

### Assignment 3 Chunlei Zhou

**Q1**

**(a)**

$$\begin{aligned} \text{OR}(\text{abuse}==1|\text{mothalc}==1) &= \frac{\text{odds}(\text{abuse}==1|\text{mothalc}==1)}{\text{odds}(\text{abuse}==1|\text{mothalc}==0)} \\ &= \frac{p(\text{abuse}==1|\text{mothalc}==1)/(1-p(\text{abuse}==1|\text{mothalc}==1))}{p(\text{abuse}==1|\text{mothalc}==0)/(1-p(\text{abuse}==1|\text{mothalc}==0))} \end{aligned}$$

It is obvious that the odds ratio  $\text{OR}(\text{abuse}==1|\text{mothalc}==1)$  does not depend on the value of  $\text{fathalc}$ .

$$\begin{aligned} \text{OR}(\text{abuse}==1|\{\text{mothalc}==1 \text{ and } \text{fathalc}==1\}, \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\}) \\ &= \frac{\text{odds}(\text{abuse}==1|\text{mothalc}==1 \text{ and } \text{fathalc}==1)}{\text{odds}(\text{abuse}==1|\text{mothalc}==0 \text{ and } \text{fathalc}==0)} \\ &= \frac{p(\text{mothalc}==1 \text{ and } \text{fathalc}==1|\text{abuse}==1)/p(\text{mothalc}==1 \text{ and } \text{fathalc}==1|\text{abuse}==0)}{p(\text{mothalc}==0 \text{ and } \text{fathalc}==0|\text{abuse}==1)/p(\text{mothalc}==0 \text{ and } \text{fathalc}==0|\text{abuse}==0)} \\ &= \frac{p(\text{mothalc}==1|\text{abuse}==1)p(\text{mothalc}==0|\text{abuse}==0)}{p(\text{mothalc}==1|\text{abuse}==0)p(\text{mothalc}==0|\text{abuse}==1)} \\ &\quad * \frac{p(\text{fathalc}==0|\text{abuse}==0)p(\text{fathalc}==0|\text{abuse}==1)}{p(\text{fathalc}==1|\text{abuse}==0)p(\text{fathalc}==1|\text{abuse}==1)} \\ &= \text{OR}(\text{abuse}==1|\text{mothalc}==1) * \text{OR}(\text{abuse}==1|\text{fathalc}==1) \end{aligned}$$

**(b)**

There are three cases for

$$\text{OR}(\text{abuse}==1|\{\text{mothalc}==1 \text{ or } \text{fathalc}==1\}, \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\})$$

1<sup>st</sup>:

$$\text{OR}_1 = \text{OR}(\text{abuse}==1|\{\text{mothalc}==1 \text{ and } \text{fathalc}==1\}, \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\})$$

Denote

$$\{\text{mothalc}==1 \text{ and } \text{fathalc}==1\} \text{ as } A \text{ and } \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\} \text{ as } B$$

Then we get

$$\eta_A = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2$$

And

$$\eta_B = \beta_0$$

So we get

$$OR_1 = e^{\eta_A - \eta_B} = e^{\beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2}$$

Since  $X_1$  and  $X_2$  are identity metrics so we can have

$$OR_1 = e^{\beta_1 + \beta_2 + \beta_3} = e^{0.63}$$

So we have

$$\beta_1 + \beta_2 + \beta_3 = 0.63$$

2<sup>nd</sup>:

$$OR_2 = OR(\text{abuse}==1 | \{\text{mothalc}==1 \text{ and } \text{fathalc}==0\}, \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\})$$

Denote

$$\{\text{mothalc}==1 \text{ and } \text{fathalc}==0\} \text{ as A and } \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\} \text{ as B}$$

Then we get

$$\eta_A = \beta_0 + \beta_1 X_1$$

And

$$\eta_B = \beta_0$$

So we get

$$OR_2 = e^{\eta_A - \eta_B} = e^{\beta_1 X_1}$$

Since  $X_1$  is an identity metrics so we can have

$$OR_2 = e^{\beta_1} = e^{0.63}$$

So we have

$$\beta_1 = 0.63$$

3<sup>rd</sup>:

$$OR_3 = OR(\text{abuse}==1 | \{\text{mothalc}==0 \text{ and } \text{fathalc}==1\}, \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\})$$

Denote

$$\{\text{mothalc}==0 \text{ and } \text{fathalc}==1\} \text{ as A and } \{\text{mothalc}==0 \text{ and } \text{fathalc}==0\} \text{ as B}$$

Then we get

$$\eta_A = \beta_0 + \beta_2 X_2$$

And

$$\eta_B = \beta_0$$

So we get

$$OR_2 = e^{\eta_A - \eta_B} = e^{\beta_2 X_2}$$

Since  $X_2$  is an identity metrics so we can have

$$OR_2 = e^{\beta_2} = e^{0.63}$$

So we have

$$\beta_2 = 0.63$$

It is easy to get that

$$\beta_3 = -0.63$$

(c)

(i)

		2.5 %	97.5 %
(Intercept)	0.0984387	0.09126191	0.1061799
mothalc	1.4697085	1.10368787	1.9571141
fathalc	1.6551649	1.40537665	1.9493498

The odds ratios  $OR(\text{abuse} == 1 \text{ j } \text{mothalc} == 1)$  and  $OR(\text{abuse} == 1 \text{ j } \text{fathalc} == 1)$  from model (1), with approximate 95% confidence intervals are 1.469709 and 1.655165 respectively.

(ii)

Call:

```
glm(formula = yabuse ~ mothalc * fathalc, family = "binomial",  
     data = alcohol)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.6518	-0.4332	-0.4332	-0.4332	2.1966

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-2.31877	0.03898	-59.491	< 2e-16 ***
mothalc	0.39696	0.19933	1.992	0.0464 *
fathalc	0.50618	0.08742	5.791	7.02e-09 ***
mothalc:fathalc	-0.02558	0.29282	-0.087	0.9304

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 6349.8 on 9821 degrees of freedom  
Residual deviance: 6303.8 on 9818 degrees of freedom  
AIC: 6311.8

Number of Fisher Scoring iterations: 5

Based on the summary, there is no evidence that  $\beta_3 \neq 0$ .

**Q2**

**(a)**

$$\text{Odds}(\text{Document is Authentic} \mid X) = \frac{P(X \mid \text{Document is Authentic})}{p(X \mid \text{Document is Forged})} * \frac{P(\text{Document is Authentic})}{1 - p(\text{Document is Authentic})}$$

$$P(X \mid \text{Document is Authentic}) = \frac{n!}{n_1!n_5!n_9!} * p_{A1}^{n_1} p_{A5}^{n_5} p_{A9}^{n_9}$$

$$P(X \mid \text{Document is Forged}) = \frac{n!}{n_1!n_5!n_9!} * p_{F1}^{n_1} p_{F5}^{n_5} p_{F9}^{n_9}$$

So,

$$\text{Odds}(\text{Document is Authentic} \mid X) = \frac{\frac{n!}{n_1!n_5!n_9!} * p_{A1}^{n_1} p_{A5}^{n_5} p_{A9}^{n_9}}{\frac{n!}{n_1!n_5!n_9!} * p_{F1}^{n_1} p_{F5}^{n_5} p_{F9}^{n_9}} * \frac{\pi_A}{1 - \pi_A} = \frac{p_{A1}^{n_1} p_{A5}^{n_5} p_{A9}^{n_9}}{p_{F1}^{n_1} p_{F5}^{n_5} p_{F9}^{n_9}} * \frac{\pi_A}{1 - \pi_A}$$

**(b)**

$$\log[\text{Odds}(\text{Document is Authentic} \mid X)]$$

$$= n_1 \log p_{A1} + n_5 \log p_{A5} + n_9 \log p_{A9} - n_1 \log p_{F1} - n_5 \log p_{F5} - n_9 \log p_{F9} + \log \pi_A - \log(1 - \pi_A)$$

$$= n_1 \log(p_{A1}/p_{F1}) + n_5 \log(p_{A5}/p_{F5}) + n_9 \log(p_{A9}/p_{F9}) + \log\left(\frac{\pi_A}{1 - \pi_A}\right)$$

So,

$$a = \log(p_{A1}/p_{F1}); b = \log(p_{A5}/p_{F5}); c = \log(p_{A9}/p_{F9}); d = \log\left(\frac{\pi_A}{1 - \pi_A}\right).$$

**(c)**

We have

$$p_{A1} = 0.72, p_{A5} = 0.17, p_{A9} = 0.11, p_{F1} = p_{F5} = p_{F9} = 1/3.$$

Given

$$X = (n_1, n_5, n_9) = (7, 5, 8) \text{ and } \pi_A = \frac{1}{2}$$

We can get

$$\text{Odds}(\text{Document is Authentic} \mid X) = \frac{p_{A1}^{n_1} p_{A5}^{n_5} p_{A9}^{n_9}}{p_{F1}^{n_1} p_{F5}^{n_5} p_{F9}^{n_9}} * \frac{\pi_A}{1 - \pi_A} = 1.09 * 10^{-3}$$

Since

$$\text{Odds}(\text{Document is Authentic} \mid X) = \frac{P(\text{Document is Authentic} \mid X)}{1 - P(\text{Document is Authentic} \mid X)}$$

We can have

$$P(\text{Document is Authentic} \mid X) = 1.09 * 10^{-3}$$

$$P(\text{Document is Forged} \mid X) = 1 - P(\text{Document is Authentic} \mid X) = 0.9989$$

### Q3

(a)

Report the P-values:

```
[1] "npreg : P-value = 3.20684823540759e-06"  
[1] "glu : P-value = 2.01071777376076e-28"  
[1] "bp : P-value = 1.02775734163267e-05"  
[1] "skin : P-value = 2.26930291623554e-09"  
[1] "bmi : P-value = 1.02588422797911e-11"  
[1] "ped : P-value = 7.38261979346636e-08"  
[1] "age : P-value = 3.45737086135437e-17"
```

This suggests that the full model should be used to build an accurate classifier.

(b)

Please see code.

(c)

$$CE = (n_{21} + n_{12}) / (n_{11} + n_{12} + n_{21} + n_{22})$$

$$\text{sens} = n_{22} / (n_{22} + n_{12})$$

$$\text{spec} = n_{11} / (n_{11} + n_{21})$$

Please refer to the code for the function.

(d)

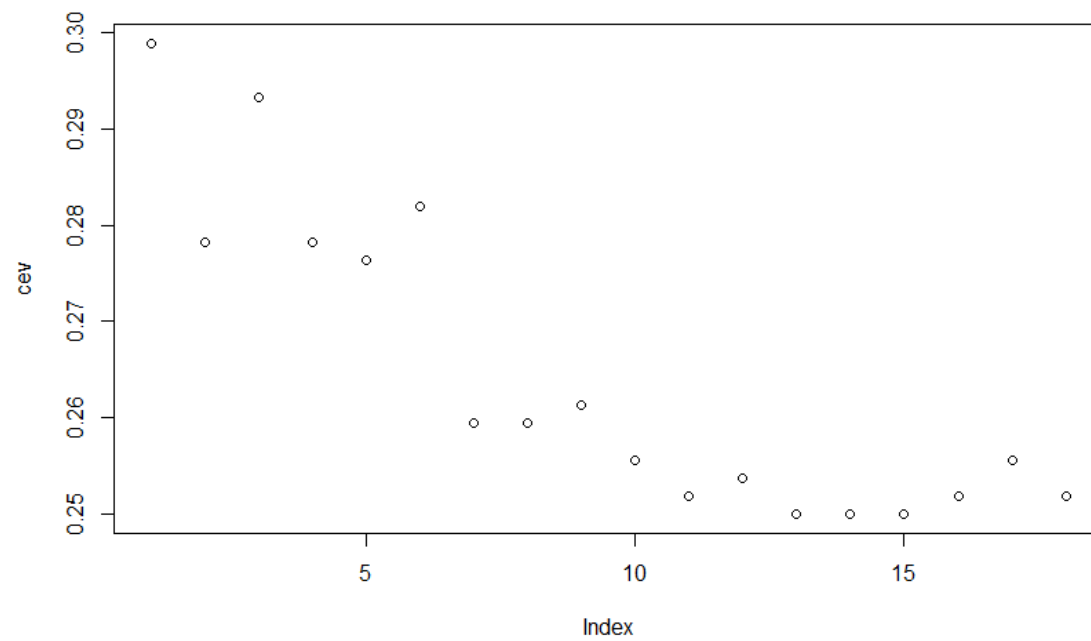
After building the KNN classifier accordingly, I found the K that minimized CE is 25.  
The minimum CE = 0.25.

(ii)

#### Summary Statistics

	CE	sens	spec
KNN	0.2500000	0.4576271	0.8957746
LDA	0.2105263	0.5988701	0.8845070
QDA	0.2124060	0.6384181	0.8270270

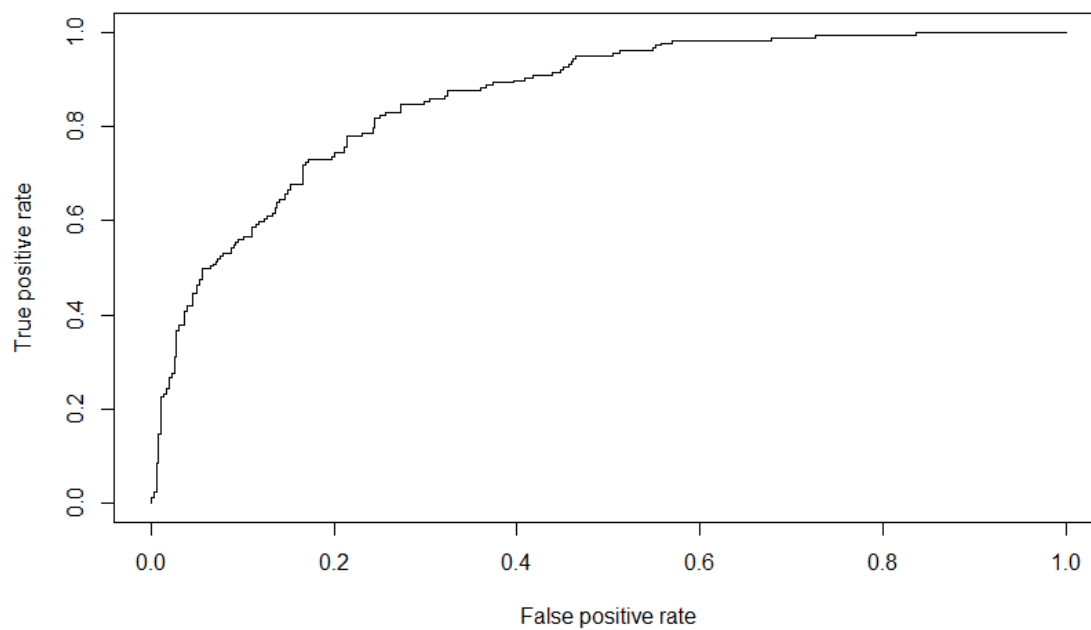
Based on the table, it can be claimed that there is no single model better than others.



(e)

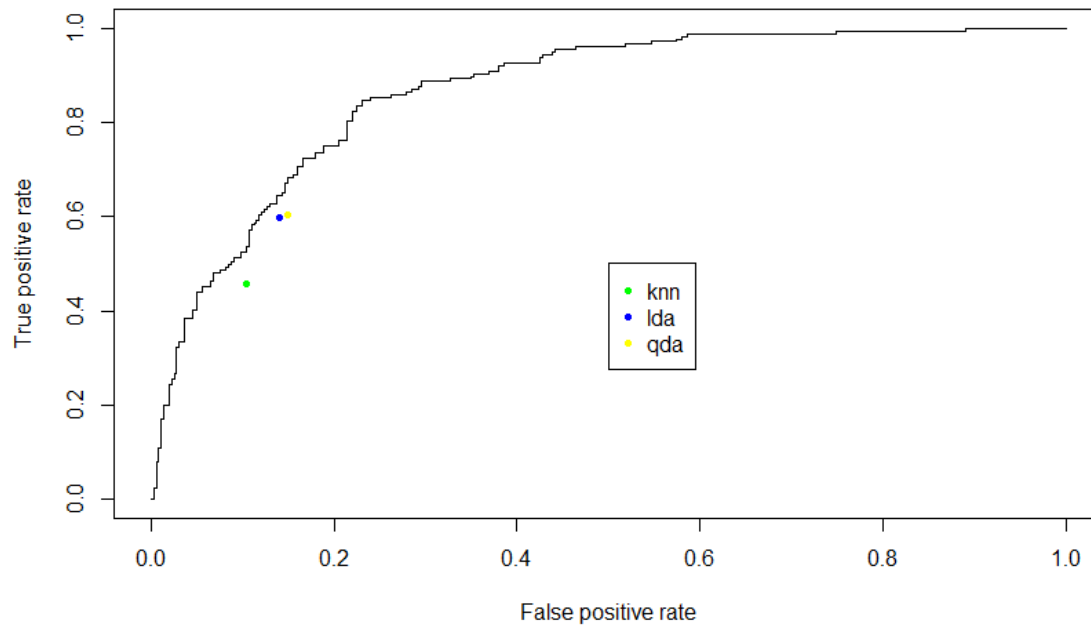
(i) Please refer to the code.

(ii) ROC curve



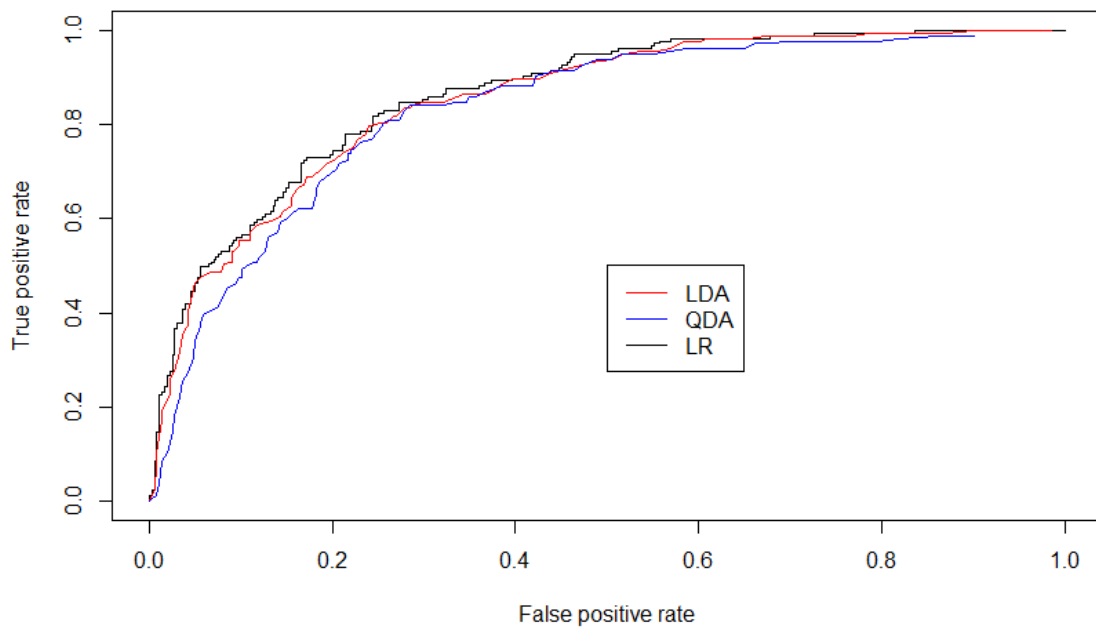


**(iii) ROC curve**



The plot shows a trade-off and based so that can help programmers to choose the model that make the most efficient use of the data.

**(iv) Neither LDA nor QDA seems preferable.**

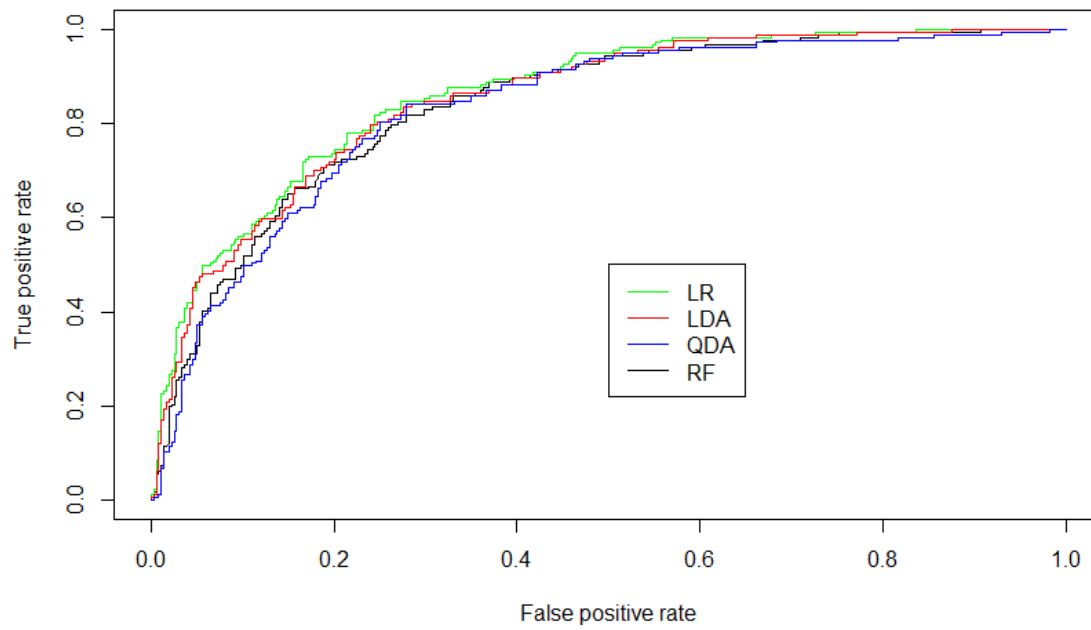


#### Q4

(a)

Please refer to the code.

(b)



This form of classifier does not offer any advantage over LDA or QDA for this application.