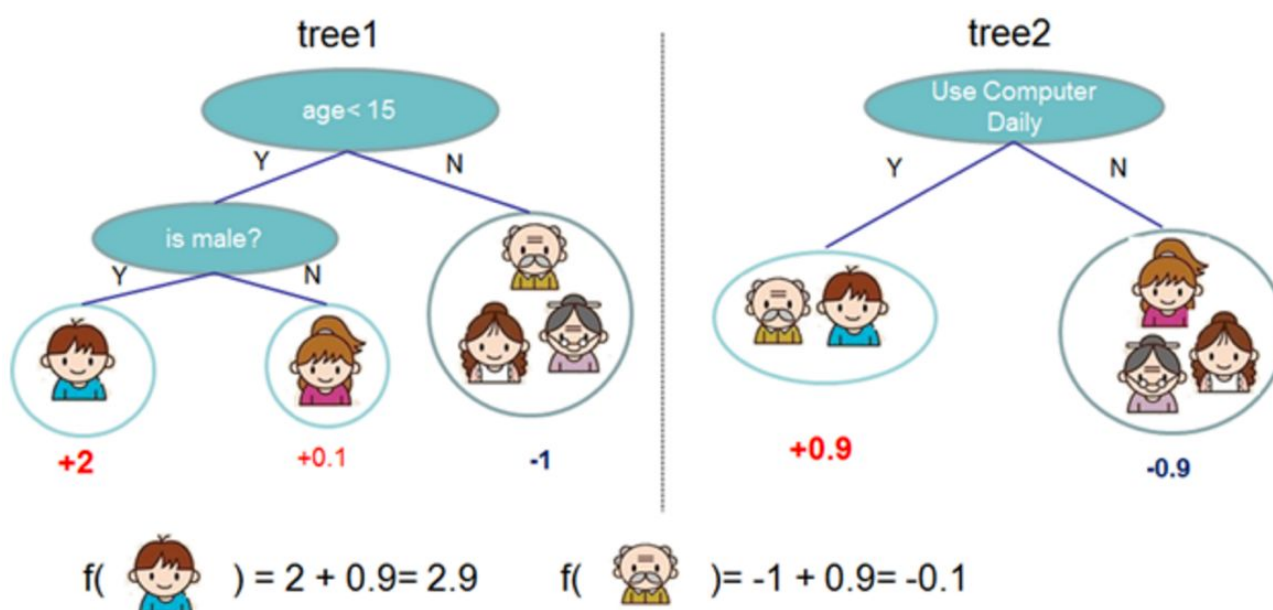


Fig. 7 Tree Ensemble Model. The final prediction for a given example is the sum of predictions from each tree

We are using 5-fold cross validation to estimate accuracy. This will split our dataset 5 parts, train on 4 and test on 1 and repeat for all combinations of train-test splits. Also, we are using the metric of accuracy to evaluate models. This is a ratio of the number of correctly predicted instances in divided by the total number of instances in the dataset multiplied by 100 to give a percentage (e.g. 95% accurate). We will be using the scoring variable when we run build and evaluate each model next.

Benchmark Model

Below table highlights performances of various models that were tried with their accuracies and errors. The standard Extreme Gradient Boosting (XGB) model with default parameters yields 91.4% accuracy on training data. So will consider it as benchmark and try to beat the benchmark with hyperparameter turning.

Algorithms	Accuracy	Std. Error
Logistic Regression	0.909611	0.005315
CART	0.888800	0.005325
Random Forest	0.912525	0.006937
AdaBoost	0.913218	0.007852
XGBoost	0.914640	0.007852

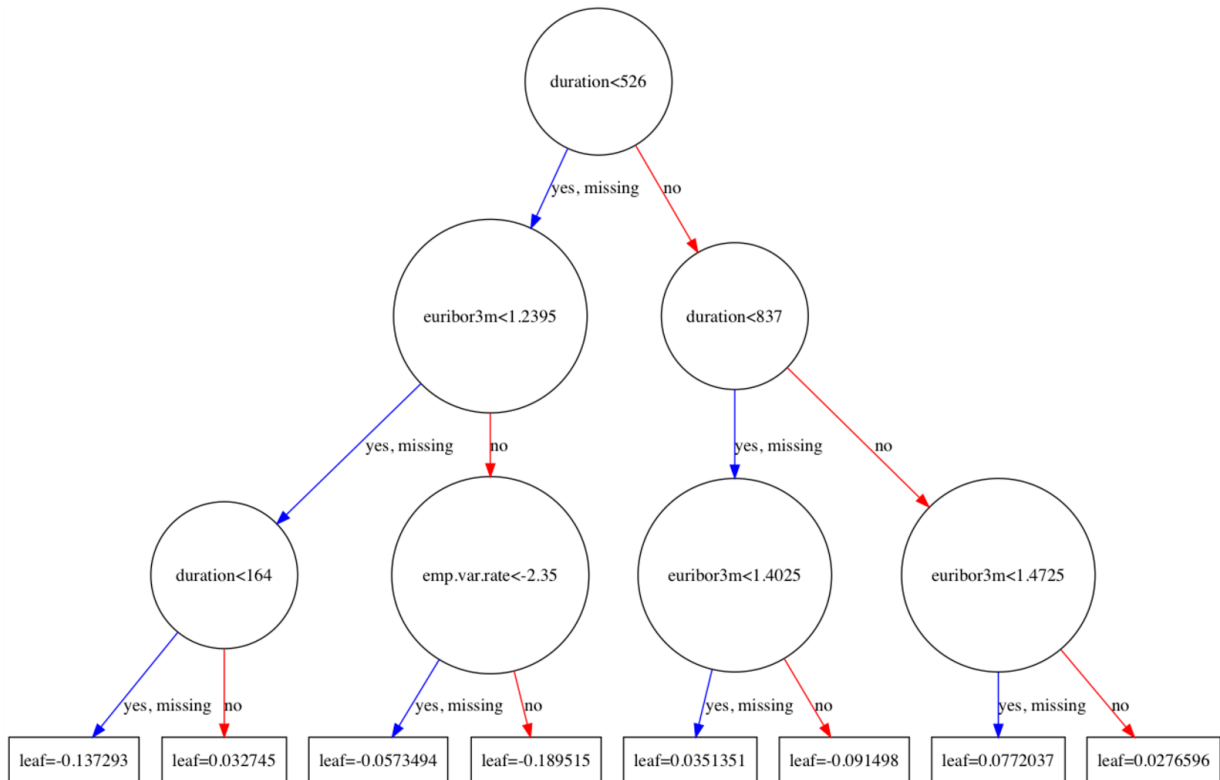


Fig. 8 Decision Tree

III. Methodology

Data Preprocessing

We will prepare the data by splitting feature and target/label columns and also check for quality of given data and perform data cleaning. To check if the model I created is any good, I will split the data into `training` and `validation` sets to check the accuracy of the best model. We will split the given `training` data in two ,70% of which will be used to train our models and 30% we will hold back as a `validation` set.

There are several non-numeric columns that need to be converted. Many of them are simply yes/no, e.g. housing. These can be reasonably converted into 1/0 (binary) values. Other columns, like profession and marital, have more than two values, and are known as categorical variables. The recommended way to handle such a column is to create as many columns as possible values (e.g. profession_admin, profession_blue-collar, etc.), and assign a 1 to

one of them and 0 to all others. These generated columns are sometimes called dummy variables, and we will use the `pandas.get_dummies()` function to perform this transformation.

Several Data preprocessing steps like preprocessing feature columns, identifying feature and target columns, data cleaning and creating training and validation data splits were followed and can be referenced for details in attached jupyter notebook.

Implementation

The project follows typical predictive analytics hierarchy as shown in Fig. 9:

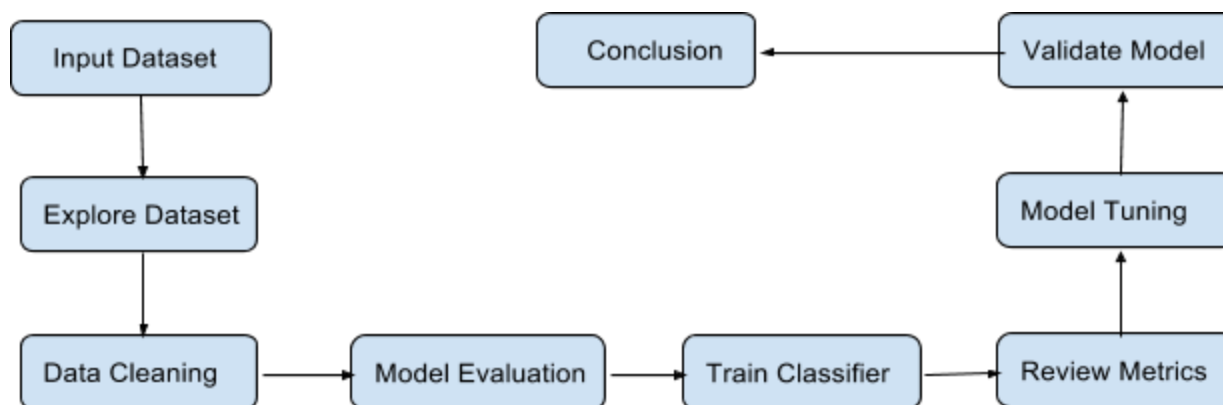


Fig. 9 Project Flow

Following the direction of the arrow as shown, with the dataset we chose, we first performed exploratory analysis by seeing distribution of some of the feature that seems to be relevant for the problem we plan to solve. We also did a check of correlations using scatterplot matrix in this process. We then performed data cleaning steps to identify columns with highest amount of null values and make a decision if we want to fill the missing values or completely drop the columns from the analysis. We processed the data through a pre-processing function to identify features that have values like 'yes', 'no', 'unknown' and convert them to 1, 0, `np.nan`. For categorical features we converted them to dummy variables so that the dataset becomes fully numeric and it then is easy to process it through our chosen algorithm for meaningful outcome.

Before we jump into the model evaluation, we split our full dataset into training and validation splits with a 70:30 ratio and then perform another check to make sure the splits have equal percent of responses to avoid deviations in the classification iterations. We used tree based algorithms to perform a check on accuracy and standard error metrics to pick a model that performs best. Typically Random Forest and XGBoost are well known to perform at best in usual supervised machine learning applications and we did confirm this by observing that XGBoost indeed has highest accuracy. One key observation we made was the accuracy of XGBoost sklearn wrapper was better compared to the XGBoost from dmlc library.

Once we benchmarked the accuracy of untuned XGBoost model, we then performed tuning of XGBoost model's hyperparameters in batches using GridSearchCV. This way we could keep track of best outcomes of each parameter and apply it to the next batch.

Refinement

We will prepare the data by splitting feature and target/label columns and also check for quality of given data and perform data cleaning. To check if the model we created is any good, we will split the data into training and validation sets to check the accuracy of the best model. We will split the given training data in two ,70% of which will be used to train our models and 30% we will hold back as a validation set.

We performed hyper tuning of parameters of XGBoost wrapper from scikit-learn library and the parameters tuned were shown in below table:

Parameter	Description	Values Tested	Best Value
max_depth	maximum depth of a tree to control overfitting	(3,4,5,6,7,8,9)	5
min_child_weight	minimum sum of weights of all observations required in a child	(2,3,4,5,6,7)	3
gamma	minimum loss reduction required to make a split	(0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.9,1.0)	0.1
reg_alpha	L1 regularization term on weights	(0,2,4)	2
reg_lambda	L2 regularization term on weights	(1,3,5)	5
subsample	Subsample ratio of the training instance	(0.5,1.0)	1
colsample_bytree	Subsample ratio of columns when constructing each tree	(0.5,1.0)	1
colsample_bylevel	Subsample ratio of columns for each split, in each level	(0.5,1.0)	1

Complete code of the simulation can be found here:

<https://github.com/tarunparmar/Udacity-ML-Capstone/blob/master/CapstoneProject.ipynb>

One refinement we could do was to improve recall rate for minority class labels (1 in our case) as the dataset was imbalanced. This may reduce the overall accuracy but the overall performance of model in predicting the outcome would be more robust.

IV. Results

Model Evaluation and Validation

We applied and compared XGB model out of the box vs the rest of tree based models. The metrics we used are calculated using sklearn wrapper so can be trusted for the model performance. Our end goal was to have a tuned model that could beat the untuned benchmark which it did by a very fine margin. So the solution described below is satisfactory to our initial expectations. We generated a final model with above list of tuned parameters. The output of this tuned model came just about 0.2% higher in accuracy vs the untuned model. Code snippet of final model shown below:

```
# Select best model
# Make predictions on validation dataset using tuned parameters
tuned_model = XGBClassifier(silent=True, nthread=-1, max_delta_step=0.7, seed=0, objective='reg:linear',
                           max_depth=5, min_child_weight=3, gamma=0.1, reg_alpha=2, reg_lambda=5, subsample=1,
                           colsample_bytree=1, colsample_bylevel=1)
tuned_fit = tuned_model.fit(X_train, y_train)
tuned_pred = tuned_model.predict(X_validation)

print("Accuracy Score: ", accuracy_score(y_validation, tuned_pred))
print("-----")
print("Confusion Matrix: \n", confusion_matrix(y_validation, tuned_pred))
print("-----")
print("Classification Report: \n", classification_report(y_validation, tuned_pred))
```

Accuracy Score: 0.920288095816

Confusion Matrix:

```
[[10615  350]
 [ 635  757]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.94	0.97	0.96	10965
1	0.68	0.54	0.61	1392
avg / total	0.91	0.92	0.92	12357

The recall rate is quite low as we discussed in above section. This can be improved by providing more weight to the positive labels `1`, this in turn decreases the overall accuracy but the model performance becomes robust to classify the outcome labels.

Justification

There is room for improvement on the final results, the tuned final model made no significant improvement over the untuned model. There could be more ways we could improve the score, particularly by selecting a subset of features using the ranking obtained from observing feature importance ranking and then performing the same exercise we did as describe above. But this will be out of scope of this report for now.

	Benchmark Model (Untuned XGBClassifier Model)	Final Model (Tuned XGBClassifier Model)
Accuracy Score	0.9184	0.9202

V. Conclusion

Free-Form Visualization

The importance of each feature was visualized with the following plot:

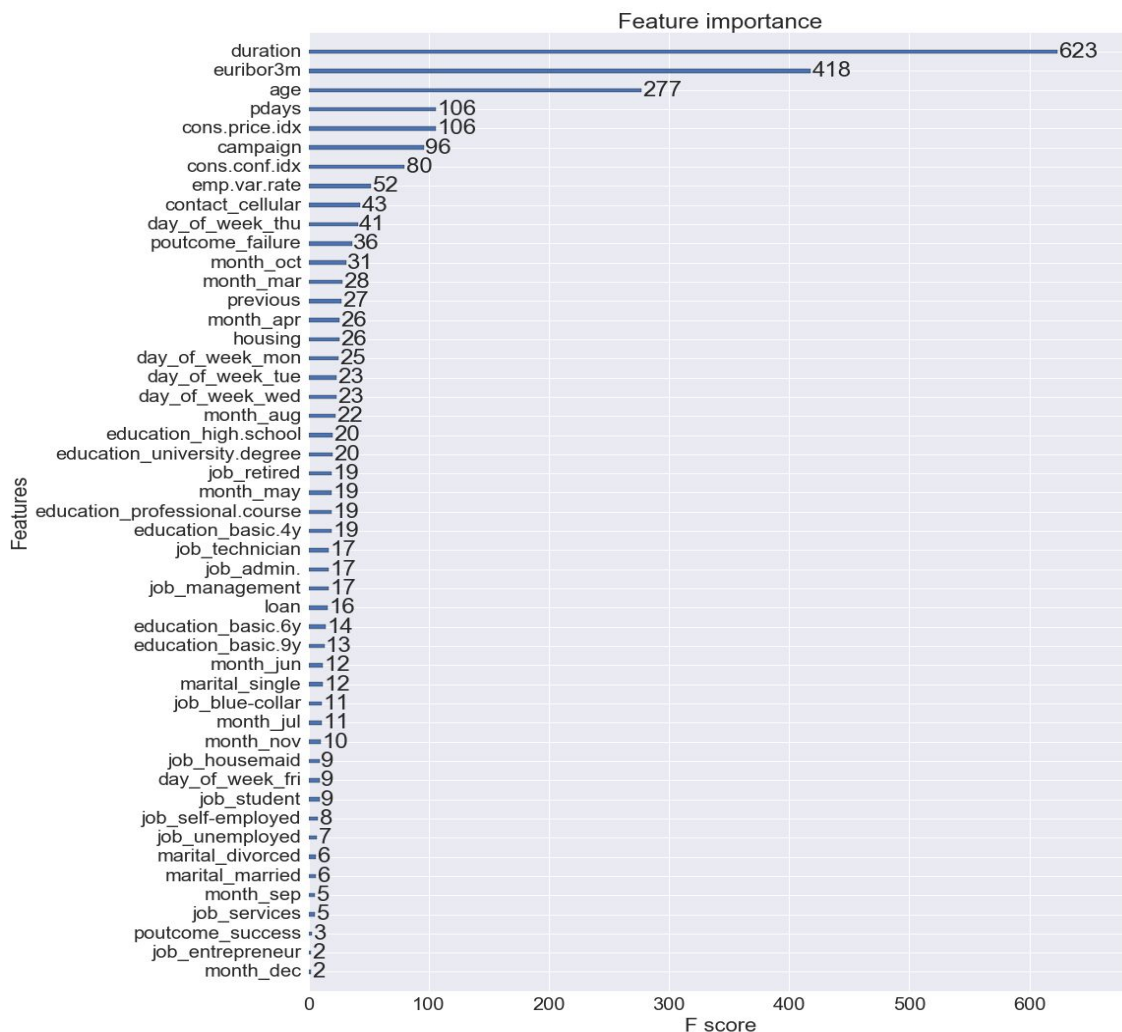


Fig.10 Important Feature ranking

Reflection

The most important and time consuming part of the problem was data cleansing and processing as there were multiple columns with null values. Once the data was prepared and ready, the next challenge was to pick an algorithm that could be best suited for the problem we choose to solve. With experience and

prior knowledge and also the outcome of accuracy on training data, we observed that XGBoost performed the best out of others.

If we look at the important predictors in their order of importance, as per the tree shown in Fig. 10, they are: duration, euribor3m, pdays, age, cons.price.idx, campaign, emp.var.rate, contact_cellular, and so on. We know characteristics of an entity are not in control of bank or any company manager as they are individual characteristics. Removing them from important predictors leaves us with duration, pdays, previous, housing, contact, and loan as option. Among these duration can be enhanced and in this case it's the most important predictor as well. Based on data analyzed and domain following are few suggestions for improving duration.

- Outbound call creates negative attitude towards bank due to the intrusion of privacy. Bank should decrease the outbound call rate and use inbound calls for cross-selling intelligently to increase the duration of the call. Agents may pitch about profits of term deposit for a particular client during inbound calls.
- Bank should have a clear promotion and they should put the value proposition up front – waiving fees, offering something free, or promoting a bundled service at a discounted price increases response rates versus just informing customers about a product. The value proposition should be right up front, with a clear call to action.
- “Duration” has positive effect on people saying “yes”. This is because the longer the conversations on the phone, the higher interest the customer will show to the term deposit. The Bank ought to focus on the potential clients who have significant call duration and moreover who have reacted emphatically amid the past campaign.

Improvements

One of the improvements we could do was to perform tuning to improve recall rate to improve overall prediction performance of the model. The other improvement could be to work with a subset of features that are high in variable importance pareto. Another key aspect would be to review metrics like log-loss to understand how quick the model is able to tune.

References

- [1] Ou, C., Liu, C., Huang, J. and Zhong, N. ‘One Data mining for direct marketing’, Springer-Verlag Berlin Heidelberg, pp. 491–498., 2003.
- [2] http://en.wikipedia.org/wiki/Direct_marketing. Wikipedia has a tool to generate citations for particular articles related to direct marketing.
- [3] O'guinn, Thomas.” Advertising and Integrated Brand Promotion”. Oxford Oxfordshire: Oxford University Press. p. 625. ISBN 978-0-324-56862-2. , 2008.

- [4] Petrison, L. A., Blattberg, R. C. and Wang, P. 'Database marketing: Past present, and future', Journal of Direct Marketing, 11, 4, 109–125, 1997.
- [5] Eniafe Festus Ayetiran, "A Data Mining-Based Response Model for Target Selection in Direct Marketing", I.J.Information Technology and Computer Science, 2012, 1, 9-18.
- [6] <https://archive.ics.uci.edu/ml/datasets/bank+marketing>