

Final Project: Loan Approval

Data 621 Business Analytics and Data Mining

Aadi Kalloo, Nathan Lim, Asher Meyers, Daniel Smilowitz, Logan Thomson

Due July 21, 2016

Contents

Abstract	1
Introduction	2
Literature Review	2
Methodology	3
Results	7
Discussion	9
References	11
Appendix – R Code	12

Abstract

Regression and classification techniques are used to predict loan approval status of customers of a fictional bank given fourteen predictors. Data was obtained from a “data challenge”, and the given variables were supplemented with transformations. Models were trained on a dataset of 614 cases, and subsequently used to predict the binary outcome of 367 test cases. Techniques utilized include logistic regression, ridge regression, the Adaboost algorithm, and a neural network. Analyses were performed to determine the optimal number of iterations and optimal number of nodes for the Adaboost algorithm and neural network, respectively. Cross validation results were not linearly co-variant with accuracy scores on the test data set. Using tools supplied by the data challenge, accuracy scores were used to compare models and the Adaboost algorithm was found to produce the best results. The optimal model created yielded an accuracy of 0.8056.

Keywords: loan prediction, classification, logistic regression, adaboost

Introduction

Between 2009-2014, approximately 95 million mortgage loan applications were filed in the United States alone (Home Mortgage Disclosure Act, 2015). When further loan types, including education and personal loans, are considered, it becomes even clearer that the ability to accurately identify reliable borrowers is a high priority for banks and other financial institutions today. Likewise, consumers are significantly affected as well. Surveys and analyses performed by Lusardi & Scheresberg (2013) has shown that individuals of poor financial literacy, and subsequently poor credit history, are declined by banks for loans and mortgages, and turn to high-cost borrowing (i.e. payday loans). Thus, it becomes clear that the factors that determine which individuals are dependable borrowers is information important to both financial institutions and consumers alike.

The dataset of interest contains information about customers of a hypothetical bank and was obtained from a data challenge hosted at AnalyticsVidhya. This data set and challenge were chosen due to its relevance to the current loan-centric economy outlined above, and to address the need for more transparent analyses of optimal methods for financially-related prediction models.

Literature Review

Although predictors can vary by loan or income level, studies investigating loan approval predictors have found that factors such as Marital Status, Loan Duration, and Number of Dependents can serve as statistically significant predictors of loan approval status (Agbemava et al., 2016). However, as loan approval criteria can vary greatly between financial institutions there is a significant lack of published materials investigating the effects that various factors may play on loan approval status. Furthermore, the majority of financial institutions use proprietary algorithms and statistical methods in order to determine loan approvals.

It should be noted that ratio of training set cases to test set cases is 1.67:1. This naturally leads to markedly decreased algorithm training than the 4:1 ratio that is normally used in prediction tasks (Dobin & Simon, 2011).

Of significant interest to us were the concepts of additive logistic regression or “boosting” as outlined by Friedman et al. (2012). The original concept of boosting was introduced by Schapire (1990) and describes a method of increasing the performance of “weak” classifiers by combining many of them into a more powerful model (Friedman et al., 2012). The authors of this paper introduce the Adaboost algorithm, which is used and analyzed in detail for this project. The Adaboost algorithm is characterized by the use of weights on the target variable, increasing the weight on each iteration for misclassified values (Friedman et al., 2012).

Given that this project was centered around a data challenge, it was also of particular significance to us to research data mining and machine learning methods that could potentially provide greater accuracy scores, despite being outside the scope of IS621 Data Mining. The primary method of prediction in this arena is the neural network. While the foundational ideas of artificial neural network date as far back as the 1940s (Warren & Pitts, 1943), the procedures used by the **neuralnet** package are based on Riedmiller and Braun (1993). As a brief overview, the **neuralnet** package calculates the following function as described by Gunther & Fritsch (2010):

$$o(x) = f(w_0 + \sum_{j=1}^J w_j \cdot f(w_{0j} + \sum_{i=1}^n w_{ij}x_i))$$

where w_0 is the intercept of the output neuron, w_{0j} is the intercept of the j^{th} hidden neuron, and w_j is the synaptic weight. This mathematical model can be expressed visually, and is depicted in the following figure:

Given the fact that loan prediction has traditionally been a field of interest for those who can benefit from financial gain (i.e. big banks, as opposed to academia), research and methodological studies in this area is sparse. Most of the current literature in this area has been focused on financial solvency issues, namely bankruptcy. It has been found previously that logistic regression, neural networks, and classification trees are the three most commonly used data mining methods in the prediction of bankruptcy (Olson et al., 2012). As a means of extending this work to the arena of personal loans, this project serves to compare some of the most commonly used methods of prediction in loan outcomes. We believe that the information and methods investigated here can trickle down to a

smaller scale and be beneficial to smaller economies such as microlending, where loans tend to be oriented around civilian populations instead of financial institutions (Conlin, 1999).

Methodology

The dataset has 614 rows (each representing a customer) and 11 predictor variables. Furthermore, there is 1 identification variable, and 1 response variable: **Loan_Status**, a binary categorical variable representing whether each customer has been approved for a loan given their credentials. The goal of the data challenge and of this investigation is to predict the approval of loan applications based on the predictors.

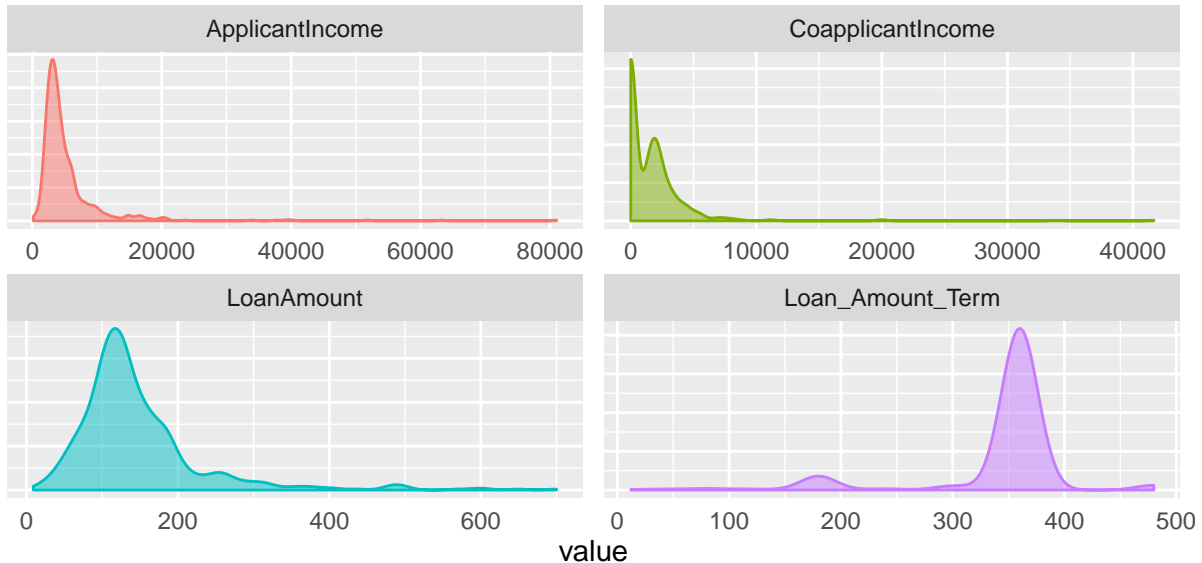
Data Exploration

The classes of the predictor and response variables are presented below:

	Class	Levels/Range
Gender	numeric	0 - 1
Married	numeric	0 - 3
Dependents	numeric	0 - 1
Education	numeric	0 - 1
Self_Employed	integer	150 - 81000
ApplicantIncome	numeric	0 - 41667
CoapplicantIncome	integer	9 - 700
LoanAmount	integer	12 - 480
Loan_Amount_Term	integer	0 - 1
Credit_History	factor	3
Property_Area	numeric	0 - 1
Loan_Status	numeric	0 - 1

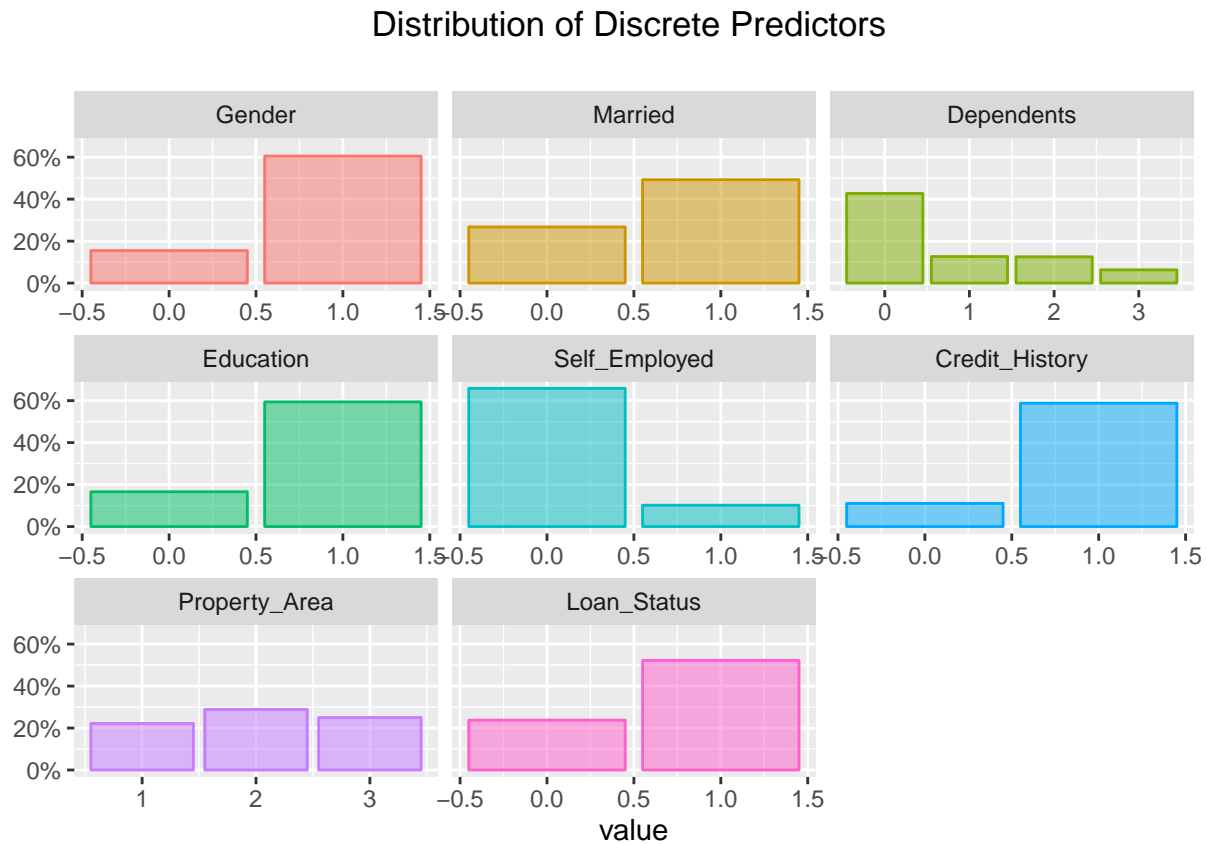
It can be seen that four of the predictors are continuous variables; the remaining 7 predictors, as well as the response, are categorical variables, with all but one of these variables (**Dependents**) being binary. The distributions of the continuous predictors are presented below:

Distribution of Continuous Predictors



It is clear that each of the continuous variables exhibits a strong level of skewness – **ApplicantIncome**, **CoapplicantIncome**, and **LoanAmount** are right-skewed, while **Loan_Amount_Term** is left-skewed. It can be seen that the peak in the distribution of **Loan_Amount_Term** is located at 360 – this corresponds to 30 years, the typical length of a mortgage.

The distributions of the discrete predictors, as well as the target, are presented below:



The last figure in the chart shows that roughly 75% of loan applicants are approved for a loan.

With the exception of **Property_Area**, each of the discrete predictors exhibits a dominant category. The distributions indicate that, with all other variables held constant, typical applicants:

- Are male
- Are married
- Have no children
- Are graduates
- Are not self-employed
- Meet the firm's credit requirements

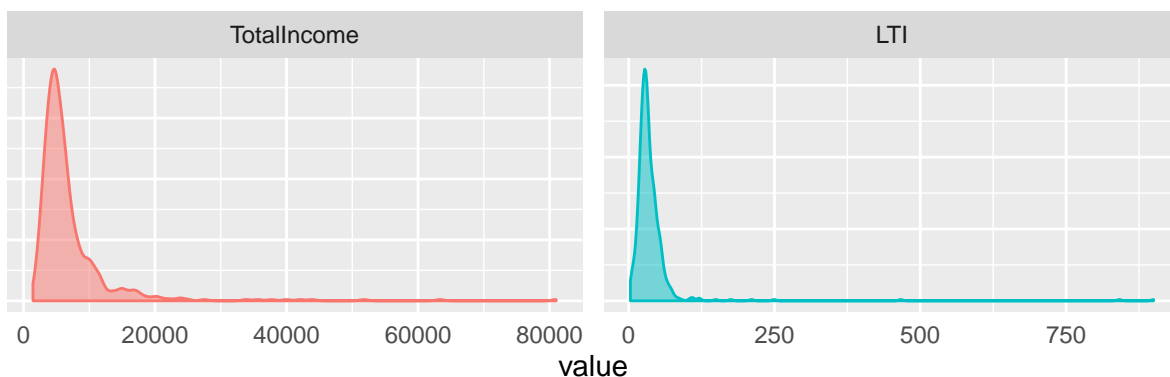
Inspection of the raw data set reveals that there are a number of missing values; these are addressed prior to the development of predictive models.

Data Manipulation

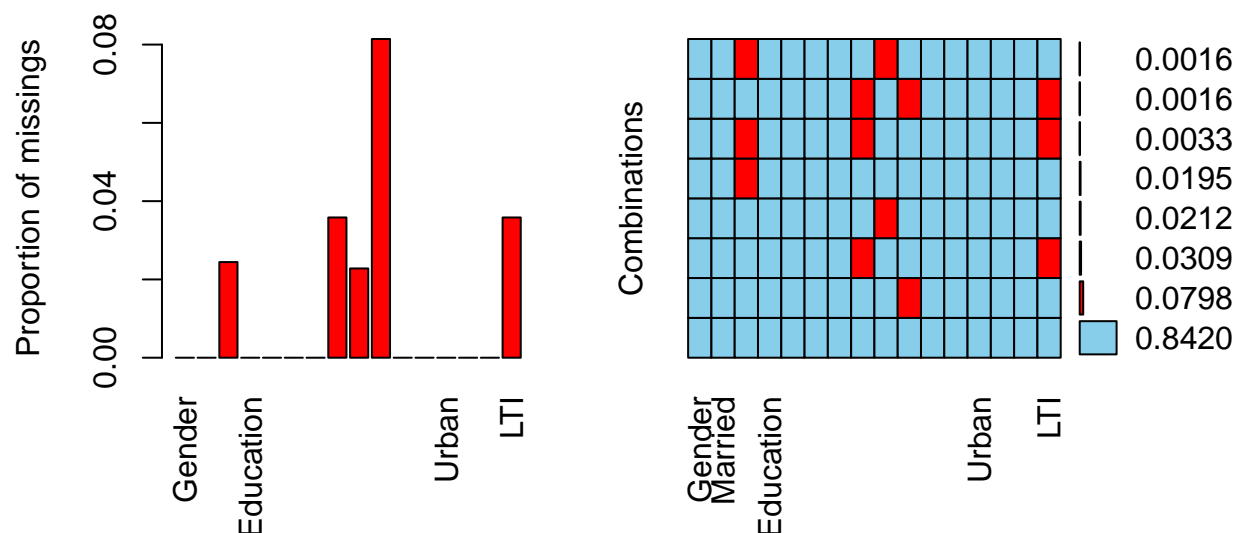
The supplied training data set originally has 11 predictor variables, six of which are character vectors which needed to be converted to binary categorical variables. The response variable, stated as a simple “Y/N” was also converted into the 0 or 1 format for modelling. As it was deemed not significant and also not chosen by any automated model, **Rural** is excluded from the **Property_Area** predictor, leaving only **Urban** and **SemiUrban**, which are turned into separate binary categorical predictors. **Property_Area** is then excluded from the training dataset.

All other variables in the data set were numeric, or already in a binary format. Given the included variables, two new features were created. **TotalIncome** is the sum of applicant and co-applicant incomes, and **LTI** (Loan to Income Ratio) is the loan amount ($\times 1000$) divided by the **ApplicantIncome** variable. The distributions of these created variables are shown below; they exhibit a similar degree of right-skewness as the three continuous predictors used to create them.

Distribution of Created Predictors



Missing values were also present in the data, either in the form of empty characters (“”), or NAs. Four of the predictors, **Gender**, **Married**, **Dependents**, and **Self_Employed** contained empty characters. Being only a small proportion of the cases ($< 1\%$ to 8%), these were in a way imputed when the predictors were converted into binary categorical variables. Four other variables, **LoanAmount**, **Loan_Amount_Term**, **Credit_History**, and **LTI** all contained small numbers of NAs. These are either dropped from the dataset, or imputed with the median, depending on the model. Alternative imputation methods, such as filling with a weighted “coin toss” value yielded less accurate results in the models. A representation of the proportion of missing values is presented below:



Model Creation

In total, six models were created and used for prediction of the target variable. A tool was supplied by AnalyticsVidhya that reported accuracy (% of cases correctly predicted) upon submission of the test data predictions.

Logistic regression is a very common method of classifying data (Thomas, 2000). In logistic regression, a linear regression is used where the dependent variable is a non-linear function of the probability of the event occurring. In this project, the purpose of the logistic regression method is to classify cases as either “Approved for loan” or “Denied”. The general logistic regression model can be represented as:

$$\ln \frac{p_i}{1 - p_i} = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n + \varepsilon$$

A few different logistic regression models were developed for this project. The first was a full model using all variables. This model made use of the `glm` function. This model was created as a basic point of comparison. This model produced an AIC of 935.

The second model created utilized both the `glm` function and the `leaps` package in order to find the best combination/subset of variables that minimizes the Bayesian Information Criterion (BIC). The “best subsets” approach is a method of automated variable selection that seeks to find the best predictors of the target variable. The `regsubsets` function as part of the `leaps` package found the `Married`, `Credit_History` and `SemiUrban` variables to be the best subset of predictors for this data set. The AIC reported for this model was 518.

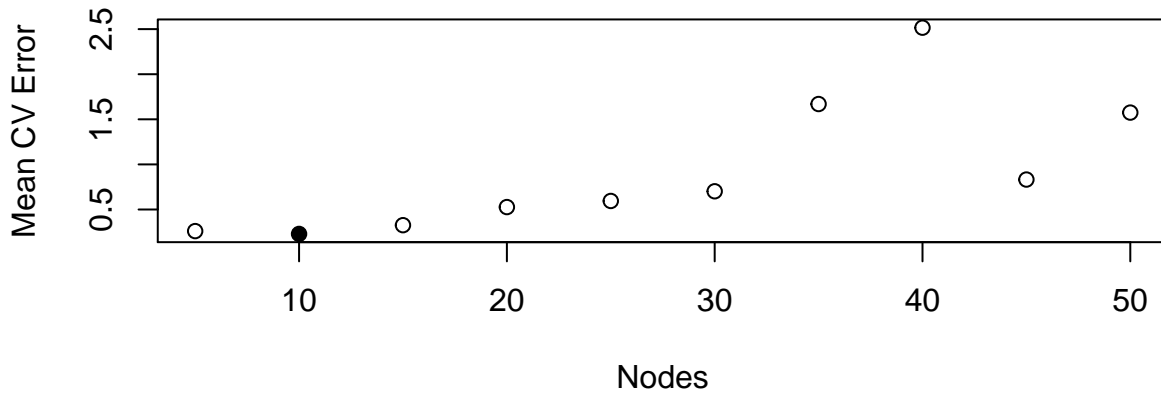
Ridge regression methods were combined with logistic regression to create a “Logistic Ridge” model, the third of this project, using the `ridge` library. As this library is no longer actively in development, it was manually obtained and installed. Like Ridge Regression, Logistic Ridge Regression introduces a tuning parameter as a penalty to avoid overfitting (Cule & De Iorio, 2012). This is useful as the ridge regression method aims to reduce variance and produce lower mean squared error for higher-dimensional data such as that used in the full model for these data. The ridge penalty used in the `ridge` package as outlined by Cule & De Iorio (2012) is given as:

$$k = \frac{p}{\hat{\beta}'\hat{\beta}}$$

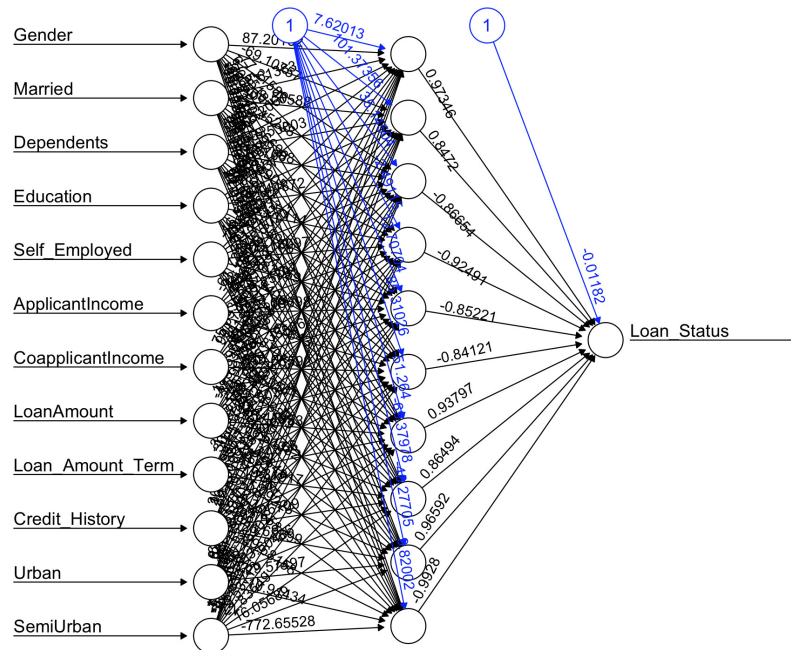
The Adaboost algorithm was used for this project in two forms. The first was through the use of the `ada` package. For this second implementation, the effect of the number of iterations of the algorithm on the final accuracy of the test set was analyzed with iterations in the range of 50-525 with steps of 25. It is hypothesized that accuracy will gradually increase with iterations to a peak, then gradually decrease as overfitting starts to become a prominent issue. The second was a manually coded implementation based on the algorithmic outline provided by Friedman et al. (2000). A manual approach was used for the purposes of comparing the academically published to the packaged algorithm. Given that these algorithms work by combining basic models, it is believed that they will perform better than a single logistic regression model alone.

A simple neural network method was pursued. The neural network approach was performed using normalized data with a single hidden layer. The number of nodes to use in the hidden layer was determined using a 10-fold Cross Validation Approach. A plot of the mean CV error against the number of hidden nodes is shown below:

Determination of Number of Nodes for neuralnet



From the above data, a network of 10 nodes was chosen due to its minimal cross validation error. This resulted in a network where the input variables were each sent to the 10 hidden nodes where weights were applied based on the variable's significance, and then sent to an output node that used its own weight to output the target variable classification. The illustration of the network used for prediction of the target variable is included here:



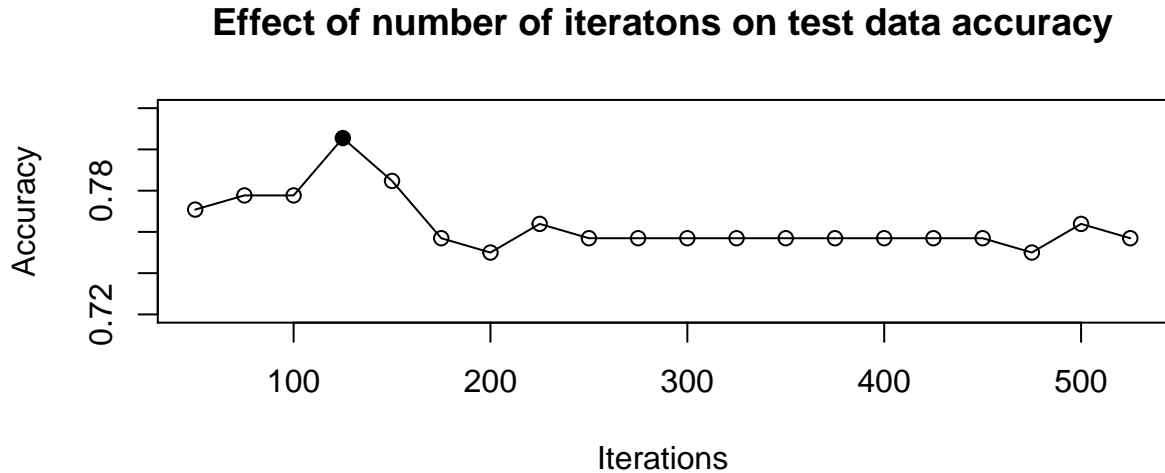
Results

Accuracy

The results obtained are summarized as follows:

Model	Accuracy
1 ada package	0.8055
2 Best Subsets	0.7778
3 Logistic Ridge	0.7708
5 Full Model, Poisson	0.7708
4 Full Model, Binomial	0.7638
6 Neural Network (neuralnet package)	0.7431
7 Manually implemented adaboost algorithm	0.7153

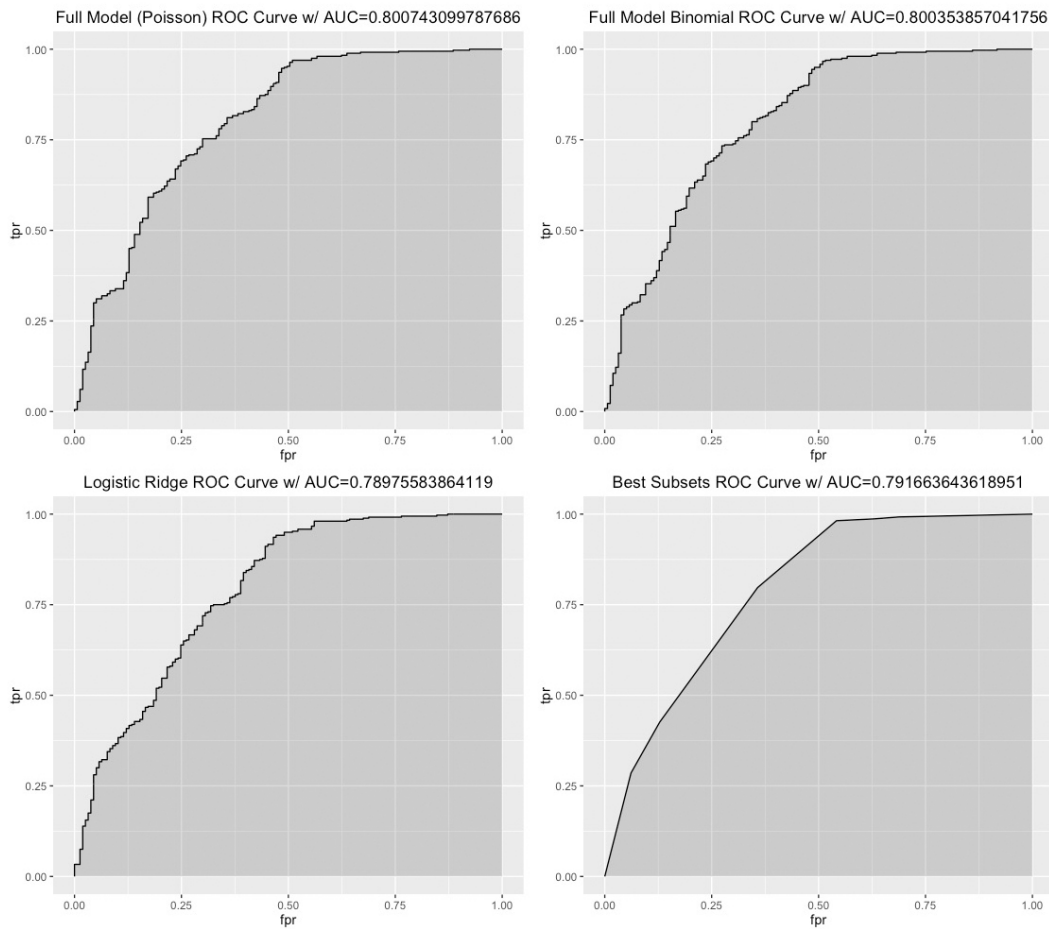
It can be seen that the best results were achieved through use of the `ada` package. The following figure illustrates the results of an analysis of test data accuracy against iterations of the `ada` function:



Optimal accuracy was attained at 125 iterations, with a steady decline and plateau thereafter. These results were in keeping with our hypothesis, however it was expected that the peak accuracy would occur at a higher iteration.

ROC Curves

ROC curves for the linear regression models were created and are displayed here.



The ROC curves shown here illustrate that the performance of the four binary classifiers represented do not vary greatly from one another – all classifiers show an Area Under the Curve (AUC) value of approximately 0.8. However, as shown in the table above, the accuracy of these classifiers vary greatly on the test data with a maximum and minimum of the four listed here being 0.7778 and 0.7638, respectively. While this may seem to be a difference of “only” 0.014 on test data accuracy, this translates to a differences of approximately 80 places in the data challenge rankings. This is a potent example of the importance of seemingly miniscule improvements in prediction models.

10-Fold Cross-Validation

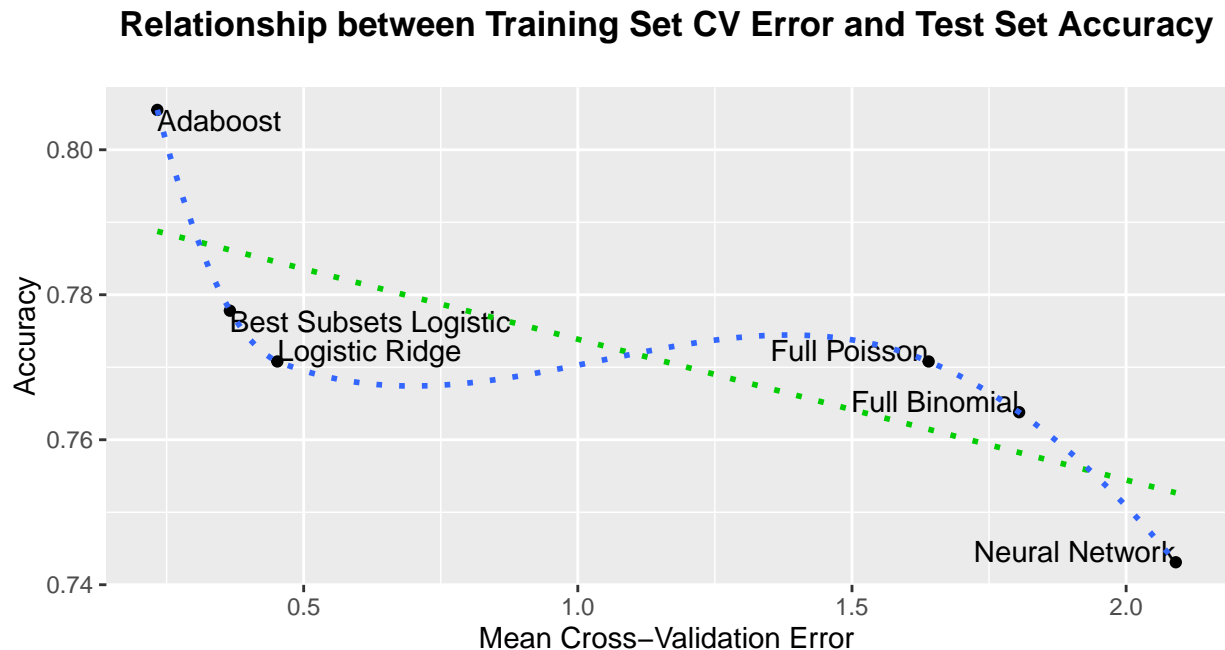
The 10-Fold Cross-Validation approach was used again to determine which model produced the lowest error on training data. This method was chosen as it allows all observations to be used for both training and validation, while also allowing each model to be tested on unseen data. These results are compared in the Discussion section below to the accuracy on test set data. A summary of the mean CV error for each model is shown here:

Model	Mean CV Error
Logistic Ridge Regression Model	0.2331
Neural Network	0.3654
Additive (Boosted) Logistic Regression Model	0.452
Best Subsets Logistic Regression Model	1.64
Full Logistic Regression Model	1.805
Full Regression Model, Poisson	2.091

Discussion

Many data mining methods were explored for this project. It was determined that the additive logistic regression algorithm (**ada** package) produced the best results, with the best subsets logistic regression coming in second. It is believed that the advantage of the **ada** package is due to its ability to combine a number of simpler models. Our hypothesis regarding the manually implemented Adaboost algorithm was incorrect, as it resulted in lower accuracy than the best subsets regression model.

Here the relationship between the Cross Validation error on training data and accuracy on the test data is shown:



While this plot does indicate a general negative relationship between mean CV error and test data accuracy (i.e. as CV error increases, test accuracy decreases), the sample is too small to determine whether or not this data is linear or follows a higher order relationship. That said, the associated regression data for the linear best fit in the plot above is shown here:

	Estimate	Std. Error	t value	Pr(> t)
CV	-0.01942	0.007334	-2.647	0.05714
(Intercept)	0.7933	0.009799	80.96	1.395e-07

Table 5: Fitting linear model: Accuracy \sim CV

Observations	Residual Std. Error	R^2	Adjusted R^2
6	0.01368	0.6366	0.5458

This model shows that the slope given by CV, while not quite significant yet ($p = 0.057$), does have the potential to become significant upon the growth of the sample.

The performance of the best subsets logistic regression model was surprising given its reliance on only three variables (**Married**, **Credit_History** and **SemiUrban**) for prediction of the target variable, **Loan_Status**. Despite this strong performance, the best subsets logistic regression model alone resulted in an overall leaderboard rank of 161. The most significant predictors found here differs from those found by Agbemava et al. (2016). This is most likely due to the population studied – higher income individuals in this project as opposed to lower income cases.

Using our best model, the typical applicant approved for a loan was a married male with no dependents, not self-employed, earning approximately \$4000 per month, and requested a loan of about \$125000. These results are similar to that seen in the training data set and leans towards higher earners with a lower debt-to-income ratio when compared to the average applicant.

Limitations of this project include the requirement to work with only the data supplied by the data challenge – the given 14 predictor variables. It is most likely that stronger models could be built upon the acquisition of more data. However, the data provided reflects the realistic scenario of having limited data from a loan application, particularly in the situation where an individual may be a new customer of the bank. In subsequent trials and analyses, this project would benefit greatly from further experimentation with more advanced additive logistic regression models and with neural network architectures.

At the time of writing this report, we are currently tied for 8th place out of a total of 268 participants. Our results show that even foundational data mining and machine learning techniques can produce powerful prediction models and quite accurate results.

References

- Agbemava, E., Nyarko, I. K., Adade, T. C., & Bediako, A. K. (2016). Logistic Regression Analysis Of Predictors Of Loan Defaults By Customers Of Non-Traditional Banks In Ghana. *European Scientific Journal*, 12(1).
- Conlin, M. (1999). Peer group micro-lending programs in Canada and the United States. *Journal of Development Economics*, 60(1), 249-269.
- Cule, E., & De Iorio, M. (2012). A semi-automatic method to guide the choice of ridge parameter in ridge regression. arXiv preprint arXiv:1205.0686.
- Dobbin, K. K., & Simon, R. M. (2011). Optimally splitting cases for training and testing high dimensional classifiers. *BMC medical genomics*, 4(1), 1.
- Günther, F., & Fritsch, S. (2010). neuralnet: Training of neural networks. *The R journal*, 2(1), 30-38.
- “Home Mortgage Disclosure Act.” Consumer Financial Protection Bureau: An Official Website of the United States Government. Consumer Financial Protection Bureau, n.d. Web. 17 July 2016. <http://www.consumerfinance.gov/data-research/hmda/>.
- Lusardi, A., & Scheresberg, C. D. B. (2013). Financial literacy and high-cost borrowing in the United States (No. w18969). National Bureau of Economic Research.
- Olson, D. L., Delen, D., & Meng, Y. (2012). Comparative analysis of data mining methods for bankruptcy prediction. *Decision Support Systems*, 52(2), 464-473.
- Thomas, L. C. (2000). A survey of credit and behavioural scoring: forecasting financial risk of lending to consumers. *International journal of forecasting*, 16(2), 149-172.

Appendix – R Code

```
##### setup #####
library(knitr)
opts_chunk$set(echo = FALSE, warning = FALSE, message = FALSE,
               comment = NA, fig.align = "center", error = TRUE)

library(glmnet)
library(leaps)
library(pROC)
library(car)
library(MASS)
library(ROCR)
library(ggplot2)
library(stringr)
library(dplyr)
library(reshape2)
library(vcd)
library(VIM)
library(pander)
library(tidyr)
library(e1071)

train <- read.csv("https://github.com/dsmilo/DATA621/raw/master/Final-Project/loan_train.csv",
                  stringsAsFactors = FALSE)
test = read.csv("https://github.com/dsmilo/DATA621/raw/master/Final-Project/loan_test.csv",
                 stringsAsFactors = FALSE)
test2 = test

predict.regsubsets <- function(object, newdata, id, ...) {
  form <- as.formula(object$call[[2]])
  mat <- model.matrix(form, newdata)
  coefi <- coef(object, id = id)
  xvars <- names(coefi)
  mat[, xvars] %*% coefi
}

plotROC <- function(model, ndata, gtruth, name = "") {
  prob <- predict(model, newdata = ndata, type = "response")
  pred <- prediction(prob, gtruth)
  perf <- performance(pred, measure = "tpr", x.measure = "fpr")
  auc <- performance(pred, measure = "auc")
  auc <- auc@y.values[[1]]
  roc.data <- data.frame(fpr = unlist(perf@x.values), tpr = unlist(perf@y.values),
                        model = "GLM")
  ggplot(roc.data, aes(x = fpr, ymin = 0, ymax = tpr)) + geom_ribbon(alpha = 0.2) +
    geom_line(aes(y = tpr)) + ggtitle(paste0(name, " ROC Curve w/ AUC=",
      (auc)))
}

fillwithmedian <- function(x) {
```

```

    median_val = median(x, na.rm = TRUE)
    x[is.na(x)] = median_val
    return(x)
}

#####

#### literature review ####
library(ReadImages) #Had to install manually
library(grid)
download.file("https://cs231n.github.io/assets/n1/neural_net2.jpeg",
             "nn1.jpeg")
img <- read.jpeg("nn1.jpeg")
grid.raster(img)
#####

#### data exploration and manipulation ####
train3 = train
test2 = test

train3$Gender = ifelse(train$Gender == "Male", 1, 0)
train3$Married = ifelse(train$Married == "Yes", 1, 0)
train3$Dependents[train3$Dependents == "3+"] = 3
train3$Dependents = as.numeric(train3$Dependents)
train3$Education = ifelse(train$Education == "Graduate", 1, 0)
train3$Self_Employed = ifelse(train$Self_Employed == "Yes", 1,
                               0)

train3$Urban = ifelse(train$Property_Area == "Urban", 1, 0)
train3$SemiUrban = ifelse(train$Property_Area == "Semiurban",
                           1, 0)
train3$TotalIncome = train3$ApplicantIncome + train3$CoapplicantIncome
train3$LTI = (train3$LoanAmount * 1000)/train3$ApplicantIncome #Loan to income ratio

train3$Loan_Status = ifelse(train$Loan_Status == "Y", 1, 0)

test2$Gender = ifelse(test$Gender == "Male", 1, 0)
test2$Married = ifelse(test$Married == "Yes", 1, 0)
test2$Dependents[test2$Dependents == "3+"] = 3
test2$Dependents = as.numeric(test2$Dependents)
test2$Education = ifelse(test$Education == "Graduate", 1, 0)
test2$Self_Employed = ifelse(test$Self_Employed == "Yes", 1,
                               0)

test2$Urban = ifelse(test$Property_Area == "Urban", 1, 0)
test2$SemiUrban = ifelse(test$Property_Area == "Semiurban", 1,
                           0)
test2$TotalIncome = test2$ApplicantIncome + test2$CoapplicantIncome
test2$LTI = (test2$LoanAmount * 1000)/test2$ApplicantIncome #Loan to income ratio

```

```

train4 = train3[complete.cases(train3), ]

train = train3[, -1]
train2 = train4
test = test2

train_var = train
train_var$Property_Area = as.factor(train_var$Property_Area)
var_class <- data.frame(Class = rep(NA, ncol(train_var) - 4),
  Levels = rep(NA, ncol(train_var) - 4), stringsAsFactors = FALSE,
  check.names = FALSE, row.names = names(train_var)[1:(ncol(train_var) -
    4)])
names(var_class) = c("Class", "Levels/Range")

for (i in 1:(ncol(train_var) - 3)) {
  var_class[i - 1, 1] <- class(train_var[, i])
  val = length(levels(train_var[, i]))
  if (val == 0) {
    val = paste0(floor(min(train_var[, i], na.rm = TRUE)),
      " - ", max(train_var[, i], na.rm = TRUE))
  }
  var_class[i - 1, 2] <- val
}
pander(var_class)

# continuous variables
cont_vars <- train_var %>% dplyr::select(c(ApplicantIncome, CoapplicantIncome,
  LoanAmount, Loan_Amount_Term))

melted <- melt(cont_vars)

ggplot(melted, aes(value)) + geom_density(aes(fill = variable,
  col = variable), alpha = 0.5, show.legend = FALSE) + facet_wrap(~variable,
  scales = "free") + scale_y_continuous("", labels = NULL) +
  ggtitle("Distribution of Continuous Predictors \n") + theme(axis.ticks = element_blank())

# discrete variables
train_var$Property_Area <- as.numeric(train_var$Property_Area)
poisson_vars <- train_var %>% dplyr::select(-c(ApplicantIncome,
  CoapplicantIncome, LoanAmount, Loan_Amount_Term, TotalIncome,
  LTI, Urban, SemiUrban))

melted <- melt(poisson_vars)

library(scales)
ggplot(melted, aes(value, y = 6 * ..count../sum(..count..))) +
  geom_bar(aes(fill = variable, col = variable), alpha = 0.5,
    show.legend = FALSE) + facet_wrap(~variable, scales = "free_x") +
  scale_y_continuous("", labels = percent) + ggtitle("Distribution of Discrete Predictors \n")

# created variables

```

```

created_vars <- train_var %>% dplyr::select(c(TotalIncome, LTI))

melted <- melt(created_vars)

ggplot(melted, aes(value)) + geom_density(aes(fill = variable,
  col = variable), alpha = 0.5, show.legend = FALSE) + facet_wrap(~variable,
  scales = "free") + scale_y_continuous("", labels = NULL) +
  ggtitle("Distribution of Created Predictors \n") + theme(axis.ticks = element_blank())

# missing values
aggr(train_var, plot = TRUE, sortVars = FALSE, numbers = TRUE,
  cex.lab = TRUE)
#####

##### Model Creation ##### Full Model
glmflagfull <- glm(Loan_Status ~ . - Property_Area, data = train,
  family = poisson, control = list(maxit = 100))
summary(glmflagfull)
vif(glmflagfull) #Check for collinearity, VIF > 10
predglmflagfull <- data.frame(class = train$Loan_Status, logit = predict(glmflagfull,
  train))

plotROC(model = glmflagfull, ndata = train, gtruth = train$Loan_Status,
  name = "Full Model (Poisson)")

a = predict(glmflagfull, test, type = "response")
testdf = data.frame(Loan_ID = test$Loan_ID, Loan_Status = a)
testdf$Loan_Status[is.na(testdf$Loan_Status)] = median(testdf$Loan_Status,
  na.rm = TRUE)
testdf$Loan_Status = ifelse(testdf$Loan_Status > 0.5, "Y", "N")

write.csv(testdf, paste0("finalresults-", as.numeric(Sys.time()),
  ".csv"))
#####

# Best Subsets
regfit.full <- regsubsets(Loan_Status ~ . - Property_Area - Loan_ID,
  data = train3, nvmax = 17)
summary(regfit.full)
par(mar = c(1, 1, 1, 1))
par(mfrow = c(1, 2))
plot(regfit.full, scale = "bic", main = "Predictor Variables vs. BIC")
reg.summary <- summary(regfit.full)
reg.summary$bic
plot(reg.summary$bic, xlab = "Number of Predictors", ylab = "BIC",
  type = "l", main = "Best Subset Selection Using BIC")
minbic <- which.min(reg.summary$bic)
points(minbic, reg.summary$bic[minbic], col = "brown", cex = 2,
  pch = 20)
coef(regfit.full, minbic)
var_names = names(coef(regfit.full, minbic))[2:length(names(coef(regfit.full,

```

```

    minbic)))]
length(var_names)

Model_toEval = paste0("glm(Loan_Status ~ ", paste(var_names,
    collapse = " + "), ", data = train, family = binomial(link='logit'))")

bestsubset21 = eval(parse(text = Model_toEval))

summary(bestsubset21)

plotROC(model = bestsubset21, ndata = train, gtruth = train$Loan_Status,
    name = "Best Subsets")

a = predict(bestsubset21, test, type = "response")
testdf = data.frame(Loan_ID = test$Loan_ID, Loan_Status = a)
testdf$Loan_Status[is.na(testdf$Loan_Status)] = median(testdf$Loan_Status,
    na.rm = TRUE)
testdf$Loan_Status = ifelse(testdf$Loan_Status > 0.5, "Y", "N")
#####

# Logistic Ridge
library(ridge) #installed from source
train$Loan_Status = as.numeric(as.character(train$Loan_Status))
train_lr <- data.frame(lapply(train[, -11], fillwithmedian))
train_lr$Property_Area = train$Property_Area

lr2 = logisticRidge(Loan_Status ~ . - Property_Area, data = train_lr)

plotROC(model = lr2, ndata = train, gtruth = train$Loan_Status,
    name = "Logistic Ridge")

a = predict(lr2, test2, type = "response")
testdf = data.frame(Loan_ID = test$Loan_ID, Loan_Status = a)
testdf$Loan_Status[is.na(testdf$Loan_Status)] = median(testdf$Loan_Status,
    na.rm = TRUE)
testdf$Loan_Status = ifelse(testdf$Loan_Status > 0.5, "Y", "N")
write.csv(testdf, paste0("finalresults-", as.numeric(Sys.time()),
    ".csv"))
#####

# Full model binomial
glmflagfull <- glm(Loan_Status ~ . - Property_Area, data = train,
    family = binomial, control = list(maxit = 100))
summary(glmflagfull)
vif(glmflagfull) #Check for collinearity, VIF > 10
predglmflagfull <- data.frame(class = train$Loan_Status, logit = predict(glmflagfull,
    train))

plotROC(model = glmflagfull, ndata = train, gtruth = train$Loan_Status,
    name = "Full Model Binomial")

a = predict(glmflagfull, test, type = "response")

```



```

testdf = data.frame(Loan_ID = test$Loan_ID, Loan_Status = a)
testdf$Loan_Status[is.na(testdf$Loan_Status)] = median(testdf$Loan_Status,
  na.rm = TRUE)
testdf$Loan_Status = ifelse(testdf$Loan_Status > 0.5, "Y", "N")
write.csv(testdf, paste0("finalresults-", as.numeric(Sys.time()),
  ".csv"))
#####

# Neural Net Determination of nodes for neuralnet
k = 10
set.seed(1306)
folds = sample(1:k, nrow(train), replace = TRUE)
cv.errors1 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors2 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors3 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors4 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors5 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors6 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors7 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors8 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))

train_2 <- data.frame(lapply(train[, -11], fillwithmedian))
train_2$Property_Area = train$Property_Area
scaled = data.frame(lapply(scaled, fillwithmedian))

num_layers = 0

f1 <- as.formula(paste("Loan_Status ~", paste(n[!n %in% "Loan_Status"],
  collapse = " + ")))
for (j in 1:k) {

  num_layers = num_layers + 5
  print(num_layers)
  dat = scaled[folds != j, ]
  dat = dat[complete.cases(dat), ]
  model5 <- neuralnet(f1, data = dat, hidden = c(num_layers),
    linear.output = TRUE, stepmax = 1e+07)

  for (i in 1:10) {
    print(i)
    f = train_2[folds == j, ]
    f = f[complete.cases(f), ]
    dat = scaled[folds == j, ]
    dat2 = dat[complete.cases(dat), -11]
    pred5 = neuralnet::compute(model5, dat2)
    netresult = pred5$net.result[, 1]
    sc = scaled$Loan_Status[folds == j]
    sc = sc[complete.cases(sc)]
    cv.errors5[j, i] = mean((sc - netresult)^2, na.rm = TRUE)
  }
}

```

```

cv_means = rowMeans(cv.errors5)
node_counts = seq(5, 50, 5)
nn_df = data.frame(Nodes = node_counts, CV_Error = cv_means)
names(nn_df) = c("Nodes", "Mean CV Error")
write.csv(nn_df, "~/Downloads/nn_df.csv", row.names = FALSE)

nn_df = read.csv("https://raw.githubusercontent.com/aadikaloo/AadiMSDA/master/IS621-Data-Mining/nn_c
plot(x = nn_df$Nodes, y = nn_df$Mean.CV.Error, xlab = "Nodes",
     ylab = "Mean CV Error", main = "Determination of Number of Nodes for neuralnet",
     pch = ifelse(nn_df$Nodes == 10, 19, 1))
#####

## Neural network
train4 = train3
train4 = train4[complete.cases(train4), ]
train4 = train4[, -c(1, 12)]
maxs <- apply(train4, 2, max)
mins <- apply(train4, 2, min)
scaled <- as.data.frame(scale(train4, center = mins, scale = maxs -
  mins))
n <- names(scaled)
f <- as.formula(paste("Loan_Status ~", paste(n[!n %in% "Loan_Status"],
  collapse = " + ")))
nn = neuralnet(f, data = scaled, hidden = c(10), linear.output = TRUE,
  stepmax = 1e+06)
plot(nn)

max2 <- function(x) {
  return(max(x, na.rm = TRUE))
}

min2 <- function(x) {
  return(min(x, na.rm = TRUE))
}

test <- data.frame(lapply(test, fillwithmedian))
test = test2
test = test[, -c(1, 12)]
test <- data.frame(lapply(test, fillwithmedian))

maxs.t <- apply(test, 2, max2)
mins.t <- apply(test, 2, min2)
scaled.t <- as.data.frame(scale(test, center = mins.t, scale = maxs.t -
  mins.t))
scaled.t = as.data.frame(lapply(scaled.t, fillwithmedian))
pr.test <- neuralnet::compute(nn, scaled.t)
pr.t_ <- pr.test$net.result * (max(train4$Loan_Status) - min(train4$Loan_Status)) +
  min(train4$Loan_Status)
pr_nn.t = ifelse(pr.t_ > 0.5, "Y", "N")

test1 = read.csv("https://github.com/dsmilo/DATA621/raw/master/Final-Project/loan_test.csv",
  stringsAsFactors = FALSE)

```

```

df = data.frame(Loan_ID = test1$Loan_ID, Loan_Status = as.data.frame(as.vector(pr_nn.t)))
names(df) = c("Loan_ID", "Loan_Status")
sum(is.na(df[, 2])) #hopefully 0
write.csv(df, paste0("nn-", as.numeric(Sys.time()), ".csv"),
  row.names = FALSE)

download.file("https://github.com/aadikaloo/AadiMSDA/raw/master/IS621-Data-Mining/nn2.jpg",
  "nn2.jpeg")
img <- read.jpeg("nn2.jpeg")
grid.raster(img)
#####

# ada package analysis
model1 = list()

for (i in 1:20) {
  iters = 25 + 25 * i
  model1[[i]] = ada(formula = Loan_Status ~ ., data = train_x1,
    iter = iters, loss = "logistic")
  results = as.numeric(as.character(predict(model1[[i]], test)))
  results = ifelse(results == 1, "Y", "N")
  testdf = data.frame(Loan_ID = test2$Loan_ID, Loan_Status = results)
  write.csv(testdf, paste0("finalresults-", as.numeric(Sys.time()),
    " - ", iters, ".csv"), row.names = FALSE)
}

ada_results = read.csv("https://raw.githubusercontent.com/aadikaloo/AadiMSDA/master/IS621-Data-Mining/ada_results.csv",
  header = FALSE)
names(ada_results) = c("Iterations", "Accuracy")
#####

# Adaboost manual implementation
wts = list()

models = list()
alpha = list()
error_rate = list()
final = 0
predictions = list()
total = 12

train_x <- data.frame(lapply(train[-c(1, 2), -11], fillwithmedian))
test_x = data.frame(lapply(test2[, -c(1, 12)], fillwithmedian))

k = total
set.seed(1306)
folds = sample(1:k, nrow(train_x), replace = TRUE)
folds = rep(1:12, nrow(train_x)/12)

wts[[1]] = rep(1/length(train_x1[, 1]), length(train_x1[, 1]))

```

```

Y_hat_run = list()
Y_hat_run[[1]] = rep(1, length(train_x1[, 1]))

for (i in 1:total) {
  train_x1 = train_x

  models[[i]] = glm(Loan_Status ~ ., data = train_x1, family = binomial,
    weights = wts[[i]])
  Y_hat = predict(models[[i]], train_x1, type = "response")
  Y_hat_run = Y_hat_run[[i]] * Y_hat
  Y_hat = ifelse(Y_hat > 0.5, 1, 0)
  error_rate[[i]] = sum(wts[[i]] * (train_x1$Loan_Status !=
    Y_hat))/length(Y_hat)
  alpha[[i]] = (1/2) * log((1 - error_rate[[i]])/error_rate[[i]])
  loan_status = ifelse(train_x1$Loan_Status > 0, 1, -1)
  Y_hat = ifelse(Y_hat == 1, 1, -1)
  wts[[i + 1]] = wts[[i]] * exp(-alpha[[i]] * loan_status *
    Y_hat)
  wts[[i + 1]] = wts[[i + 1]]/sum(wts[[i + 1]])
}

for (i in 1:total) {
  predictions[[i]] = predict(models[[i]], test_x, type = "response")
  predictions[[i]] = alpha[[i]] * predictions[[i]]
}

toEval = "predictions[[1]]"

for (i in 2:4) {
  toEval = paste0(toEval, " * predictions[[" , i, "]]")
}

prop = table(train$Loan_Status)[[1]]/table(train$Loan_Status)[[2]]
new = Y_hat_run[order(Y_hat_run)]
pred2 = ifelse((rank(Y_hat_run)) > prop * length(test[, 1]),
  1, 0)

results = eval(parse(text = toEval))
results = ifelse(results < error_rate[[1]], "N", "Y")

testdf = data.frame(Loan_ID = test2$Loan_ID, Loan_Status = results)

write.csv(testdf, paste0("finalresults-", as.numeric(Sys.time()),
  ".csv"), row.names = FALSE)
#####

##### Results #####
data_results = data.frame(Model = c("ada package", "Best Subsets",
  "Logistic Ridge", "Full Model, Binomial", "Full Model, Poisson",
  "Neural Network (neuralnet package)", "Manually implemented adaboost algorithm"),
  Accuracy = c("0.8055", "0.7778", "0.7708", "0.7638", "0.7708",

```

```

      "0.7431", "0.7153"), stringsAsFactors = FALSE)
data_results2 = data_results[order(data_results$Accuracy, decreasing = TRUE),
  ]
kable(data_results2)

ada_results = read.csv("https://raw.githubusercontent.com/aadikaloo/AadiMSDA/master/IS621-Data-Mining/ada_results.csv",
  header = FALSE)
names(ada_results) = c("Iterations", "Accuracy")
plot(x = ada_results$Iterations, y = ada_results$Accuracy, xlab = "Iterations",
  ylab = "Accuracy", ylim = c(0.72, 0.82), main = "Effect of number of iterations on test data accuracy",
  pch = ifelse(ada_results$Accuracy == 0.8055, 19, 1))
lines(smooth.spline(x = ada_results$Iterations, y = ada_results$Accuracy))

download.file("https://raw.githubusercontent.com/aadikaloo/AadiMSDA/master/IS621-Data-Mining/roc2.jpg",
  "roc.jpeg")
img <- read.jpeg("roc.jpeg")
grid.raster(img)

# 10-fold CV
k = 10
set.seed(1306)
folds = sample(1:k, nrow(train), replace = TRUE)
cv.errors1 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors2 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors3 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors4 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors5 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors6 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors7 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))
cv.errors8 = matrix(NA, k, 10, dimnames = list(NULL, paste(1:10)))

train_2 <- data.frame(lapply(train[, -11], fillwithmedian))
train_2$Property_Area = train$Property_Area
scaled = data.frame(lapply(scaled, fillwithmedian))

f1 <- as.formula(paste("Loan_Status ~", paste(n[!n %in% "Loan_Status"],
  collapse = " + ")))
for (j in 1:k) {
  print(j)

  model1 <- glm(Loan_Status ~ . - Property_Area, data = train_2[folds !=
    j, ], family = poisson, control = list(maxit = 100))
  model2 <- glm(formula = Loan_Status ~ Married + Credit_History +
    SemiUrban, family = binomial(link = "logit"), data = train_2[folds !=
    j, ])
  model3 <- logisticRidge(Loan_Status ~ . - Property_Area,
    data = train_2[folds != j, ])
  model4 <- glm(Loan_Status ~ . - Property_Area, data = train_2[folds !=
    j, ], family = binomial, control = list(maxit = 100))
  dat = scaled[folds != j, ]
  dat = dat[complete.cases(dat), ]
  model5 <- neuralnet(f1, data = dat, hidden = c(5, 3), linear.output = TRUE,

```

```

    stepmax = 1e+07)
model6 <- ada(formula = Loan_Status ~ ., data = train_x1[folds !=
  j, ], iter = 125, loss = "logistic")

for (i in 1:10) {
  print(i)
  f = train_2[folds == j, ]
  f = f[complete.cases(f), ]

  pred1 = predict(model1, f, id = i)
  cv.errors1[j, i] = mean((train_2$Loan_Status[folds ==
    j] - pred1)^2, na.rm = TRUE)

  pred2 = predict(model2, f, id = i)
  cv.errors2[j, i] = mean((train_2$Loan_Status[folds ==
    j] - pred2)^2, na.rm = TRUE)

  pred3 = predict(model3, f, id = i)
  cv.errors3[j, i] = mean((train_2$Loan_Status[folds ==
    j] - pred3)^2, na.rm = TRUE)

  pred4 = predict(model4, f, id = i)
  cv.errors4[j, i] = mean((train_2$Loan_Status[folds ==
    j] - pred4)^2, na.rm = TRUE)

  dat = scaled[folds == j, ]
  dat = dat[complete.cases(dat), -11]
  pred5 = neuralnet::compute(model5, dat2)
  netresult = pred5$net.result[, 1]
  sc = scaled$Loan_Status[folds == j]
  sc = sc[complete.cases(sc)]
  cv.errors5[j, i] = mean((sc - netresult)^2, na.rm = TRUE)

  pred6 = predict(model6, train_x1[folds == j, ], id = i,
    type = "probs")
  cv.errors6[j, i] = mean((train_2$Loan_Status[folds ==
    j] - pred6)^2, na.rm = TRUE)

}

}

mean.cv.errors1 <- apply(cv.errors1, 2, mean)
mean.cv.errors2 <- apply(cv.errors2, 2, mean)
mean.cv.errors3 <- apply(cv.errors3, 2, mean)
mean.cv.errors4 <- apply(cv.errors4, 2, mean)
mean.cv.errors5 <- apply(cv.errors5, 2, mean)
mean.cv.errors6 <- apply(cv.errors6, 2, mean)

all.cv.error = data.frame(mean(mean.cv.errors1), mean(mean.cv.errors2),
  mean(mean.cv.errors3), mean(mean.cv.errors4), mean(mean.cv.errors5),
  mean(mean.cv.errors6))

```

```

names(all.cv.error) = c("Full Regression Model, Poisson", "Best Subsets Logistic Regression Model",
  "Logistic Ridge Regression Model", "Full Logistic Regression Model",
  "Neural Network", "Additive (Boosted) Logistic Regression Model")
all.cv.error = t(all.cv.error)
names(all.cv.error) = c("Model", "Mean CV Error")

mean_cv_results = read.csv("https://github.com/aadikaloo/AadiMSDA/raw/master/IS621-Data-Mining/allcv
  check.names = FALSE)
mean_cv = mean_cv_results[order(mean_cv_results$"Mean CV Error"),
  ]
row.names(mean_cv) = NULL
pander(mean_cv)
#####

#### discussion ####
acc_cv <- data.frame(Accuracy = as.numeric(data_results2$Accuracy[1:6]),
  CV = mean_cv$"Mean CV Error")
rownames(acc_cv) <- c("Adaboost", "Best Subsets Logistic", "Logistic Ridge",
  "Full Poisson", "Full Binomial", "Neural Network")

ggplot(acc_cv, aes(x = CV, y = Accuracy)) + geom_point() + geom_text(aes(label = rownames(acc_cv)),
  vjust = "inward", hjust = "inward", check_overlap = TRUE) +
  geom_smooth(alpha = 0.25, method = glm, fill = 3, col = 3,
    se = FALSE, lty = 3) + scale_x_continuous("Mean Cross-Validation Error") +
  ggtitle("Relationship between Training Set CV Error and Test Set Accuracy\n") +
  theme(plot.title = element_text(lineheight = 0.8, face = "bold")) +
  geom_smooth(alpha = 0.25, method = "loess", se = FALSE, lwd = 1,
    lty = 3)

alm = lm(Accuracy ~ CV, data = acc_cv)
pander(summary(alm))
#####

```