

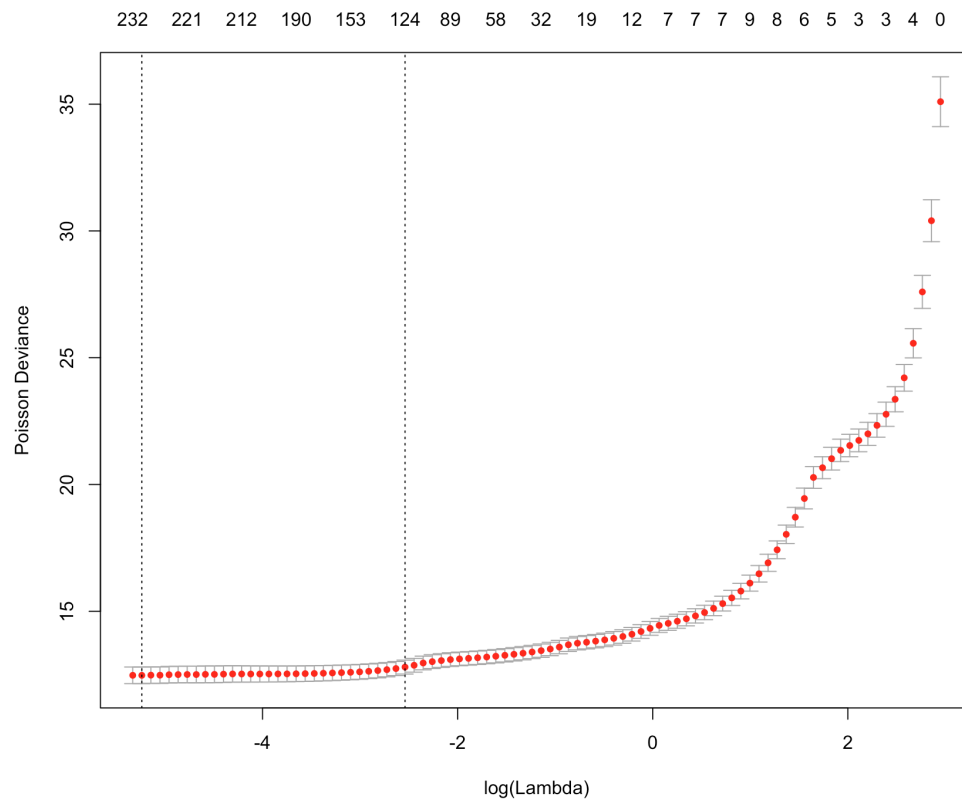
CS498 AML HW7

Huiyun Wu (hwu63)

Yidi Yang (yyang160)

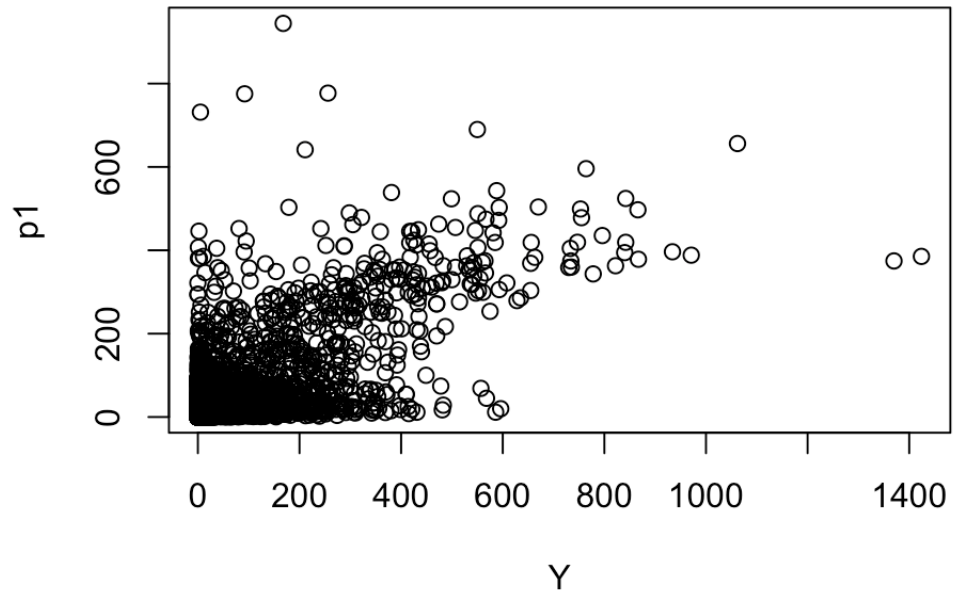
12.3

1. Plot of the cross-validated deviance of the model.



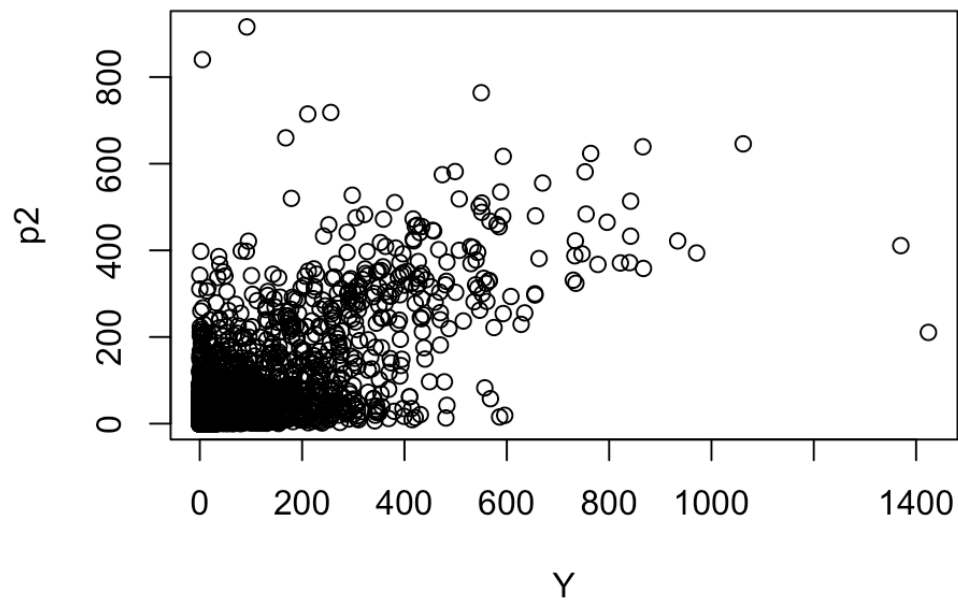
2. Scatter plots of true values vs predicted values **on training data**.

a. Lambda = "**lambda.1se**" = 0.07892024.



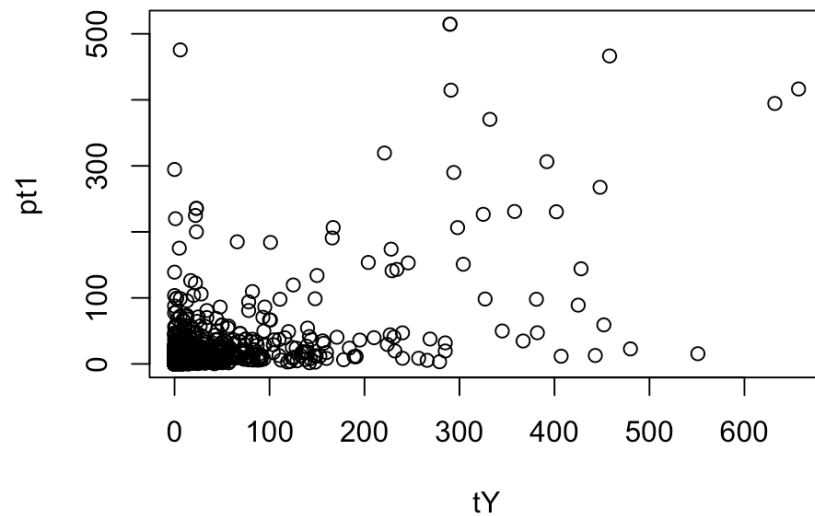
b.

c. Lambda = "**lambda.min**" = 0.005314608.

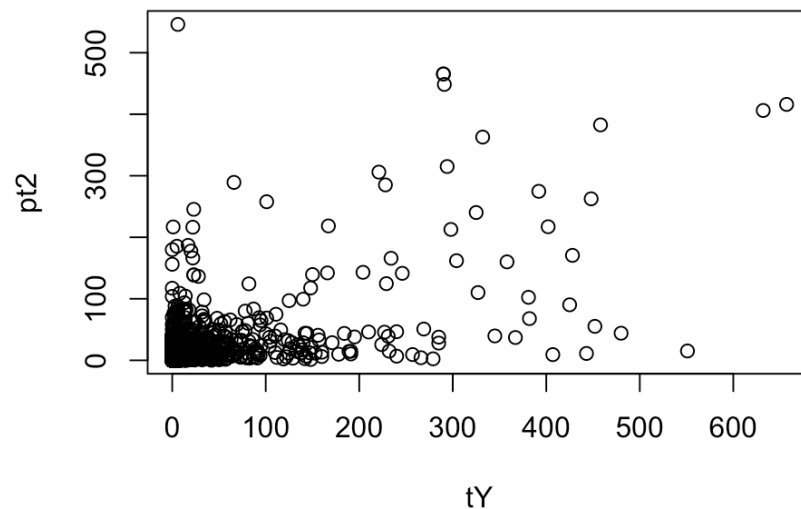


3. Scatter plots of true values vs predicted values **on testing data**.

a. Lambda = "**lambda.1se**" = 0.07892024.



b. Lambda = "**lambda.min**" = 0.005314608.



Performance of the model:

The model performs not ideally. According to the plots above, the scales of the true values and the predicted ones differ and there are many obvious points that are predicted wrong.

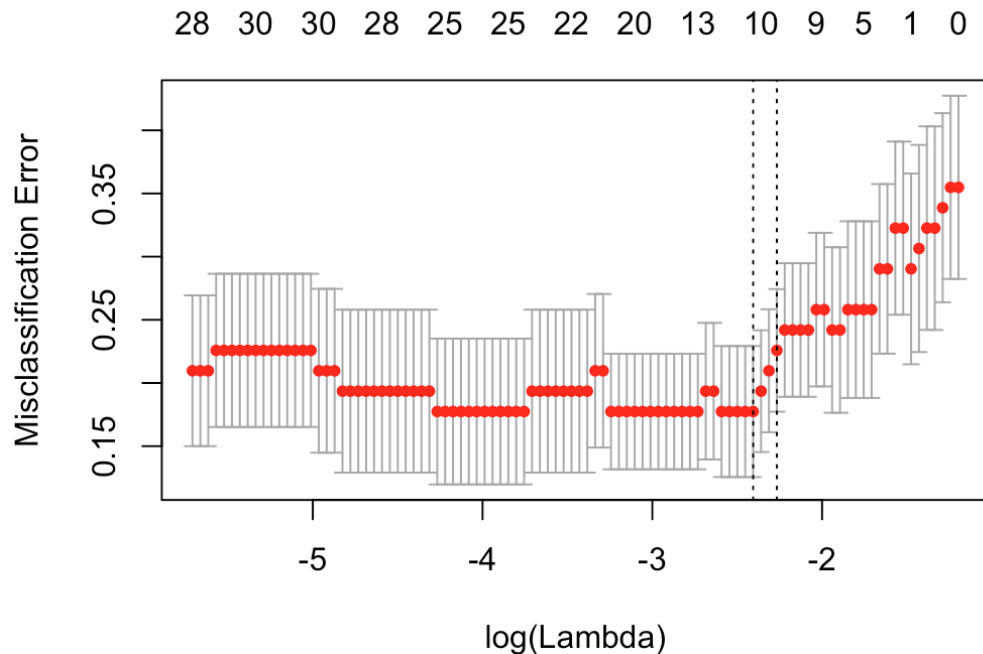
Why is the model difficult?

The true values contain too many zeros so that the training data has bias and the dependent variable is less relevant to the features. Those that has non-zero values are sparse and look like outliers on the plots, which may also influence the reliability of the model.

12.4

4. Plot of the (mis)classification error of the model against the regularization variable.

(family = "binomial", type.measure = "class")

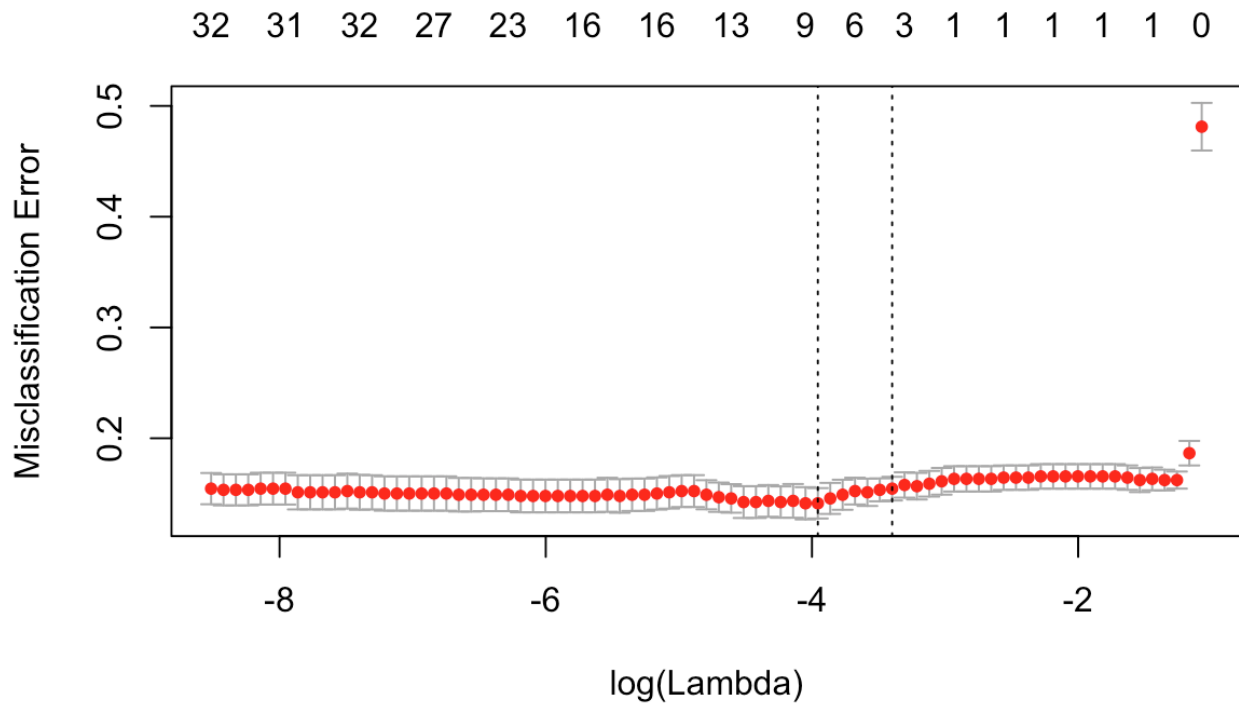


Compare with the baseline prediction:

The common type in this case is tumor. The accuracy of baseline prediction is 0.6451613, with 40 tumor samples and 22 normal ones. The accuracy of the logistic model using $\lambda = 0.07485259$ (λ_{\min}) is 0.9354839. The logistic model has significantly higher accuracy than the baseline prediction model does.

5. Plot of the classification error of the model against the regularization variable.

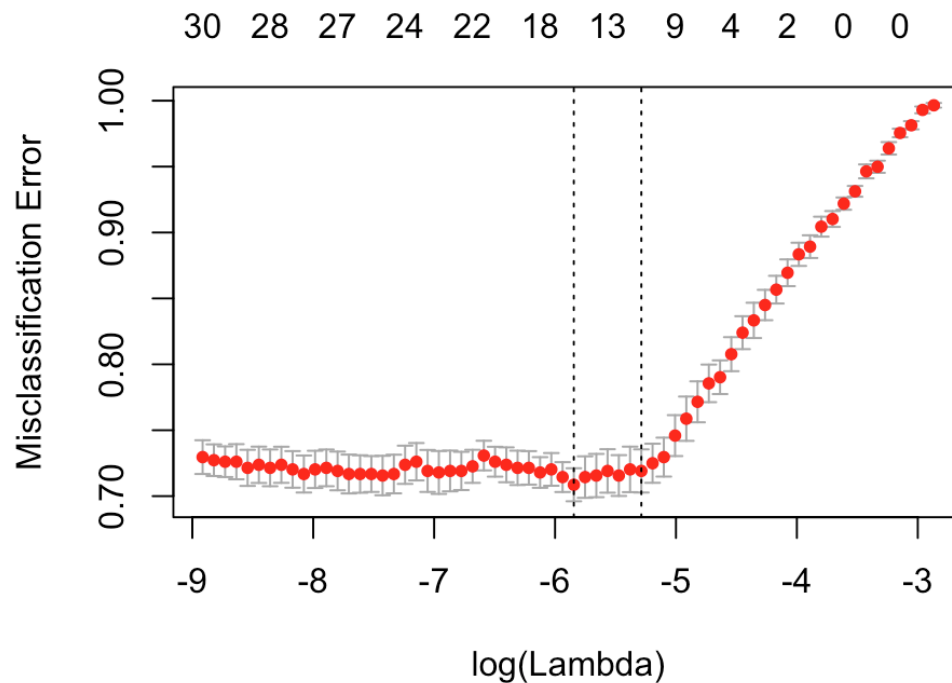
(family = "binomial", alpha=1, type.measure = "class")



Compare with the baseline prediction:

The common type in this case is male. The accuracy of baseline prediction is 0.50883. The accuracy of the logistic model with lasso using $\lambda = 0.01202526$ (lambda.min) is 0.8509934. The logistic model has significantly higher accuracy than the baseline prediction model does.

6. Plot of the classification error of the model against the regularization variable.



Compare with the baseline prediction:

Since there are multiple common types, we choose “BXD100” at random to do the comparison. The accuracy of baseline prediction using class “BXD100” (predicting all the instances as BXD100) is 0.02331002. The accuracy of the logistic model with lasso **using lambda = 0.01202526 (lambda.min)** is 0.6130536. The logistic model has significantly higher accuracy than the baseline.

7. Code Screenshots.

```
1 library(glmnet)
2 library(readr)
3 # -----12.3-----
4 blogData <- read_csv("~/Downloads/BlogFeedback/blogData_train.csv", col_names = FALSE)
5 X <- as.matrix(blogData[,c(1:280)])
6 Y <- as.matrix(blogData[,281])
7 cv <- cv.glmnet(X, Y, family="poisson")
8 plot(cv)
9
10 p1 <- predict(cv, X, s='lambda.1se', type='response')
11 p2 <- predict(cv, X, s='lambda.min', type = 'response')
12 plot(Y, p1)
13 plot(Y, p2)
14
15 setwd("~/Downloads/BlogFeedback/")
16 testfiles <- list.files(pattern = "*00.csv")
17 testdata <- lapply(testfiles, function(i){
18   read.csv(i, header=FALSE)
19 })
20 library(data.table)
21 testdata <- rbindlist(testdata)
22 tX <- as.matrix(testdata[,c(1:280)])
23 tY <- as.matrix(testdata[,281])
24 pt1 <- predict(cv, tX, s='lambda.1se', type='response')
25 pt2 <- predict(cv, tX, s='lambda.min', type='response')
26 plot(tY, pt1)
27 plot(tY, pt2)
28 # -----12.4-----
29 genedata <- read.table("~/Desktop/CS498AML/HW7/gene.txt", quote = "\"")
30 tissues <- read.table("~/Desktop/CS498AML/HW7/tissues.txt", quote="\"")
31 genedata <- as.matrix(t(genedata))
32 tissues[tissues > 0] = 1
33 tissues[tissues < 0] = 0
34 table(tissues)
35 tissues <- as.matrix(tissues)
36 logis <- cv.glmnet(genedata, tissues, type.measure = "class", family="binomial")
37 plot(logis)
38 blacc <- 40/62
39 print(blacc)
40 results <- predict(logis, genedata, s="lambda.min")
41 results[results > mean(results)] = 1
42 results[results < mean(results)] = 0
43 mean(results==tissues)
```

(Continued)

```

44 #-----12.5-----
45 micedata <- read_csv("~/Downloads/Crusio1.csv", col_names = TRUE)
46 micedata <- micedata[,c(2,4:41)]
47 micedata <- micedata[complete.cases(micedata),]
48 micedataX <- as.matrix(micedata[,c(2:39)])
49 micedataY <- as.matrix(micedata[,1])
50
51 micedataY[micedataY == 'f'] = 0
52 micedataY[micedataY == 'm'] = 1
53 table(micedataY)
54 logislasso <- cv.glmnet(micedataX, micedataY, type.measure = "class", alpha = 1, family="binomial")
55 plot(logislasso)
56 gender <- predict(logislasso, micedataX)
57 gender <- ifelse(gender > mean(gender), 1, 0)
58 mean(gender==micedataY,alpha=1)
59
60 print(logislasso$lambda.min)
61
62 micedata <- read_csv("~/Downloads/Crusio1.csv", col_names = TRUE)
63 micedata <- micedata[,c(1,4:41)]
64 micedata <- micedata[complete.cases(micedata),]
65 micedata <- micedata[micedata$strain %in% names(table(micedata$strain))[table(micedata$strain) > 9], ]
66
67 micedataY <- as.matrix(micedata[,1])
68 micedataX <- as.matrix(micedata[,c(2:39)])
69 micedataY <- as.matrix(micedataY[micedataY %in% names(table(micedataY))[table(micedataY) > 9], ])
70
71 logislasso <- cv.glmnet(micedataX, micedataY, type.measure = "class", alpha= 1, family="multinomial")
72 plot(logislasso)
73 mean(predict(logislasso, micedataX, s="lambda.min", type="class") == micedataY, alpha=1)
74 table(micedataY)
75 print(logislasso$lambda.min)

```

(End)