

# LASSO\_NN

Cavalry

10/29/2020

## Step 0: Load Libraries

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --
## v ggplot2 3.3.2      v purrr  0.3.4
## v tibble  3.0.4      v dplyr  1.0.2
## v tidyr   1.1.2      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.0
## Warning: package 'ggplot2' was built under R version 3.6.2
## Warning: package 'tibble' was built under R version 3.6.2
## Warning: package 'tidyr' was built under R version 3.6.2
## Warning: package 'readr' was built under R version 3.6.2
## Warning: package 'purrr' was built under R version 3.6.2
## Warning: package 'dplyr' was built under R version 3.6.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(readr)
library(caret)

## Loading required package: lattice
## Warning: package 'lattice' was built under R version 3.6.2
##
## Attaching package: 'caret'
## The following object is masked from 'package:purrr':
##
##     lift

library(WeightedROC)
source("./LASSO_NN.R")

## Warning: package 'glmnet' was built under R version 3.6.2
## Loading required package: Matrix
##
## Attaching package: 'Matrix'
```

```
## The following objects are masked from 'package:tidyr':
##
##   expand, pack, unpack
## Loaded glmnet 4.0-2
##
## Attaching package: 'neuralnet'
## The following object is masked from 'package:dplyr':
##
##   compute
## Warning: package 'e1071' was built under R version 3.6.2
source("./svm.R")
```

## step 1: Load Data

```
datatrain = read_csv("../data/feature_train_balanced_data.csv")

##
## -- Column specification -----
## cols(
##   .default = col_double()
## )
## i Use `spec()` for the full column specifications.
datatest = read_csv("../data/feature_test_data.csv")

##
## -- Column specification -----
## cols(
##   .default = col_double()
## )
## i Use `spec()` for the full column specifications.
location = which(colnames(datatrain) == "HEALTH")
trainx = datatrain[,-location]
trainy = datatrain$HEALTH
testx = datatest[, -location]
testy = datatest$HEALTH

weight_test <- rep(NA, length(testy))
for (v in unique(testy)){
  weight_test[testy == v] = 0.5 * length(testy) / length(testy[testy == v])
}
```

## Step 2: LASSO method

```
begin <- Sys.time()
lasso <- lassoReg(trainx, trainy)
summary(lasso)

##           Length Class      Mode
## a0           1    -none-  numeric
```

```
## beta      50      dgCMatrix S4
## df         1      -none-    numeric
## dim        2      -none-    numeric
## lambda     1      -none-    numeric
## dev.ratio  1      -none-    numeric
## nulldev    1      -none-    numeric
## npasses    1      -none-    numeric
## jerr        1      -none-    numeric
## offset     1      -none-    logical
## call       7      -none-    call
## nobs       1      -none-    numeric
```

```
coef <- rbind("(intercept)" = lasso$a0, as.data.frame(as.matrix(lasso$beta))) %>%
  dplyr::arrange(desc(abs(s0))) %>% rename(Coef = s0)
```

```
coef
```

```
##              Coef
## (intercept)  0.5664298638
## FS1_1        -0.1225947683
## fpl          0.0738409785
## FWB2_1       -0.0631319263
## FS2_1        -0.0590656855
## FS1_4        -0.0576781116
## SWB_2        0.0527978563
## FWB2_4       -0.0464642620
## FS1_2        -0.0432104251
## FWB1_1       0.0400702432
## ACT1_1       0.0390507433
## FS1_7        -0.0318541502
## SWB_1        0.0314262665
## SUBKNOWL1    0.0310169724
## FWB1_3       -0.0301354835
## FS1_5        -0.0266557270
## ASK1_1       -0.0229151142
## FSscore      0.0227865827
## FS1_3        -0.0208277030
## FWB1_2       0.0197213286
## SUBNUMERACY1 0.0191392218
## FS2_3        0.0189430641
## FINGOALS     0.0185804675
## MANAGE1_4    -0.0169795874
## PROPPLAN_3   0.0151763375
## PROPPLAN_1   -0.0143315117
## FWB1_6       -0.0130012505
## SWB_3        0.0127521195
## ASK1_2       0.0126977696
## MANAGE1_3    -0.0124272348
## FWBscore     -0.0123473225
## FWB2_3       -0.0106533231
## sample       -0.0104712910
## FWB1_4       0.0095153457
## GOALCONF     0.0093415322
## CHANGEABLE   -0.0087591672
## FS1_6        0.0086364512
```

```
## LMscore      -0.0085541812
## FWB2_2       0.0078657332
## MANAGE1_2    -0.0076115008
## MANAGE1_1     0.0072534828
## AUTOMATED_1  -0.0056257656
## ACT1_2       0.0045198728
## FS2_2        -0.0036105505
## FWB1_5       -0.0034181454
## SUBNUMERACY2 -0.0025427767
## SAVEHABIT    0.0021405228
## FRUGALITY    -0.0012310856
## AUTOMATED_2  0.0009062519
## PROPPLAN_4   0.0003389113
## PROPPLAN_2   0.0000000000

pred <- lassoPred(lasso, testx)
pred <- ifelse(pred > mean(pred), 1, 0)

cat("MSE is", mean((pred - testy)^2), ". Accuracy is", mean(pred==testy), ".")

## MSE is 0.3652283 . Accuracy is 0.6347717 .
confusionMatrix(factor(pred), factor(testy))$byClass

##           Sensitivity      Specificity      Pos Pred Value
##           0.8016194      0.6042899      0.2701228
##           Neg Pred Value      Precision      Recall
##           0.9434180      0.2701228      0.8016194
##           F1      Prevalence      Detection Rate
##           0.4040816      0.1544715      0.1238274
## Detection Prevalence      Balanced Accuracy
##           0.4584115      0.7029547

tpr.fpr <- WeightedROC(pred, testy, weight_test)
auc.log <- WeightedAUC(tpr.fpr)

cat("AUC is", auc.log)

## AUC is 0.7029547

end <- Sys.time()
timedif <- end - begin
cat("Time for running", "LASSO", "is", timedif)

## Time for running LASSO is 0.5473251
timedif

## Time difference of 0.5473251 secs
```

## Step 3: Support Vector Machine

### 3.1 SVM Tune

```
begin <- Sys.time()
source("./svm.R")
# Tune SVM
set.seed(2020)
```

```

opt.svm <- svm_tune(scale(as.matrix(trainx)), trainy)
bestgamma = opt.svm$best.parameters$gamma
bestcost = opt.svm$best.parameters$cost

# runs for about 21 minutes.
# tuned result is gamma=0.001, cost = 0.21
bestgamma; bestcost

end <- Sys.time()
timedif <- end - begin
cat("Time for tuning","SVM","is", timedif)
timedif

```

### 3.2 SVM Train

```

begin <- Sys.time()
#svm_fit <- svm_train(scale(as.matrix(trainx)), trainy, bestgamma, bestcost)
svm_fit <- svm_train(scale(as.matrix(trainx)), trainy, 0.001, 0.21)

```

### 3.3 SVM Test

```

svm_pred <- svm_test(svm_fit, scale(as.matrix(testx)))

svm.pred <- ifelse(svm_pred > mean(svm_pred), 1, 0)

# Calculate Accuracy
cat("MSE is",mean((svm.pred-testy)^2),
    "and accuracy is",mean(svm.pred==testy))

```

```
## MSE is 0.355222 and accuracy is 0.644778
```

```
confusionMatrix(factor(svm.pred), factor(testy))$byClass
```

```
##          Sensitivity          Specificity      Pos Pred Value
##          0.7773279          0.6205621          0.2723404
##      Neg Pred Value          Precision          Recall
##          0.9384787          0.2723404          0.7773279
##              F1          Prevalence      Detection Rate
##          0.4033613          0.1544715          0.1200750
## Detection Prevalence      Balanced Accuracy
##          0.4409006          0.6989450

```

```

tpr.fpr <- WeightedROC(svm.pred, testy, weight_test)
auc.svm <- WeightedAUC(tpr.fpr)
cat("AUC is",auc.svm)

```

```
## AUC is 0.698945
```

```

end <- Sys.time()
timedif <- end - begin
cat("Time for running","SVM","is", timedif)

```

```
## Time for running SVM is 3.09303
```

```
timedif
```

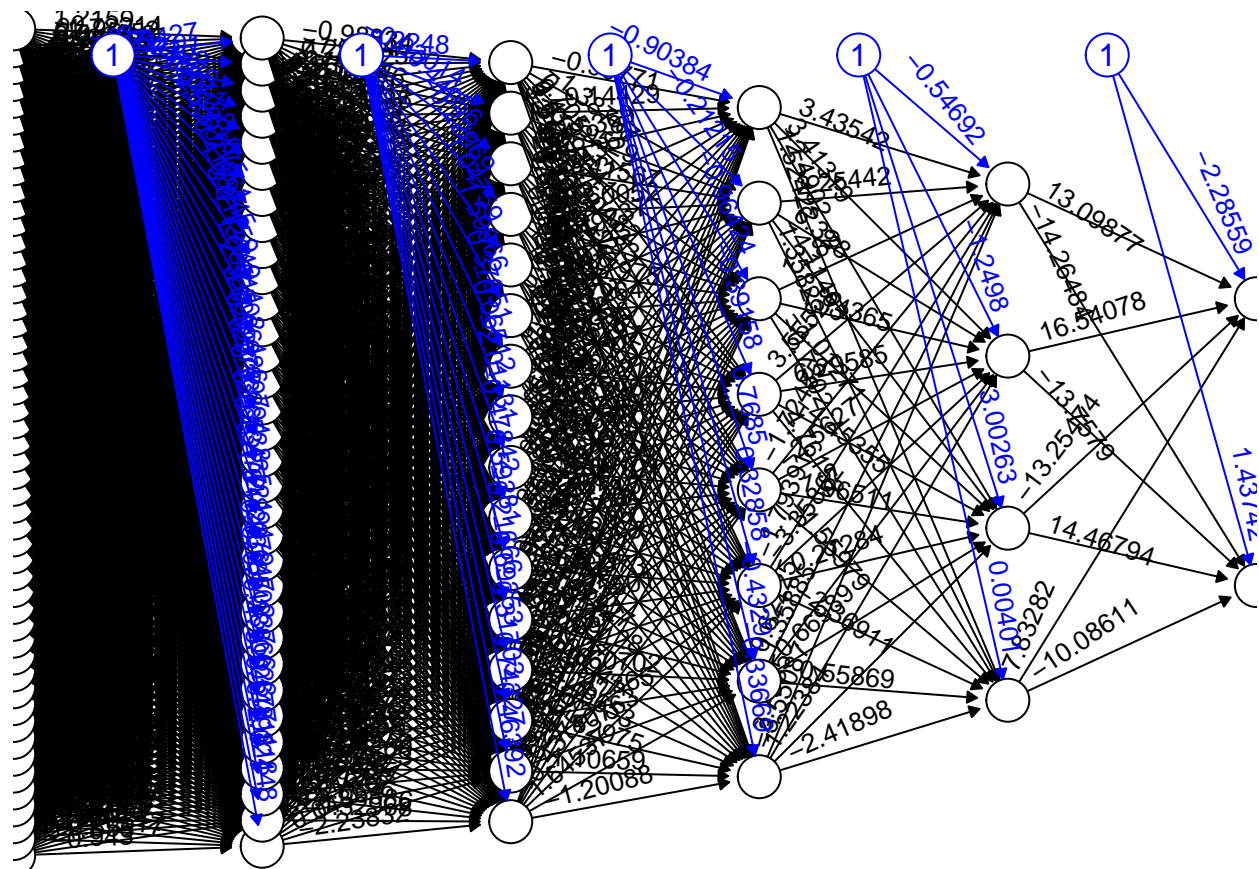
```
## Time difference of 3.09303 mins
```

## Step 4 Neural Net

```
source("./LASSO_NN.R")
begin <- Sys.time()
nndata <- data.frame(
  scale(as.matrix(trainx)), HEALTH = trainy)

nnfit <- nnReg(nndata, hd = c(32, 16, 8, 4))

plot(nnfit, rep = "best")
```



```
end <- Sys.time()
timedif <- end - begin
cat("Time for training", "neuralnet", "is", timedif)
```

```
## Time for training neuralnet is 22.09293
```

```
timedif
```

```
## Time difference of 22.09293 secs
```

```
begin <- Sys.time()
pred_nn <- predict(nnfit, testx)
summary(pred_nn)
```

```
##           V1           V2
## Min.      :0.0000002   Min.      :0.000000
## 1st Qu.:0.9931049     1st Qu.:0.000000
## Median :1.0000000     Median :0.000000
## Mean      :0.8748955   Mean      :0.131950
## 3rd Qu.:1.0000000     3rd Qu.:0.009879
## Max.      :1.0000000   Max.      :1.000000

nn.prediction <- round(pred_nn[,1])
summary(nn.prediction)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.0000  1.0000  1.0000  0.8987  1.0000  1.0000

confusionMatrix(factor(nn.prediction), factor(testy))$byClass

##           Sensitivity      Specificity      Pos Pred Value
##           0.10526316      0.89940828      0.16049383
##      Neg Pred Value      Precision      Recall
##           0.84620738      0.16049383      0.10526316
##           F1      Prevalence      Detection Rate
##           0.12713936      0.15447154      0.01626016
## Detection Prevalence      Balanced Accuracy
##           0.10131332      0.50233572

tpr.fpr <- WeightedROC(nn.prediction, testy, weight_test)
auc.nn <- WeightedAUC(tpr.fpr)

cat("AUC is", auc.nn)

## AUC is 0.5023357

end <- Sys.time()
timedif <- end - begin
cat("Time for predicting using","neuralnet","is", timedif)

## Time for predicting using neuralnet is 0.01925993
timedif

## Time difference of 0.01925993 secs
```

## Step 5 Naive Bayes

```
begin <- Sys.time()

nbfit = nbReg(datatrain)
nbpred = predict(nbfit, testx)

mean((as.numeric(nbpred)-1 - testy)^2)

## [1] 0.3327079

confusionMatrix(nbpred, factor(testy))$byClass

##           Sensitivity      Specificity      Pos Pred Value
##           0.6761134      0.6656805      0.2697900
##      Neg Pred Value      Precision      Recall
```

```
##          0.9183673          0.2697900          0.6761134
##          F1          Prevalence          Detection Rate
##          0.3856813          0.1544715          0.1044403
## Detection Prevalence    Balanced Accuracy
##          0.3871169          0.6708969
```

```
tpr.fpr <- WeightedROC(as.numeric(nbpred)-1, testy, weight_test)
auc.nb <- WeightedAUC(tpr.fpr)

cat("AUC is", auc.nb)
```

```
## AUC is 0.6708969
```

```
end <- Sys.time()
timedif <- end - begin
cat("Time for running", "Naive Bayies", "is", timedif)
```

```
## Time for running Naive Bayies is 0.836715
timedif
```

```
## Time difference of 0.836715 secs
```