# Tutorial 3

Q1) refer to data uploaded in table2_7 that measured from types of insects :

(a) Ch . Concinna "A" and

(b) Ch . Heikertlingeri . "B"

It shows two variables :

• $X_1$: first the width of the first joint

• $X_2$: the width for the second joint .

(الجدول التالي يبين عملية تشريح حشرات الشاتوكتيميا لعشرين من ذكور الخنافس الصغيرة حيث أن المتغيرات هي عرض المفصل (الكاحل) الأول والثاني )

Find

1- Find the estimated Fisher's linear discriminant function.
2- Classify the new insect with observation (194 , 124 )

Solution:

1-

```
> setwd("C:/Users/Rad16/OneDrive/سطح المكتب/stat339")
> data <- read.csv(file="table2_7.csv", header=TRUE, sep=";")
> data
   Type  X1  X2
1     A 191 131
2     A 185 134
3     A 200 137
4     A 173 127
5     A 171 128
6     A 160 118
7     A 188 134
8     A 186 129
9     A 174 131
10    A 163 115
11    B 186 107
12    B 211 122
13    B 201 144
14    B 242 131
15    B 184 108
16    B 211 118
17    B 217 122
18    B 223 127
19    B 208 125
20    B 199 124
>
> data1 <- data[1:10,2:3] #data of A insect
> data2 <- data[11:20,2:3] #data of B insect
>
> n1 <- nrow(data1)
> n1
[1] 10
> n2 <- nrow(data2)
> n2
[1] 10
```

```
C:/Users/Rad16/OneDrive/سطح المكتب/stat339/
> xbar1 <- colMeans(data1)
> xbar1
   X1    X2
179.1 128.4
> xbar2 <- colMeans(data2)
> xbar2
   X1    X2
208.2 122.8
>
> s1<- cov(data1)#var each joint and cov jint with other for insect A
> s2 <- cov(data2)
>
> spooled <- (((n1-1)*s1)+((n2-1)*s2)) / (n1+n2-2) #since equal variance assumption
>
> spooled
         X1        X2
X1 231.25000 87.33333
X2  87.33333 81.88889
>
> inv.spooled <- solve(spooled)
> #or we can find inverse by Generalized Inverse of a Matrix
> library(MASS)
> inv.spooled <- ginv(spooled)
> inv.spooled
            [,1]         [,2]
[1,]  0.007240593 -0.007721989
[2,] -0.007721989  0.020447060
>
> #The cutoff point to determine group membership of the observation vector is then
 found
> mhat <- ((xbar1-xbar2)%*%inv.spooled%*%(xbar1+xbar2)) / (2)
> mhat
          [,1]
[1,] -6.571125
>
```

2-

```
[1,] -6.571125
>
> y <- (xbar1-xbar2)%*%inv.spooled
> y
            [,1]      [,2]
[1,] -0.2539444 0.3392134
> x0 <- c(194, 124)
> y0 <- (xbar1-xbar2)%*%inv.spooled%*%x0
> y0
          [,1]
[1,] -7.202747
>
>
> if (y0 >= mhat){
+             print ("A")
+       } else {
+             print ("B")
+  }
[1] "B"
```

we classify it to insect B since Y<m

Q2) upload "iris" data frame with 150 cases (rows) and 5 variables (columns) where

It shows two variables :

• $X_1$: iris where its Species :setosa

• $X_2$: iris where its Species : versicolor

Find

1. the estimated Fisher's linear discriminant function.
2. Classify the new observation of row #40

Sepal.Length Sepal.Width Petal.Length Petal.Width

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width |
|---|---|---|---|---|
| 40 | 5.1 | 3.4 | 1.5 | 0.2 |

Solution:

```
> library(MASS) #Load package 'MASS' to call iris data
> data(iris)
> str(iris)
'data.frame':   150 obs. of  5 variables:
 $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
 $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
 $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
 $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
 $ Species     : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1
...
> iris
    Sepal.Length Sepal.Width Petal.Length Petal.Width   Species
1            5.1         3.5          1.4         0.2    setosa
2            4.9         3.0          1.4         0.2    setosa
3            4.7         3.2          1.3         0.2    setosa
4            4.6         3.1          1.5         0.2    setosa
5            5.0         3.6          1.4         0.2    setosa
6            5.4         3.9          1.7         0.4    setosa
7            4.6         3.4          1.4         0.3    setosa
8            5.0         3.4          1.5         0.2    setosa
9            4.4         2.9          1.4         0.2    setosa
10           4.9         3.1          1.5         0.1    setosa
11           5.4         3.7          1.5         0.2    setosa
12           4.8         3.4          1.6         0.2    setosa
13           4.8         3.0          1.4         0.1    setosa
14           4.3         3.0          1.1         0.1    setosa
15           5.8         4.0          1.2         0.2    setosa
16           5.7         4.4          1.5         0.4    setosa
17           5.4         3.9          1.3         0.4    setosa
18           5.1         3.5          1.4         0.3    setosa
19           5.7         3.8          1.7         0.3    setosa
20           5.1         3.8          1.5         0.3    setosa
21           5.4         3.4          1.7         0.2    setosa
```

```
> data <- iris[-101:-150, ]
> data <- iris[1:100, ] #just to focus on first 100 obs.
>
> data1 <- data[data$Species == "setosa",][ ,1:4]
> data1 <- data[1:50,1:4]
> data2 <- data[data$Species == "versicolor",][ ,1:4]
> data2 <- data[51:100,1:4]
>
> n1 <- nrow(data1)
> n2 <- nrow(data2)
> xbar1 <- colMeans(data1)
> xbar2 <- colMeans(data2)
> s1 <- cov(data1)
> s2 <- cov(data2)
> spooled <- (((n1-1)*s1)+((n2-1)*s2)) / (n1+n2-2)
> inv.spooled <- solve(spooled)
> inv.spooled
             Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length    11.634235   -6.552973    -7.998430    3.884391
Sepal.Width     -6.552973   14.236847     3.274256  -10.853906
Petal.Length    -7.998430    3.274256    21.497513  -26.658191
Petal.Width      3.884391  -10.853906   -26.658191   87.666138
> mhat<- ((xbar1-xbar2)%*%inv.spooled%*%(xbar1+xbar2)) / (2)
> mhat
          [,1]
[1,] -13.96174
> x0 <- t(as.matrix(data[40,1:4])) #assume that it is 40th rows
> y0 <- (xbar1-xbar2)%*%inv.spooled%*%x0
> y0
         40
[1,] 38.02906
> group <- ifelse(y0 >= mhat, "setosa", "versicolor")
> group
      40
[1,] "setosa"
```

we classify it to "setosa" species since y>m as it is actual be.

Q3)According to slide 87 about "Longely" dataset in R that describes 7 economic variables observed from 1947 to 1962 used to predict the number of people employed yearly (n=16).

  a. Fit classical multiple linear regression and ridge regression
  b. Obtain the estimated regression coefficient. Which technique provide the smallest coefficients?
  c. Perform LOOCV to get best lambda?

Solution:

First we call the data and split it into independent and dependent variable:



**a.** Fit classical multiple linear regression to find its Coefficients.

Fit ridge regression to find its Coefficients.





**b.**

- Increasing lambda increases the shrinking of the coefficients.

- If the sample size is very small compared to the number of covariates, estimation is not efficient and therefore we might not get the desirable shrinking.