

## STAT 318 Assignment 2

---

Student Name : Xiao Meng

Student ID : 88211337

---

1.

With the fitting logistic regression model, when a student get a GPA value  $\geq 7$ ,

$$\ln\left(\frac{p(Y=1)}{1-p(Y=1)}\right) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$$

(a)

If they study for 5 hours and attend 36 classes, which means  $X_1 = 5$ ,  $X_2 = 36$ . With the given estimated coefficients  $\hat{\beta}_0 = -16$ ,  $\hat{\beta}_1 = 1.4$ ,  $\hat{\beta}_3 = 0.3$ , we can get

$$\begin{aligned} p(Y=1) &= \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2}} \\ &= \frac{e^{-16 + 1.4 \times 5 + 0.3 \times 36}}{1 + e^{-16 + 1.4 \times 5 + 0.3 \times 36}} \\ &= 0.8518(4dp) \end{aligned}$$

(b)

As above, when  $p(Y=1) = 0.5$ ,  $X_2 = 18$ ,

$$\begin{aligned} X_1 &= \frac{\ln\left(\frac{p(Y=1)}{1-p(Y=1)}\right) - \hat{\beta}_0 - \hat{\beta}_2 X_2}{\hat{\beta}_1} \\ &= \frac{0 - (-16) - 0.3 \times 18}{1.4} \\ &= 7.57(4dp) \end{aligned}$$

---

2.

(a)

With the training data, the multiple logistic regression based on predictors  $X_1$ ,  $X_3$  is

$$\ln\left(\frac{p(X)}{1-p(X)}\right) = 0.2204 - 1.3149X_1 - 0.2174X_3$$

As the p-value of variables are both much less than 0.001, which represents that both variables are statistically significant. Meanwhile, the coefficients of variables are negative shows that when value of predictors increases, logit value of this logistic regression model will decrease. It will more likely to be a genuine banknote.

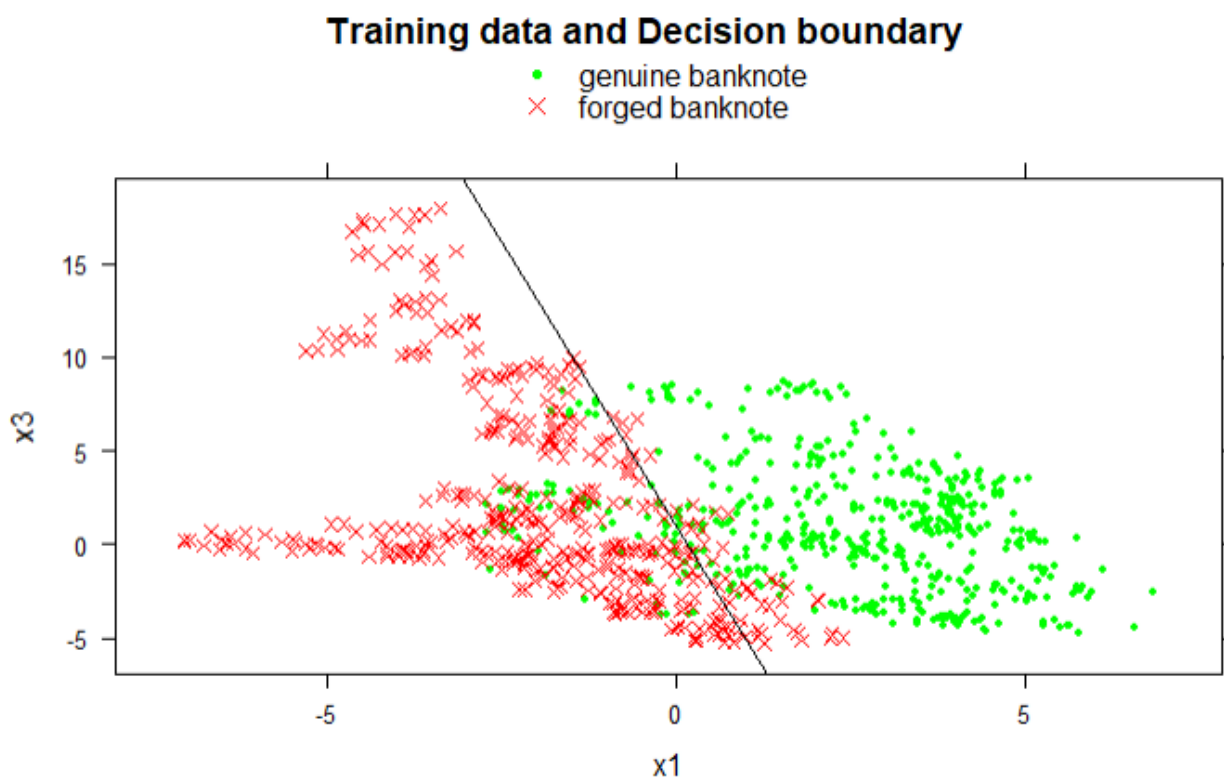
(b)

Using the Bootstrap and 1000 replicates, the estimated standard errors for  $\hat{\beta}_1$  is 0.0894(4dp),  $\hat{\beta}_2$  is 0.0267(4dp).

	original	bias	std.error
t1*	0.2204101	0.002915229	0.12143542
t2*	-1.3148902	-0.010448191	0.08944161
t3*	-0.2173841	-0.001463668	0.02672804

(c) i

Plot the training data and the decision boundary for  $\theta = 0.5$ .

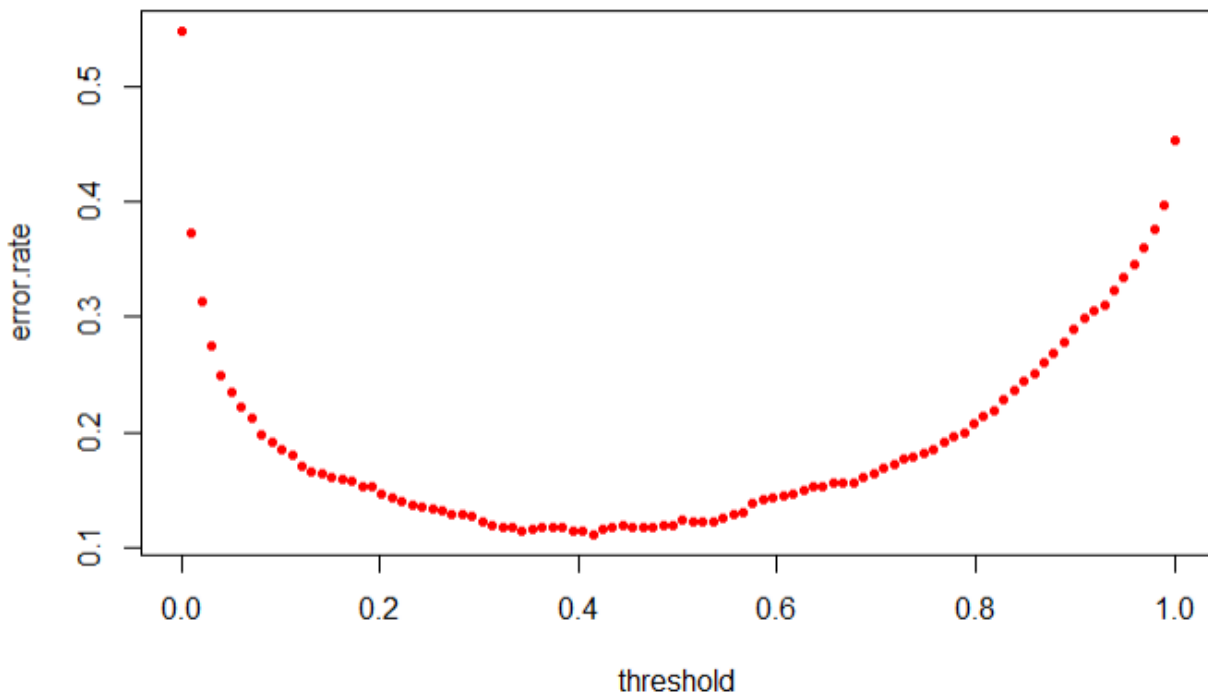


(c) ii

glm.pred	genuine banknote	forged banknote
genuine banknote	204	24
forged banknote	32	152

The accuracy rate is  $86.41\% \left( \frac{204+152}{412} \right)$ . From the confusion matrix we can see, there are 24 forged banknotes incorrectly assigned to genuine banknotes, while this error rate is  $13.64\% \left( \frac{24}{24+152} \right)$ . For the actual situation, this error rate is very important, which might lead some serious consequences.

(c) iii



From the figure above we can see when the threshold is 0.42, there will be the minimal training error rate 11.15%.

glm.pred	genuine banknote	forged banknote
genuine banknote	200	13
forged banknote	36	163

Applying this best threshold to testing data, the accuracy rate increases to 88.11%. Meanwhile the incorrectly predicted forged banknotes decreases to 13 and this error rate falls to 7.39%. The new threshold gives a good results.

3.

Model	Training error	Testing error
Logistic Regresion	12.08%	13.59%
LDA	12.08%	13.35%
QDA	11.46%	11.17%

(a)

Fitting an LDA model with training data, the training error is 12.08% and testing error is 13.35%.

(b)

Fitting an QDA model with training data, the training error is 11.46% and testing error is 11.17%.

(c)

From the table above we can see, among these three models, QDA gives the best performance which has both the lowest training error and testing error. While the errors of Logistic Regression and LDA are quite similar.

As LDA regression is based on the assumption that each classification shares the same variance. However the variance of these two classification is 8.17 and 18.92, which is quite different. This reason might make LDA perform worse in this data set.

From figure of training data and the decision boundary, we know these two classifications have some overlapping points, which might make the boundary not to be linear. Therefore, the Logistic regression model has a higher error rate than QDA. Besides, the QDA does not make assumptions about the same variance of each classification. So in this data set, I recommend the QDA which performs the best.

---

4.

As  $\pi_0 = 0.4$ ,  $\pi_1 = 1 - \pi_0 = 0.6$ , the decision boundary is

$$\begin{aligned}\pi_0 f_0(x) &= \pi_1 f_1(x) \\ 0.4 \times \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{8}x^2} &= 0.6 \times \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{8}(x-2)^2} \\ e^{-\frac{1}{8}x^2 + \frac{1}{8}(x-2)^2} &= \frac{3}{2} \\ e^{\frac{1-x}{2}} &= \frac{3}{2} \\ x &= 1 - 2\ln\left(\frac{3}{2}\right) \\ x &= 0.1891(4dp)\end{aligned}$$

The Bayes error rate is

$$\begin{aligned}&1 - \left( \int_{-\infty}^{0.1891} \pi_0 f_0(x) + \int_{0.1891}^{\infty} \pi_1 f_1(x) \right) \\ &= 1 - \left( \int_{-\infty}^{0.1891} 0.4 \times \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{8}x^2} + \int_{0.1891}^{\infty} 0.6 \times \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{8}(x-2)^2} \right) \\ &= 1 - (0.2151 + 0.4904) \\ &= 0.2945(4dp)\end{aligned}$$

---

## Appendix R Code

```
1  #=====
2  #----- 1(a) -----
3  #=====
4  beta0 = -16
5  beta1 = 1.4
6  beta2 = 0.3
7  x1 = 5
8  x2 = 36
9  y = (exp(beta0 + beta1 * x1 + beta2 * x2))/(1 + exp(beta0 + beta1 * x1 + beta2 *
10 x2))
11 y
12
13 #=====
14 #----- 1(b) -----
15 #=====
16 y = 0.5
17 x2 = 18
18 (log(y / (1 - y)) - beta0 - beta2 * x2)/beta1
19
20
21 #=====
22 #----- 2(a) -----
23 #=====
24 #fit the training datas from BankTrain.csv with logistic model
25 #Banktrain = read.csv("/Users/mac/Desktop/Computer Science/Data
Mining/Assignment/A2/BankTrain.csv")
26
27 Banktrain = read.csv("E:/文档/UC/318/Assignment/A2/BankTrain.csv")
28 glm.fit = glm(y ~ x1 + x3 , data = Banktrain, family = binomial)
29 summary(glm.fit)
30
31
32 #=====
33 #----- 2(b) -----
34 #=====
35 library(boot)
36 set.seed(2)
37 #create a new function to fit the model with a single row
38 boot.fn = function(data, index){
39   return (coef(glm(y ~ x1 + x3, data = Banktrain, family = binomial, subset =
index)))
40 }
41
42 #estimate the standard errors for coefficients using the bootstrap
43 boot(data = Banktrain, statistic = boot.fn, R = 1000)
```

```

44
45 #compare to the standard errors calculating by the model
46 summary(glm(y ~ x1 + x3, data = Banktrain, family = binomial))
47
48
49 #=====
50 #----- 2(c) -----
51 #=====
52 library(lattice)
53 # parameters of boundary when threshold = 0.5
54 slope = -coef(glm.fit)[2] / coef(glm.fit)[3]
55 intercept = -coef(glm.fit)[1] / coef(glm.fit)[3]
56
57 xyplot(x3 ~ x1, data = Banktrain, groups = y, pch = c(20, 4), col = c("green",
"red"), main = 'Training data and Decision boundary',
58         key=list(space='top',
59                 points = list(pch = c(20, 4), col = c('green', 'red')),
60                 text = list(lab = c('genuine banknote', 'forged banknote'))),
61         panel = function(...){
62             panel.xyplot(...)
63             panel.abline(intercept, slope)}}
64
65
66 #=====
67 #----- 2(d) -----
68 #=====
69 #predict probabilities with testing datas from BankTest.csv
70 Banktest = read.csv("E:/文档/UC/318/Assignment/A2/BankTest.csv")
71
72 glm.probs = predict(glm.fit, Banktest, type = "response")
73
74 #create a vector with the results of predictions from logistic model
75 glm.pred = rep("0", 412)
76 glm.pred[glm.probs > .5] = "1"
77
78 #create a matrix to classify how the predictions perform
79 table(glm.pred, Banktest$y)
80 mean(glm.pred == Banktest$y)
81
82
83 #=====
84 #----- 2(e) -----
85 #=====
86 #predict probabilities with testing datas from BankTest.csv
87
88 glm.probs.train = predict(glm.fit, Banktrain, type = "response")
89
90 #create a vector with the results of predictions from logistic model
91 threshold = seq(0, 1, length = 100)
92 error.rate = list()
93
94 for(x in threshold){
95     glm.pred.train = rep("0", 960)

```

```

96   glm.pred.train[glm.probs.train > x] = "1"
97   error.training = mean(glm.pred.train != Banktrain$y)
98   error.rate = c(error.rate, list(error.training))
99 }
100 min(unlist(error.rate))
101 match(c(min(unlist(error.rate))), error.rate)
102
103 plot(threshold, error.rate, pch = 20, col = 'red')
104
105 glm.probs.test = predict(glm.fit, Banktest, type = "response")
106 glm.pred.test = rep("0", 412)
107 glm.pred.test[glm.probs.test > .42] = "1"
108
109 #create a matrix to classify how the predictions perform
110 table(glm.pred.test, Banktest$y)
111 mean(glm.pred.test == Banktest$y)
112
113
114
115 #=====
116 #----- 3(a) -----
117 #=====
118 #LDA model
119 library(MASS)
120 lda.fit = lda(y ~ x1 + x3, data = Banktrain)
121
122 #predict classification of training data and training error
123 lda.class.train = predict(lda.fit, Banktrain)$class
124 mean(lda.class.train != Banktrain$y)
125
126
127 #predict classification of test data and test error
128 lda.class.test = predict(lda.fit, Banktest)$class
129 mean(lda.class.test != Banktest$y)
130
131
132 #=====
133 #----- 3(b) -----
134 #=====
135 #QDA model
136 qda.fit = qda(y ~ x1 + x3, data = Banktrain)
137
138 #predict classification of training data and training error
139 qda.class.train = predict(qda.fit, Banktrain)$class
140 mean(qda.class.train != Banktrain$y)
141
142 #predict classification of test data and test error
143 qda.class.test = predict(qda.fit, Banktest)$class
144 mean(qda.class.test != Banktest$y)
145
146 #=====
147 #----- 3(c) -----

```

```

148 #=====
149 #Compute training error and test error of logistic model
150
151 #train error
152 glm.probs.train1 = predict(glm.fit, Banktrain, type = "response")
153 glm.pred.train1 = rep("0", 960)
154 glm.pred.train1[glm.probs.train1 > .5] = "1"
155 mean(glm.pred.train1 != Banktrain$y)
156
157 #test error
158 glm.probs.test1 = predict(glm.fit, Banktest, type = "response")
159 glm.pred.test1 = rep("0", 412)
160 glm.pred.test1[glm.probs.test1 > .5] = "1"
161 mean(glm.pred.test1 != Banktest$y)
162
163 sd(Banktrain$x1)^2
164 sd(Banktrain$x3)^2
165
166
167 #=====
168 #----- 4 -----
169 #=====
170 fx = function(x, pi0, pi1){return (pi0 * dnorm(x, 0, 2) - pi1 * dnorm(x, 2, 2))}
171 root = uniroot(fx, c(-5, 5), pi0 = 0.4, pi1 = 0.6, tol = 0.0001)
172 boundary = root$root
173 boundary
174 pi0 = 0.4
175 pi1 = 1 - pi0
176 fx0 = function(x) pi0 * dnorm(x, 0, 2)
177 fx1 = function(x) pi1 * dnorm(x, 2, 2)
178 1 - (integrate(fx0, -Inf, boundary)$value + integrate(fx1, boundary, Inf)$value)
179

```