

SEIS 736 Machine Learning
Project Report
December 9, 2015

Give Me Some Credit

Team Members: Princewill Eneh, Suraj Jois, Xiaozhao Song, Yalei Yan, Ruogu Wang

Table of Content

• Introduction	3
• Dataset	3
• Methodology	4
• Tools	4
• Import Data	5
Import our datasets.	
• Preliminary Analysis & Data Exploration	6
Check for Colinearity, Skewness and Missing Data	
• Stratified Sampling	6
Perform Stratified sampling to prevent low precision in our model	
• Data Cleansing Process	
Remove outliers	
• Imputation	7
Impute missing Data	
• Normalization	7
Make our data unitless	
• Regularization	8
Find the best predictors	
• Algorithms	12
Logistic Regression with Cross validation	
Linear Discriminant Analysis	
Support Vector Machine	
Artificial Neural Network	
Random Forest	
• Evaluate our Results	23
Evaluate result using Confusion Matrix, ROC, AUC etc	
• Performance	25
Compare to SAP HANA	
• Lessons Learned and Challenges	25
• Conclusion	26
• Reference	26

Introduction & Objective

Banks play a crucial role in market economies. They decide who can get finance and on what terms and can make or break investment decisions. For markets and society to function, individuals and companies need access to credit.

Credit scoring algorithms, which make a guess at the probability of default, are the method banks use to determine whether or not a loan should be granted. This project requires us to improve on the state of the art in credit scoring, by predicting the probability that somebody will experience financial distress in the next two years.

The goal of this project is to build a model that performs better than 77% accuracy which is the specified threshold to beat.

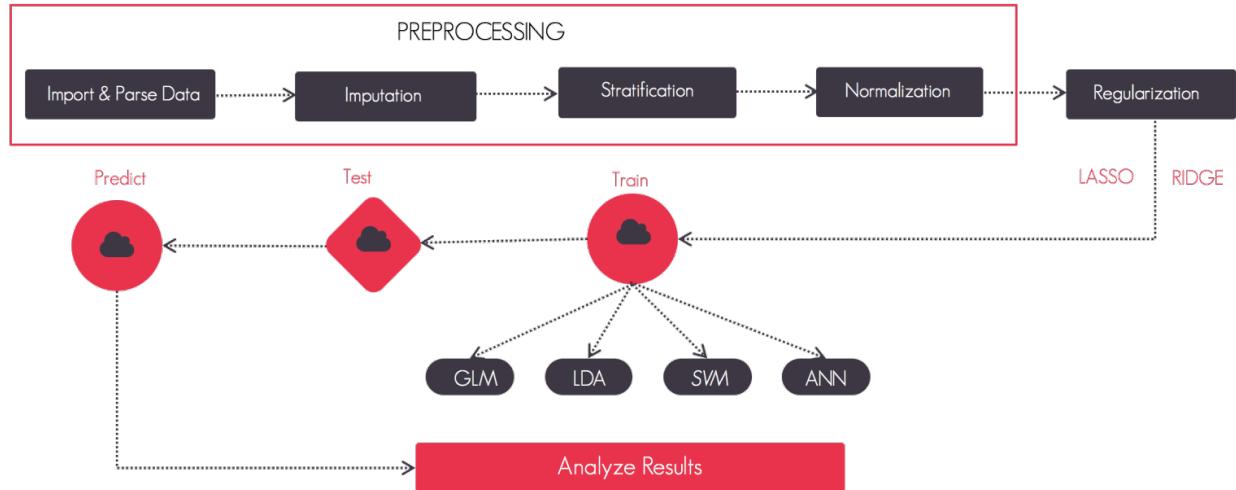
Datasets

- Our dataset consists of two csv files. One for training and the other for testing. The csv files contain about 150,000 records of customers for training and 100,000 for testing.
- Below is the data dictionary of the dataset.

Variable Name	Description	Type
SeriousDlqin2yrs	Person experienced 90 days past due delinquency or worse	Y/N
RevolvingUtilizationOfUnsecuredLines	Total balance on credit cards and personal lines of credit except real estate and no installment debt like car loans divided by the sum of credit limits	percentage
age	Age of borrower in years	integer
NumberOfTime30-59DaysPastDueNotWorse	Number of times borrower has been 30-59 days past due but no worse in the last 2 years.	integer
DebtRatio	Monthly debt payments, alimony, living costs divided by monthly gross income	percentage
MonthlyIncome	Monthly income	real
NumberOfOpenCreditLinesAndLoans	Number of Open loans (installment like car loan or mortgage) and Lines of credit (e.g. credit cards)	integer
NumberOfTimes90DaysLate	Number of times borrower has been 90 days or more past due.	integer
NumberRealEstateLoansOrLines	Number of mortgage and real estate loans including home equity lines of credit	integer
NumberOfTime60-89DaysPastDueNotWorse	Number of times borrower has been 60-89 days past due but no worse in the last 2 years.	integer
NumberOfDependents	Number of dependents in family excluding themselves (spouse, children etc.)	integer

Methodology

We plan to use several Machine learning techniques and methodologies. Below is the plan of action for achieving the projects set goal.



Tools

The primary tool used for this project was **R**. We chose R because it is a free flexible statistical analysis toolkit. All of the standard data analysis tools are built right into the R language: from accessing data in various formats, to data manipulation (transforms, merges, aggregations, etc.), to traditional and modern statistical models (regression, ANOVA, GLM, tree models, etc.). All are included in an object-oriented framework that makes it easy to programmatically extract out and combine just the information you need from the results, rather than having to cut-and-paste from a static report.

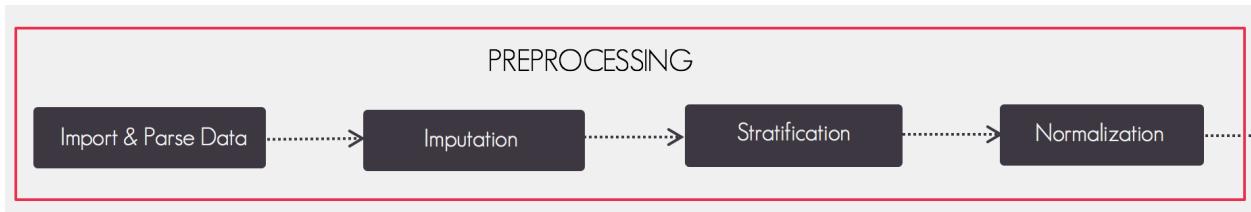
We also used a data wrangling tool called **Trifecta Wrangler** (<https://www.trifecta.com>). The tool was helpful in us cleaning up the data as Data wrangling is the critical first step of the analysis process – and one that has traditionally consumed a significant amount of time and effort. Successful analysis relies upon accurate, well-structured data that has been formatted for the specific needs of the task at hand. In fact, wrangling is as much a part of the data analysis process as the final results. So we employed this tool to visualize and understand our data before any data analysis is performed.



Finally, we wanted to see how some of the Algorithms performed when used on an in memory database like **SAP HANA**. We ran our dataset on AWS against Decision Tree and BPNN algorithms.

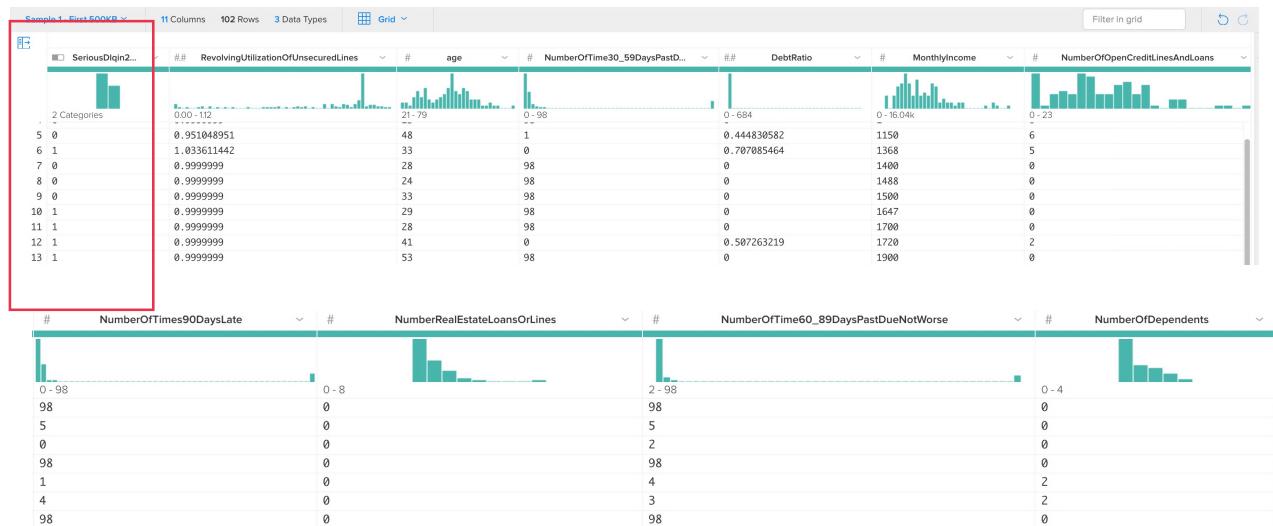
Preprocessing

The preprocessing phase of this project was critical as we wanted to have the best data for the predictive analysis we wished to perform. Some of the techniques employed include Imputation, Stratification and Normalization.



Import Data

We plan to understand our dataset as much as possible therefore before working on it, we need to know the distribution of each attribute. The first task is to import the raw data into Trifecta Wrangler. Once the data is imported we can see the distribution of all attribute. Below is the preview of the distribution.



We then import the data into R and get a summary of the data. The summary indicates we have 150000 observations and 12 variables.

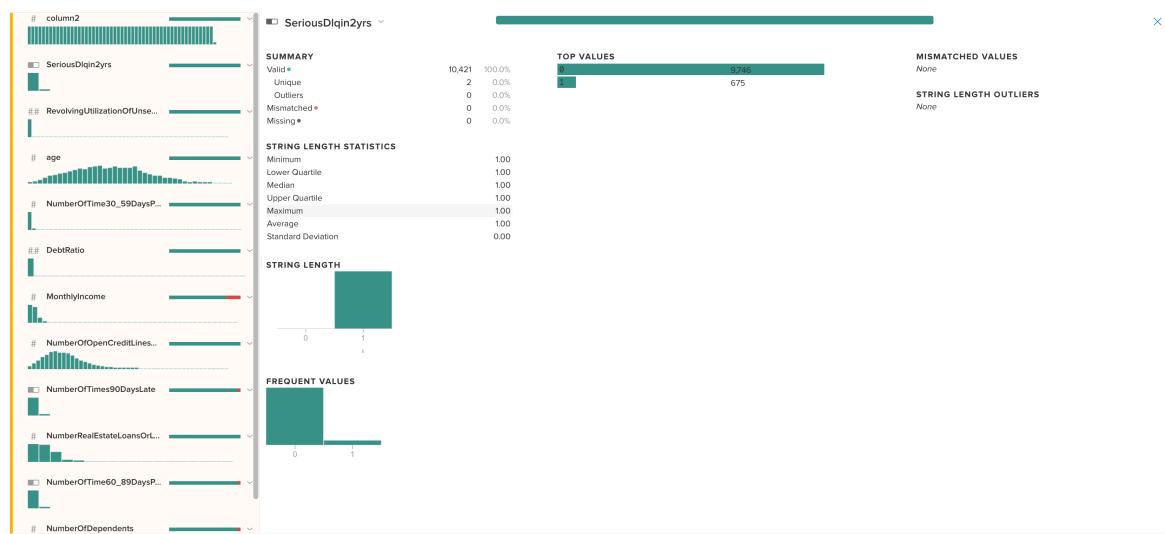
Data	
train	150000 obs. of 12 variables

Data Exploration

Next we checked for colinearity, skewness and missing data since we know that the presence of those will result in problematic model .

Stratified Sampling

Looking at our dataset we noticed that the training data had about 93% of 0s and 7% of the ones. Using this dataset, the way it is to train our data resulted in a good accuracy but precision was poor and the model could not generalize to our test data.



To combat this issue, we decided to stratify our training data and make the 1's and 0's the same size that is, 50% for 1's and 50% for 0's. This will ensure that our model is robust and can generalize to any test data.

Data Cleansing Process

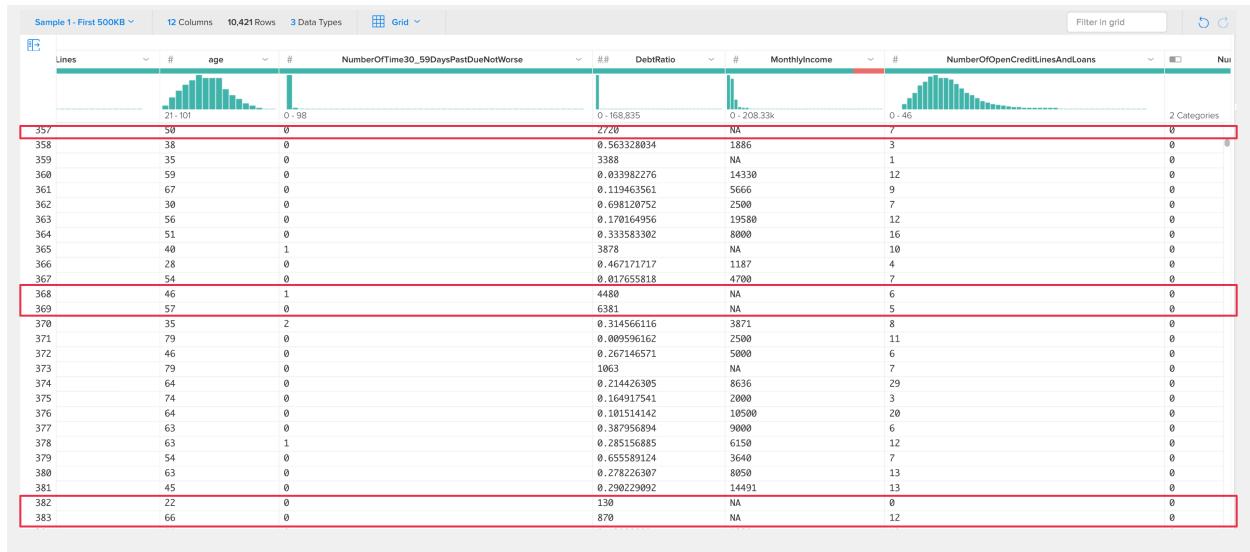
As is always the case with most datasets, the problem of missing values is one that always needs to be resolved for us to generate models from GLM, to SVM etc.

Imputation

Algorithms like logistic regression, Neural network, Support vector machine etc. would not work if there are missing values in the table. So to solve this, we had two choices;

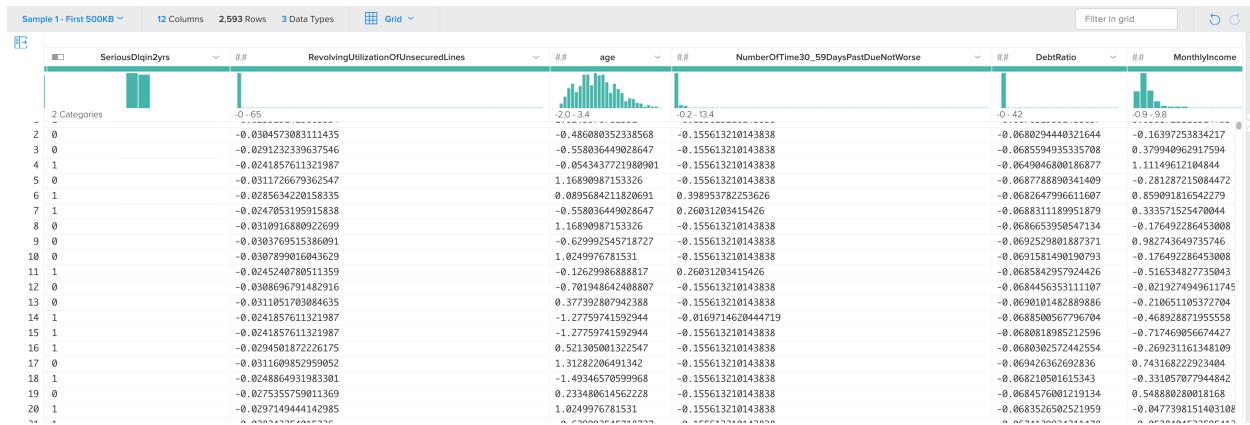
1. Impute the missing values
2. Remove rows with missing values in them

Each choice had its pros and cons. For choice one we will be retaining all the rows in our dataset which gives us more rows to train with. But doing this takes a lot of time and processing power. Choice two was the quickest but meant we will be removing about 30,000 records from our 150,000 record dataset. We decided to go with this choice as we concluded that 120,000 records for training is still a decent amount of records.



Normalization

Looking at the imported data, we notice a vast range of scale and unit differences in the attributes. Since algorithms like SVM require normalized data, we normalized the training data. Below is a preview of the result.



Regularization

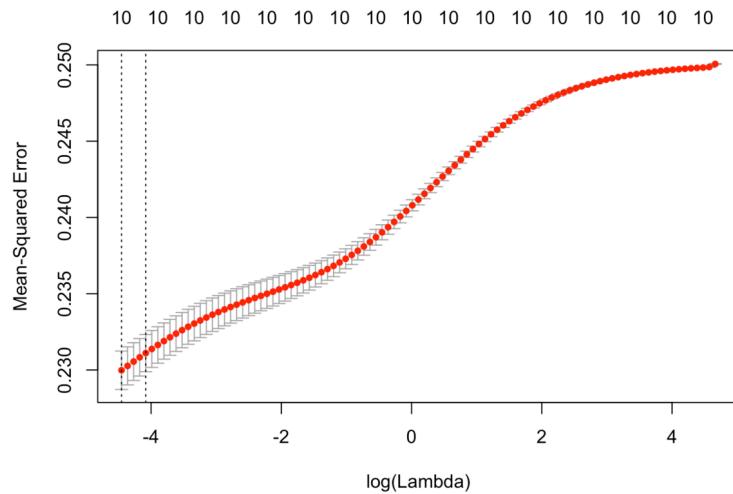
Regularization helps us find the best predictors and prevent overfitting. We wanted to make sure that our model does not overfit. We perform regularization using lasso and ridge. First we start with Ridge.

RIDGE

```
1 x.tr <- model.matrix(SeriousDlqin2yrs ~ RevolvingUtilizationOfUnsecuredLines + age + NumberOfTime30.59DaysPastDu  
eNotWorse  
2           + DebtRatio + MonthlyIncome + NumberOfOpenCreditLinesAndLoans + NumberOfTimes90DaysLate +  
3           NumberRealEstateLoansOrLines + NumberOfTime60.89DaysPastDueNotWorse  
4           + NumberOfDependents, data = training_data)[, -1]  
5 y.tr <- training_data$SeriousDlqin2yrs  
6  
7 x.val <- model.matrix(SeriousDlqin2yrs ~ RevolvingUtilizationOfUnsecuredLines + age + NumberOfTime30.59DaysPastD  
ueNotWorse  
8           + DebtRatio + MonthlyIncome + NumberOfOpenCreditLinesAndLoans + NumberOfTimes90DaysLate +  
9           NumberRealEstateLoansOrLines + NumberOfTime60.89DaysPastDueNotWorse  
10          + NumberOfDependents, data = testing_data)[, -1]  
11 y.val <- testing_data$SeriousDlqin2yrs  
12  
13  
14  
15  
16  
17
```

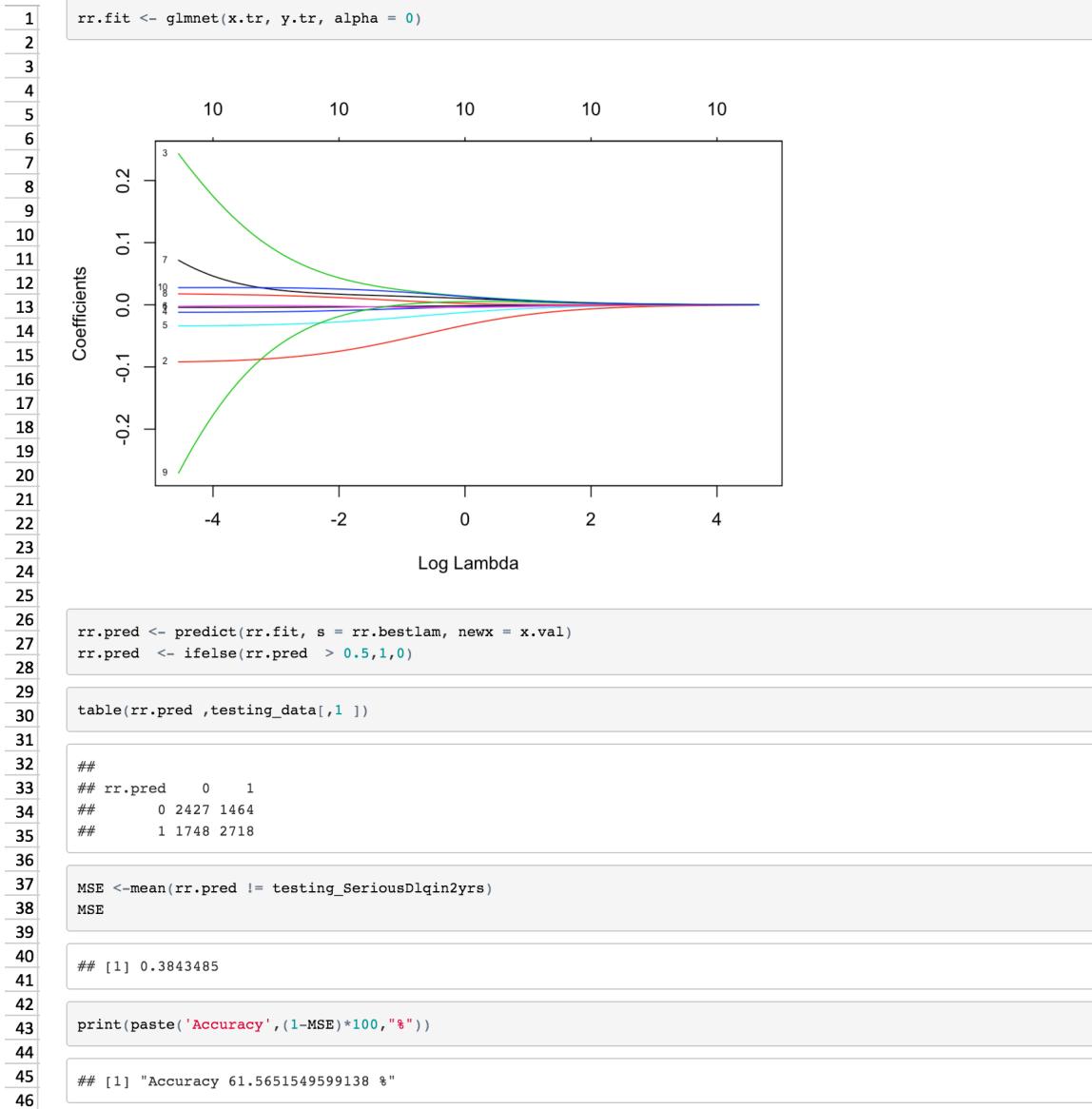
Cross Validation to obtain best lambda

```
18 set.seed(10)  
19 rr.cv <- cv.glmnet(x.tr, y.tr, alpha = 0)
```



First we turn our training data table into a Matrix (*line 1 above*) so we can perform regularization. Next we set seed (*line 19 above*) so we can repeat the process and get the same result. On line 20 we perform a cross validation to get the best Lambda for our ridge regularization. The Graph above denotes that our best lambda is between **-4.5** and **-4**.

Next we try to find the best predictors. Line 1 on the next page shows how this is done in R. The graph on the next page shows that almost all the predictors converged at **0** at the same time. We expected this might happen as an initial logistic regression model we generated should very significant p-value for all the **10** independent variables.



We now try to predict our dependent variable using the model we generated after applying the best lambda and eliminating not so useful predicators according to ridge regularization (line 26-28).

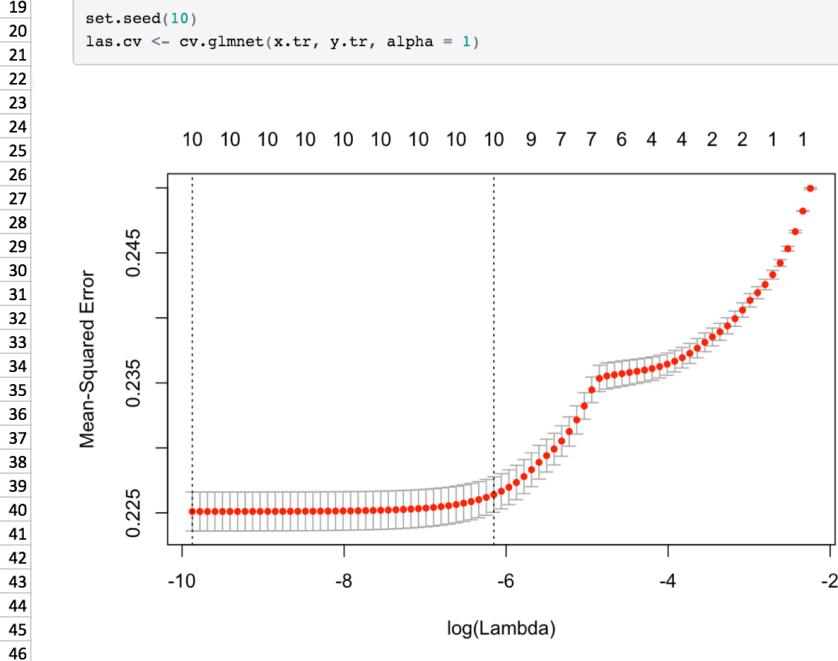
Now, it's time to see how ridge regularization performed. We compare the predicted output to the actual output from our testing data. The confusion matrix on line 33 shows a high **FP** and **FN** values. The accuracy is **61%**.

We anticipated the above result as we mentioned earlier. Almost all the independent variables are useful in predicting our independent variable. Therefore, the use of ridge regularization resulted in an underfitting model.

LASSO

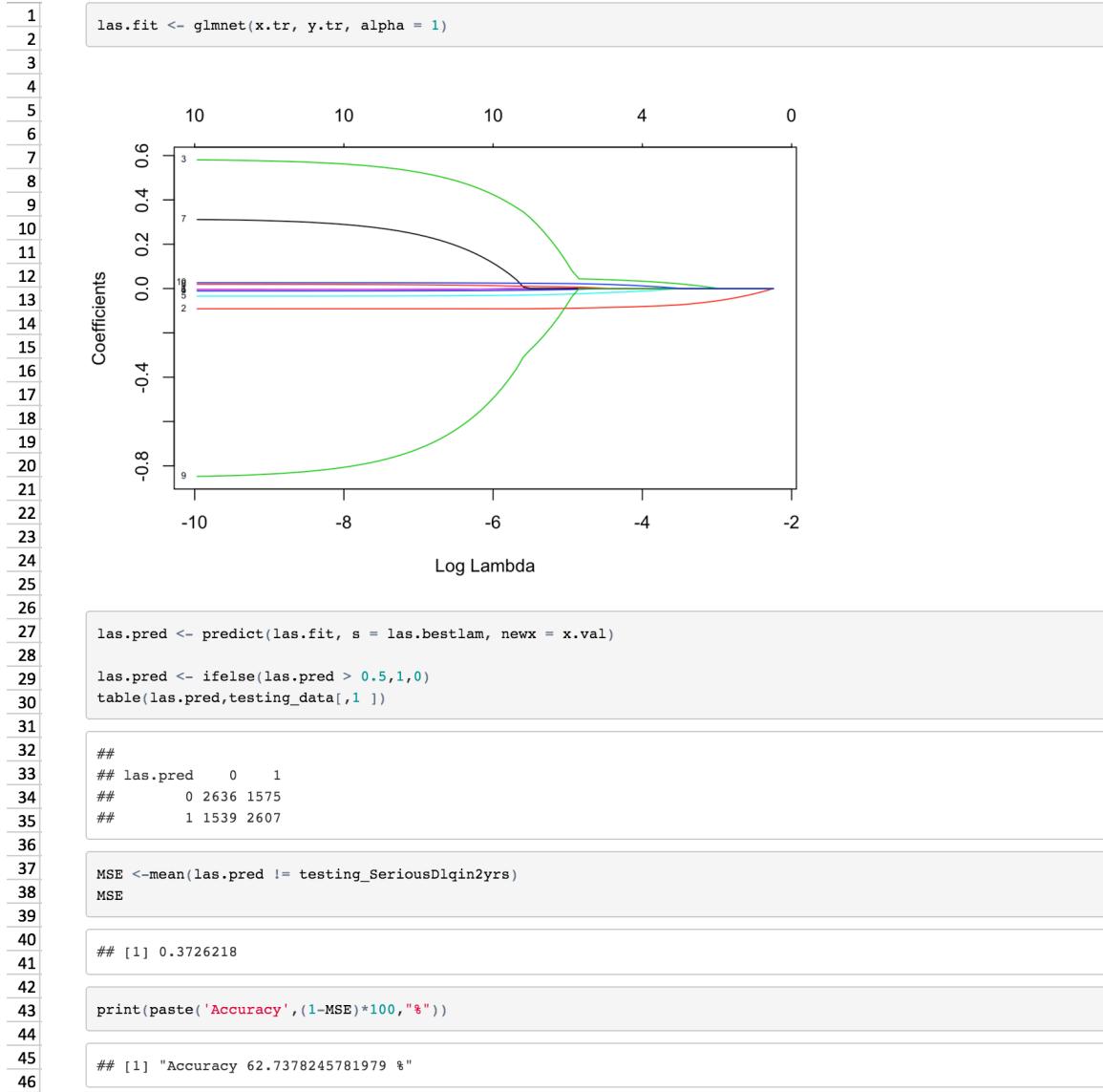
```
1 x.tr <- model.matrix(SeriousDlqin2yrs ~ RevolvingUtilizationOfUnsecuredLines + age + NumberOfTime30.59DaysPastDu  
2 eNotWorse  
3 + DebtRatio + MonthlyIncome +NumberOfOpenCreditLinesAndLoans + NumberOfTimes90DaysLate +  
4 NumberRealEstateLoansOrLines + NumberOfTime60.89DaysPastDueNotWorse  
5 + NumberOfDependents, data = training_data)[, -1]  
6 y.tr <- training_data$SeriousDlqin2yrs  
7  
8 x.val <- model.matrix(SeriousDlqin2yrs ~ RevolvingUtilizationOfUnsecuredLines + age + NumberOfTime30.59DaysPastDu  
9 ueNotWorse  
10 + DebtRatio + MonthlyIncome +NumberOfOpenCreditLinesAndLoans + NumberOfTimes90DaysLate +  
11 NumberRealEstateLoansOrLines + NumberOfTime60.89DaysPastDueNotWorse  
12 + NumberOfDependents, data = testing_data)[, -1]  
13 y.val <- testing_data$SeriousDlqin2yrs  
14  
15  
16
```

CV to obtain best lambda



Like ridge we also have to turn our training data table into a Matrix (*line 1 above*) so we can perform lasso regularization. Next we set seed (*line 19 above*) so we can repeat the process and get the same result. On line 20 we perform a cross validation to get the best Lambda for our ridge regularization. The Graph above denotes that our best lambda is between **-9.8** and **-6.2**.

Next we try to find the best predictors. Line 1 on the next page shows how this is done in R. The graph on the next page shows that almost all the predictors converged at **0** at the same time except for predictor **2** and **3**. The coefficients converge to **0** at log Lambda of **-6** which is earlier than ridge which was around **2**. We expected this might happen as lasso converges coefficients to zero faster than ridge.



We now try to predict our dependent variable using the model we generated after applying the best lambda and eliminating the not so useful predictors according to lasso regularization (line 27).

Now, it's time to see how lasso regularization performed. We compare the predicted output to the actual output from our testing data. The confusion matrix on line 33 shows a high FP and FN rates. The accuracy is 62%.

We anticipated the above result as we mentioned earlier. Almost all the independent variables are useful in predicting our independent variable. Therefore, the use of lasso regularization resulted in an underfitting model.

Although lasso regularization was fractionally better than ridge in this case, it's still not require for our particular problem. As we will demonstrate later, the result we got using both ridge and lasso regularization were not significantly better than when we did not use them.

Algorithms

Logistic Regression

We choose Logistic Regression as one of the classification algorithm. As mentioned in the tools section, we use R as our tool to get the logistic regression result. The commands don't have any big difference from the MATLAB commands.

The below graphs are for model results.

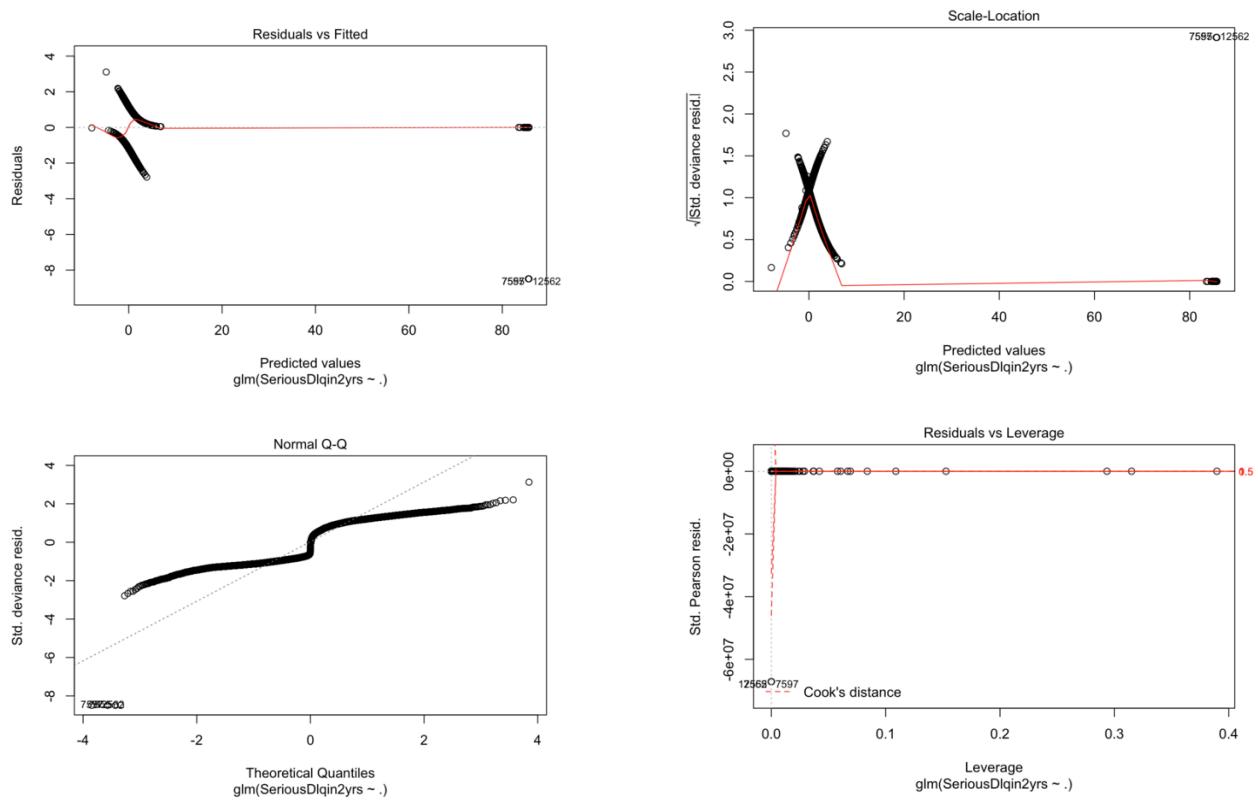
```
1 glmdl <- glm(SeriousDlqin2yrs ~ ., family=binomial(link='logit'), data=training_data)
2
3 ## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
4
5 summary(glmdl)
6
7
8
9 ## Call:
10 ## glm(formula = SeriousDlqin2yrs ~ ., family = binomial(link = "logit"),
11 ##       data = training_data)
12 ##
13 ## Deviance Residuals:
14 ##      Min        1Q     Median        3Q       Max
15 ## -8.4904   -1.0161   -0.2558   1.0741   3.1066
16 ##
17 ## Coefficients:
18 ##                               Estimate Std. Error z value Pr(>|z|)
19 ## (Intercept)               0.49780  0.03893 12.789 < 2e-16
20 ## RevolvingUtilizationOfUnsecuredLines -0.01224  0.02210 -0.554 0.579748
21 ## age                      -0.37683  0.02539 -14.841 < 2e-16
22 ## NumberOfTime30.59DaysPastDueNotWorse 3.62786  0.20618 17.596 < 2e-16
23 ## DebtRatio                 -0.06161  0.02726 -2.260 0.023799
24 ## MonthlyIncome              -0.20013  0.03543 -5.648 1.62e-08
25 ## NumberOfOpenCreditLinesAndLoans 0.01395  0.02793  0.500 0.617396
26 ## NumberOfTimes90DaysLate    3.18999  0.26769 11.917 < 2e-16
27 ## NumberOfRealEstateLoansOrLines 0.11153  0.02835  3.934 8.35e-05
28 ## NumberOfTime60.89DaysPastDueNotWorse -0.53621  0.32775 -1.636 0.101827
29 ## NumberOfDependents          0.09467  0.02452  3.861 0.000113
30 ##
31 ## (Intercept)                ***
32 ## RevolvingUtilizationOfUnsecuredLines
33 ## age                      ***
34 ## NumberOfTime30.59DaysPastDueNotWorse ***
35 ## DebtRatio                  *
36 ## MonthlyIncome                ***
37 ## NumberOfOpenCreditLinesAndLoans
38 ## NumberOfTimes90DaysLate     ***
39 ## NumberOfRealEstateLoansOrLines ***
40 ## NumberOfTime60.89DaysPastDueNotWorse
41 ## NumberOfDependents          ***
42 ##
43 ## ---
44 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
45 ##
46 ## (Dispersion parameter for binomial family taken to be 1)
47 ##
48 ## Null deviance: 11585.3 on 8356 degrees of freedom
49 ## Residual deviance: 9917.1 on 8346 degrees of freedom
50 ## AIC: 9939.1
51 ##
52 ## Number of Fisher Scoring iterations: 8
```

In line 2 we use the **logit** family of binomial classification and specify our training data. We also indicate that our dependent variable is **SeriouslyDlqin2yrs** and the rest are independent variable

We run this line to generate the model and the resulting model is summarized.

Line 17 to 21 shows us the coefficients of the independent variable as well as their respective p-values. R has a nice way of showing which p-values are significant. (***) means a very significant p-value or a p-value less than 0.001. As you can tell the All but two of the predictors have a significant p-value.

A plot of the residual vs fitted value, Scale-location, NormalQ_Q and Residual vs leverage of the resulting model can be seen below.



Next we try to predict using the model that we generated in the previous page. Line 1 show we are using the response type of Response to generate the probabilities of each instance belonging to the 1 or 0. The variable **GlmDlpredict** now has a vector of all the predicted probabilities.

```

1 glmdlpredict <- predict(glmmdl,testing_data,type='response')
2 glmdlpredict <- ifelse(glmdlpredict > 0.5,1,0)
3
4 table(glmdlpredict,testing_data[,1])
5
6 ##
7 ## glmdlpredict 0 1
8 ## 0 3182 1525
9 ## 1 993 2657
10
11 MSE <- mean(glmdlpredict != testing_SeriousDlqin2yrs)
12 MSE
13
14 ## [1] 0.3013043
15
16 print(paste('Accuracy',(1-MSE)*100,"%"))
17
18 ## [1] "Accuracy 69.8695704200072 %"
19
20 sqrt(mean((glmdlpredict - testing_data$SeriousDlqin2yrs)^2))
21
22 ## [1] 0.5489119
23
24 confusionMatrix(glmdlpredict, testing_SeriousDlqin2yrs)
25
26 ## Confusion Matrix and Statistics
27 ##
28 ##             Reference
29 ## Prediction 0 1
30 ##          0 3182 1525
31 ##          1 993 2657
32 ##          Accuracy : 0.6987
33 ##                     95% CI : (0.6887, 0.7085)
34 ##          No Information Rate : 0.5004
35 ##          P-Value [Acc > NIR] : < 2.2e-16
36 ##          Kappa : 0.3975
37 ##          Mcnemar's Test P-Value : < 2.2e-16
38 ##          Sensitivity : 0.7622
39 ##          Specificity : 0.6353
40 ##          Pos Pred Value : 0.6760
41 ##          Neg Pred Value : 0.7279
42 ##          Prevalence : 0.4996
43 ##          Detection Rate : 0.3808
44 ##          Detection Prevalence : 0.5632
45 ##          Balanced Accuracy : 0.6987
46 ##          'Positive' Class : 0
##
```

The next step on line two is to assign 1's and 0's to those probabilities with probabilities >0.5 assigned a 1 class and those less than 0.5 assigned as 0 class.

Once this is done we then compare this result to our testing data line 4 and get the MSE which is 0.301 in this case.

We then plot the confusion matrix (line 24) to see how good our model performed against our test data. We will provide and in-depth review of this result in the result and Analysis section of this report.

Linear Discriminant Analysis

The purpose of linear discriminant analysis (LDA) is to find the linear combinations of the original variables (10 in our dataset) that gives the best possible separation between the groups (default/ not default) in our data set.

If we want to separate the records by to give credit or not, the records come from 2 different variants, so the number of groups (G) is 2, and the number of variables is 10. The maximum number of useful discriminant functions that can separate the records is the minimum of G-1, and so in this case it is the minimum of 1 to 10.

```
1 ldamdl<-lda(SeriousDlqin2yrs ~., data = training_data)
2 ldamdl
3
4 ## Call:
5 ## lda(SeriousDlqin2yrs ~ ., data = training_data)
6 ##
7 ## Prior probabilities of groups:
8 ##      0         1
9 ## 0.5004188 0.4995812
10 ##
11 ## Group means:
12 ##   RevolvingUtilizationOfUnsecuredLines      age
13 ## 0          0.014683021 0.2230021
14 ## 1          -0.009784796 -0.2055636
15 ##   NumberofTime30.59DaysPastDueNotWorse  DebtRatio MonthlyIncome
16 ## 0          -0.1142063 0.01536177 0.07281932
17 ## 1          0.0859057 -0.01843612 -0.08549153
18 ##   NumberOfOpenCreditLinesAndLoans NumberofTimes90DaysLate
19 ## 0          0.04646429 -0.10069598
20 ## 1          -0.04778837 0.07623902
21 ##   NumberRealEstateLoansOrLines NumberofTime60.89DaysPastDueNotWorse
22 ## 0          0.016906134 -0.08432083
23 ## 1          0.004875361 0.05711844
24 ##   NumberOfDependents
25 ## 0          -0.08777666
26 ## 1          0.08663385
27 ##
28 ## Coefficients of linear discriminants:
29 ##                               LD1
30 ## RevolvingUtilizationOfUnsecuredLines -0.02546800
31 ## age                                -0.60302787
32 ## NumberofTime30.59DaysPastDueNotWorse 3.88982613
33 ## DebtRatio                            -0.07301906
34 ## MonthlyIncome                         -0.22669692
35 ## NumberOfOpenCreditLinesAndLoans       -0.03388304
36 ## NumberofTimes90DaysLate              2.10898015
37 ## NumberRealEstateLoansOrLines        0.12768022
38 ## NumberofTime60.89DaysPastDueNotWorse -5.70389196
39 ## NumberOfDependents                  0.17499667
```

We can carry out a linear discriminant analysis using the “lda()” function from the R “MASS” package. To use this function, we first need to install the “MASS” R package.

To calculate the separation achieved by discriminant function, we first need to calculate the value of each discriminant function, by substituting the variables’ values into the linear combination for the discriminant function (eg. $-0.025*V2 - 0.603*V3 + 3.889*V4 - 0.073*V5 - 0.226*V6 - 0.033*V7 - 1.661*V8 + 2.108*V9 + 0.127*V10 - 5.703*V11 + 0.174*V12$ for the first discriminant function).

```

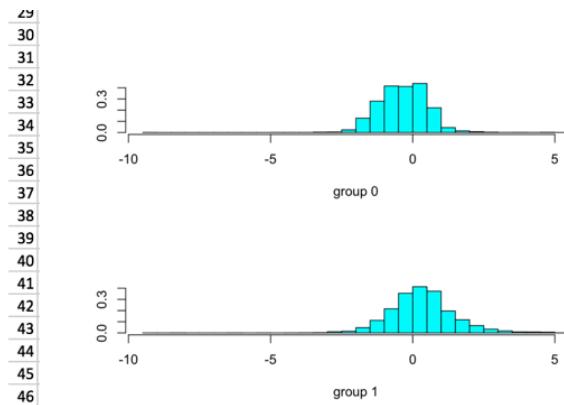
1 ldamdlpredict<-predict(ldamdl, newdata = testing_data[,c(2,3,4,5,6,7,8,9,10,11)])$class
2
3
4 table(ldamdlpredict,testing_data[,1])
5
6
7 ##
8 ## ldamdlpredict      0      1
9 ##                  0 2640 1577
10 ##                 1 1535 2605

```

We use the predict function in R to do the prediction against test data. We then build the confusionMatrix between ldamdlpredict and Testing_SeriousDlq dataset. To get the above confusion matrix results.

There are $1535+1577=3122$ samples that are misclassified, out of $(1535+1577 +2640+2605=)$ 8537 samples, therefore the accuracy is 0.627. The misclassification rate is quite high, and therefore the accuracy of the allocation rule is low.

To get a more accurate idea of how well the discriminant functions separates the groups, we would need to see a stacked histogram of the values for the two using some unseen “test set”, that is, using a set of data that was not used to calculate the linear discriminant function.



A nice way of displaying the results of a linear discriminant analysis (LDA) is to make a stacked histogram of the values of the discriminant function for the samples from different groups

We can do this using the “ldahist()” function in R. For example, to make a stacked histogram of the first discriminant function’s values for samples of the 2 different groups, we type:

```
> ldahist(data = testing_data.lda.values$x[2,3,4,5,6,7,8,9,10,11], g= $class)
```

We can see from the histogram that groups/classes 1 and 2 are not well separated by the discriminant functions, since the values for the zero class are between -4 and +3, while the values for class 1 are between -3 and 4, and so there is overlap in values. This clearly makes it challenging for the model to give good prediction accuracy.

Support Vector Machine

As a powerful tool of pattern recognition, support vector machine (SVM) was used to build a prediction model using the training data, and pass the test data to evaluate the quality of the prediction model with the ultimate goal of using the model to forecast the likelihood of experiencing serious financial distress for potential loaners.

```
1  svmmdl <- ksvm(SeriousDlqin2yrs ~ ., data=training_data, type = "C-svc", kernel = "rbfdot",kpar = list(sigma =
2  0.1), C = 10, prob.model = TRUE)
3
4  ## line search fails -1.01618 -0.08577839 1.823325e-05 7.463505e-06 -1.64527e-08 -9.743073e-09 -3.727038e-13
5
6  svmmdl
7
8  ## Support Vector Machine object of class "ksvm"
9  ##
10 ## SV type: C-svc (classification)
11 ## parameter : cost C = 10
12 ## Gaussian Radial Basis kernel function.
13 ## Hyperparameter : sigma = 0.1
14 ## Number of Support Vectors : 5150
15 ##
16 ## Objective Function Value : -48507.17
17 ## Training error : 0.249252
18 ## Probability model included.
```

In this project, we adopted kernel function of SVM , which is available in the kernlab package in R, to avoid the curse of dimensionality. More specially, we used the **rbfdot** kernel function, which is also referred to as Radial Basis kernel “Gaussian”. **C-svc** was passed as a parameter to do C classification.

Other parameters used in model building include hyper-parameter **sigma=0.1** (inverse kernel width for rbfdot, **C=10** (cost of constraints violation), and **prob.model=TRUE**(for calculating class probabilities).

```

1 svmmdlpredict <- predict(svmmdl,testing_data[,c(2,3,4,5,6,7,8,9,10,11)])
2 table(svmmdlpredict,testing_data[,1])
3
4
5 ##
6 ## svmmdlpredict    0     1
7 ##                 0 3572 1570
8 ##                 1  603 2612
9
10 MSE <-mean(svmmdlpredict != testing_SeriousDlqin2yrs)
11 MSE
12
13 ## [1] 0.2600215
14
15 print(paste('Accuracy',(1-MSE)*100,"%"))
16
17 ## [1] "Accuracy 73.9978461170276 %"
18
19
20 confusionMatrix(svmmdlpredict, testing_SeriousDlqin2yrs)
21
22
23 ## Confusion Matrix and Statistics
24 ##
25 ##             Reference
26 ## Prediction   0     1
27 ##                 0 3572 1570
28 ##                 1  603 2612
29 ##
30 ##                 Accuracy : 0.74
31 ##                           95% CI : (0.7304, 0.7494)
32 ##                           No Information Rate : 0.5004
33 ##                           P-Value [Acc > NIR] : < 2.2e-16
34 ##
35 ##                           Kappa : 0.4801
36 ## Mcnemar's Test P-Value : < 2.2e-16
37 ##
38 ##                           Sensitivity : 0.8556
39 ##                           Specificity : 0.6246
40 ##                           Pos Pred Value : 0.6947
41 ##                           Neg Pred Value : 0.8124
42 ##                           Prevalence : 0.4996
43 ##                           Detection Rate : 0.4274
44 ##                           Detection Prevalence : 0.6153
45 ##                           Balanced Accuracy : 0.7401
46 ##
47 ## 'Positive' Class : 0
48 ##

```

As a result, 5150 support vectors (out of 25071 training data points) were identified. Support vectors are the data points that lie closest to the decision surface. The range of the number of support vectors can be wide, depending on how much slack we allow and the complexity of our model. The percentage of support vectors is 20.5% for our training data, suggesting that the possibility for our model to overfit the train data is low.

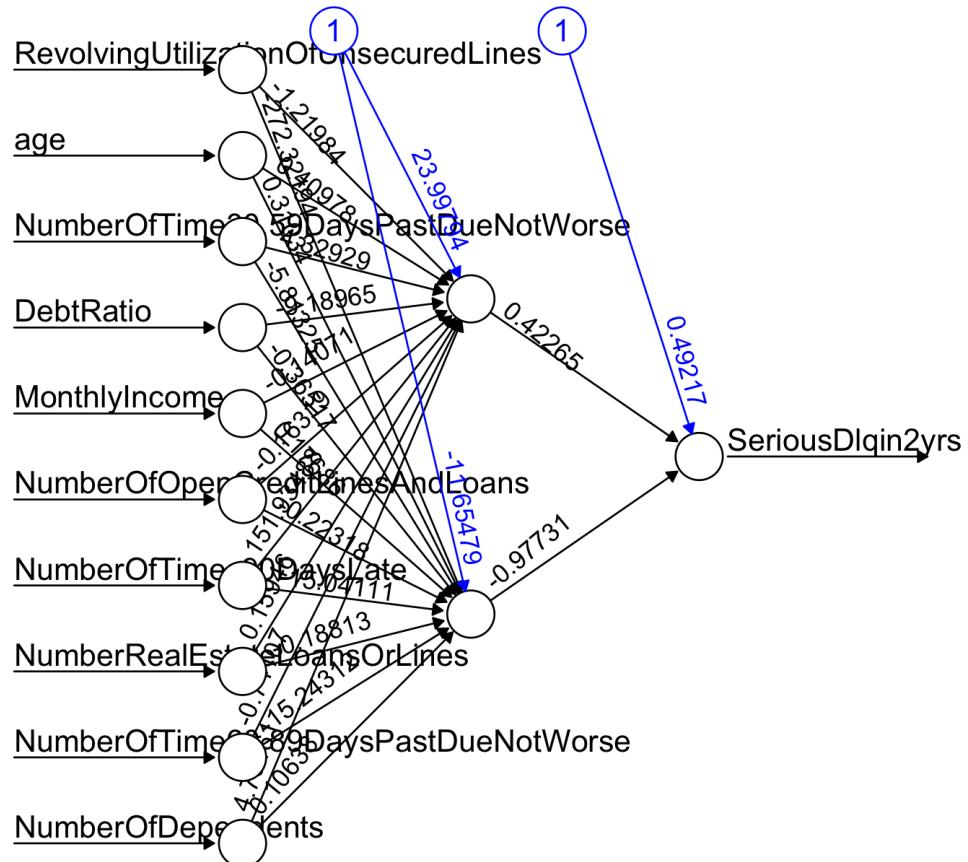
Using the test data (a total of 8357 records), the accuracy of our SVM model is 0.74, with a total of 2173 records being misclassified. In SVM, the accuracy depends on the trade-off between the complexity of a model and a large-margin that misclassifies some of the training data in the interest of better generalization. We will explore the results in details in the Analysis section.

Neural Network

For neural network we proceed the same way we did for all the other algorithms. First we train using our training data. Some of the parameters in the training function include hidden and threshold. **Hidden** refers to a vector of integers specifying the number of hidden neurons (vertices) in each layer. **Threshold** refers to a numeric value specifying the threshold for the partial derivatives of the error function as stopping criteria.

```
nnmdl <- neuralnet(SeriousDlqin2yrs ~ RevolvingUtilizationOfUnsecuredLines + age + NumberOfTime30.59DaysPastDueNotWorse  
+ DebtRatio + MonthlyIncome + NumberOfOpenCreditLinesAndLoans + NumberOfTimes90DaysLate +  
NumberOfRealEstateLoansOrLines + NumberOfTime60.89DaysPastDueNotWorse  
+ NumberOfDependents, data = training_data, hidden=2, threshold=0.01)
```

This model will be used as a parameter in our ANN predict function to predict the class of each of the dependent variable.



The image above represents the plot of the **nnmdl** model that was generated using our training dataset.

```

1 nnmdlpredict <- compute(nnmdl, testing_data[,c(2,3,4,5,6,7,8,9,10,11)])
2
3 results <- data.frame(actual = testing_SeriousDlqin2yrs, prediction = nnmdlpredict$net.result)
4 nnmdlpredict <- round(results$prediction)
5
6 table(nnmdlpredict, testing_data[, 1])
7
8 ##
9 ## nnmdlpredict      0      1
10 ##             0 3234  973
11 ##             1  941 3209
12
13 MSE <-mean(nnmdlpredict != testing_SeriousDlqin2yrs)
14 MSE
15
16 ## [1] 0.2290295561
17
18 print(paste('Accuracy', (1-MSE)*100, "%"))
19
20 ## [1] "Accuracy 77.0970443939213 %"
21
22
23 confusionMatrix(nnmdlpredict, testing_SeriousDlqin2yrs)
24
25 ## Confusion Matrix and Statistics
26 ##
27 ##          Reference
28 ## Prediction      0      1
29 ##             0 3234  973
30 ##             1  941 3209
31 ##          Accuracy : 0.7709704
32 ## 95% CI : (0.7618078, 0.7799438)
33 ## No Information Rate : 0.5004188
34 ## P-Value [Acc > NIR] : < 0.00000000000000022
35 ##          Kappa : 0.5419435
36 ## McNemar's Test P-Value : 0.4785833
37 ##          Sensitivity : 0.7746108
38 ##          Specificity : 0.7673362
39 ##          Pos Pred Value : 0.7687188
40 ##          Neg Pred Value : 0.7732530
41 ##          Prevalence : 0.4995812
42 ##          Detection Rate : 0.3869810
43 ## Detection Prevalence : 0.5034103
44 ##          Balanced Accuracy : 0.7709735
45 ##          'Positive' Class : 0
46 ##

```

On line 1 of the code above, we predict the dependent variable in our test set using the model generated in the previous page. Once the prediction is complete we generate a confusion matrix in line 27. Neural network give us an Accuracy of 77.1%(line 31) the highest of the algorithms so far. We will discuss this further in our result analysis section.

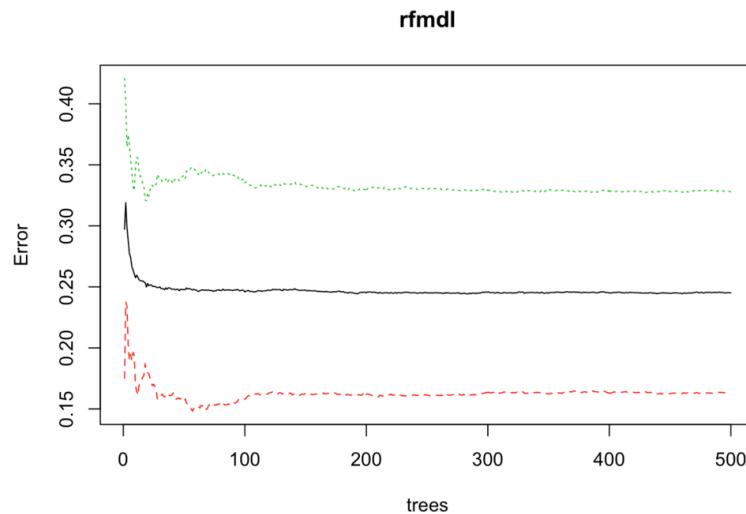
Random Forest

So far we've been working with algorithms that we've talk about in class. But out of curiosity we decided to test out some other highly talked about algorithm in the data science community.

Random forest is one of such algorithm. It's a notion of the general technique of random decision forests that are an ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of overfitting to their training set (Wikipedia).

```
set.seed(100)
rfmdl <- randomForest(training_data[,-c(1,2,7,12)], factor(training_data$SeriousDlqin2yrs),
    sampsize=1000, do.trace=TRUE, importance=TRUE, ntree=500, forest=TRUE)
```

```
plot(rfmdl)
```



First we set a seed so we can replicate our result. Then we build our model. Some of the parameters in the training function include sampsize, do.trace, importance, ntree, and forest.

sampsize	Size(s) of sample to draw. For classification, if sampsize is a vector of the length the number of strata, then sampling is stratified by strata, and the elements of sampsize indicate the numbers to be drawn from the strata.
importance	Should importance of predictors be assessed?
localImp	Should casewise importance measure be computed? (Setting this to TRUE will override importance.)
do.trace	If set to TRUE, give a more verbose output as randomForest is run. If set to some integer, then running output is printed for every do.trace trees.

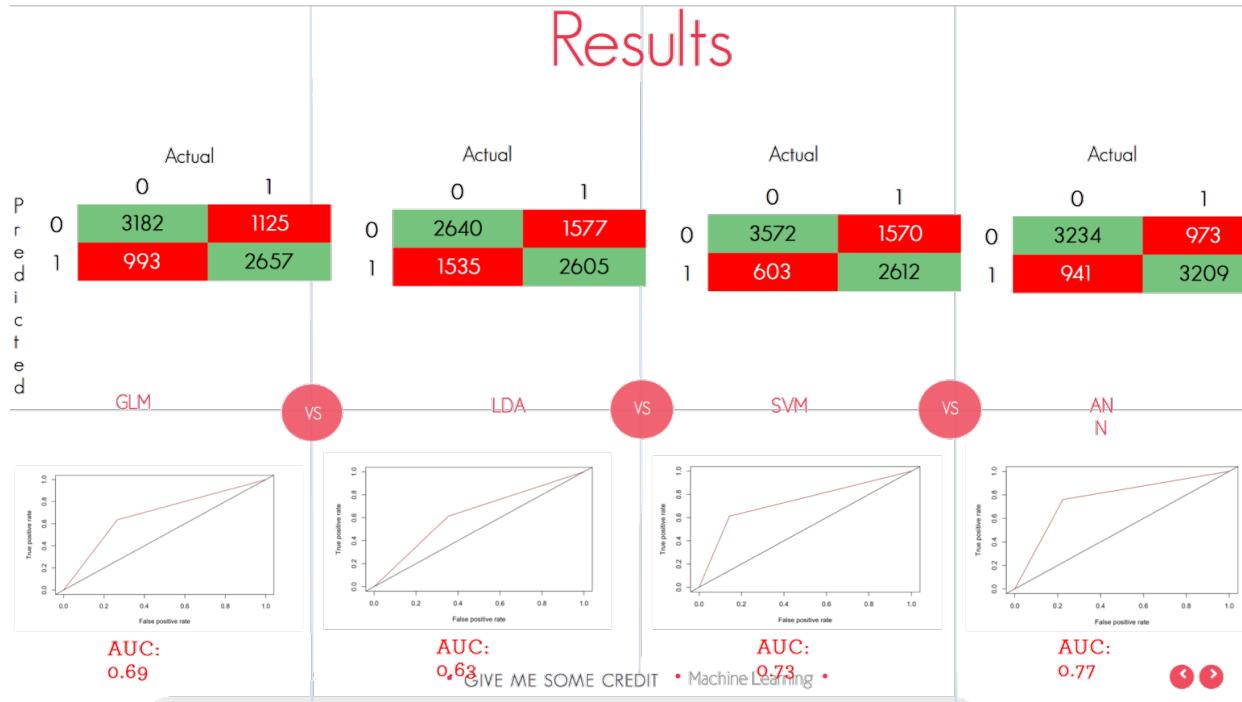
```
1 rfmdlpredict <- data.frame(SeriousDlqin2yrs=predict(rfmdl,testing_data[,-c(1,2,7,12)],type="prob")[,2])
2 rfmdlpredict <- ifelse(rfmdlpredict > 0.5,1,0)
3
4 table(rfmdlpredict,testing_data[,1])
5
6 ##
7 ## rfmdlpredict     0      1
8 ##              0 3420 1401
9 ##              1  726 2810
10
11 MSE <-mean(rfmdlpredict != testing_SeriousDlqin2yrs)
12 MSE
13
14 ## [1] 0.2545172
15
16
17 print(paste('Accuracy',(1-MSE)*100,"%"))
18
19 ## [1] "Accuracy 74.5482828766304 %"
20
```

On line 1 on the code above, we predict the dependent variable in our test set using the model generated in the previous page. Once the prediction is complete we generate a confusion matrix . Neural network give us an Accuracy of 74.5%(line 19 below) a little less accurate than neural network. We will discuss this further in our result analysis section.

```
1 confusionMatrix(rfmdlpredict, testing_SeriousDlqin2yrs)
2
3
4 ## Confusion Matrix and Statistics
5 ## Reference
6 ##           0      1
7 ##      0 3420 1401
8 ##      1  726 2810
9 ##
10 ##           Accuracy : 0.7455
11 ##                  95% CI : (0.736, 0.7548)
12 ##      No Information Rate : 0.5039
13 ##      P-Value [Acc > NIR] : < 2.2e-16
14 ##
15 ##           Kappa : 0.4916
16 ##  Mcnemar's Test P-Value : < 2.2e-16
17 ##
18 ##           Sensitivity : 0.8249
19 ##           Specificity : 0.6673
20 ##           Pos Pred Value : 0.7094
21 ##           Neg Pred Value : 0.7947
22 ##           Prevalence : 0.4961
23 ##           Detection Rate : 0.4092
24 ##           Detection Prevalence : 0.5769
25 ##           Balanced Accuracy : 0.7461
26 ##
27 ##           'Positive' Class : 0
28 ##
29
30 rfmdlpredict<- as.numeric(rfmdlpredict)
31
32 rrfroc<-prediction(rfmdlpredict, testing_SeriousDlqin2yrs, label.ordering = NULL)
33
34 rrfroc.perf <- performance(rrfroc, measure = "tpr", x.measure = "fpr")
35 plot(rrfroc.perf, col = "dark red")
36 abline(a=0, b= 1)
```

Results Evaluation and Algorithm comparison

Now that we've generated all the confusion matrix and the corresponding ROC curves, we will now evaluate the result of these algorithm. A snapshot of the result all 4 algorithms can be seen below



A grand view of all the results indicated that **LDA** was the worst performing algorithm. We suspect the reason could be that Linear Discriminant Analysis (LDA) is most commonly used as dimensionality reduction technique in the pre-processing step for pattern-classification and machine learning applications. And since our data set is not normally distributed, LDA found it difficult to correctly classify some instances. It should be mentioned that LDA assumes normal distributed data, features that are statistically independent, and identical covariance matrices for every class.

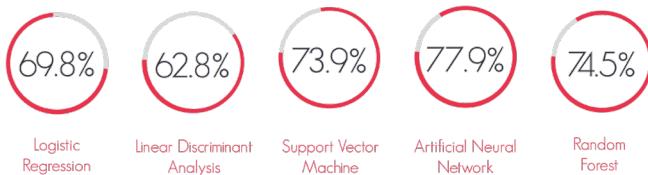
Logistic Regression with an AUC of 0.69 is the second to the last performing in terms of accuracy. Unlike LDA Logistic regression has several advantages over discriminant analysis some which are that it is more robust: the independent variables don't have to be normally distributed, or have equal variance in each group and It does not assume a linear relationship between the independent variable and the dependent variable. Unfortunately, the advantages of logistic regression come at a cost: it requires much more data to achieve stable, meaningful results. With standard regression, and LDA, typically 20 data points per predictor is considered the lower bound. For logistic regression, at least 50 data points per predictor is necessary to achieve stable results. This is not the case with our dataset which had about 10 predictors.

Third on the list is **SVM**, which had an AUC of 0.73. SVM was good in this case because SVM has a regularization parameter, which makes the user think about avoiding over-fitting. it also uses the kernel trick, so you can build in expert knowledge about the problem via engineering the kernel.

Finally, we examine the result of ANN. **Artificial Neural Network** performed the best with an AUC of 0.77. As in most cases it's usually difficult to explain why ANN gives the best results but our inclination is that you'll find that finite ANNs with even one hidden layer can approximate a huge class of functions. This property turns out to be fairly robust relative to the exact choice of network: by this we mean, you don't have to have the choice of network exactly right in order to approximate a particular function - there may be a network which is ideal for that function, but you can still do quite well with a similar network (of sufficient complexity).

With a new appreciation for the computational power of neural networks, we can return to the problem of learning real-world patterns like the one we have. We use ANN in situation were we suspect there is a rule governing how this a value changes as a function of the input. However, the scenario is very complicated, and you are not sure how to model it with an equation or an algorithm, this is where we ask ANN to model the rule for us. This in most cases returns very good results. This we suspect is the case here.

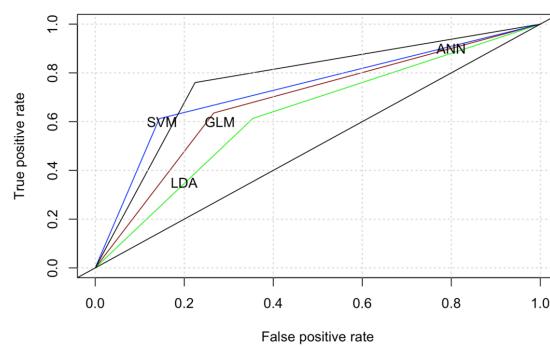
Below is a snapshot of the Accuracy of all the algorithms.



The ROC shows also the performance of all 4 algorithms with ANN occupying the most area and LDA doing not as well.

plot all ROCS

```
plot( glmroc.perf, col = "dark red")
plot(ladaroc.perf, add = TRUE, col = "green")
plot(svmroc.perf, add = TRUE, col = "blue")
plot(nnroc.perf, add = TRUE, col = "black")
grid()
abline(a=0, b= 1)
text(.8,.9,"ANN")
text(.28,.6,"GLM")
text(.15,.6,"SVM")
text(.2,.35,"LDA")
```

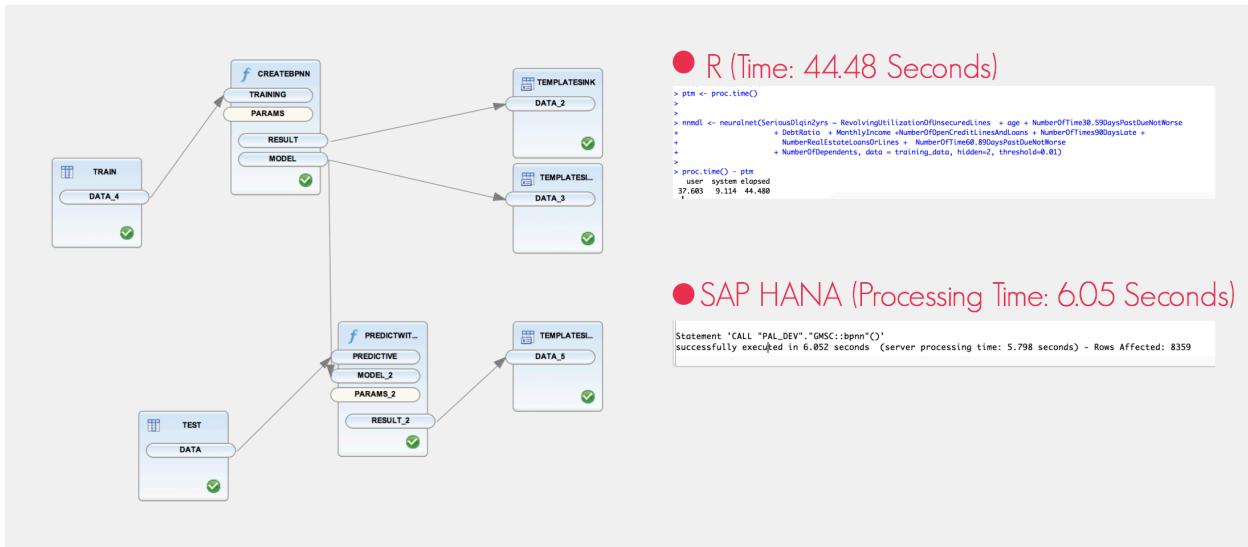


Performance Comparison with SAP HANA

We wanted to see how the algorithm performed when ran in an in memory database like SAP HANA. To this we wanted to compare the algorithm that took the longest time to run in R and compare it with the run on SAP HANA.

ANN averaged about 40 seconds when ran on our dataset in R environment. The longest time of all the other algorithms. This dropped down to about 6 seconds when ran on SAP HANA.

We can easily conclude that given the computational complexity of ANN, we would recommend running ANN on an in memory database especially if your ANN structure has several layers and neurons.



Lessons Learned and Challenges

There were several challenges that we encountered and lessons we learnt throughout this project but these three were particularly impactful.

Preprocessing and wrangling the data: Preprocessing is an essential part of predictive analysis. This part of the project took the most time and was very critical when it came time to use the algorithms. From normalization to stratification and outlier elimination it played a key role in making sure we are not feeding garbage to our algorithms.

Handling the Skew of the data: Skew in data can impact the result you get. This we discovered earlier in our wrangling process. Skew data can produce results with high accuracy and low precision. A lesson we've learnt and tend to look out for in future work with datasets.

Choosing the right parameters for the algorithms: Choosing the right parameters can impact your results especially for parametric algorithms like the ones we examined in this project. We did not explore a lot of the parameters present for most of the algorithms. We suspect that a more thorough examination and use of these would improve our model significantly.

Conclusion

In conclusion, we were able to improve fractionally by .9% the predictive accuracy of the credit scoring sample model, which had an accuracy of 77% by using Artificial Neural Network. ANN performed the best out of all the other algorithm for our particular use case.

References

G. C. Cawley and N. L. C. Talbot, Over-fitting in model selection and subsequent selection bias in performance evaluation, Journal of Machine Learning Research, 2010. Research, vol. 11, pp. 2079-2107, July 2010. (pdf)

<https://victorfang.wordpress.com/2011/05/10/advantages-and-disadvantages-of-logistic-regression/>

<https://www.quora.com/Why-do-Artificial-Neural-Networks-work-the-way-they-do>