

A Project report on

*A data driven approach to
predict the success of bank
telemarketing.*

In Foundation of Financial Data Sciences (FE- 582- GROUP 9)

By

NIKHIL NIRHALE

VITEJ BARI

SENUJ JAIN

DEEPANSHU TYAGI

TABLE OF CONTENTS

1. ABSTRACT.....	
2. INTRODUCTION.....	
3. PACKAGES	
1. RATTLE.....	
2. RMINER.....	
4. DATA MINING METHODS.....	
1. LOGISTIC REGRESSION.....	
2. DECISION TREES.....	
3. SVM.....	
4. ARTIFICIAL NEURAL NETWORK.....	
5. APPLICATION OF DATA MINING METHODS TO THE DATA.....	
1. SCREENSHOTS OF THE CODE.....	
6. PREDICTION ANALYSIS AND OBSERVATIONS.....	
7. CONCLUSIONS.....	
8. REFERENCES.....	

ABSTRACT

We propose a data mining (DM) approach to predict the success of telemarketing calls for selling bank long-term deposits. A Portuguese retail bank was addressed, with data collected from 2008 to 2013, thus including the effects of the recent financial crisis. We analysed a large set of 150 features related with bank client, product and social-economic attributes. A semi-automatic feature selection was explored in the modelling phase, performed with the data prior to July 2012 and that allowed to select a reduced set of 22 features. We are comparing four DM models: logistic regression, decision trees (DTs), neural network (NN) and support vector machine. We will be selecting the best model for this data using advanced metrics using the Rminer Package.

INTRODUCTION

Marketing selling campaigns constitute a typical strategy to enhance business. Companies use direct marketing when targeting segments of customers by contacting them to meet a specific goal. Centralizing customer remote interactions in a contact centre eases operational management of campaigns.

Technology enables rethinking marketing by focusing on maximizing customer lifetime value through the evaluation of available information and customer metrics, thus allowing us to build longer and tighter relations in alignment with business demand.

Decision support systems (DSSs) use information technology to support managerial decision making. There are several DSSs sub-fields, such as personal and intelligent DSSs. Personal DSSs are related with small-scale systems that support a decision task of one manager, while intelligent DSSs use artificial intelligence techniques to support decisions. Another related DSS concept is Business Intelligence (BI), which is an umbrella term that includes information technologies, such as data warehouses and data mining (DM), to support decision making using business data. DM can play a key role in personal and intelligent DSSs, allowing the semi-automatic extraction of explanatory and predictive knowledge from raw data. Classification is the most common DM task and the goal is to build a data driven model that learns an unknown underlying function that maps several input variables, which characterize an item (e.g., bank client), with one labelled output target (e.g., type of bank deposit sell: “failure” or “success”).

In this project, we propose a personal and intelligent DSS that can automatically predict the result of a phone call to sell long term deposits by using a DM approach. Such DSS is valuable to assist managers in prioritizing and selecting the next customers to be

contacted during bank marketing campaigns. For instance, by using a Lift analysis that analyses the probability of success and leaves to managers only the decision on how many customers to contact. Consequently, the time and costs of such campaigns, would be reduced. Also, by performing fewer and more effective phone calls, client stress and intrusiveness would be diminished.

This research focus on targeting through telemarketing phone calls to sell long-term deposits. Within a campaign, the human agents execute phone calls to a list of clients to sell the deposit (outbound) or, if meanwhile the client calls the contact-centre for any other reason, he is asked to subscribe the deposit (inbound). Thus, the result is a binary unsuccessful or successful contact.

This study considers real data collected from a Portuguese retail bank, from May 2008 to June 2013, in a total of 49,732 phone contacts. For evaluation purposes, a time ordered split was initially performed, where the records were divided into training (four years) and test data (one year). The training data is used for feature and model selection and includes all contacts executed up to June 2012, in a total of 45,211 examples. The test data is used for measuring the prediction capabilities of the selected data-driven model, including the most recent 4521 contacts, from July 2012 to June 2013.

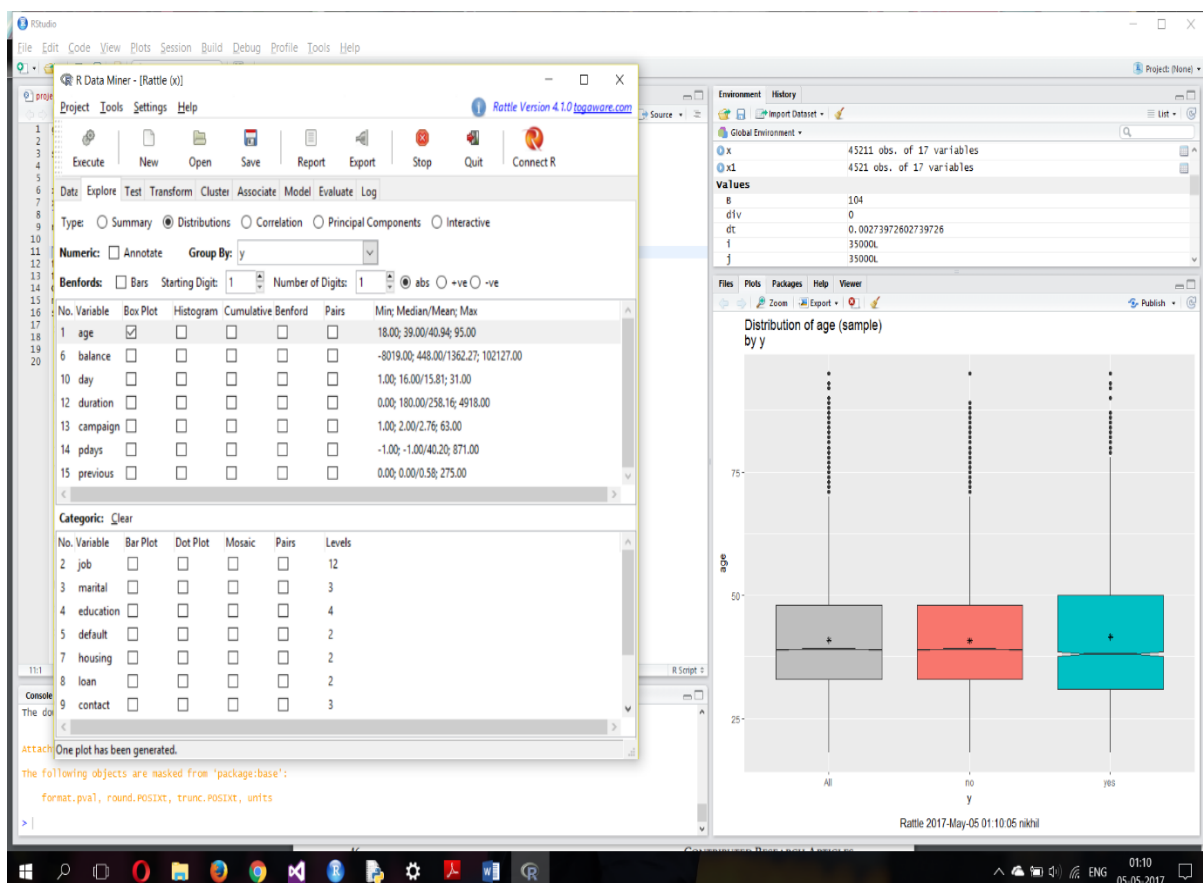
DATA SUMMARY AND EDA ALREADY INCLUDED IN PROGRESS REPORT.

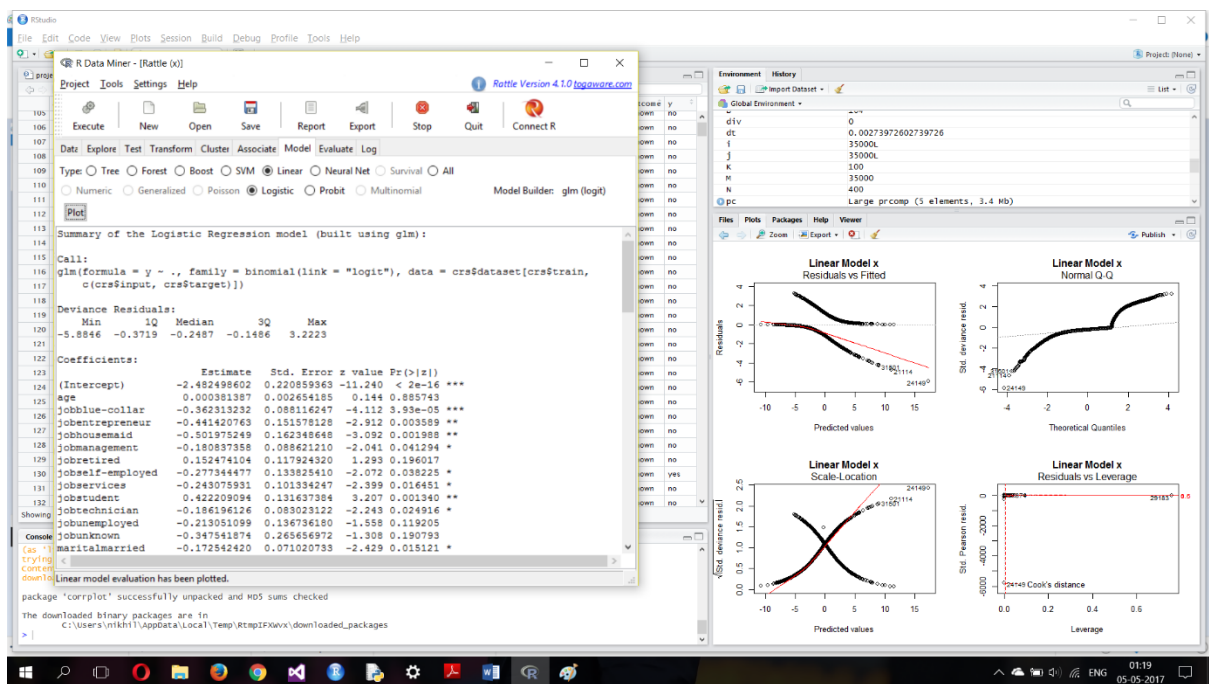
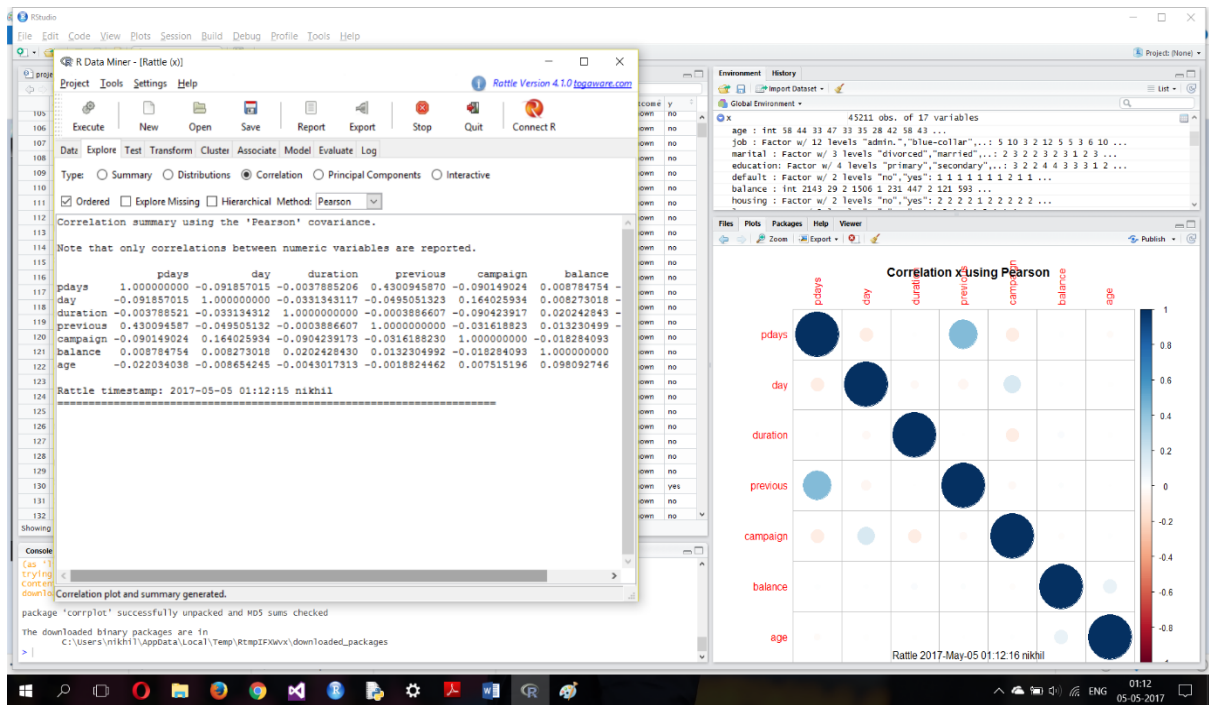
PACKAGES

1.RATTLE

Rattle (the R Analytical Tool to Learn Easily) is a graphical data mining application written in and providing a pathway into (Williams,2009b). It has been developed specifically to ease the transition from basic data mining, as necessarily offered by GUIs, to sophisticated data analyses using a powerful statistical language.

In our Project we have used this package to conduct some basic EDA to look for outliers and missing values. We have conducted this analysis in our Progress report to conclude that our data does not include any missing values and effect of outliers is specifically not very relevant in this case. Below are some screenshots of the RATTLE GUI.





2.RMINER

This package facilitates the use of data mining algorithms in classification and regression (including time series forecasting) tasks by presenting a short and coherent set of functions. It is a very powerful package and simplifies the data mining task by many bounds.

For e.g. To Fit a SVM or NN model on a data set the data needs to scale to mean of 0 and standard deviation of 1. In many R packages this must be done manually by writing relevant set of codes. But Rminer does this for us inherently, and is done specifying a parameter to choose the scaling.

Set Of Functions provided by Rminer :-

CasesSeries-Create a training set (data. Frame) from a time series using a sliding window.

Crossvaldata-Computes k-fold cross validation for rminer models.

Delevels-Reduce (delete) or replace levels from a factor variable (useful for preprocessing datasets).

Fit-Fit a supervised data mining model (classification or regression) model

Holdout-Computes indexes for holdout data split into training and test sets.

Importance- Measure input importance (including sensitivity analysis) given a supervised data mining model.

Imputation-Missing data imputation (e.g. substitution by value or hotdeck method).

Lforecast-Compute long term forecasts.

Mgraph-Mining graph function

Mining-Powerful function that trains and tests a particular fit model under several runs and a given validation method

Mmetric-Compute classification or regression error metrics.

Mparheuristic-Function that returns a list of searching (hyper)parameters for a particular classification or regression model

predict.fit-predict method for fit objects (rminer)

rminer-Internal-Internal rminer Functions

sa_fri1-Synthetic regression and classification datasets for measuring input importance of supervised learning models

savemining- Load/save into a file the result of a fit (model) or mining functions.

sin1reg-sin1 regression dataset

vecplot-VEC plot function (to use in conjunction with Importance function).

We would be using the Functions Highlighted in Red.

DATA MINING METHODS

1.LOGISTIC REGRESSION.

Logistic regression is a method for fitting a regression curve, $y = f(x)$, when y is a categorical variable. The typical use of this model is predicting y given a set of predictors x . The predictors can be continuous, categorical or a mix of both.

Unlike actual regression, logistic regression does not try to predict the value of a numeric variable given a set of inputs. Instead, the output is a probability that the given input point belongs to a certain class.

Logistic Regression Model seeks to:

- model the probability of an event occurring depending on the values of the independent variables, which can be categorical or numerical.
- estimate the probability that an event occurs for a randomly selected observation versus the probability that an event does not occur.
- predict the effect of a series of variables on a binary response variables.
- classify observations by estimating the probability that an observation in a particular category.

While using our data, we divide the data in training set and testing set. Then we model our Logistic regression model using the training data set. This model considers all the independent variables which are responsible for the dependent variables. Then we predict the dependent variables and check its accuracy with the testing data.

The better the prediction the more the accuracy we get from modelling the Logistic Regression.

The central premise of Logistic Regression is the assumption that your input space can be separated into two nice 'regions', one for each class, by a linear boundary. By Linear boundary we mean that, for two dimensions, it's a straight line and no curving. For three

dimensions, it's a plane. And so on. But for this to make sense, the data points must be separable into the two aforementioned regions by a linear boundary. If your data points do satisfy this constraint, they are said to be linear-separable. Then the model depending on the independent variables, decided the placement of the dependent variable and predict the data.

2. DECISION TREES.

Decision tree algorithm partitions the data samples into two or more subsets so that the samples within each subset are more homogeneous than in the previous subset. Decision trees are of two main types, one is Classification tree, and the other is Regression tree. For our dataset, we used Classification tree. In the classification setting, the classification error rate is used as a criterion for making the binary splits.

DECISION-TREE FOR THIS DATA.

main features that characterize maximum margin algorithm: a non-linear function is learned by linear learning machine mapping into high dimensional kernel induced feature space. The capacity of the system is controlled by parameters that do not depend on the dimensionality of feature space.

In SVM regression, the input \mathbf{x} is first mapped onto a m -dimensional feature space using some fixed (nonlinear) mapping, and then a linear model is constructed in this feature space. Using mathematical notation, the linear model (in the feature space) $f(\mathbf{x}, \mathbf{w})$ is given by

$$f(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^m w_j g_j(\mathbf{x}) + b$$

where $g_j(\mathbf{x}), j = 1, \dots, m$ denotes a set of nonlinear transformations, and b is the “bias” term. Often the data are assumed to be zero mean (this can be achieved by preprocessing), so the bias term is dropped.

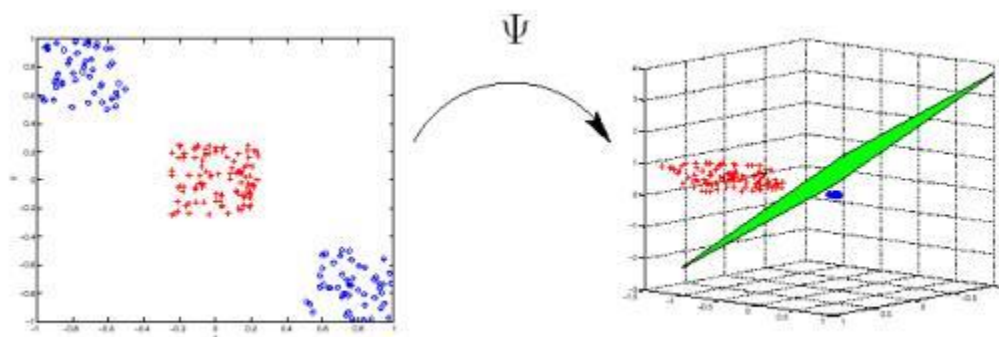


Figure. Non-linear mapping of input examples into high dimensional feature space.

The quality of estimation is measured by the loss function $L(y, f(\mathbf{x}, \mathbf{w}))$. SVM regression uses a new type of loss function called ε -insensitive loss function proposed by Vapnik :

$$L_{\varepsilon}(y, f(\mathbf{x}, \mathbf{w})) = \begin{cases} 0 & \text{if } |y - f(\mathbf{x}, \mathbf{w})| \leq \varepsilon \\ |y - f(\mathbf{x}, \mathbf{w})| - \varepsilon & \text{otherwise} \end{cases}$$

4. ARTIFICIAL NEURAL NETWORK.

Artificial neural networks (ANNs) or connectionist systems are a computational model used in machine learning, computer science and other research disciplines, which is based on a large collection of connected simple units called artificial neurons, loosely analogous to axons in a biological brain. Connections between neurons carry an activation signal of varying strength.

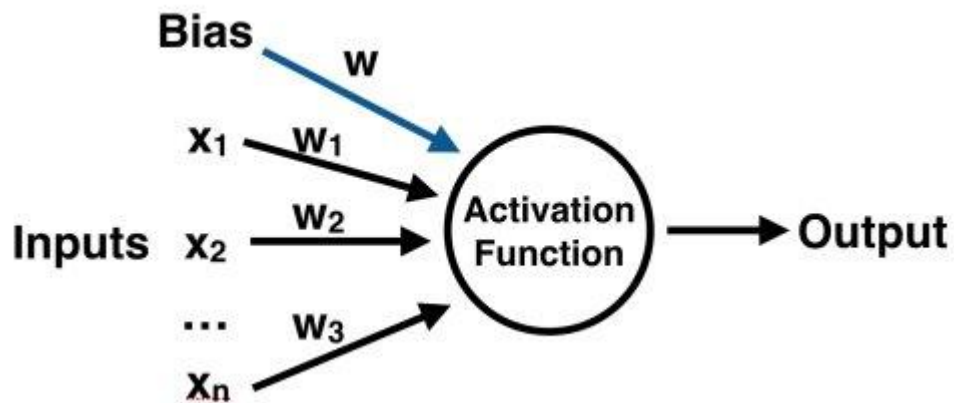
Such systems can be trained from examples, rather than explicitly programmed, and excel in areas where the solution or feature detection is difficult to express in a traditional computer program.

The signals and state of artificial neurons are real numbers, typically between 0 and 1. There may be a threshold function or limiting function on each connection and on the unit itself, such that the signal must surpass the limit before propagating. Back propagation is the use of forward stimulation to modify connection weights, and is sometimes done to train the network using known correct outputs. However, the success is unpredictable: after training, some systems are good at solving problems while others are not. Training typically requires several thousand cycles of interaction.

Neural Networks are a machine learning framework that attempts to mimic the learning pattern of natural biological neural networks. Biological neural networks have interconnected neurons with dendrites that receive inputs, then based on these inputs they produce an output signal through an axon to another neuron. We will try to mimic this process through the use of Artificial Neural Networks (ANN), which we will just refer to as neural networks from now on. The process of creating a neural network begins with the most basic form, a single perceptron.

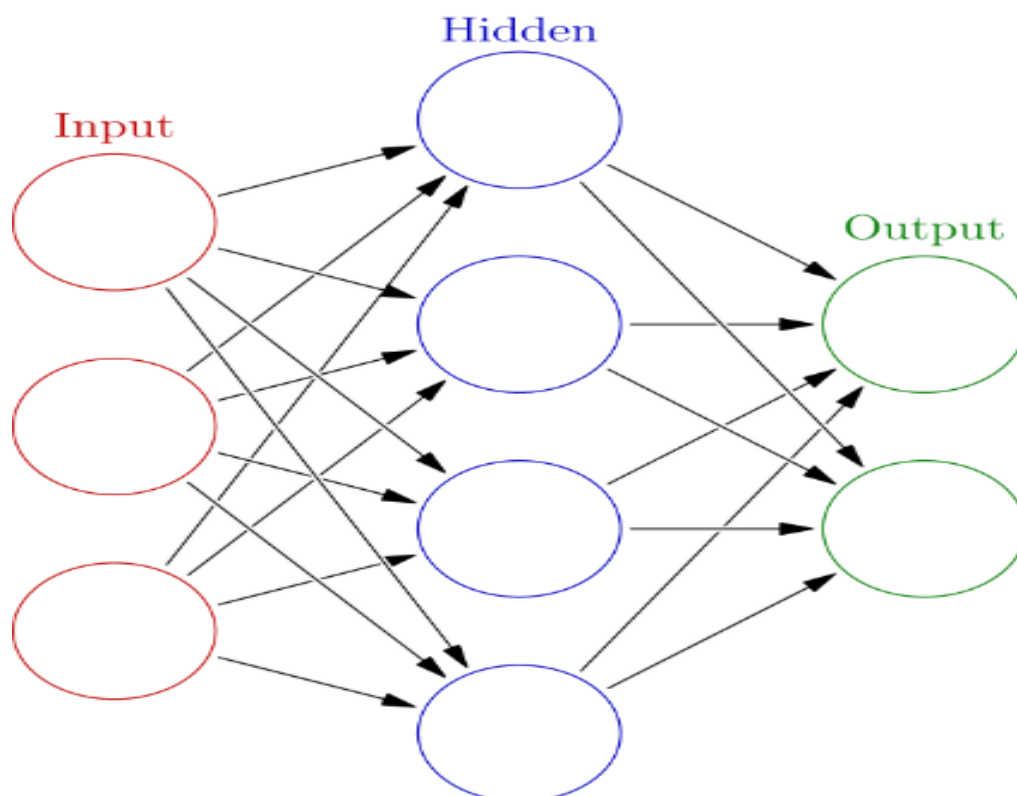
To create a neural network, we simply begin to add layers of perceptron's together, creating a multi-layer perceptron model of a

neural network. You'll have an input layer which directly takes in your feature inputs and an output layer which will create the resulting outputs. Any layers in between are known as hidden layers because they don't directly "see" the feature inputs or outputs.



A PERCEPTRON MODEL (ABOVE)

A NEURAL NETWORK WITH MULTIPLE PERCEPTRONS (BELOW)



APPLICATION OF DATA MINING METHODS TO THE DATA

SCREENSHOTS OF THE CODE

```
1 getwd()
2 rm(list=ls())
3 setwd("F:/Study/Financial data science/project")
4
5
6 x<-read.csv("bank-full.csv",sep = ";")
7 x1<-read.csv("bank.csv",sep=";")
8
9 library(rattle)
10 library(rminer)
11
12 fit_1<-fit(y~.,data = x,model = "mlp",task="p",scale="inputs")
13 fit_1_1<-fit(y~.,data = x,model = "mlp",task="c",scale="inputs")
14 f_NN<-predict(fit_1,newdata=x1)
15 f_NN_1<-predict(fit_1_1,newdata=x1)
16 CF_NN<-data.frame(x1$y,nn=f_NN)
17 CF_NN_1<-data.frame(x1$y,nn=f_NN_1)
18 table(real=CF_NN_1[,1],NN=CF_NN_1[,2])
19 savemodel(fit_1,"NN")
20
21
22 fit_2<-fit(y~.,data = x,model = "svm",task="p",scale="inputs")
23 fit_2_2<-fit(y~.,data = x,model = "svm",task="c",scale="inputs")
24 f_SVM<-predict(fit_2,newdata=x1)
25 f_SVM_2<-predict(fit_2_2,newdata=x1)
26 CF_SVM<-data.frame(x1$y,svm=f_SVM)
27 CF_SVM_2<-data.frame(x1$y,svm=f_SVM_2)
28 table(real=CF_SVM_2[,1],SVM=CF_SVM_2[,2])
29 savemodel(fit_2,"SVM")
30
31
32 fit_3<-fit(y~.,data = x,model = "lr",task="p")
33 fit_3_3<-fit(y~.,data = x,model = "lr",task="c")
34 f_LR<-predict(fit_3,newdata=x1)
35 f_LR_3<-predict(fit_3_3,newdata=x1)
36 CF_LR<-data.frame(x1$y,lr=f_LR)
37 CF_LR_3<-data.frame(x1$y,lr=f_LR_3)
38 table(real=CF_LR_3[,1],LR=CF_LR_3[,2])
39 savemodel(fit_3,"LR")
40
41
42 fit_4<-fit(y~.,data = x,model = "dt",task="p")
43 fit_4_4<-fit(y~.,data = x,model = "dt",task="c")
44 f_DT<-predict(fit_4,newdata=x1)
45 f_DT_4<-predict(fit_4_4,newdata=x1)
46 CF_DT<-data.frame(x1$y,dt=f_DT)
47 CF_DT_4<-data.frame(x1$y,dt=f_DT_4)
48 savemodel(fit_4,"DT")
49 table(real=CF_DT_4[,1],DT=CF_DT_4[,2])
50
51
52
53 #-----ROC-----
54 ROC_NN<-mmetric(CF_NN[,1],CF_NN[,c(2,3)],"ROC")
55 AUC_NN<-ROC_NN$roc$auc
56 ROC_NN<-ROC_NN$roc$roc
```



```

project.R % Final Results % imp_table % x1 % x %
Source on Save
53 #-----ROC-----
54 ROC_NN<-mmetric(CF_NN[,1],CF_NN[,c(2,3)],"ROC")
55 AUC_NN<-ROC_NN$roc$auc
56 ROC_NN<-ROC_NN$roc$roc
57 plot(ROC_NN[,1],ROC_NN[,2],type="l",col="green",xlab="FPR",ylab="TPR")
58 ROC_SVM<-mmetric(CF_SVM[,1],CF_SVM[,c(2,3)],"ROC")
59 AUC_SVM<-ROC_SVM$roc$auc
60 ROC_SVM<-ROC_SVM$roc$roc
61 lines(ROC_SVM[,1],ROC_SVM[,2],type="l",col="red")
62 ROC_LR<-mmetric(CF_LR[,1],CF_LR[,c(2,3)],"ROC")
63 AUC_LR<-ROC_LR$roc$auc
64 ROC_LR<-ROC_LR$roc$roc
65 lines(ROC_LR[,1],ROC_LR[,2],type="l",col="blue")
66 ROC_DT<-mmetric(CF_DT[,1],CF_DT[,c(2,3)],"ROC")
67 AUC_DT<-ROC_DT$roc$auc
68 ROC_DT<-ROC_DT$roc$roc
69 lines(ROC_DT[,1],ROC_DT[,2],type="l",col="orange")
70 abline(0,1)
71 legend(0.6,0.4,c("NN","SVM","LR","DT"),fill = c("green","red","blue","orange"))
72 #-----LIFT-----
73 LIFT_NN<-mmetric(CF_NN[,1],CF_NN[,c(2,3)],"LIFT")
74 AUC_LIFT_NN<-LIFT_NN$lift$area
75 LIFT_NN<-LIFT_NN$lift$lift
76 plot(LIFT_NN[,1],LIFT_NN[,2],type="l",col="green",xlab="SAMPLE SIZE",ylab="RESPONSES")
77 LIFT_SVM<-mmetric(CF_SVM[,1],CF_SVM[,c(2,3)],"LIFT")
78 AUC_LIFT_SVM<-LIFT_SVM$lift$area
79 LIFT_SVM<-LIFT_SVM$lift$lift
80 lines(LIFT_SVM[,1],LIFT_SVM[,2],type="l",col="red")
81 LIFT_LR<-mmetric(CF_LR[,1],CF_LR[,c(2,3)],"LIFT")
82 AUC_LIFT_LR<-LIFT_LR$lift$area
83 LIFT_LR<-LIFT_LR$lift$lift
84 lines(LIFT_LR[,1],LIFT_LR[,2],type="l",col="blue")
85 LIFT_DT<-mmetric(CF_DT[,1],CF_DT[,c(2,3)],"LIFT")
86 AUC_LIFT_DT<-LIFT_DT$lift$area
87 LIFT_DT<-LIFT_DT$lift$lift
88 lines(LIFT_DT[,1],LIFT_DT[,2],type="l",col="orange")
89 abline(0,1)
90 legend(0.6,0.4,c("NN","SVM","LR","DT"),fill = c("green","red","blue","orange"))
91
92 #-----IMPORTANCE ANALYSIS-----
93 IMP_NN<-Importance(fit_1,data = x,method = "SA")
94 IMP_NN_DSA<-Importance(fit_1,data = x,method = "DSA")
95 IMP_DT<-Importance(fit_4,data = x,method = "SA")
96 IMP_SVM<-Importance(fit_2,data = x,method = "SA")
97 IMP_LR<-Importance(fit_3,data = x,method = "SA")
98 vecplot(IMP_NN_DSA,sort = "decreasing",graph = "VEC",xlab = "value range of attributes",ylab = "pr
99 sum(IMP_DT$value)
100 sum(IMP_SVM$value)
101 sum(IMP_LR$value)
102 IMP_NN_DSA$data
103 sum(IMP_NN$value)
104 IMP_NN$data
105 imp_table<-cbind.data.frame(features=colnames(x),relative_importance=IMP_NN$imp,SA=IMP_NN$value)
106 write.csv(imp_table,"relativeimportanceandSA.csv")
107
108

```

PREDICTION ANALYSIS AND OBSERVATIONS

CONFUSION TABLES

NN CONFUSION TABLE

	NN	
REAL	NO	YES
NO	3856	144
YES	230	291

ACCURACY: 91.72%

SVM CONFUSION TABLE

	SVM	
REAL	NO	YES
NO	3936	64
YES	319	202

ACCURACY: 91.71%

DT CONFUSION TABLE

	DT	
REAL	NO	YES
NO	3896	104
YES	348	173

ACCURACY: 90%

LR CONFUSION TABLE

	LR	
REAL	NO	YES
NO	3912	88
YES	352	169

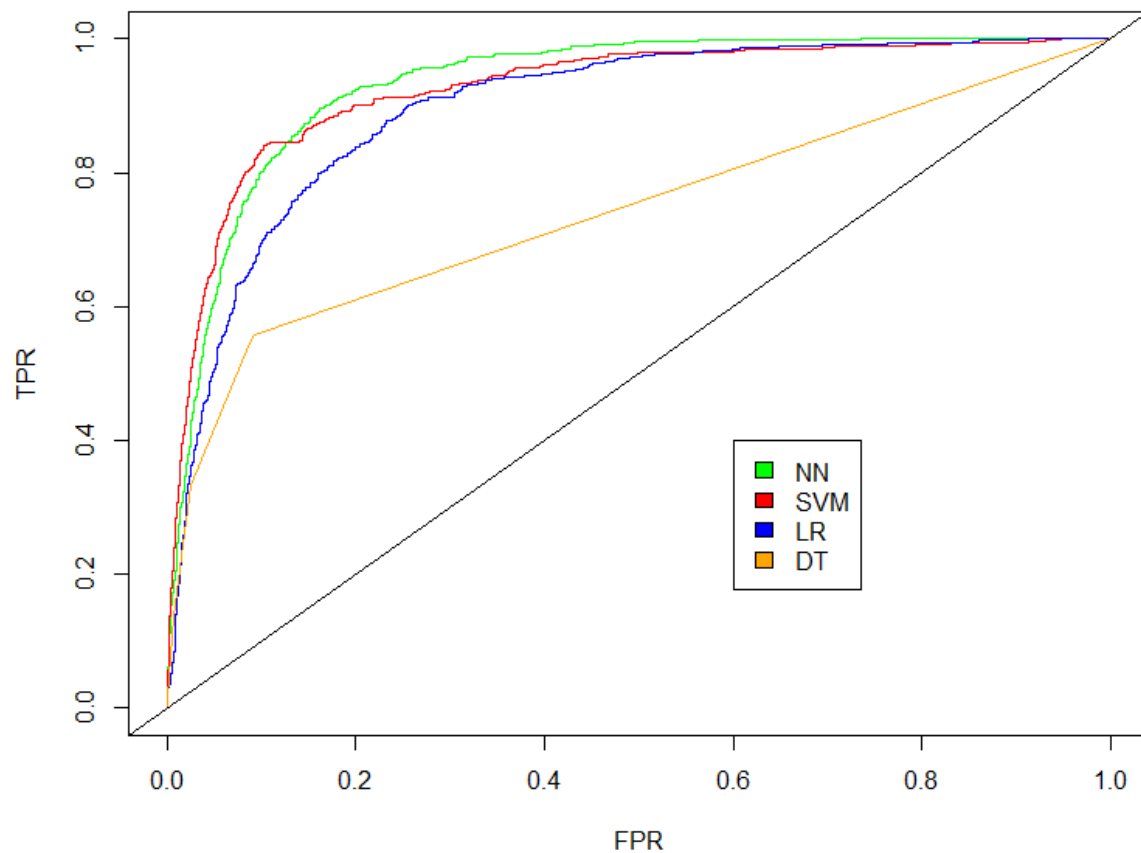
ACCURACY: 90%

As we could see that the accuracy rates are shown by the confusion tables are quite good for all the methods, which creates ambiguity. We know that DT and LR models are very rigid for classifying such uncorrelated and nonlinear data. We are going to use some advanced metrics to highlight this fact.

ROC AND LIFT CURVE ANALYSIS

1. ROC (Receiver Operating characteristics) CURVE

The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The true-positive rate is also known as sensitivity, recall or probability of detection in machine learning. The false-positive rate is also known as the fall-out or probability of false alarm and can be calculated as $(1 - \text{specificity})$. The ROC curve is thus the sensitivity as a function of fall-out.



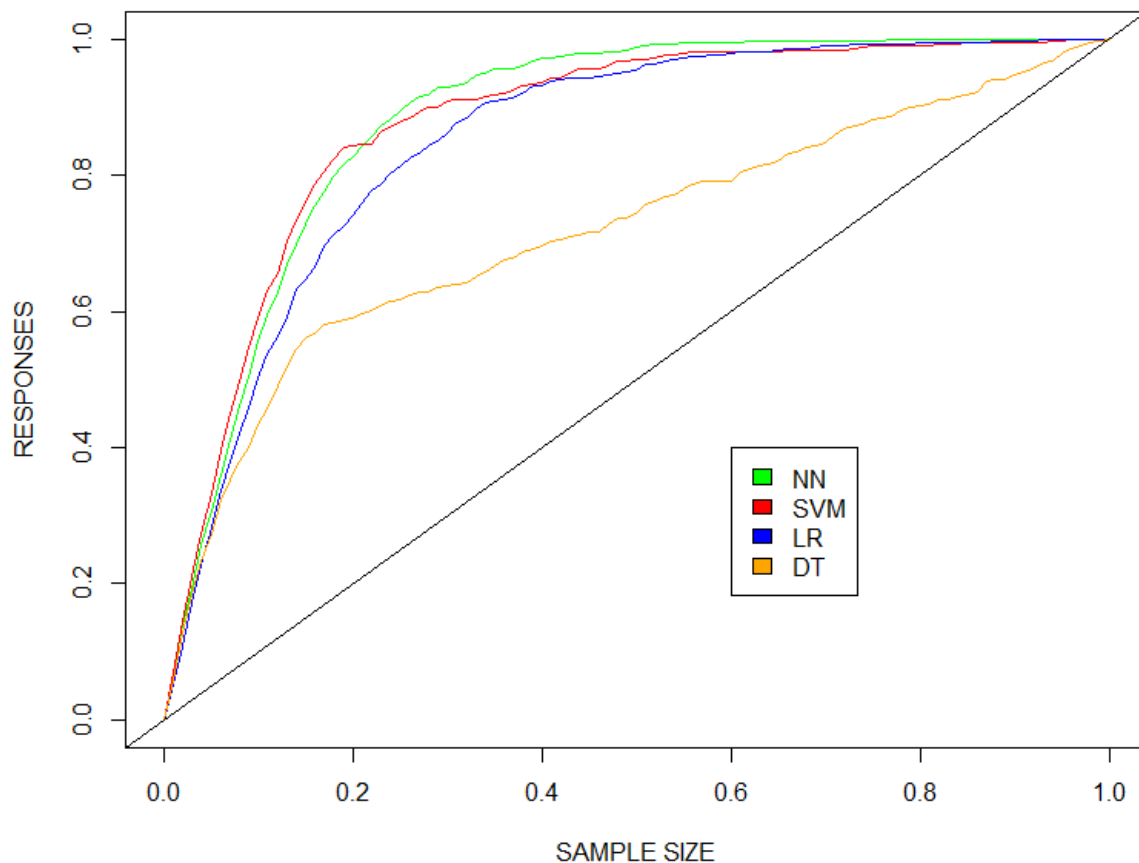
AREA UNDER THE CURVE FOR ALL MODELS

MODEL	NN	SVM	LR	DT
AUC	0.9325019	0.9250221	0.8969952	0.741935

2. LIFT CURVE

In data mining and association rule learning, lift is a measure of the performance of a targeting model (association rule) at predicting or classifying cases as having an enhanced response (with respect to the population as a whole), measured against a random choice targeting model. A targeting model is doing a good job if the response within the target is much better than

the average for the population as a whole. Lift is simply the ratio of these values: target response divided by average response.



AREA UNDER THE CURVE FOR ALL MODELS

MODEL	NN	SVM	LR	DT
AUC	0.882524	0.8759981	0.8510077	0.7177063

Looking the ROC AND LIFT metrics we could conclude that the NN model is the best classification method for this model. But to cement this fact from the input point of view we conduct sensitivity analysis.

3. SENSITIVITY AND IMPORTANCE ANALYSIS

1. RELATIVE IMPORTANCE TABLE

FEAUTURES	RELATIVE_IMPORTANCE
AGE	0.019478
JOB	0.012609
MARITAL	0.003256
EDUCATION	0.008671
DEFAULT	0.00224
BALANCE	0.004295
HOUSING	0.088887
LOAN	7.01E-05
CONTACT	0.002457
DAY	0.053672
MONTH	0.160348
DURATION	0.130027
CAMPAIGN	0.006945
PDAYS	0.005068
PREVIOUS	0.409708
POUTCOME	0.092268

Relative importance as the name suggest gives relative attribute weights for the model. And this means it creates forcefull attribute co-relations, which says that the model creates a perfect fit for the attributes as the sum is 1, but which never is the case. That is why we study non-relative Sensitivity analysis.

2. RELATIVE IMPORTANCE TABLE WITH SENSITIVITY ANALYSIS

FEATURES	RELATIVE IMPOTANCE	SENSITIVITY ANALYSIS
AGE	0.019478	0.01614
JOB	0.012609	0.010448
MARITAL	0.003256	0.002698
EDUCATION	0.008671	0.007185
DEFAULT	0.00224	0.001856
BALANCE	0.004295	0.003559
HOUSING	0.088887	0.073654
LOAN	7.01E-05	5.81E-05
CONTACT	0.002457	0.002036
DAY	0.053672	0.044474
MONTH	0.160348	0.132868
DURATION	0.130027	0.107743
CAMPAIGN	0.006945	0.005755
PDAYS	0.005068	0.004199
PREVIOUS	0.409708	0.339494
POUTCOME	0.092268	0.076455
SUM	1	0.8286241

SVM SA SUM- 0.8082777

LR SA SUM- 0.7590074

DT SA SUM- 0.2792326

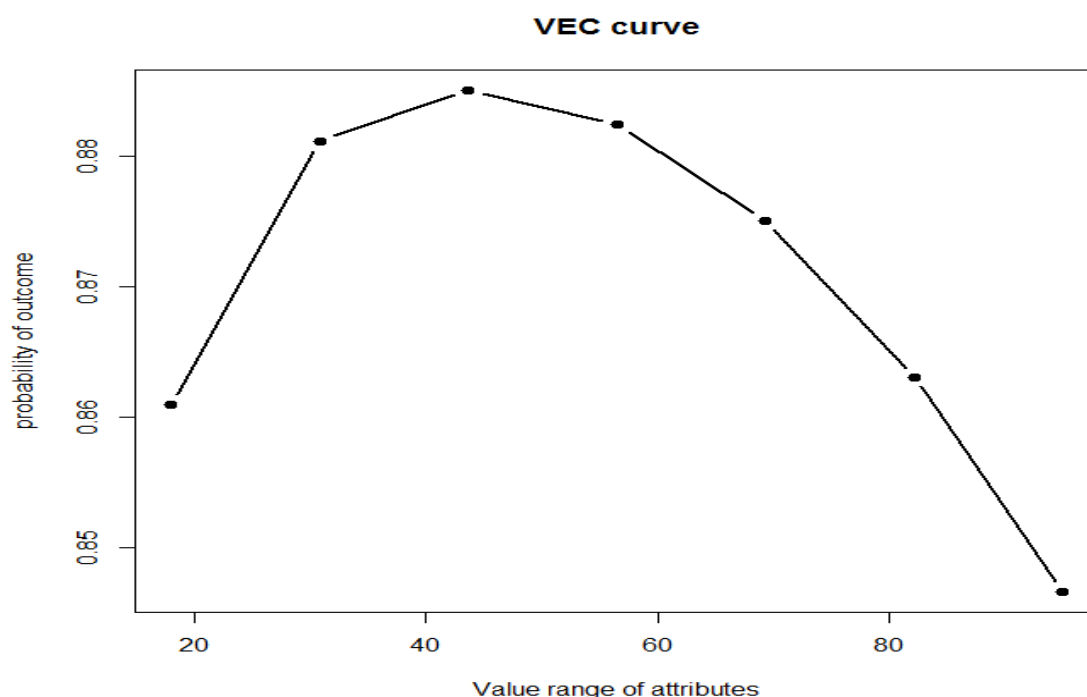
As we could see the SA sum is the also the highest for NN model i.e. if we keep all the attributes zero ,even then there is a chance that we will get an outcome of this model. The lower this chance the better the model fits the data.

And from the remaining observations DT is most surely not a good fit for our data because of the very low sensitivity sum.

Hence looking at the two tables above which reflect the feature sensitivity of the NN model, we could conclude which factors are to be considered while making the calls. For, example previously contacted customers who said yes are highly likely to give a positive response to this telemarketing drive and consequently secure the bank deposit with a positive outcome.

4. VEC-VARIABLE EFFECTIVE CURVE

The VEC curve shown above shows the probability of outcome for the top 6 features



CONCLUSIONS

The obtained results are credible for the banking domain and provide valuable knowledge for the telemarketing campaign manager. In our observations we concluded a good classification model and applied sensitivity analysis to get the features or attributes which effect the outcome the most.

The result of this is that the company could give valuable time and importance to a customer base which is largely going to affect their business. Such analysis could lead to huge savings in time and money for the relevant institution. Which is why institutions have started giving increasing importance to the power of data analytics, which is the stepping stone to AI.

REFERENCES

<https://www.wikipedia.org/>

<https://cran.r-project.org/web/packages/rminer/rminer.pdf>

<https://cran.r-project.org/web/packages/rattle/rattle.pdf>

[\[Moro et al., 2014\] S. Moro, P. Cortez and P. Rita. A Data-Driven Approach to Predict the Success of Bank Telemarketing. Decision Support Systems, Elsevier, 62:22-31, June 2014](#)

<http://www.kdnuggets.com/2016/08/beginners-guide-neural-networks-r.html>