

STK1110 - Statistiske metoder og dataanalyse høsten 2022

Løsningsforslag til ekstraoppgave 3.1

Ekstraoppgave 3.1

a) Vi finner først den kumulative fordelingen til Y_i -ene:

$$F_Y(y; \theta) = P(Y_i \leq y) = P(g(X_i) \leq y) = P(h(g(X_i)) \leq h(y)) = P(X_i \leq h(y)) = F_X(h(y); \theta).$$

Vi derivere og finner at tettheten til Y_i -ene er gitt ved:

$$f_Y(y; \theta) = F'_Y(y; \theta) = F'_X(h(y); \theta)h'(y) = f_X(h(y); \theta)h'(y)$$

Likelihooden basert på X_i -ene er nå

$$\mathcal{L}_X(\theta) = \prod_{i=1}^n f_X(x_i; \theta),$$

og likelihooden basert på Y_i -ene er

$$\mathcal{L}_Y(\theta) = \prod_{i=1}^n f_Y(y_i; \theta) = \prod_{i=1}^n f_X(h(y_i); \theta)h'(y_i) = \left(\prod_{i=1}^n f_X(x_i; \theta) \right) \left(\prod_{i=1}^n h'(y_i) \right) = \left(\prod_{i=1}^n h'(y_i) \right) \mathcal{L}_X(\theta).$$

Vi merker oss at $\prod_{i=1}^n h'(y_i)$ er en konstant med hensyn til parameteren θ . Derfor er de to likelihoodene $\mathcal{L}_Y(\theta)$ og $\mathcal{L}_X(\theta)$ proporsjonale. Det betyr at θ som maksimerer $\mathcal{L}_X(\theta)$ også er den θ som maksimerer $\mathcal{L}_Y(\theta)$. Det vil si at maximum likelihood estimatoren basert på X_i -ene og maximum likelihood estimatoren basert på Y_i -er vil være den samme.

b) For normalfordelingen har vi at $E[X] = \mu$ og $E[X^2] = V(X) + (E[X])^2 = \sigma^2 + \mu^2$.

Momentestimatorene er derfor gitt som løsningene av likningene

$$\begin{aligned} \bar{X} &= \hat{\mu} \\ \frac{1}{n} \sum_{i=1}^n X_i^2 &= \hat{\sigma}^2 + \hat{\mu}^2 \end{aligned}$$

Det gir

$$\begin{aligned} \hat{\mu} &= \bar{X} \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \end{aligned}$$

c) Fra eksempel 7.18 i læreboka har vi at maximum likelihood estimatorene $\hat{\mu}$ og $\hat{\sigma}^2$ er lik momentestimatorene fra punkt b). I punkt a) har vi vist at strengt voksende transformasjoner av X_i -ene

ikke endrer maximum likelihood estimatorene. Altså er maximum likelihood estimatorene for μ og σ^2 basert på Y_i -ene også lik momentestimatorene fra punkt b).

d) Vi har at $Y_i = \exp(X_i)$. Derfor er

$$\begin{aligned} E[Y_i] &= E[e^{X_i}] = M_X(1) = \exp[\mu + \sigma^2/2] \\ E[Y_i^2] &= E[e^{2X_i}] = M_X(2) = \exp[2\mu + 2\sigma^2] \end{aligned}$$

e) Ved å bruke resultatet i punkt d), får vi at momentestimatorene basert på Y_i -ene er løsningen av likningene

$$\begin{aligned} \bar{Y} &= \exp[\hat{\mu} + \hat{\sigma}^2/2] \\ \frac{1}{n} \sum_{i=1}^n Y_i^2 &= \exp[2\hat{\mu} + 2\hat{\sigma}^2] \end{aligned}$$

Vi tar logaritmen på begge sider av likningene over og får likningene:

$$\begin{aligned} \ln(\bar{Y}) &= \hat{\mu} + \frac{\hat{\sigma}^2}{2} \\ \ln\left(\frac{1}{n} \sum_{i=1}^n Y_i^2\right) &= 2\hat{\mu} + 2\hat{\sigma}^2 \end{aligned}$$

Her multipliserer vi den første likningen med 2 og trekker den fra den andre likningen. Da får vi

$$\hat{\sigma}^2 = \ln\left(\frac{1}{n} \sum_{i=1}^n Y_i^2\right) - 2 \ln(\bar{Y})$$

Så setter vi dette inn i den første likningen og får:

$$\hat{\mu} = 2 \ln \bar{Y} - \frac{1}{2} \ln\left(\frac{1}{n} \sum_{i=1}^n Y_i^2\right)$$

f)

```

1 > # f)
2 > # For reproduserbarhet angit vi seed:
3 > set.seed(1)
4 > # Simuler data
5 > x = rnorm(n=10, mean=1, sd=sqrt(3))
6 > y = exp(x)
7 >
8 > # Lag en ramme for å skrive inn resultatet.
9 > result.matrix = as.data.frame(matrix(NA, nrow=2, ncol=2))
10 > colnames(result.matrix) = c("mu", "sigma.sq")
11 > rownames(result.matrix) = c("X", "Y")
12 >
13 > # Momentestimatorene basert på X.
14 > result.matrix["X", "mu"] = mean(x)
15 > result.matrix["X", "sigma.sq"] = mean(x^2) - mean(x)^2
16 >
17 > # Momentestimatorene basert på Y.
18 > result.matrix["Y", "mu"] = 2*log(mean(y)) - (1/2)*log(mean(y^2))
19 > result.matrix["Y", "sigma.sq"] = log(mean(y^2)) - 2*log(mean(y))
20 >
21 > # Vis resultatet.
22 > print(result.matrix)
23      mu sigma.sq
24 X 1.228982 1.645149
25 Y 1.460774 1.209689

```

Vi ser at momentestimatene basert på X_i -ene og Y_i -ene er forskjellige.

g)

```
1 > # g)
2 > repetition = 1000
3 >
4 > # Lag en ramme for å skrive inn resultatet.
5 > result.array = array(data=NA, dim=c(2,2,repetition))
6 > dimnames(result.array) = list(
7 +   c("X", "Y"),
8 +   c("mu", "sigma.sq"),
9 +   paste("rept.", 1:repetition, sep="")
10 + )
11 >
12 > # Gjenta simulasjonen 1000 ganger.
13 > for (i in 1:repetition) {
14 +
15 +   # Simuler data.
16 +   x = rnorm(n=10, mean=1, sd=sqrt(3))
17 +   y = exp(x)
18 +
19 +   # Momentestimatorene basert på X.
20 +   result.array["X","mu",i] = mean(x)
21 +   result.array["X","sigma.sq",i] = mean(x^2) - mean(x)^2
22 +
23 +   # Momentestimatorene basert på Y.
24 +   result.array["Y","mu",i] = 2*log(mean(y)) - (1/2)*log(mean(y^2))
25 +   result.array["Y","sigma.sq",i] = log(mean(y^2)) - 2*log(mean(y))
26 + }
27 >
28 > # Skjevhet
29 > bias.matrix = apply(X=result.array, MARGIN=1:2, FUN=mean)
30 > bias.matrix[,"mu"] = bias.matrix[,"mu"] - 1
31 > bias.matrix[,"sigma.sq"] = bias.matrix[,"sigma.sq"] - 3
32 > print(bias.matrix)
33       mu      sigma.sq
34 X -0.0122316 -0.2390005
35 Y  0.5920356 -1.8587630
36 >
37 > # Standardavvik
38 > sd.matrix = apply(X=result.array, MARGIN=1:2, FUN=sd)
39 > print(sd.matrix)
40       mu      sigma.sq
41 X 0.5602708 1.2738232
42 Y 0.6664552 0.4239948
```

For μ viser simuleringene at estimatoren basert på X_i -ene både har mindre skjevhet og standardfeil enn den basert på Y_i -ene.

For σ^2 har estimatoren basert på X_i -ene mindre skjevhet enn den basert på Y_i -ene. Men estimatoren basert på Y_i -er har mindre standardfeil enn den basert på X_i -ene.

h)

```
1 > # h)
2 > repetition = 1000
3 >
4 > # Lag en ramme for å skrive inn resultatet.
5 > result.array.2 = array(data=NA, dim=c(2,2,repetition))
6 > dimnames(result.array.2) = list(
7 +   c("X", "Y"),
8 +   c("mu", "sigma.sq"),
9 +   paste("rept.", 1:repetition, sep="")
10 + )
11 >
```

```

12 > # Gjenta simulasjonen 1000 ganger.
13 > for (i in 1:repetition) {
14 +
15 +   # Simuler data.
16 +   x = rnorm(n=1000, mean=1, sd=sqrt(3))
17 +   y = exp(x)
18 +
19 +   # Momentestimatorene basert på X.
20 +   result.array.2["X","mu",i] = mean(x)
21 +   result.array.2["X","sigma.sq",i] = mean(x^2) - mean(x)^2
22 +
23 +   # Momentestimatorene basert på Y.
24 +   result.array.2["Y","mu",i] = 2*log(mean(y)) - (1/2)*log(mean(y^2))
25 +   result.array.2["Y","sigma.sq",i] = log(mean(y^2)) - 2*log(mean(y))
26 + }
27 >
28 > # Skjevhet
29 > bias.matrix = apply(X=result.array.2, MARGIN=1:2, FUN=mean)
30 > bias.matrix[, "mu"] = bias.matrix[, "mu"] - 1
31 > bias.matrix[, "sigma.sq"] = bias.matrix[, "sigma.sq"] - 3
32 > print(bias.matrix)
33           mu      sigma.sq
34 X 0.0005411164 -0.002528499
35 Y 0.2221945841 -0.463651350
36 >
37 > # Standardavvik
38 > sd.matrix = apply(X=result.array.2, MARGIN=1:2, FUN=sd)
39 > print(sd.matrix)
40           mu      sigma.sq
41 X 0.05388164 0.1356096
42 Y 0.19488271 0.5258223

```

Vi har økt n fra 10 til 1000. For både μ og σ^2 har nå estimatorene basert på X_i -ene mindre skjevhet og standardfeil enn estimatorene basert på Y_i -ene.

i) Vi plottes histogram og qq-plott med følgende R-kode:

```

1 # i)
2 # Plot
3 pdf(file = paste("./plot/n10.pdf", sep=""), width = 6, height = 12)
4 par(mfrow=c(4,2))
5
6 # Histogram
7 hist(result.array["X","mu",], breaks=150, font.main=1)
8 # Display the correct value.
9 abline(v=1, col="red", lty=2, lwd=1)
10 # QQ plot
11 qqnorm(result.array["X","mu",])
12 qqline(result.array["X","mu",])
13
14 # Histogram
15 hist(result.array["Y","mu",], breaks=150, font.main=1)
16 # Display the correct value.
17 abline(v=1, col="red", lty=2, lwd=1)
18 # QQ plot
19 qqnorm(result.array["Y","mu",])
20 qqline(result.array["Y","mu",])
21
22 # Histogram
23 hist(result.array["X","sigma.sq",], breaks=150, font.main=1)
24 # Display the correct value.
25 abline(v=3, col="red", lty=2, lwd=1)
26 # QQ plot
27 qqnorm(result.array["X","sigma.sq",])
28 qqline(result.array["X","sigma.sq",])
29
30 # Histogram
31 hist(result.array["Y","sigma.sq",], breaks=150, font.main=1)
32 # Display the correct value.

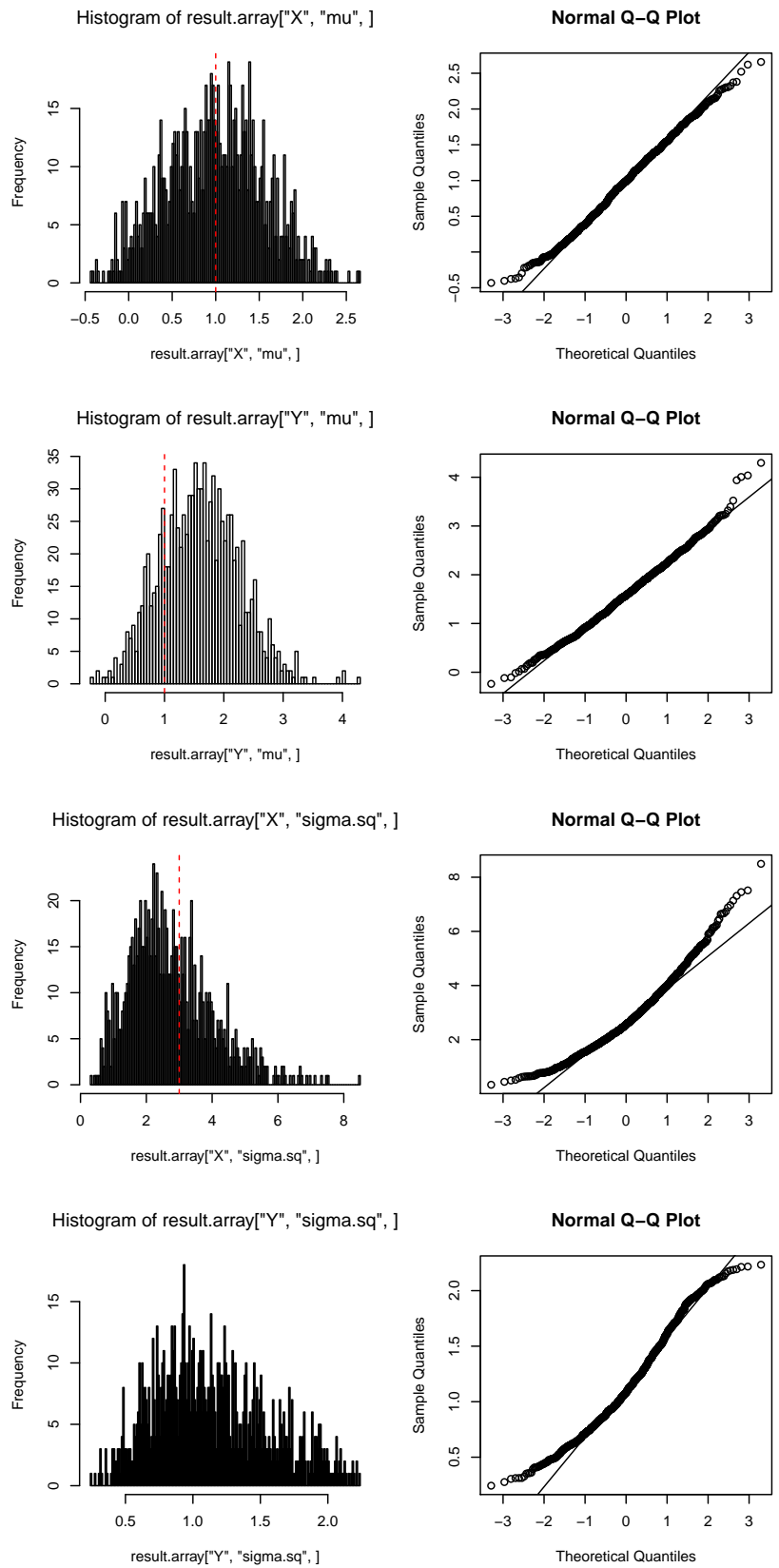
```

```

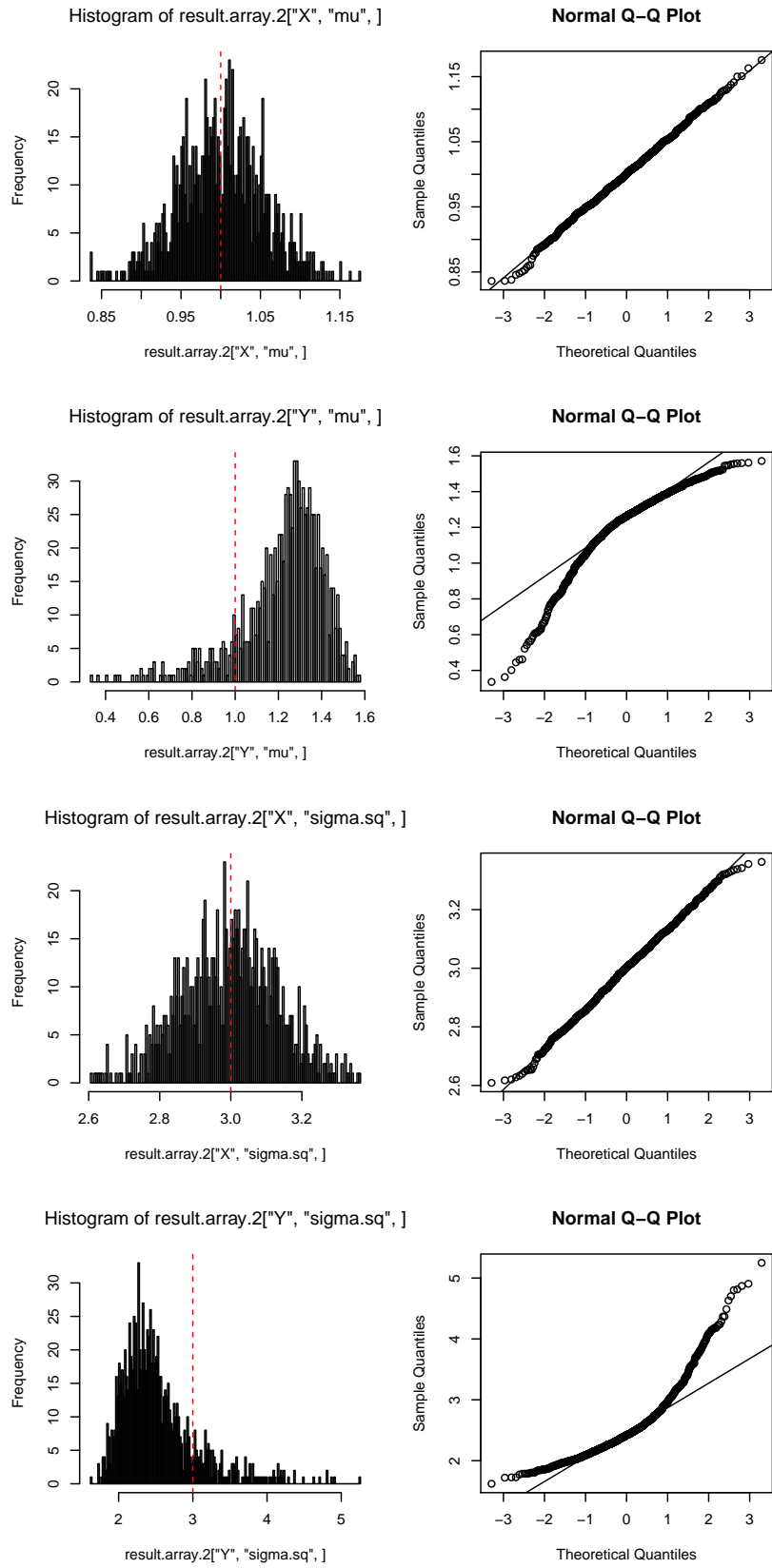
33 abline(v=3, col="red", lty=2, lwd=1)
34 # QQ plot
35 qqnorm(result.array["Y", "sigma.sq",])
36 qqline(result.array["Y", "sigma.sq",])
37
38 dev.off()
39
40
41 # Plot
42 pdf(file = paste("./plot/n1000.pdf", sep=""), width = 6, height = 12)
43 par(mfrow=c(4,2))
44
45 # Histogram
46 hist(result.array.2["X", "mu",], breaks=150, font.main=1)
47 # Display the correct value.
48 abline(v=1, col="red", lty=2, lwd=1)
49 # QQ plot
50 qqnorm(result.array.2["X", "mu",])
51 qqline(result.array.2["X", "mu",])
52
53 # Histogram
54 hist(result.array.2["Y", "mu",], breaks=150, font.main=1)
55 # Display the correct value.
56 abline(v=1, col="red", lty=2, lwd=1)
57 # QQ plot
58 qqnorm(result.array.2["Y", "mu",])
59 qqline(result.array.2["Y", "mu",])
60
61 # Histogram
62 hist(result.array.2["X", "sigma.sq",], breaks=150, font.main=1)
63 # Display the correct value.
64 abline(v=3, col="red", lty=2, lwd=1)
65 # QQ plot
66 qqnorm(result.array.2["X", "sigma.sq",])
67 qqline(result.array.2["X", "sigma.sq",])
68
69 # Histogram
70 hist(result.array.2["Y", "sigma.sq",], breaks=150, font.main=1)
71 # Display the correct value.
72 abline(v=3, col="red", lty=2, lwd=1)
73 # QQ plot
74 qqnorm(result.array.2["Y", "sigma.sq",])
75 qqline(result.array.2["Y", "sigma.sq",])
76
77 dev.off()

```

Vi kan se fra figurene på de neste sidene at estimatorene basert på X_i -ene har mindre skjevhet enn estimatorene basert på Y_i -ene. Når vi får flere observasjoner blir normaltilnærmingen bedre for estimatorene basert på X_i -ene.



Figur 1: Resultat fra ekstraoppgave 7 (g) ($n = 10$)



Figur 2: Resultat fra ekstraoppgave 7 (h) ($n = 1000$)