

**STK1110**  
**Statistiske metoder og dataanalyse**

**OBLIG 1**

Egil Furnes  
Student: 693784

## Problem 1

a)

Vi skal utlede moment-estimatorene for  $\alpha$  og  $\gamma$ .

Vi har for  $X \sim \text{Gamma}(\alpha, \gamma)$ :

$$\mathbb{E}[X] = \frac{\alpha}{\gamma} \quad V[X] = \frac{\alpha}{\gamma^2}$$

Setter lik  $\bar{x}$  og  $s^2$ :

$$\hat{\gamma} = \frac{\bar{x}}{s^2} \quad \hat{\alpha} = \frac{\bar{x}^2}{s^2}$$

Vi kan da bruke en parameterisering med skala  $\beta = 1/\gamma$ :

$$\hat{\alpha} = \frac{\bar{x}}{s^2} \quad \hat{\beta} = \frac{s^2}{\bar{x}}$$

```

1 x <- scan("data/forsikringskrav.txt")
2 xb <- mean(x)
3 s2 <- mean((x-xb)^2)
4 alpha <- xb^2/s2
5 gamma <- xb/s2
6 beta <- s2/xb

```

```

1 > xb
2 [1] 24.13847
3 > s2
4 [1] 836.3836
5 > alpha
6 [1] 0.696649
7 > gamma
8 [1] 0.02886053
9 > beta
10 [1] 34.6494

```

Følgende har vi funnet:

$$\bar{x} = 24.14 \quad s^2 = 836.38 \quad \hat{\alpha} = 0.69 \quad \hat{\gamma} = 0.0289 \quad \hat{\beta} = 34.65$$

b)

For i.i.d.  $x_1, \dots, x_n \sim \text{Gamma}(\alpha, \gamma)$  med tethet  $f(x; \alpha, \gamma) = \frac{\gamma^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\gamma x}$   $\alpha, \gamma > 0$  er log-likelihood-funksjonen:

$$l(\alpha, \gamma) = n\alpha \log \gamma - n \log \Gamma(\alpha) + (\alpha - 1) \sum_{i=1}^n \log x_i - \gamma \sum_{i=1}^n x_i$$

Vi har også  $\hat{\alpha} = \bar{x}^2/s^2$  og  $\hat{\gamma} = \bar{x}/s^2$  som gir:

```

1 n <- length(x)
2 logL <- n*alpha*log(gamma) - n*lgamma(alpha) + (alpha-1)*sum(log(x))
      - gamma*sum(x)

```

```

1 > n
2 [1] 6377
3 > logL
4 [1] -27264.51

```

Med dette har vi funnet at log-likelihood verdien er:

$$\log L = -27264.51$$

c)

Fra tidligere har vi uttrykket for  $\log L$

$$l(\alpha, \gamma) = n\alpha \log \gamma - n \log \Gamma(\alpha) + (\alpha - 1) \sum_{i=1}^n \log x_i - \gamma \sum_{i=1}^n x_i$$

Nå kan vi derivere med hensyn på  $\gamma$  og sette lik 0 for å finne maksimum likelihood.

$$\begin{aligned} \frac{\partial l}{\partial \gamma} &= \frac{n\alpha}{\gamma} - \sum_{i=1}^n x_i \\ \Rightarrow \hat{\gamma}(\alpha) &= \frac{n\alpha}{\sum_i x_i} = \frac{\alpha}{\bar{x}} \end{aligned}$$

For å sammenligne med moment-estimatoren for  $\gamma$  for gitt  $\alpha$ :

$$\bar{x} = \frac{\alpha}{\gamma} \quad \Rightarrow \quad \gamma = \frac{\alpha}{\bar{x}}$$

Som er samme uttrykk som

$$\hat{\gamma}_{MLE}(\alpha) = \hat{\gamma}_{MOM}(\alpha) = \frac{\alpha}{\bar{x}}$$

Kan også sjekke at det er et maksimum med

$$\frac{\partial^2 l}{\partial \gamma^2} = -\frac{n\alpha}{\gamma^2} < 0 \quad \text{for } \alpha, \gamma > 0$$

d)

For denne funksjonen bør vi bruke  $\log(\alpha)$  er for å sikre at  $\alpha > 0$  uten å måtte begrense optimeringen på noe vis.

```

1 alpha0 <- xb^2/s2
2 gamma0 <- xb/s2
3
4 negloglikgamma <- function(logalpha, x = x)
5 {
6   n <- length(x)
7   alpha <- exp(logalpha)
8   gamma <- alpha / mean(x) # MLE for gamma gitt alpha
9   logL <- n*alpha*log(gamma) - n*lgamma(alpha) +
10    (alpha - 1)*sum(log(x)) - gamma*sum(x)
11   -logL                         # returnér negativ logL
12 }
13
14 fit.ml <- optim(log(alpha0), negloglikgamma, x = x, method = "BFGS")
15
16 alpha_mle <- exp(fit.ml$par)
17 gamma_mle <- alpha_mle / xb
18 logL_mle <- -fit.ml$value

```

```

1 > alpha_mle
2 [1] 1.388346
3 > gamma_mle
4 [1] 0.05751589
5 > logL_mle
6 [1] -26488.28
7 > logL
8 [1] -27264.51

```

Her finner vi likelihood-verdiene som ser rimelige ut i forhold til den vi fant i b).

$$\log L = -27264.51 \quad \log L_{MLE} = -26488.28$$

e)

```

1 set.seed(1705)
2 B <- 2000
3 boot <- replicate(B, {
4   xb_ <- mean(xb <- sample(x, n, TRUE))
5   s2_ <- mean((xb - xb_)^2)
6   a0 <- xb_^2 / s2_
7   a <- exp(optim(log(a0), negloglikgamma, x = xb, method = "BFGS")$par)
8   c(alpha = a, gamma = a / xb_)
9 })
10
11 alpha_se <- sd(boot["alpha", ])
12 gamma_se <- sd(boot["gamma", ])
13 alpha_ci <- quantile(boot["alpha", ], c(.025, .975))
14 gamma_ci <- quantile(boot["gamma", ], c(.025, .975))

```

```

1 > alpha_se
2 [1] 0.03080641
3 > gamma_se
4 [1] 0.002002218
5 > alpha_ci
6      2.5%    97.5%
7 1.331434 1.449643
8 > gamma_ci
9      2.5%    97.5%
10 0.05383500 0.06154386

```

Med dette finner vi følgende 95% konfidens-intervaller for  $\alpha$  og  $\gamma$

$$\begin{aligned}\hat{\alpha}_{CI} &\in [1.331, 1.450] \\ \hat{\gamma}_{CI} &\in [0.0538, 0.0615]\end{aligned}$$

f)

```

1 set.seed(1705)
2 B <- 2000
3 mu_hat <- mean(x)
4 mu_boot <- replicate(B, mean(sample(x, n, TRUE)))
5 mu_ci <- quantile(mu_boot, c(.025, .975))
6
7 mu_hat
8 mu_ci
9
10 # decisions via CI inclusion of 25
11 rej_05 <- !(25 >= mu_ci[1] & 25 <= mu_ci[2])
12 rej_01 <- !(25 >= quantile(mu_boot, c(.005, .995))[1] &
13                  25 <= quantile(mu_boot, c(.005, .995))[2])
14 rej_05
15 rej_01

```

```

1 > rej_05
2 2.5%
3 TRUE
4 > rej_01
5 0.5%
6 FALSE

```

Vi ser det at `rej_05` gir *TRUE* men `rej_01` gir *FALSE*, og dermed er det naturlig å tro at P-veriden for testen er et sted  $0.01 < p < 0.05$ .

## Problem 2

a)

Vi finner et 95% konfidensintervall

$$\bar{X} \pm t_{0.975,14} \frac{S}{\sqrt{15}}$$

```

1 x <- c(525,587,547,558,591,531,571,551,566,622,561,502,556,565,562)
2 n <- length(x)
3 xb <- mean(x)
4 s <- sd(x)
5 tcrit <- qt(.975, df = n-1)
6 lo <- xb - tcrit * s / sqrt(n)
7 hi <- xb + tcrit * s / sqrt(n)

```

```

1 > xb
2 [1] 559.6667
3 > s
4 [1] 28.55988
5 > lo
6 [1] 543.8507
7 > hi
8 [1] 575.4826

```

Med dette finner vi

$$\bar{x} \approx 559.67, \quad S \approx 28.56, \quad \mu \in [543.85, 575.48]$$

b)

I denne oppgaven genererer jeg 1000 datasett med størrelse  $n = 15$  bestående av stokastiske variabler  $X_1, \dots, X_{15} \stackrel{\text{uif}}{\sim} N(558, 30^2)$  med `rnorm()`.

```

1 set.seed(1705)
2 B <- 10000
3 n <- 15
4 mu <- 558
5 sigma <- 30
6
7 cov <- replicate(B, {
8   y <- rnorm(n, mu, sigma)
9   xb <- mean(y); s <- sd(y)
10  tcrit <- qt(.975, df = n - 1)
11  lo <- xb - tcrit * s / sqrt(n); hi <- xb + tcrit * s / sqrt(n)
12  (mu >= lo) & (mu <= hi)
13 })
14

```

```

15 cover <- mean(cov)
16 se_mc <- sqrt(cover * (1 - cover) / B)

```

Fra utskriften nedenfor finner vi at simuleringen ga en dekning som nærmer seg veldig det teoretiske konfidensintervallet  $0.9534 \approx 0.95$ . Dette er å forvente i en Monte Carlo simulering, ettersom vi har introdusert ikke-reduserbart støy med en standardfeil på omtrent 0.0021.

```

1 > cover
2 [1] 0.9534
3 > se_mc
4 [1] 0.002107805

```

Andelen  $\approx 95\%$  av disse simuleringene inneholder den sanne verdien  $\mu = 558$  som stemmer overens med forventet antocyaninnhold for blåbær.

c)

I denne deloppgaven benytter vi heller det tilnærmede utvalget

$$\left( \bar{X} - 1.96 \frac{S}{\sqrt{15}}, \quad \bar{X} + 1.96 \frac{S}{\sqrt{15}} \right)$$

med

$$\bar{X} = \frac{1}{15} \sum_{i=1}^{15} X_i \quad S^2 = \frac{1}{15-1} \sum_{i=1}^{15} (X_i - \bar{X})^2$$

```

1 set.seed(1705)
2 cov_z <- replicate(B, {
3   y <- rnorm(n, mu, sigma)
4   xb <- mean(y)
5   s <- sd(y)
6   lo <- xb - 1.96 * s / sqrt(n)
7   hi <- xb + 1.96 * s / sqrt(n)
8   (mu >= lo) & (mu <= hi)
9 })
10 cover_z <- mean(cov_z)
11 se_mc_z <- sqrt(cover_z * (1 - cover_z) / B)

```

Med denne fremgangsmåten finner vi at 93.34% av intervallene inneholder  $\mu = 558$ , i motsetning til forrige metode som ga 95.34%. Med dette kan vi konkludere at å bruke 1.96 heller enn  $t_{0.975,14}$  blir et for smalt intervall.

```

1 > cover_z; se_mc_z
2 [1] 0.9334
3 [1] 0.00249328

```

**d)**

Nå ser vi heller på  $\sigma$  og hvor mange intervaller som inneholder  $\sigma = 30$ .

```

1 set.seed(1705)
2 cov_sigma <- replicate(B, {
3   y <- rnorm(n, mu, sigma)
4   s <- sd(y)
5   lo <- sqrt((n - 1) * s^2 / qchisq(.975, df = n - 1))
6   hi <- sqrt((n - 1) * s^2 / qchisq(.025, df = n - 1))
7   (sigma >= lo) & (sigma <= hi)
8 })
9 cover_sigma <- mean(cov_sigma)
10 se_mc_sigma <- sqrt(cover_sigma * (1 - cover_sigma) / B)

```

Med denne fremgangsmåten finner vi at 94.89% av intervallene inneholder  $\sigma = 30$  med et standardavvik på omrent 0.002.

```

1 > cover_sigma; se_mc_sigma
2 [1] 0.9489
3 [1] 0.002202017

```

**e)**

Nå gjentar vi oppgave b) med de nye datasettene, med antakelse om at  $Z_1, \dots, Z_{15} \stackrel{\text{uif}}{\sim} t_7$ .

```

1 set.seed(1705)
2 cov_t7 <- replicate(B, {
3   z <- rt(n, df = 7)
4   y <- mu + sigma * z
5   xb <- mean(y)
6   s <- sd(y)
7   tcrit <- qt(.975, df = n - 1)
8   lo <- xb - tcrit * s / sqrt(n); hi <- xb + tcrit * s / sqrt(n)
9   (mu >= lo) & (mu <= hi)
10 })
11 cover_t7 <- mean(cov_t7)
12 se_mc_t7 <- sqrt(cover_t7 * (1 - cover_t7) / B)

```

Med dette finner vi at dekningen er 95.31% altså tilbake på statistisk signifikansnivå 5%. Vi ser det at selv om data ikke er normalfordelte her, men heller  $t_7$ -fordelt med tyngre haler er middelverdien relativt robust allerede for  $n = 15$ .

```

1 > cover_t7; se_mc_t7
2 [1] 0.9531
3 [1] 0.002114247

```

f)

Vi lager nå et konfidensintervall for  $\tilde{\sigma}$  og undersøke andelen som inneholder denne gitt

$$V(Z_i) = \frac{7}{7 - 2}$$

$$\tilde{\sigma}^2 = V(X_i) = V(\mu + \sigma Z_i) = \sigma^2 V(Z_i) = 1.4\sigma^2$$

```

1 set.seed(1705)
2 tilde_sigma <- sqrt(1.4) * sigma
3
4 cov_t7_sigma <- replicate(B, {
5   y <- mu + sigma * rt(n, df = 7)
6   s <- sd(y)
7   lo <- sqrt((n - 1) * s^2 / qchisq(.975, df = n - 1))
8   hi <- sqrt((n - 1) * s^2 / qchisq(.025, df = n - 1))
9   (tilde_sigma >= lo) & (tilde_sigma <= hi)
10 })
11
12 cover_t7_sigma <- mean(cov_t7_sigma)
13 se_mc_t7_sigma <- sqrt(cover_t7_sigma * (1 - cover_t7_sigma) / B)

```

I dette tilfellet får vi dekning for  $\tilde{\sigma}$  på 89.16%  $\ll 95\%$  og samtidig standardavvik på 0.0031. Nå ser vi langt dårligere dekning enn for i d). Dette kan være siden  $t_7$  ikke er fordelt som  $\chi^2$  og intervallene blir igjen for smale.

```

1 > cover_t7_sigma; se_mc_t7_sigma
2 [1] 0.8916
3 [1] 0.003108849

```