

EST ASSIGNMENT - 3 (WE)

Problem 1

1. When the response variable has finite number of outcomes (values). In other words, they are categorical in nature, logistic regression is used
eg: Red, Blue, Green, True / False.

Linear regression is not applicable since it can produce responses below 0 and beyond 1, i.e., it is not bounded between 0 & 1. ~~This is not~~
In other words it produces continuous response which is not suitable for handling categorical responses.

2. Logistic regression is used to model probability of an outcome based on the observations.

The relationship between independent variables & probability:

If $X = (X_1, X_2 \dots X_p)$, then

$$P(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

where $P(X)$ is the probability obtained for the given independent variables.

3. Odds refers to the ratio between the total number of favorable outcomes to the ~~not~~ total number of unfavorable outcomes.

If ' p ' is the probability of an event occurring, then $1-p$ is the probability of that event not occurring.

Then the Odds is given by
 $\text{Odds} = (\frac{P}{1-P})$

Consider, probability of rain to be 0.8, then the probability of not raining is 0.2. The odds of it rain on the particular day is $0.8/0.2 = 4$. In terms of ration it is 4:1, thus we can say the chances of raining is 4 times than not raining.

4. To Equation 4.2 :-

$$P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Eqn 4.4 :-

$$\log \frac{P(x)}{1 - P(x)} = \beta_0 + \beta_1 x$$

~~Note~~ In consistent with definitions of the text book,

Consider eqn 4.2,

$$P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

taking log,

~~$$\log P(x) = \beta_0 + \beta_1 x$$~~

~~$$\frac{P(x)}{1 - P(x)} = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$~~

~~$$\frac{P(x)}{1 - P(x)} = \frac{e^{\beta_0 + \beta_1 x}}{1}$$~~

taking log, we get

~~$$\log \left(\frac{P(x)}{1 - P(x)} \right) = \log(e^{\beta_0 + \beta_1 x}) = \beta_0 + \beta_1 x.$$~~

The above equation is 4.4. hence equivalent to 4.4

$$P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$1 - P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$1 - P(x) = \frac{1}{1 + e^{\beta_0 + \beta_1 x}}$$

assuming $P(x) \neq 0$, divide the above eqn by $P(x)$
we get,

$$\frac{1 - P(x)}{P(x)} = \frac{1}{1 + e^{\beta_0 + \beta_1 x}} \cdot \frac{1}{P(x)}$$

$$\frac{1 - P(x)}{P(x)} = \frac{1}{1 + e^{\beta_0 + \beta_1 x}} \cdot \frac{1 + e^{\beta_0 + \beta_1 x}}{e^{\beta_0 + \beta_1 x}}$$

$$\therefore \frac{P(x)}{1 - P(x)} = e^{\beta_0 + \beta_1 x}$$

taking log, we get

$$\log\left(\frac{P(x)}{1 - P(x)}\right) = \log(e^{\beta_0 + \beta_1 x})$$

$$\log\left(\frac{P(x)}{1 - P(x)}\right) = \beta_0 + \beta_1 x. \quad \text{--- (1)}$$

the above equation is equivalent to 4.2.

(b) logit is defined as the log of odds i.e. $\log\left(\frac{P(x)}{1 - P(x)}\right)$

from (1) which is equivalent to logistic

regression $P(x) = e^{\beta_0 + \beta_1 x} / (1 + e^{\beta_0 + \beta_1 x})$. We
can say that logistic regression has a
logit that is linearly in x .

$$5) \text{Odds}(x) = \frac{p(x)}{1-p(x)} = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 - \beta_i x_i)$$

$$\text{Odd}(x_i) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 - \dots + \beta_i x_i)$$

$$\text{Odd}(x_i + \Delta) = \exp(\beta_0 + \beta_1 (x_1) + \beta_2 x_2 + \dots - \beta_i (x_i + \Delta))$$

Hence

$$\frac{\text{Odd}(x_i + \Delta)}{\text{Odd}(x_i)} = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 - \dots - \beta_i x_i + \beta_i \Delta)}{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 - \dots - \beta_i x_i)}$$

$$= \frac{\exp(\beta_0 + \beta_1 x_1 - \dots - \beta_i x_i + \beta_i \Delta - \beta_i x_i - \beta_i x_i)}{1}$$

$$\boxed{\frac{\text{Odd}(x_i + \Delta)}{\text{Odd}(x_i)} = \exp(\beta_i \Delta)}$$

$$6) p(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)} \quad \text{If } p(x) = 0.5, \text{ then}$$

$$\frac{1}{2} = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

$$1 = \exp(\beta_0 + \beta_1 x)$$

$$0 = \beta_0 + \beta_1 x$$

$$\therefore x = -\frac{\beta_0}{\beta_1}$$

By choosing $x = -\frac{\beta_0}{\beta_1}$, we get probability to be 0.5.

The probability of 0.5 means, the chance of Y occurrence is 0.5. If the problem is binary, then we can say the occurrence of Y is no better than happening by chance.

7) ~~Multinomial log~~

The conditional probability and mult logistic regression for k -classes can be given by

$$P(Y=1|x) = \frac{\exp(\beta_{01} + \beta_1 x)}{1 + \sum_{j=1}^{k-1} \exp(\beta_{j0} + \beta_j x)} \quad p=1$$

$$P(Y=2|x) = \frac{\exp(\beta_{02} + \beta_2 x)}{1 + \sum_{j=1}^{k-1} \exp(\beta_{j0} + \beta_j x)} \quad k=2$$

$$P(Y=k-1|x) = \frac{\exp(\beta_{k-10} + \beta_{k-1} x)}{1 + \sum_{j=1}^{k-1} \exp(\beta_{j0} + \beta_j x)} \quad k=k-1$$

~~Other model uses last class index~~

~~$P(Y=k|x) = \frac{\exp(\beta_{k0} + \beta_k x)}{1 + \sum_{j=1}^{k-1} \exp(\beta_{j0} + \beta_j x)}$~~

~~Q this model use~~

Since, the model uses last class in the denominator we get,

$$P(Y=k|x) = \frac{1}{1 + \sum_{j=1}^{k-1} \exp(\beta_{j0} + \beta_j x)}$$

clearly, $\sum_{j=1}^k P(Y=j|x) = 1$. Thus satisfy the probability.

Taking log odds, and using the result from ~~Q~~ problem 1.5, i.e.,

~~$\frac{\text{odd}(x_i + \Delta)}{\text{odd}(x_i)} = \exp(\beta_i \Delta)$~~

We get model of the form,

$$\log \left(\frac{P(Y=1|x)}{P(Y=k|x)} \right) = \beta_{10} + \beta_1 x$$

$$\log \left(\frac{P(Y=2|x)}{P(Y=k|x)} \right) = \beta_{20} + \beta_2 x$$

$$\log \left(\frac{P(Y=k-1|x)}{P(Y=k|x)} \right) = \beta_{k-10} + \beta_{k-1} x$$

let the entire parameter set be given by

$$\theta = \{\beta_{10}, \beta_1, \dots, \beta_{(k+1)}, \beta_{k+1}\}$$

then we can represent the predictions made for an input x as following

$$\hat{p}_i(x, \theta) = \frac{\exp(\hat{\beta}_{10} + \hat{\beta}_i x)}{1 + \sum_{j=1}^{k+1} (\hat{\beta}_{j0} + \hat{\beta}_j x)}$$

$$\hat{p}_{k+1}(x, \theta_k) = \frac{\exp(\hat{\beta}_{(k+1)0} + \hat{\beta}_{k+1} x)}{1 + \sum_{j=1}^k (\hat{\beta}_{j0} + \hat{\beta}_j x)}$$

then the class predicted is according to

$$\hat{f}(x) = \arg \max_{j=1, \dots, k} \hat{p}_j(x, \theta_j)$$

Problem 2

(1) Linear Discriminant Analysis

Let us consider π_k as a prior probability that an observation that is chosen randomly belongs to k^{th} class.

Let us consider a density function $f_k(x)$ which denotes a particular observation has come from ~~the~~ class K .

$$\therefore f_k(x) = \Pr(x=x | Y=k)$$

Now, we know Bayes Theorem states as below:-

$$\Pr(Y=k | x=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)} \rightarrow ①$$

which means probability that x belongs to K^{th} class by total probability.

We have $P_k(x)$ which is posterior probability, which is equal to $\Pr(Y=k | x=x)$ and substitute the values of π_k and $f_k(x)$

$$\text{i.e } P_k(x) = \Pr(Y=k | x=x)$$

Before we estimate $f_k(x)$ we assume it is normal/Gaussian
So, the normal density function takes the form

$$f_k(x) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right)$$

where, σ_k^2 is the variance of k^{th} class

μ_k is the mean of k^{th} class:

In this case, we assume for LDA, all variances are same

$$\text{i.e } \sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2$$

Hence, we consider it as σ^2

Substitute the values of $f_k(x)$ in eqn ①, we get.

$$\begin{aligned} P_k(x) &= \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)} \\ &= \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x-\mu_l)^2\right)} \end{aligned}$$

We assumed all variance to be σ^2 , so ...

$$P_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}.$$

In the above equation, we take all term which do not vary across k as c because they do not contribute in maximisation.

$$\therefore \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x^2) + \mu_k^2 - 2x\mu_k\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}$$

$$\therefore \text{Here, } c \text{ value} = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x^2)\right)}{\sum_{l=1}^K \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}$$

So, we have:

$$P_k(x) = C \cdot \pi_k \exp\left(-\frac{1}{2\sigma^2}(\mu_k^2 - 2x\mu_k)\right)$$

Take log on both sides,

$$\begin{aligned} \log P_k(x) &= \log(c \cdot \pi_k \cdot \exp(-\frac{1}{2\sigma^2}(\mu_k^2 - 2x\mu_k))) \\ &= \log(c) + \log \pi_k + \log \exp(-\frac{1}{2\sigma^2}(\mu_k^2 - 2x\mu_k)) \\ &= \log(c) + \log \pi_k + -\frac{1}{2\sigma^2}(\mu_k^2 - 2x\mu_k) \\ &= \log(c) + \log(\pi_k) + \left(\frac{-\mu_k^2}{2\sigma^2} + \frac{2\mu_k x}{2\sigma^2}\right) \end{aligned}$$

$$\therefore \log P_k(x) = \frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k) + \underline{\log(c)}$$

Here, $\log(c)$ is a constant as it does not have any effect in maximization.

Therefore, the above equation is equivalent to the discriminant function $\delta_k(x)$.

$$\delta_k(x) = \frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

which shows, maximizing the discriminant function $\delta_k(x)$ is same as maximizing the original function $P_k(x)$ because the positive constant ($\log c$) has no effect on maximization.

(2) Quadratic discriminant Analysis.

This case is similar to that of linear discriminant analysis but here we have a class specific mean and variance. We already know the density function is given by

$$f_k(x) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right).$$

We also know the posterior probability function $P_k(x)$

$$P_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)}{\sum \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x - \mu_l)^2\right)}$$

Taking log on both sides

$$\log P_k(x) = \log \left(\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right) \right) - \log \left(\sum \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x - \mu_l)^2\right) \right)$$

$$\log P_k(x) + \log \left(\sum \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x - \mu_l)^2\right) \right) =$$

$$\log(\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)).$$

$$\Rightarrow \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi\sigma_k^2}}\right) + \log(\exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right))$$

$$\Rightarrow \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi\sigma_k^2}}\right) + \left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)$$

$$\therefore \log P_k(x) \cdot \log \left(\sum \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x - \mu_l)^2\right) \right) = \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi\sigma_k^2}}\right)$$

$$\underbrace{\log \left(\sum \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2}(x - \mu_l)^2\right) \right)}_{\text{The total probability}} + \left(-\frac{x^2}{2\sigma_k^2} + \frac{\mu_k^2}{2\sigma_k^2} - \frac{2x\mu_k}{2\sigma_k^2}\right)$$

remains constant as it does not have any effect in maximizing and let it be C.

$$\log P_k(x) \cdot C = \log(\pi_k) + \log\left(\frac{1}{\sqrt{2\pi\sigma_k^2}}\right) - \frac{x^2}{2\sigma_k^2} - \frac{\mu_k^2}{2\sigma_k^2} + \frac{x\mu_k}{\sigma_k^2}$$

The above equation is also equivalent to the ~~density~~
~~+~~ discriminant function $\delta_k(x)$.

$$\text{where, } \delta_K(x) = \log(\pi_K) + \log\left(\frac{1}{\sqrt{2\pi\sigma_K^2}}\right) + -\frac{1}{2\sigma_K^2}(x - \mu_K)^2$$

Therefore, in the equation below, x appears as a quadratic term in the classifier. which shows the case here, Bayes classifier is not linear but in fact quadratic

$$\log P_K(x) \cdot c = \log(\pi_K) + \log\left(\frac{1}{\sqrt{2\pi\sigma_K^2}}\right) \cdot \underbrace{\frac{x^2}{2\sigma_K^2}}_{\text{circled}} - \frac{\mu_K^2}{2\sigma_K^2} + \frac{x\mu_K}{\sigma_K^2}$$

therefore, there is the x^2 , hence quadratic.

3) Given that there are k iid samples, $x_i, i=1, 2 \dots k$ with a positive pairwise correlation ρ and $\text{Var}(x_i) = \sigma^2$ for $i=1, 2 \dots k$

~~consider~~ We know that by definition of mean, the mean of k iid random variables is given by

$$\bar{x} = \frac{1}{k} \sum_{i=1}^k x_i \quad \text{--- (1)}$$

~~consider~~,

$$\begin{aligned}\text{Var}(ax) &= E[(ax - a\mu)^2] \quad \text{where } \mu \text{ is the mean,} \\ &= E[a^2(x - \mu)^2] \\ &= a^2 E[(x - \mu)^2] = a^2 \text{Var}(x)\end{aligned}$$

\therefore (1) becomes,

$$\text{Var}(\bar{x}) = \text{Var}\left(\frac{1}{k} \cdot \sum_{i=1}^k x_i\right) = \frac{1}{k^2} \cdot \text{Var}\left(\sum_{i=1}^k x_i\right) \quad \text{--- (2)}$$

~~In the above equation, x_1, \dots, x_k are random var~~

Now, we can write,

$$\text{Var}\left(\sum_{i=1}^k x_i\right) = E\left(\left[\sum_{i=1}^k x_i\right]^2\right) - \left(E\left[\sum_{i=1}^k x_i\right]\right)^2$$

Note $\left(\sum_{i=1}^k a_i\right)^2 = \sum_{j=1}^k \sum_{i=1}^k a_i a_j$ viz equivalent to $(a_1 + a_2 + \dots + a_n)(a_1 + a_2 + \dots + a_n)$

$$\therefore E\left[\left(\sum_{i=1}^k x_i\right)^2\right] = E\left[\sum_{i=1}^k \sum_{j=1}^k x_i x_j\right] = \sum_{i=1}^k \sum_{j=1}^k E(x_i x_j)$$

$$\text{Hence } \left(E\left(\sum_{i=1}^k x_i\right)\right)^2 = \left(\sum_{i=1}^k E(x_i)\right)^2 = \left(\sum_{j=1}^k \sum_{i=1}^k E(x_i) \cdot E(x_j)\right)$$

$$\therefore \text{Var}\left(\sum_{i=1}^k x_i\right) = \sum_{i=1}^k \sum_{j=1}^k E(x_i x_j) - E(x_i) E(x_j) = \sum_{i=1}^k \sum_{j=1}^k \text{Cov}(x_i, x_j)$$

by definition of covariance.

Using the above result, (2) becomes,

$$\text{Var}(\bar{x}) = \frac{1}{k^2} \cdot \sum_{i=1}^k \sum_{j=1}^k \text{Cov}(x_i, x_j) \quad \text{--- (3)}$$

The above covariance can be written as

$$\text{cov}(x_i, x_j) = \begin{pmatrix} \sigma^2 & \rho\sigma^2 & \dots & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 & \dots & \vdots \\ \vdots & \ddots & \ddots & \rho\sigma^2 \\ \rho\sigma^2 & \dots & \rho\sigma^2 & \sigma^2 \end{pmatrix}$$

where ρ is pairwise correlation.

i.e. if $i=j$, then $\text{cov}(x_i, x_j) = \sigma^2$ & $\text{cov}(x_i, x_j) = \rho\sigma^2$

There are k^2 number of elements in $\text{cov}(x_i, x_j)$ { $i \neq j$ there are k elements where $i=j$. Hence the rest is given by k^2-k .

$$\therefore \text{Var}(\bar{x}) = \frac{1}{k^2} \sum_{i=1}^k \sum_{j=1}^k \text{cov}(x_i, x_j) = \frac{1}{k^2} (k\sigma^2 + (k^2-k)\rho\sigma^2)$$

$$= \frac{\sigma^2}{k} + \rho\sigma^2 - \frac{1}{k}\rho\sigma^2$$

$$= \rho\sigma^2 + \frac{\sigma^2}{k}(1-\rho)$$

Problem 4 (P, 20 Points)

1. (2P) Download and load the phoneme data set (phoneme.csv) from the course website. Split the dataset into training and test set according to the speaker column. Be sure to exclude the row number, speaker and response columns from the features.

Solution. Refer to section 4.1 in the code. ■

2. (3P) Fit an LDA model, compute and report train and test error.

Solution. Refer to section 4.2 in the code.

The train error is **0.055** and the corresponding test error is given by **0.080**. ■

3. (3P) Plot the projection of the training data onto the first two canonical coordinates of the LDA using the plot() function. Investigate the data projected on further dimensions using the *dimen* parameter.

Solution. Refer to section 4.3 in the code. Plot the projection of the training data onto the first two canonical coordinates of the LDA using the phonemes dataset is given by **Figure1**.

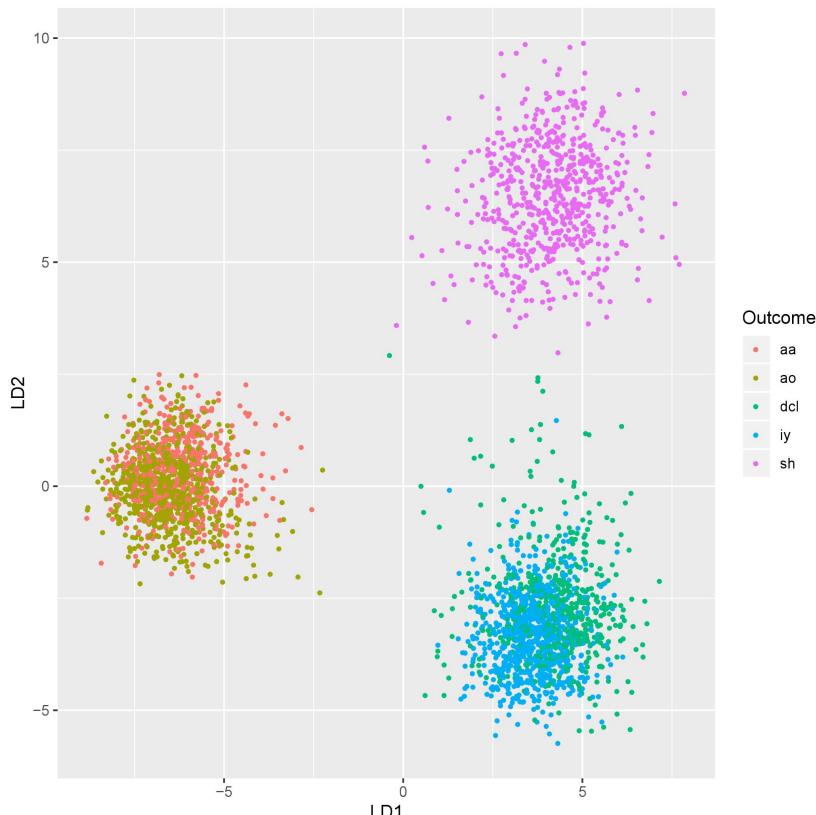


Figure 1: First two canonical coordinates of the LDA

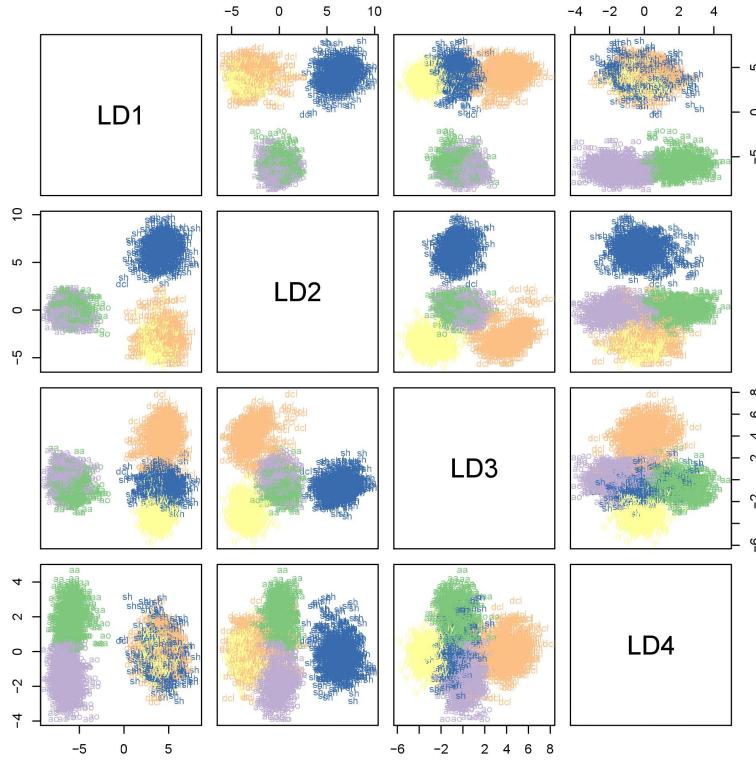


Figure 2: Plot between all the canonical coordinates of the LDA

We can see the data projected across further dimensions in **Figure2**

■

4. Select the two phonemes aa and ao. Fit an LDA model on this data set and repeat the steps done in (2).

Solution. Refer to section 4.4 in the code.

For the reduced dataset, the LDA train error is **0.106** and the test error is **0.214**.

■

5. (6P) Repeat steps (2) and (4) using QDA and report your findings. Would you prefer LDA or QDA in this example? Why?

Solution. Refer to section 4.4 in the code.

For the complete dataset, the QDA train error is **0.0** and the test error is **0.158**.

For the classes "aa" and "ao", the QDA train error is **0.0** and the test error is **0.339**.

Clearly the QDA model is overfitting to the train data and poorly generalising on the test data with respect to the QDA model. Hence, the preferred model is **LDA**.

■

6. (3P) Generate confusion matrices for the LDA and QDA model for aa and ao. Which differences can you observe between the models?

Solution. Refer to section 4.6 in the code.

The confusion matrices for the LDA and QDA model for aa and ao is given by **Table1** and **Table2** respectively.

Table 1: LDA Confusion Matrix

Pred\Actual	aa	ao
aa	121	39
ao	55	224

Table 2: QDA Confusion Matrix

Pred\Actual	aa	ao
aa	29	2
ao	147	261

From the tables it is evident that the QDA model tends to predict every observation to the class "ao", while the LDA model does fairly well in distinguishing between the two classes.

■