

# HUDM6122 Homework\_01

Chenguang Pan

Jan 28, 2023

## 0.1 Exercise 1.1

First, I made a `xlsx` version of **Table 1.1** to let R read it directly using the package ‘`readxl`’. This table is in 10x7 size. The first column is just the index of each observation, so I drop it here. Finally this dataset is in 9x7 size.

One should notice that the `sex`, `depression`, and `health` are categorical variables. The Pearson Correlation Coefficient is used for continuous rather than categorical variables. Therefore, when calculate the correlation matrix we should drop the categorical ones.

Note, there are some parameters need to be set. Since the original dataset contains missing value, I construct the correlation matrix based on all complete observations.

```
> library(readxl)
> table_11 <- read_excel("table_1.1.xlsx")
> my_data <- table_11[,c(2:7)]
> # drop the discrete vars and use only the complete observations
> my_data_cor <- round(cor(my_data[,c(2,3,6)]), use = "complete"),2)
> # the output is rounded in two decimals.
> my_data_cor
      age      IQ weight
age    1.00 -0.15 -0.12
IQ     -0.15  1.00  0.75
weight -0.12  0.75  1.00
```

## 0.2 Exercise 1.2

Fill the NA with the column's mean, and recalculate the correlation matrix.

```
> # to impute the NA with mean using a for-loop
> for (cols in c(2,3,6)) {
+   my_data[,cols][is.na(my_data[,cols])] <- mean(my_data[,cols], na.rm=T)
+ }
> # create the correlation matrix
> my_data_cor_2 <- round(cor(my_data[,c(2,3,6)]),2)
> my_data_cor_2
      age      IQ weight
age    1.00 -0.14 -0.10
IQ     -0.14  1.00  0.52
weight -0.10  0.52  1.00
```

### 0.3 Exercise 1.3

The authors may not clearly say where can readers find the original dataset. After seeing the code underlying in this book's R package `MVA`, and run this argument `demo("Ch-MVA")`, one can find the `table 1.3's` dataset information at the paragraph named `code chunk number 5`. It is in another package called `HSAUR2`. First, I draw the normal probability plots of each variable.

```
> # load the Table 1.3's original dataset
> library(HSAUR2)
> dim(pottery)
```

```
[1] 45 10
```

```
> names(pottery) # the dataset is the same
```

```
[1] "Al2O3" "Fe2O3" "MgO" "CaO" "Na2O" "K2O" "TiO2" "MnO" "BaO"
[10] "kiln"
```

```
> #layout(matrix(1:10, nc = 3))
> sapply(colnames(pottery)[1:9], function(x){
+   # use double bracket can directly return the col
+   qqnorm(pottery[[x]], main = x)
+   qqline(pottery[[x]])
+ })
```

```
$Al2O3 NULL
```

```
$Fe2O3 NULL
```

```
$MgO NULL
```

```
$CaO NULL
```

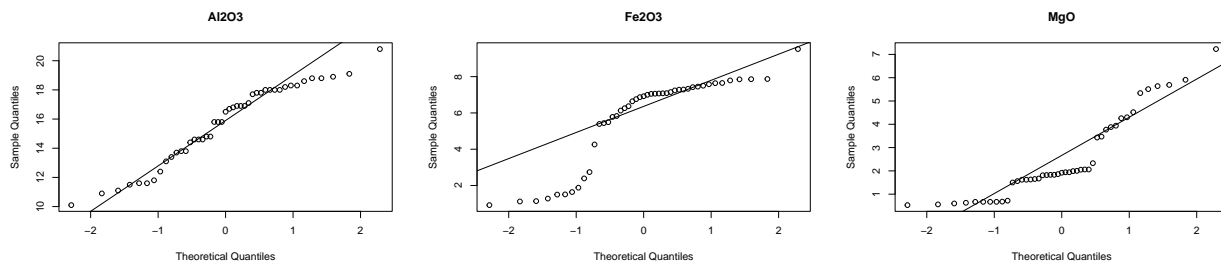
```
$Na2O NULL
```

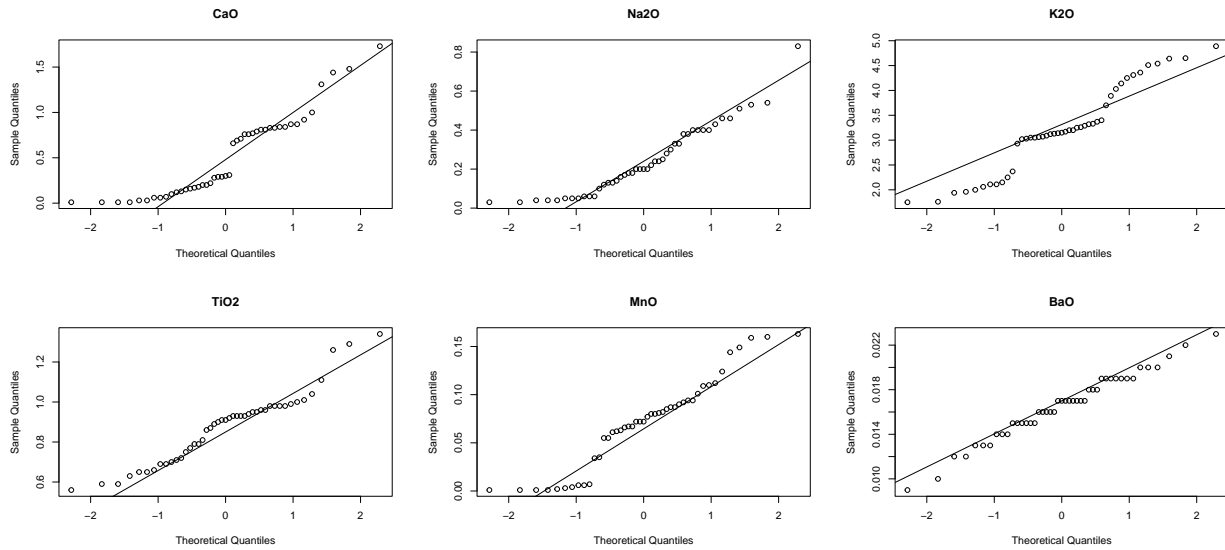
```
$K2O NULL
```

```
$TiO2 NULL
```

```
$MnO NULL
```

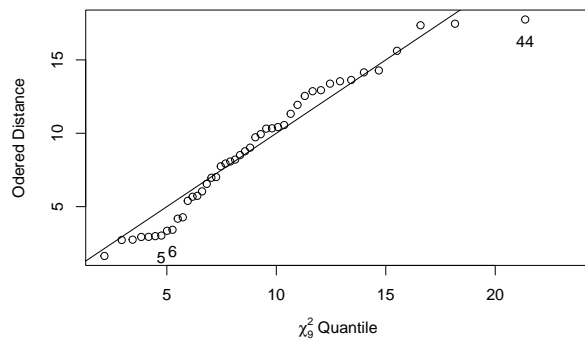
```
$BaO NULL
```





Second, I draw the chi-square plot of the data as followed.

```
> # to drop the last discrete variable
> X <- pottery[,1:9]
> # get each col's mean
> col_mean <- colMeans(X)
> # get the cov matrix
> S <- cov(X)
> # solve() function can get inverse matrix directly
> # compute the generalised distance
> d <- apply(X, 1, # array; 1= row 2 = col
+           function(X)
+             t(X - col_mean) %*% solve(S) %*% (X - col_mean))
> plot(qc <- qchisq((1:nrow(X) - 1/2)/ nrow(X), df = 9),
+      sd <- sort(d),
+      xlab = expression(paste(chi[9]^2, "Quantile")),
+      ylab = "Ordered Distance", xlim = range(qc)*c(1, 1.1))
> ous <- which(rank(abs(qc-sd), ties="random") > nrow(X) - 3)
> text(qc[ous], sd[ous]-1.5, names(ous))
> abline(a=0, b=1)
```



The overall results showed that TODO.

## 0.4 Exercise 1.4

This question can be addressed in one line using `cov2cor` function

```
> # import the cov matrix
> a_matrix<- matrix(c(3.8778, 2.8110, 3.1480, 3.5062,
+                    2.8110, 2.1210, 2.2669, 2.5690,
+                    3.1480, 2.2669, 2.6550, 2.8341,
+                    3.5062, 2.5690, 2.8341, 3.2352),4,4, byrow=T)
> round(cov2cor(a_matrix),3)
      [,1] [,2] [,3] [,4]
[1,] 1.000 0.980 0.981 0.990
[2,] 0.980 1.000 0.955 0.981
[3,] 0.981 0.955 1.000 0.967
[4,] 0.990 0.981 0.967 1.000
```

But here I provide another more specific code to achieve this goal. Could you please consider to give me extra credits? Just kidding..lol..

```
> # function to convert covariance matrix to correlation matrix
> cov2cor_test <- function(covmat) {
+   sd_vec <- sqrt(diag(covmat))
+   cormat <- covmat / outer(sd_vec, sd_vec)
+   return(cormat)
+ }
>
> round(cov2cor_test(a_matrix),3)
      [,1] [,2] [,3] [,4]
[1,] 1.000 0.980 0.981 0.990
[2,] 0.980 1.000 0.955 0.981
[3,] 0.981 0.955 1.000 0.967
[4,] 0.990 0.981 0.967 1.000
```

Finally, the results are exactly the same.

## 0.5 Exercise 1.5

Here I write a function to get the Euclidean distance matrix for any shape of matrix. Since for any  $m * n$  matrix, there number of Eu-dist matrix's columns are always  $n-1+n-2+\dots+1$ . That is, there are  $n*(n-1)/2$  columns in the Euclidean distance.

```
>
> # write a function that can get Eu-dist matrix any shape of matrix
>
> eudis_matrix<-function(a_matrix){
+   # first define a function to get the eu distance of any two vectors
+   eu_distance <- function(vec_1, vec_2){
+     eu_dist <- sqrt(sum(vec_1-vec_2)^2)
+     # this
+     return(eu_dist)
+   }
+ }
```

```

+   m <- nrow(a_matrix)
+   dist_matrix <- matrix(nrow = m, ncol = m)
+   for (i in 1:m) {
+     for (j in 1:m) {
+       dist_matrix[i,j] <- eu_distance(a_matrix[i,], a_matrix[j,])
+     }
+   }
+   return(dist_matrix)
+ }

```

The function seems good. Then, I tried it on the given dataset.

```

> # import the data
> multi_mat<- matrix(c(3,6,4,0,7,
+                      4,2,7,4,6,
+                      4,0,3,1,5,
+                      6,2,6,1,1,
+                      1,6,2,1,4,
+                      5,1,2,0,2,
+                      1,1,2,6,1,
+                      1,1,5,4,4,
+                      7,0,1,3,3,
+                      3,3,0,5,1),10,5, byrow=T)
> # create the euclidean dis-matrix on the given dataset.
> eudis_matrix(multi_mat)

```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
[1,]	0	3	7	4	6	10	9	5	6	8
[2,]	3	0	10	7	9	13	12	8	9	11
[3,]	7	10	0	3	1	3	2	2	1	1
[4,]	4	7	3	0	2	6	5	1	2	4
[5,]	6	9	1	2	0	4	3	1	0	2
[6,]	10	13	3	6	4	0	1	5	4	2
[7,]	9	12	2	5	3	1	0	4	3	1
[8,]	5	8	2	1	1	5	4	0	1	3
[9,]	6	9	1	2	0	4	3	1	0	2
[10,]	8	11	1	4	2	2	1	3	2	0

Note each row of an empirical data  $n \times p$  matrix represents an observation/vector with  $P$  dimensions. Therefore, the entry  $[i, j]$  on a  $n \times n$  Euclidean distance matrix must be a scalar rather than a vector and represents the euclidean distance between  $i$ th observation/vector and  $j$ th observation/vector.