

Multivariate Statistical Methods

Assignment 2

Allan Gholmi, Emma Wallentinsson, Rasmus Holm

2017-12-08

Question 1

```
data <- read.table("../data/T1-9.dat")
names(data) <- c("country", "100m", "200m", "400m", "800m", "1500m", "3000m", "marathon")
numeric_data <- data[, -1]

countries <- as.character(data$country)
```

a)

```
X <- as.matrix(numeric_data)
means <- colMeans(X)
covariances <- cov(X)
X_central <- X - rep(1, nrow(X)) %*% t(means)

mdist_sq <- X_central %*% solve(covariances) %*% t(X_central)
country_mdists <- diag(mdist_sq)

significance_level <- 0.1
p <- ncol(X)
quantile <- qchisq(1 - significance_level, df=p)

outliers <- country_mdists > quantile
print("Outliers without correction")
#> [1] "Outliers without correction"
countries[outliers]
#> [1] "COK" "KORN" "MEX" "PNG" "SAM"
```

No clue what the multiple-testing correction procedure refers to.

b)

The Mahalanobis takes the covariances into consideration so the distances lead to an elliptic decision boundary as opposed to the circular boundary by Euclidean distance. That indicates that North Korea is an outlier based on the covariances meaning their result does not follow the general trend.

Question 2

a)

```
bird <- read.table("../data/T5-12.DAT")

mu <- c(190, 275) #mus
x_bar <- colMeans(bird)
S <- cov(bird)
angles <- seq(0, 2 * pi, length.out=200) # make angles for circle

n <- nrow(bird)
p <- ncol(bird)

confidence_level <- 0.05

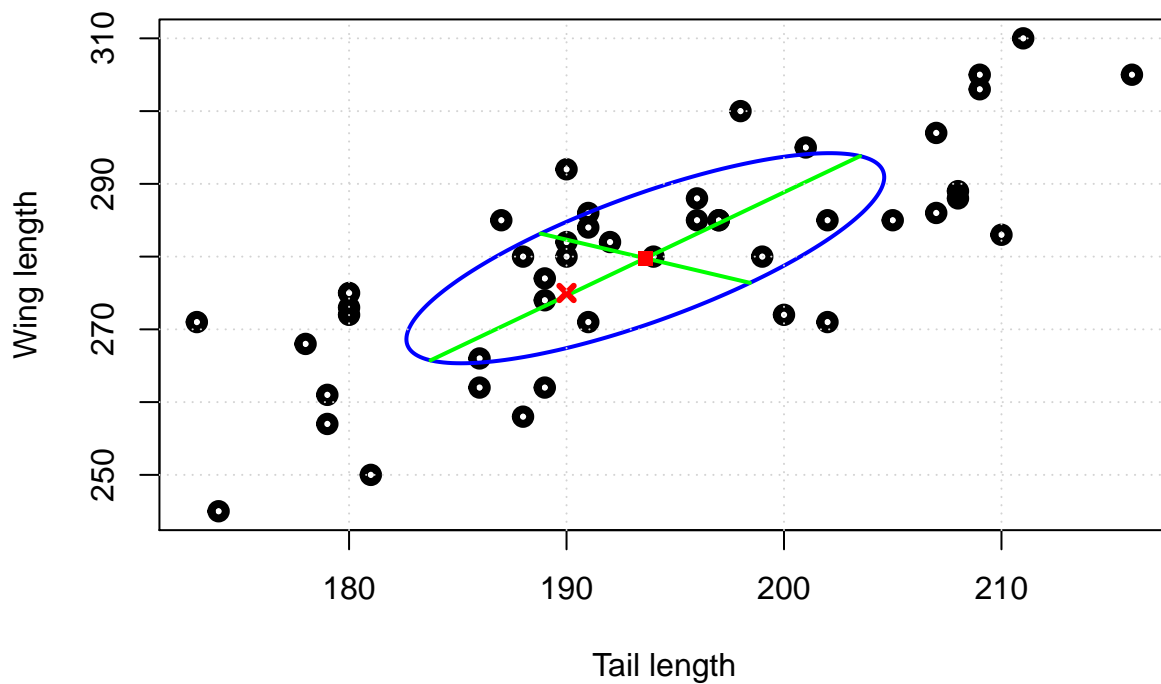
## eigenvalues and eigenvectors from covariance matrix S
eigVal <- eigen(S)$values
eigVec <- eigen(S)$vectors

#scale eevectors to unit length
scaled <- eigVec %*% diag(sqrt(eigVal)) # scale eigenvectors to length = square-root

c2 <- qchisq(1 - confidence_level, p)
c <- sqrt(c2)

xMat <- rbind(x_bar[1] + scaled[1, ], x_bar[1] - scaled[1, ])
yMat <- rbind(x_bar[2] + scaled[2, ], x_bar[2] - scaled[2, ])
ellBase <- cbind(sqrt(eigVal[1])*cos(angles), sqrt(eigVal[2])*sin(angles)) # making a circle base...

ellax <- eigVec %*% t(ellBase) # where the ellips axis goes through eigenvectors.
plot(bird, lwd="4", xlab="Tail length", ylab="Wing length")
lines((ellax+x_bar)[1, ], (ellax+x_bar)[2, ], asp=1, type="l", lwd=2, col="blue")
matlines(xMat, yMat, lty=1, lwd=2, col="green") #
points(mu[1], mu[2], pch=4, col="red", lwd=3)
grid()
points(mean(bird[,1]),mean(bird[,2]), type="p", col="red", pch=15)
```



b)

```
Tsq_offset <- sqrt(p * (n - 1) * qf(1 - confidence_level, df1=p, df2=n - p) / (n - p) * diag(S) / n)
Tsq_confidence_interval <- rbind(x_bar - Tsq_offset, x_bar + Tsq_offset)

bonferroni_offset <- sqrt(diag(S) / n) * qt(1 - confidence_level / (2 * p), df=n - 1)
bonferroni_confidence_interval <- rbind(x_bar - bonferroni_offset, x_bar + bonferroni_offset)

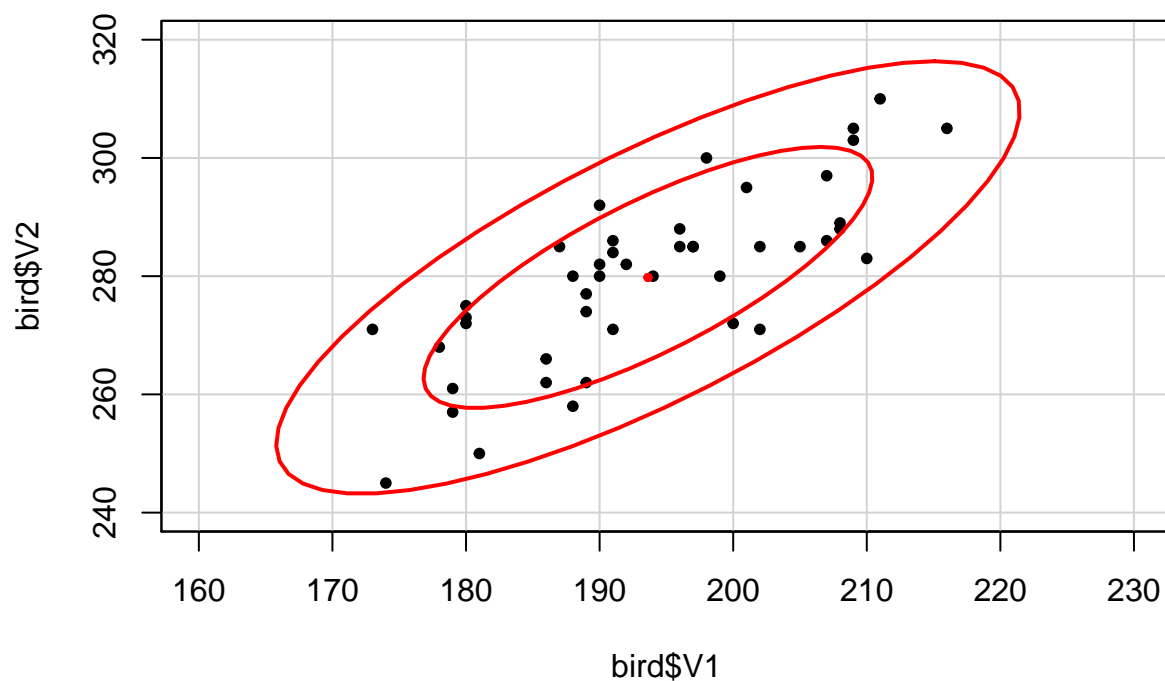
print("T-square Intervals")
#> [1] "T-square Intervals"
Tsq_confidence_interval
#>      V1      V2
#> [1,] 189.4217 274.2564
#> [2,] 197.8227 285.2992

print("Bonferroni Intervals")
#> [1] "Bonferroni Intervals"
bonferroni_confidence_interval
#>      V1      V2
#> [1,] 189.8216 274.7819
#> [2,] 197.4229 284.7736
```

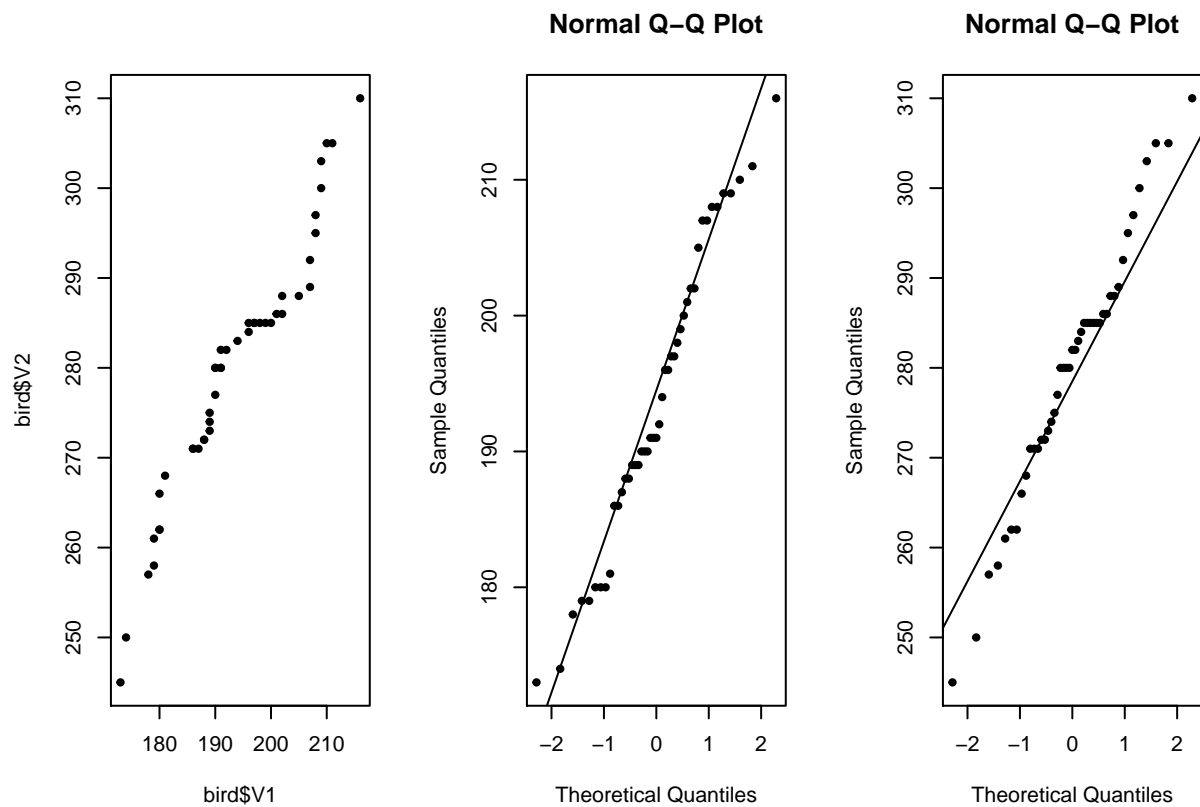
T-square test always gives wider confidence intervals since it takes the correlation between the measured variables into account. Bonferroni intervals are more precise if you are interested in the individual component means, but if you are interested in the overall data mean you should consider the T-square intervals.

c)

```
dataEllipse(x=bird$V1, y=bird$V2, pch=20, levels=c(0.68, 0.95),  
            xlim=c(160, 230), ylim=c(240, 320), center.cex=0.5)
```



```
old <- par(mfrow=c(1, 3))  
qqplot(bird$V1, bird$V2, pch=20)  
qqnorm(bird$V1, pch=20)  
qqline(bird$V1)  
qqnorm(bird$V2, pch=20)  
qqline(bird$V2)
```



```
par(old)
```

A bivariate normal distribution would be a viable population model. The qqplots do not deviate to much from the straight lines and the scatter plot shows that the points could very well have been generated from a bivariate normal distribution.

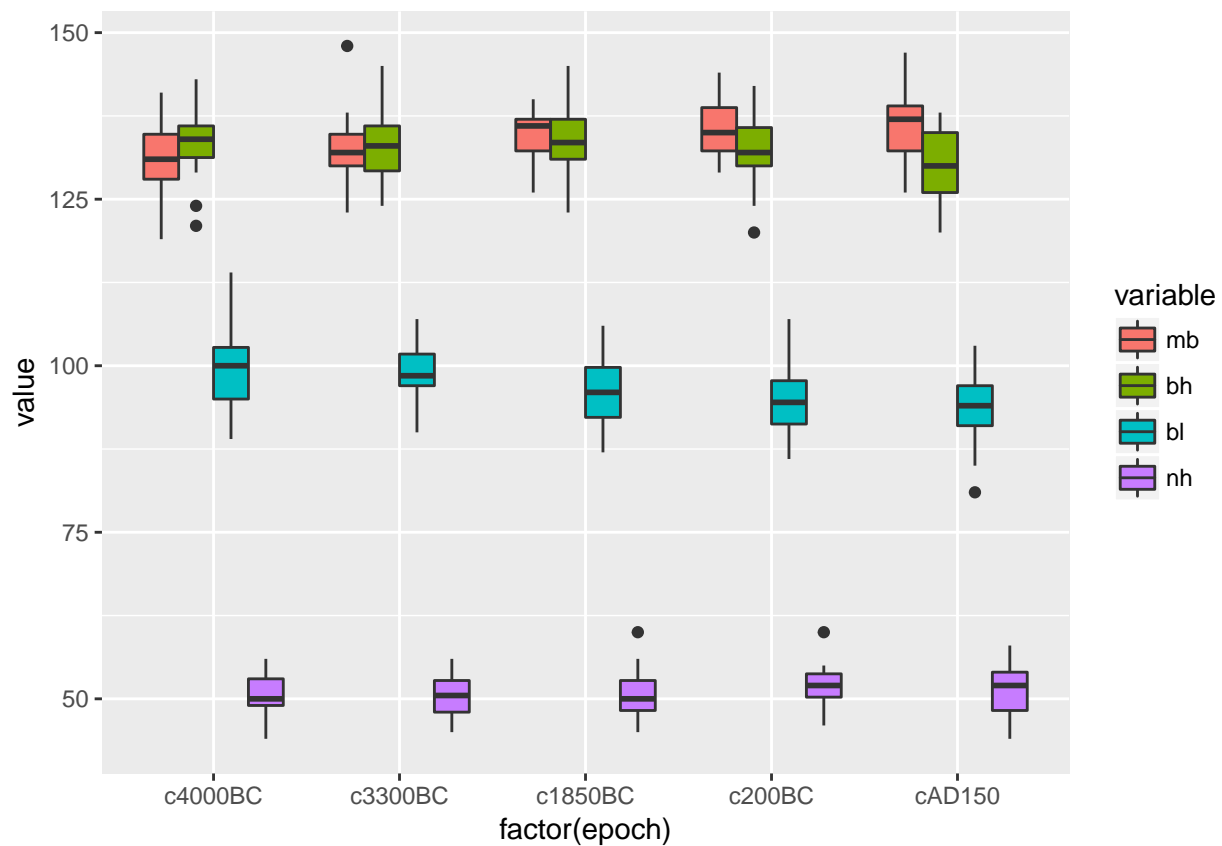
Question 3

```
library(heplots)
library(dplyr)
library(ggplot2)
library(reshape2)

data <- Skulls
numeric_data <- data[, -1]
colors <- as.numeric(data$epoch)
```

a)

```
# pairs(numeric_data, col=colors)
mm <- melt(data, id="epoch")
ggplot(mm) +
  geom_boxplot(aes(x=factor(epoch), y=value, fill=variable))
```

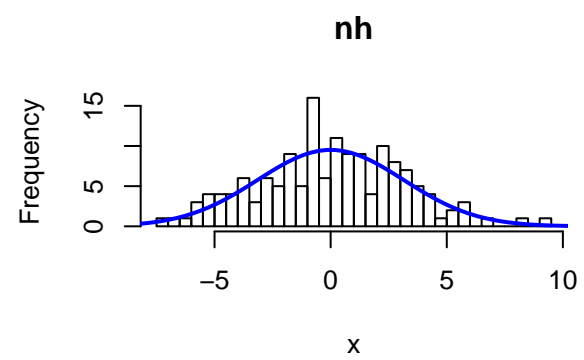
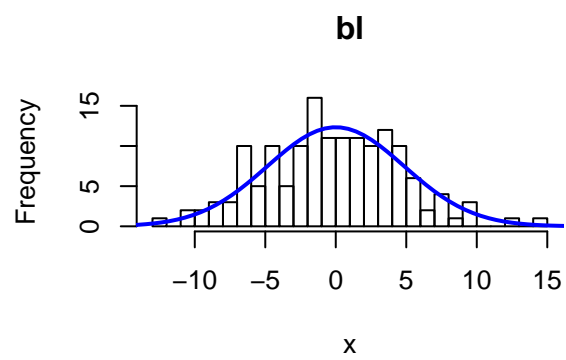
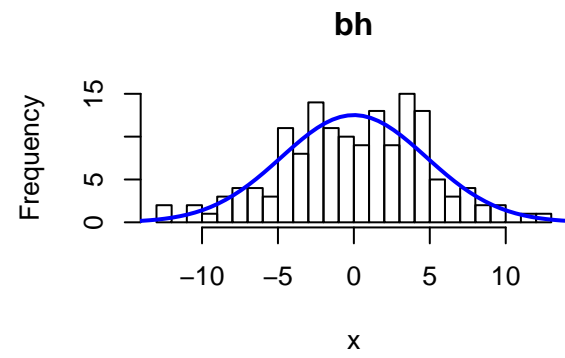
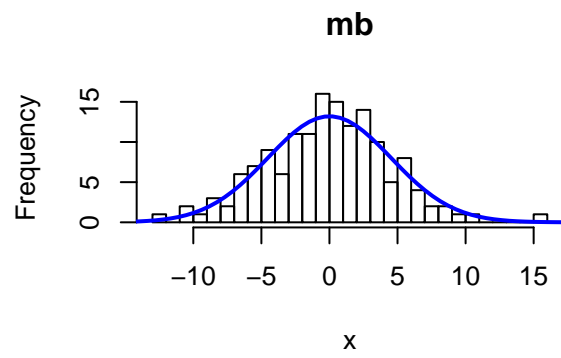


b)

```
group_means <- data %>%  
  group_by(epoch) %>%  
  summarise_all(funs(mean(., na.rm=TRUE)))  
  
fit <- manova(cbind(mb, bh, bl, nh) ~ data$epoch, data)
```

c)

```
residuals <- fit$res  
col_names <- c("mb", "bh", "bl", "nh")  
  
old <- par(mfrow=c(2, 2))  
  
for (col in 1:ncol(residuals)) {  
  x <- residuals[, col]  
  main <- col_names[col]  
  h <- hist(x, breaks=25, main=main)  
  offset <- (max(x) - min(x)) / 2  
  xfit <- seq(min(x) - offset, max(x) + offset, length = 100)  
  yfit <- dnorm(xfit, mean = mean(x), sd = sd(x))  
  yfit <- yfit * diff(h$mids[1:2]) * length(x)  
  lines(xfit, yfit, col="blue", lwd=2)  
}
```



```
par(old)
```


Appendix

Code

```
# Question 1
data <- read.table("../data/T1-9.dat")
names(data) <- c("country", "100m", "200m", "400m", "800m", "1500m", "3000m", "marathon")
numeric_data <- data[, -1]

countries <- as.character(data$country)
X <- as.matrix(numeric_data)
means <- colMeans(X)
covariances <- cov(X)
X_central <- X - rep(1, nrow(X)) %*% t(means)

mdist_sq <- X_central %*% solve(covariances) %*% t(X_central)
country_mdists <- diag(mdist_sq)

significance_level <- 0.1
p <- ncol(X)
quantile <- qchisq(1 - significance_level, df=p)

outliers <- country_mdists > quantile
print("Outliers without correction")
countries[outliers]

# Question 2

bird <- read.table("../data/T5-12.DAT")

mu <- c(190, 275) #mus
x_bar <- colMeans(bird)
S <- cov(bird)
angles <- seq(0, 2 * pi, length.out=200) # make angles for circle

n <- nrow(bird)
p <- ncol(bird)

confidence_level <- 0.05

## eigenvalues and eigenvectors from covariance matrix S
eigVal <- eigen(S)$values
eigVec <- eigen(S)$vectors

#scale evectors to unit length
scaled <- eigVec %*% diag(sqrt(eigVal)) # scale eigenvectors to length = square-root

c2 <- qchisq(1 - confidence_level, p)
c <- sqrt(c2)

xMat <- rbind(x_bar[1] + scaled[1, ], x_bar[1] - scaled[1, ])
yMat <- rbind(x_bar[2] + scaled[2, ], x_bar[2] - scaled[2, ])
```

```

ellBase <- cbind(sqrt(eigVal[1])*cos(angles), sqrt(eigVal[2])*sin(angles)) # making a circle base...

ellax <- eigVec %*% t(ellBase) # where the ellips axis goes through eigenvectors.
plot(bird, lwd="4", xlab="Tail length", ylab="Wing length")
lines((ellax+x_bar)[1, ], (ellax+x_bar)[2, ], asp=1, type="l", lwd=2, col="blue")
matlines(xMat, yMat, lty=1, lwd=2, col="green") #
points(mu[1], mu[2], pch=4, col="red", lwd=3)
grid()
points(mean(bird[,1]), mean(bird[,2]), type="p", col="red", pch=15)
Tsqr_offset <- sqrt(p * (n - 1) * qf(1 - confidence_level, df1=p, df2=n - p) / (n - p) * diag(S) / n)
Tsqr_confidence_interval <- rbind(x_bar - Tsqr_offset, x_bar + Tsqr_offset)

bonferroni_offset <- sqrt(diag(S) / n) * qt(1 - confidence_level / (2 * p), df=n - 1)
bonferroni_confidence_interval <- rbind(x_bar - bonferroni_offset, x_bar + bonferroni_offset)

print("T-square Intervals")
Tsqr_confidence_interval

print("Bonferroni Intervals")
bonferroni_confidence_interval

# Question 3
library(heplots)
library(dplyr)
library(ggplot2)
library(reshape2)

data <- Skulls
numeric_data <- data[, -1]
colors <- as.numeric(data$epoch)
# pairs(numeric_data, col=colors)
mm <- melt(data, id="epoch")
ggplot(mm) +
  geom_boxplot(aes(x=factor(epoch), y=value, fill=variable))
group_means <- data %>%
  group_by(epoch) %>%
  summarise_all(funs(mean(., na.rm=TRUE)))

fit <- manova(cbind(mb, bh, bl, nh) ~ data$epoch, data)

```