

# STAT 501 Homework 8

## Multinomial

April 27, 2018

First we standardized the data and use the standardized data in the clustering.

```
1 women <- read.table("women-track-records.dat",
2                       header=F, col.names=c("x1", "x2", "x3", "x4", "x5", "x6", "x7",
3                                               ↪ "Country"))
4 men <- read.table("men-track-records.dat",
5                   header=F, col.names=c("x1", "x2", "x3", "x4", "x5", "x6", "x7", "x8",
6                                           ↪ "Country"))
7 women.country <- women$Country
8 men.country <- men$Country
9
10 women <- women[,-ncol(women)]
11 men <- men[,-ncol(men)]
12
13 # standardize the data
14 women <- scale(women)
15 men <- scale(men)
```

- Hierarchical Clustering

Before doing hierarchical clustering, we mapped the data to the sphere. In this way the Euclidian distance is just the correlation similarity distance for the standardized the data. We adopted complete linkage in hierarchical clustering for men and women. The hierarchical tree is shown in Figure 1 and we display the clustering result with 3 clusters in Figure 2 using the 2 LDs by LDA.

```
1 # Sphere the data for correlation similarity distance
2 men.sphere<-t(apply(men,1,FUN= function(x) (x - mean(x))/sqrt(sum((x - mean(x))^2))))
3 women.sphere<-t(apply(women,1,FUN= function(x) (x - mean(x))/sqrt(sum((x -
4     ↪ mean(x))^2))))
5
6 # Use complete linkage, this is also the default method
7 hc.men <- hclust(dist(men.sphere),method="complete")
8 plot(hc.men, label = men.country,
9      main = "Complete Linkage Cluster Analysis: Men Track Records Data")
10
11 hc.women <- hclust(dist(women.sphere),method="complete")
12 plot(hc.women, label = women.country,
13      main = "Complete Linkage Cluster Analysis: Women Track Records Data")
14
15 # Use canonical discriminants to display the clusters.
16 # The first function computes linear canonical discriminants and
17 # the second function is used to plot the computed scores.
18 library(MASS)
19 # Compute canonical discriminant scores and display 2-dimensional
20 # projections of the clusters
21 hc.men.lda <- lda(men, cutree(hc.men, 3))
22 plot(hc.men.lda)
```

```

23
24 hc.women.lda <- lda(women, cutree(hc.women, 3))
25 plot(hc.women.lda)

```

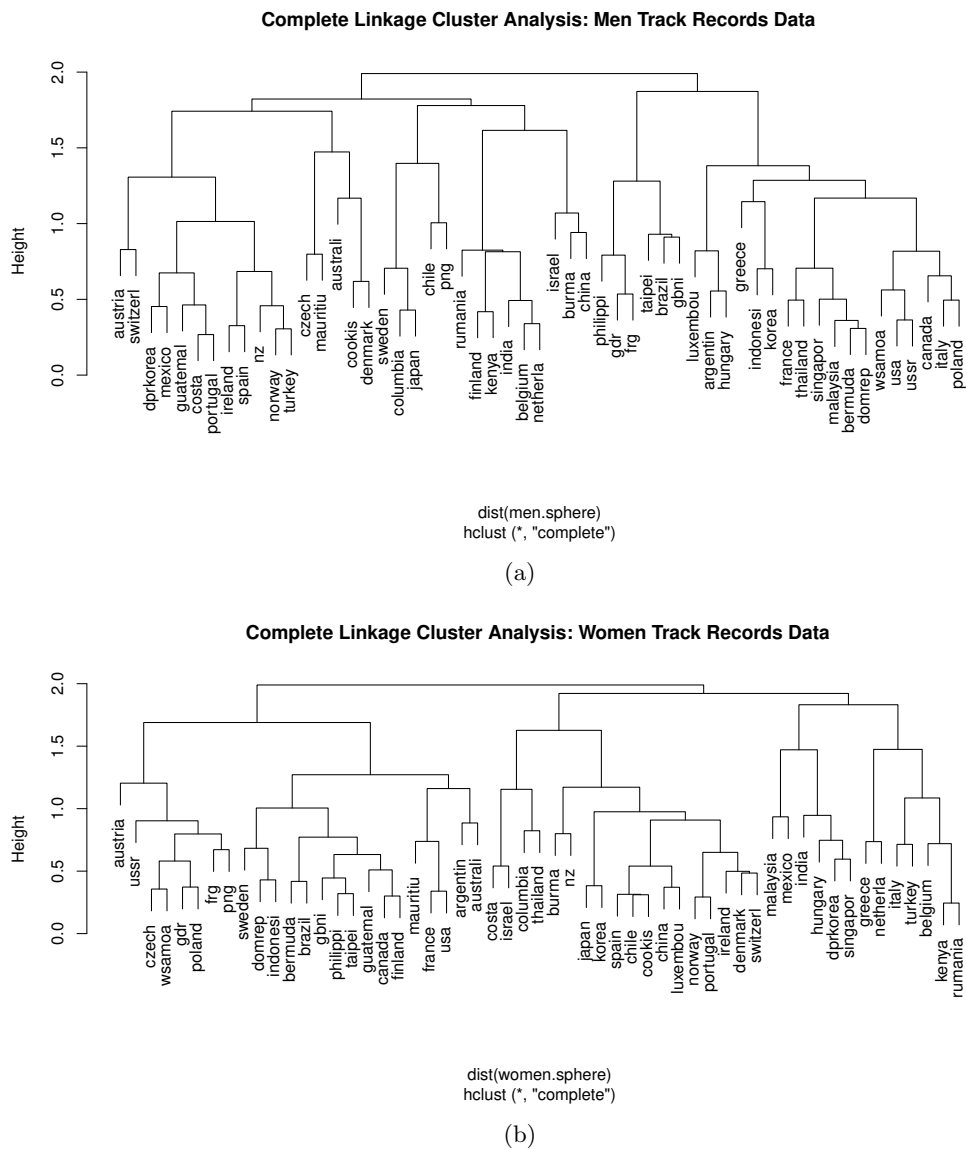


Figure 1: Hierarchical clustering tree

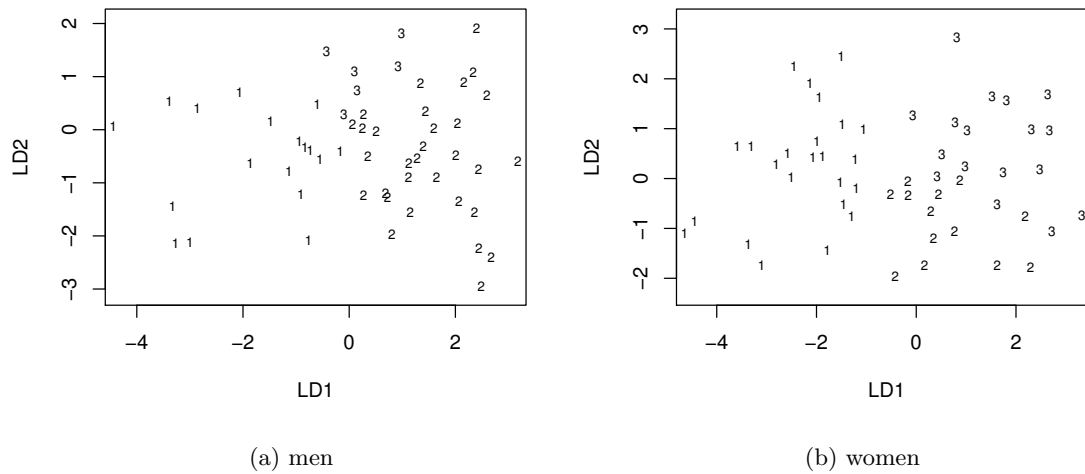


Figure 2: Hierarchical clustering results

- k-means

We use the results of hierarchical clustering with ward linkage as our initial for k-means. Then we display the results from k-means with 3 clusters using the 2 LDs from LDA in Figure 3.

```

1  # Compute K-means cluster analysis starting with results from hclust
2  kmnsinithcl <- function(x.data, nclus, ncut = nclus, hcl.tree)
3  {
4    x.hcl <- hcl.tree
5    x.cl <- cutree(x.hcl, k = ncut)
6    data.x <- data.frame(x.data, cl = x.cl)
7    means <- aggregate(. ~ cl, data = data.x, FUN = mean)
8    return(kmeans(x.data, centers= means[, -1]))
9  }
10
11 hc.men2 <- hclust(dist(men.sphere),method="ward.D2")
12 km.men <- kmnsinithcl(men, nclus = 3, ncut = 3, hcl.tree = hc.men2)
13 km.men.lda <- lda(men, km.men$cluster)
14 plot(km.men.lda)
15
16 hc.women2 <- hclust(dist(women.sphere),method="ward.D2")
17 km.women <- kmnsinithcl(women, nclus = 3, ncut = 3, hcl.tree = hc.women2)
18 km.women.lda <- lda(women, km.women$cluster)
19 plot(km.women.lda)

```

- Model based clustering

For model based clustering, we use BIC to pick the model and number of clusters. In Figure 4 we can see, for men, model VEE and 2 clusters is the best, while model VEV and 2 clusters for women is the best. The display of clustering is shown in Figure 5.

```

1  # Model based clustering based on the assumption that
2  # the data were sampled from a mixture of multivariate normal
3  # distributions with common covariance matrix
4
5  library(mclust)
6

```

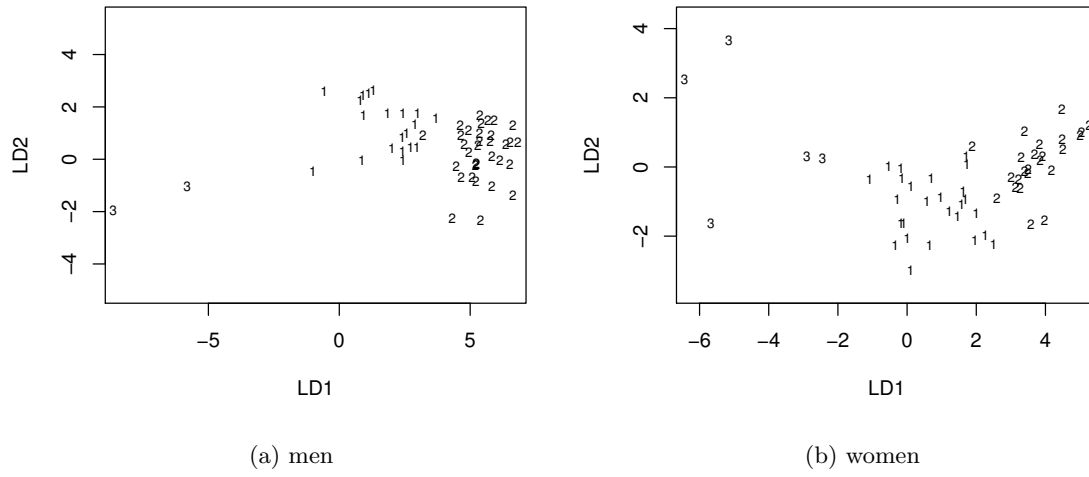


Figure 3: k-means clustering results

```

7 mcl.men <- Mclust(men)
8 plot(mcl.men$BIC)
9 plot.Mclust(mcl.men, what = "classification")
10
11 mcl.women <- Mclust(women)
12 plot(mcl.women$BIC)
13 plot.Mclust(mcl.women, what = "classification")

```

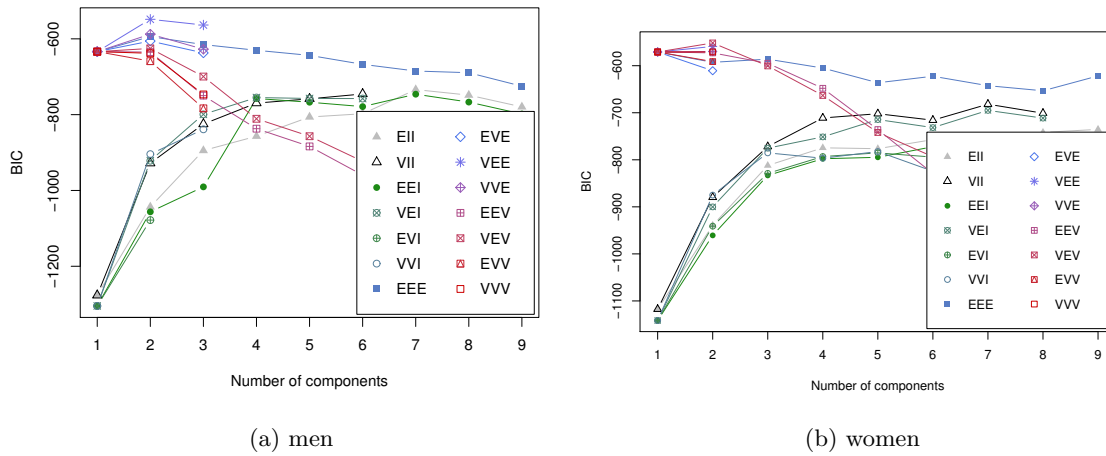
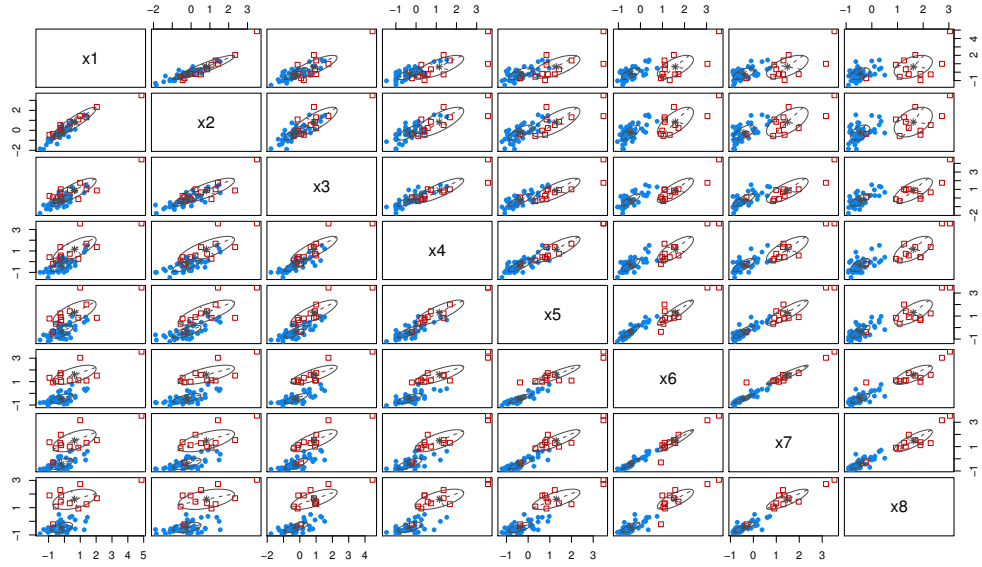
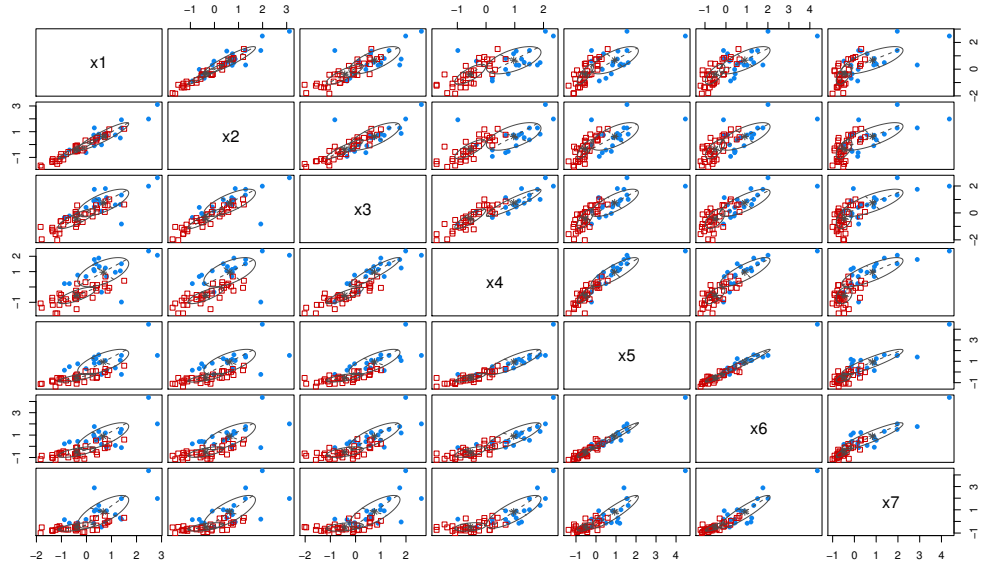


Figure 4: BIC for model selection



(a) men



(b) women

Figure 5: Model based clustering results

Comment:

The clustering results using 3 methods shows some differences in the clusters patterns between men and women. However we can see from the results using 3 methods that there might be 2 clusters for both women and men.