

# Application of Logistic Regression via packages `stats::glm` and `glmnet` for Enhancer Prediction data.

Shaurya Jauhari (Email: [shauryajauhari@gzhmu.edu.cn](mailto:shauryajauhari@gzhmu.edu.cn))

2020-03-31

This is a R Markdown document on *glmnet* and *stats::glm* packages. These packages provide functionalities to cater to logistic regression problems, amongst others. Let's begin with *glmnet* first.

## glmnet

```
install.packages("glmnet",
                 repos = "https://cran.us.r-project.org")

## Installing package into '/Users/mei/Library/R/3.6/library'
## (as 'lib' is unspecified)

##
## The downloaded binary packages are in
## /var/folders/hm/c3_fjypn62v5xh5b5ygv267m0000gn/T//RtmpBWxxif/downloaded_packages
library(glmnet)

## Loading required package: Matrix

## Loaded glmnet 3.0-2

The dataset derivation has been lengthily described here

epdata <- readRDS("../Machine_Learning_Deep_Learning/predictionData.rds")

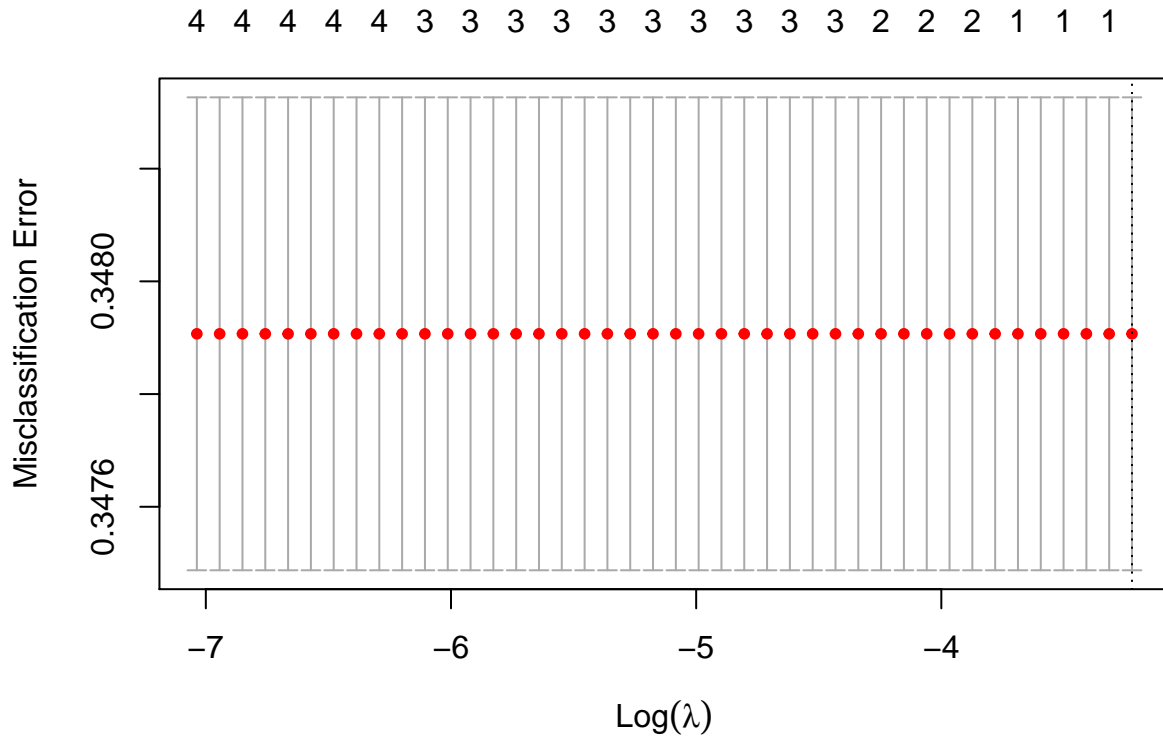
## Subsetting just features and class columns.
epdata <- epdata[,4:8]
colnames(epdata)[1:4] <- c("reads_h3k27ac", "reads_h3k4me1", "reads_h3k4me2", "reads_h3k4me3")
```

You can always seek help in the R documentation with `?`.

The class labels 0 and 1 represent “Enhancer” and “Non-Enhancer” categories, respectively.

```
set.seed(005)
cv.modelfit <- cv.glmnet(as.matrix(epdata[,1:4]),
                        epdata$class,
                        family = "binomial",
                        type.measure = "class",
                        alpha = 1,
                        nlambda = 100)

plot(cv.modelfit)
```



```
cat("There are", length(cv.modelfit$lambda),
    "lambda values in all:",
    cv.modelfit$lambda,
    ", out of which",
    min(cv.modelfit$lambda),
    "is the minimum, while",
    cv.modelfit$lambda.1se,
    "denotes the value at which the model is optimized at one standard error.")
```

```
## There are 42 lambda values in all: 0.0398659 0.03632432 0.03309737 0.03015709 0.02747802 0.02503695
```

The plot shows the models (with varying lambda values) that *glmnet* has fit, along with the misclassification error associated with each model. The first dotted line highlights the minimum misclassification error, while the second one is the highly regularized model within *1se* (one standard error).

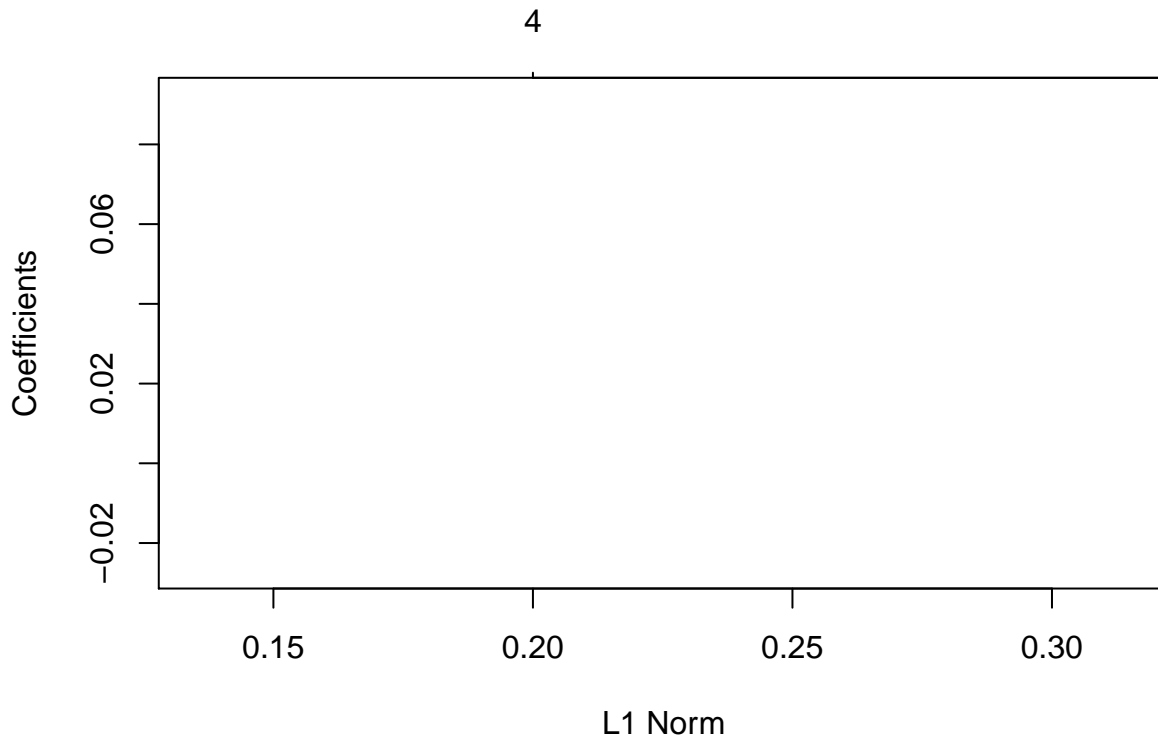
```
set.seed(2)
modelfit <- glmnet(as.matrix(epdata[,1:4]),
                  epdata$class,
                  family = "binomial",
                  alpha = 1,
                  lambda = min(cv.modelfit$lambda))
```

```
# Listing non-zero coefficients
```

```
print(modelfit$beta[,1])
```

```
## reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3
## -0.02666590 0.09195840 0.06621806 0.04035916
```

```
plot(modelfit)
```



Note

that the **features must be presented as a data matrix**, while the **response variable is a factor with two levels**. On calling the *glmnet*, we get information under 3 heads: *Df* signifies the number of non-zero coefficients from left to right, i.e. in this case coefficients for reads\_h3k27ac, reads\_h3k4me3, reads\_h3k4me2, reads\_h3k4me1; *%Dev* represents deviation; and *Lambda* represents the penalties imposed by the model. They would typically be limited to 100, but could even halt early if insufficient deviation is observed. Also, by default elastic-net (lasso+ridge) is used for regularization task by the *glmnet*, which could be set to lasso ( $\alpha=1$ ) or ridge ( $\alpha=0$ ).

```
coef(modelfit)[,1]
```

```
## (Intercept) reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3
## 0.53332669 -0.02666590 0.09195840 0.06621806 0.04035916
```

```
predict(modelfit, type="coef")
```

```
## 5 x 1 sparse Matrix of class "dgCMatrix"
##          s0
## (Intercept) 0.53332669
## reads_h3k27ac -0.02666590
## reads_h3k4me1 0.09195840
## reads_h3k4me2 0.06621806
## reads_h3k4me3 0.04035916
```

“.” here symbolizes 0.

## Exercices.

1. Try to fit the model with varying lambda, say *cv.modelfit\$lambda.1se*.
2. Try the above for  $\alpha = 0$ , i.e. ridge penalty.

3. Try the above for any value between 0 and 1; that's **elastic-net** regularization.
4. Try the above template for several available datasets at <http://archive.ics.uci.edu/ml/index.php>.

```
# Model fit with *cv.modelfit$lambda.1se*.
set.seed(5)
modelfit1 <- glmnet(as.matrix(epdata[,1:4]),
                    epdata$class,
                    family = "binomial",
                    alpha = 1,
                    lambda = cv.modelfit$lambda.1se)

# Listing non-zero coefficients

print(modelfit1$beta[,1])

## reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3
##           0           0           0           0

plot(modelfit1)

## Warning in plotCoef(x$beta, lambda = x$lambda, df = x$df, dev = x$dev.ratio, :
## No plot produced since all coefficients zero
## NULL

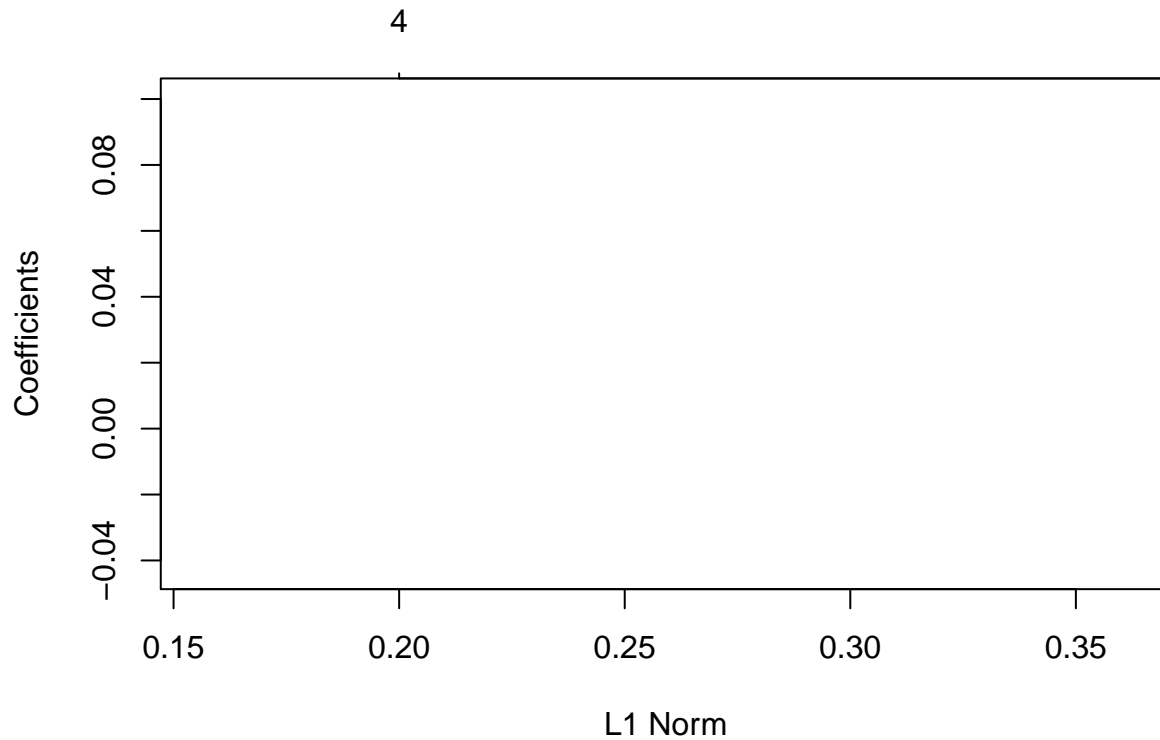
# Model fit with alpha=0.
set.seed(3)
modelfit2 <- glmnet(as.matrix(epdata[,1:4]),
                    epdata$class,
                    family = "binomial",
                    alpha = 0,
                    lambda = min(cv.modelfit$lambda))

# Listing non-zero coefficients

print(modelfit2$beta[,1])

## reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3
## -0.04296635  0.10050911  0.07047058  0.04514512

plot(modelfit2)
```



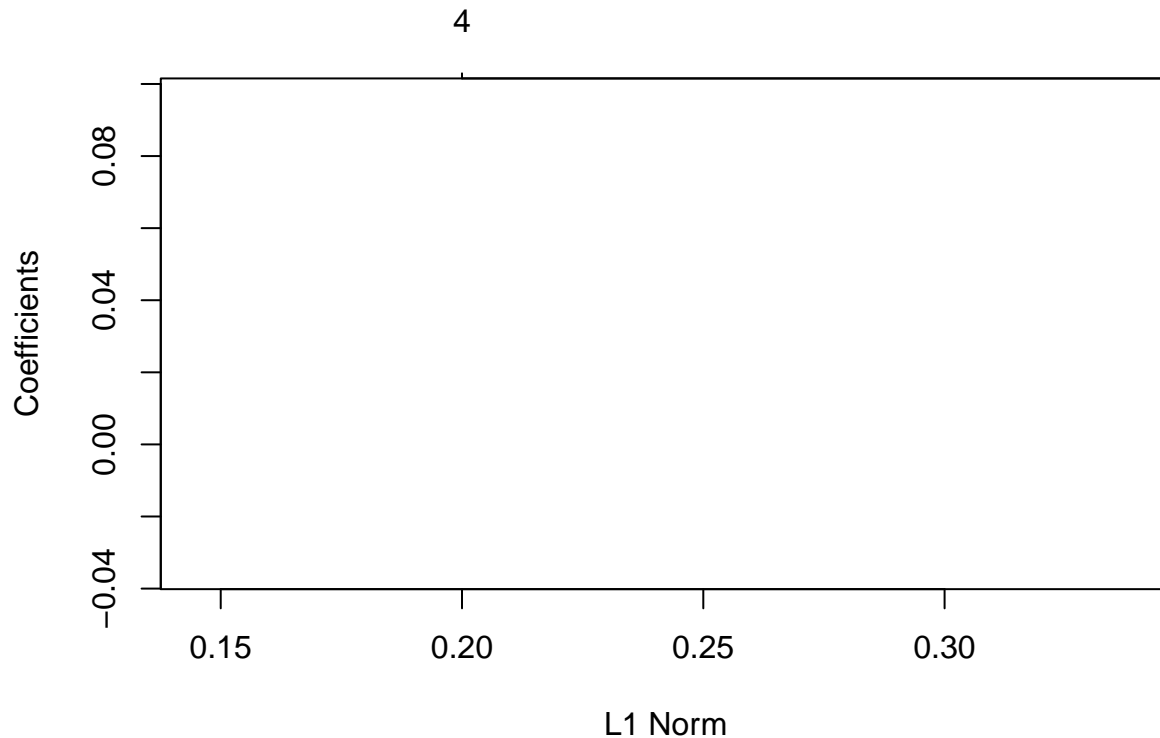
```
# Model fit with alpha=0.5.
set.seed(05)
modelfit3 <- glmnet(as.matrix(epdata[,1:4]),
                    epdata$class,
                    family = "binomial",
                    alpha = 0.5,
                    lambda = min(cv.modelfit3$lambda))

# Listing non-zero coefficients

print(modelfit3$beta[,1])

## reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3
## -0.03490984  0.09628728  0.06828960  0.04272984

plot(modelfit3)
```



Now, let's move to `stats::glm`.

```
#stats::glm()
```

The *stats* package is preloaded in R. We are particularly interested in the generalised linear models, `glm()` function. To begin, we shall customarily bifurcate our dataset into training data and testing data. The training data shall be used to build our linear model, while the testing data shall be used for its validation. Arbitrary proportions can be considered for splitting the data, however, usually 80-20 partition is reasonable.

```
set.seed(7) # for results reproducibility.
```

```
epdata$class <- as.numeric(as.factor(epdata$class))-1
part <- sample(2, nrow(epdata),
              replace = TRUE,
              prob = c(0.7,0.3))
train <- epdata[part==1,]
test <- epdata[part==2,]
cat("So, now we have",
    nrow(train),
    "training rows and",
    nrow(test),
    "testing rows")
```

```
## So, now we have 794602 training rows and 340624 testing rows
```

```
epmodel <- glm(formula = class ~ reads_h3k27ac + reads_h3k4me3 + reads_h3k4me2 + reads_h3k4me1,
               data = train,
               family = "binomial")
summary(epmodel)
```

```
##
```

```
## Call:
```

```
## glm(formula = class ~ reads_h3k27ac + reads_h3k4me3 + reads_h3k4me2 +
```

```
##      reads_h3k4me1, family = "binomial", data = train)
##
## Deviance Residuals:
##      Min        1Q      Median        3Q        Max
## -3.6695  -1.4192   0.9200   0.9502   1.1792
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.534487   0.002978  179.49  <2e-16 ***
## reads_h3k27ac -0.041838   0.003594  -11.64  <2e-16 ***
## reads_h3k4me3  0.049430   0.002312   21.38  <2e-16 ***
## reads_h3k4me2  0.065430   0.004221   15.50  <2e-16 ***
## reads_h3k4me1  0.100927   0.003602   28.02  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1026853  on 794601  degrees of freedom
## Residual deviance: 1018483  on 794597  degrees of freedom
## AIC: 1018493
##
## Number of Fisher Scoring iterations: 5
```

Here, we are taking into account all the variables as responses to the predictor variable - *Class*. Although, it can be interpreted straightforwardly, that none of the estimated coefficients of the model are statistically significant (See  $\text{Pr}(>|z|)$ ); but that's just the nature of this data, and in general terms it's better to reject all variables that have insignificant coefficients. Had we chosen to do that here, we would've left with nothing. Never mind. This demonstration is to highlight the protocol of logistic regression. Let's continue with whatever we have here, taking all.

Nonetheless, we mustn't ignore an important aspect of *multicollinearity*. Out of many ways to access that, `rms::vif()` provides an effective way to seek multicollinearity problem. `vif` stands for Variance Inflation Factor, and by norm if `vif() > 10`, we must omit the corresponding column (variable) as it does not add much to the model due to redundancy.

```
install.packages("rms",
                 repos = "https://cran.us.r-project.org")
```

```
## Installing package into '/Users/mei/Library/R/3.6/library'
## (as 'lib' is unspecified)
```

```
##
## The downloaded binary packages are in
## /var/folders/hm/c3_fjypn62v5xh5b5ygv267m0000gn/T//RtmpBWxxif/downloaded_packages
```

```
library(rms)
```

```
## Loading required package: Hmisc
## Loading required package: lattice
## Loading required package: survival
## Loading required package: Formula
## Loading required package: ggplot2
##
```

```
## Attaching package: 'Hmisc'

## The following objects are masked from 'package:base':
##
##      format.pval, units
## Loading required package: SparseM
##
## Attaching package: 'SparseM'
## The following object is masked from 'package:base':
##
##      backsolve
vif(epmodel)

## reads_h3k27ac reads_h3k4me3 reads_h3k4me2 reads_h3k4me1
##      1.724342      2.240417      3.595139      1.758284

Neither of the variables have high vif score, so they all qualify for the model.
```

```
y_train <- predict(epmodel, train, type = "response")
head(y_train)

##      12459      23570      34681      10      9121      19121
## 0.6305290 0.6305290 0.6305290 0.6331502 0.6305290 0.6305290
```

```
head(train)

##      reads_h3k27ac reads_h3k4me1 reads_h3k4me2 reads_h3k4me3 class
## 12459              0      0.00000              0              0      1
## 23570              0      0.00000              0              0      1
## 34681              0      0.00000              0              0      1
## 10                0      0.11165              0              0      1
## 9121               0      0.00000              0              0      1
## 19121              0      0.00000              0              0      1
```

These are the estimates of the class variable. To calculate the accuracy of the model we need to compare these to the original values of the response variable, 0 for “Enhancer” and 1 for “Non-Enhancer”. If you see the first observation, 0.6305290 ~ 63 % chance of being a non-enhancer, and in actuality if you look at the original data frame it is a non-enhancer. The probability (63 %) can be calculated by fitting values of coefficients in the model. Try doing that.

$$y = 0.534487 + (-0.041838 \times \text{reads\_h3k27ac}) + (0.049430 \times \text{reads\_h3k4me3}) + (0.065430 \times \text{reads\_h3k4me2}) + (0.100927 \times \text{reads\_h3k4me1})$$

```
prediction_probabilities_train <- ifelse(y_train > 0.5, 1, 0) # Probabilities to Labels conversion
confusion_matrix_train <- table(Predicted = prediction_probabilities_train, Actual = train$class)
print(confusion_matrix_train)
```

```
##      Actual
## Predicted    0      1
##           0      3      1
##           1 276447 518151
```

```
misclassification_error_train <- 1 - sum(diag(confusion_matrix_train))/sum(confusion_matrix_train)
cat("The misclassification error in train data is",
    (round(misclassification_error_train*100)), "percent")
```

```
## The misclassification error in train data is 35 percent
```



Now, we can repeat the same procedure for the test data.

```
y_test <- predict(epmodel, test, type = "response")
prediction_probabilities_test <- ifelse(y_test > 0.5, 1, 0)
confusion_matrix_test <- table(Predicted = prediction_probabilities_test, Actual = test$class)
print(confusion_matrix_test)
```

```
##           Actual
## Predicted    0     1
##           0     1     0
##           1 118502 222121
```

```
misclassification_error_test <- 1 - sum(diag(confusion_matrix_test)) / sum(confusion_matrix_test)
cat("The misclassification error in test data is",
    (round(misclassification_error_test*100)), "percent")
```

```
## The misclassification error in test data is 35 percent
```

Finally, there is also a way to ascertain if our model on the whole is statistically significant. We refer this as the Goodness-Of-Fit test.

```
overall_p <- with(epmodel,
                  pchisq(null.deviance-deviance,
                        df.null-df.residual,
                        lower.tail = FALSE))
cat("The statistical significance for the model is", overall_p, "\n")
```

```
## The statistical significance for the model is 0
```

```
cat("The confidence level for this model is",
    ((1-overall_p)*100), "percent")
```

```
## The confidence level for this model is 100 percent
```

Is this an ideal model ???