

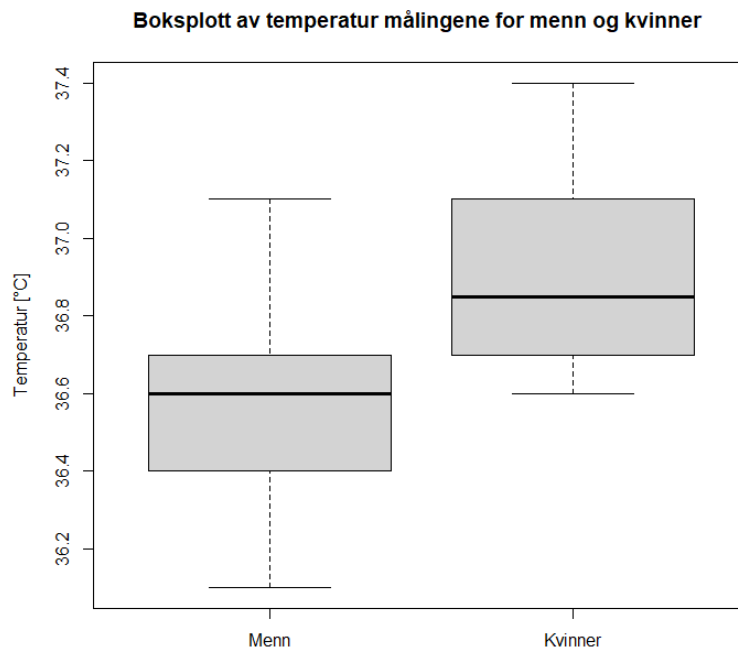
# Oblig 2 - STK1110

Kevin Alexander Aslesen

November 2022

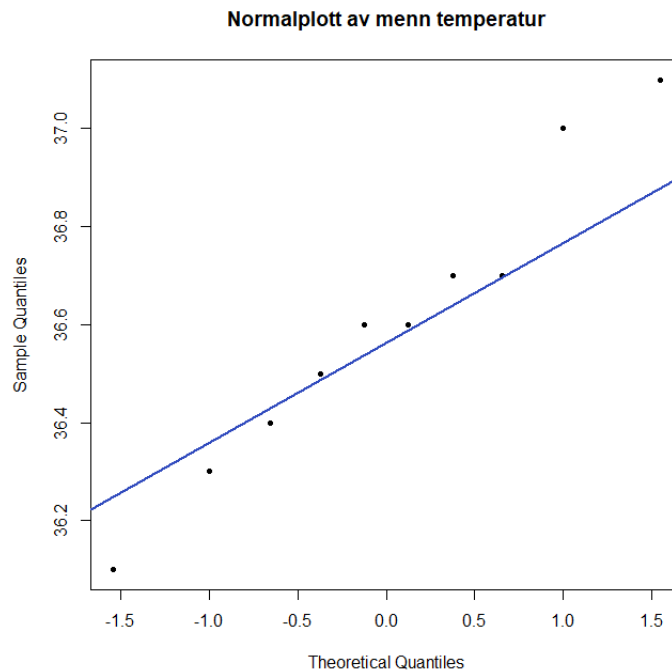
## Oppgave 1

a) BoksploTT av de målte temperaturene er vist i Figur 1. Vi ser fra boksploTT-



Figur 1: *BoksploTT for målingene av kroppstemperaturen til menn og kvinner.*

tene at kroppstemperaturene for kvinner ligger litt høyere enn for menn. Den største temperaturen for menn tilsvarer øvre kvartil for kvinnene, mens den laveste verdien for kvinne temperaturene tilsvarer gjennomsnitts temperaturen for menn. Ved å bare se på disse plottene alene så kan man fort prøve å konkludere med at kvinner har høyere kroppstemperatur enn menn. Dette skal vi se nærmere på.



Figur 2: Normalfordelingsplott for kroppstemperaturene til menn.

b) Normalfordelingsplott for målingene av temperaturen for menn og kvinner er vist hver for seg i henholdsvis Figur 2 og Figur 3. Vi ser fra begge normalplottene at vi får en ganske rett linje, som tyder på at målingene kommer fra en normalfordelt populasjon.

c) Vi vil teste  $H_0 : \mu_1 - \mu_2 = 0$  mot  $H_A : \mu_1 - \mu_2 \neq 0$  med signifikantnivå 5% ( $\alpha = 0.05$ ). Hvis variansen ikke er det samme for de to utvalgene, så benytter vi testobservatoren:

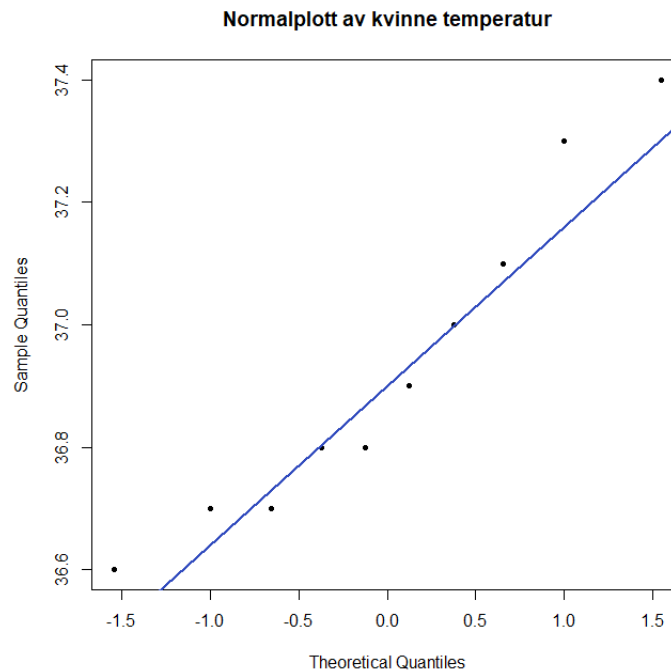
$$t_p = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{s_p^2/m + s_p^2/n}} = \frac{\bar{x} - \bar{y}}{\sqrt{s_p^2/m + s_p^2/n}}$$

hvor  $s_p^2$  er den kombinerte estimatoren for variansen  $\sigma^2$  gitt ved:

$$s_p^2 = \frac{m-1}{m+n-2} s_1^2 + \frac{n-1}{m+n-2} s_2^2$$

Siden  $m = n$  i vårt tilfelle, så reduseres dette til:

$$s_p^2 = \frac{s_1^2 + s_2^2}{2}$$



Figur 3: *Normalfordelingsplott for kroppstemperaturene til kvinner.*

hvor  $s_1$  og  $s_2$  er det målte standardavviket til temperaturen hos menn og kvinner, respectfully. Vi antar først at variansen er det samme for de to utvalgene, så vil lar  $s_p = s_1$ . Benytter da følgende kode for å finne denne testobservasjonen,  $t_p$ , gitt målingene våre:

```

1 > temp_menn = c(36.1, 36.3, 36.4, 36.6, 36.6, 36.7, 36.7, 37.0,
2   36.5, 37.1)
3 > temp_kvin = c(36.6, 36.7, 36.8, 36.8, 36.7, 37.0, 37.1, 37.3,
4   36.9, 37.4)
5 >
6 > m = length(temp_menn)
7 > n = length(temp_kvin)
8 >
9 > x.bar = mean(temp_menn)
10 > y.bar = mean(temp_kvin)
11 >
12 > s1 = sd(temp_menn)
13 > s2 = sd(temp_kvin)
14 > s_p = s1
15 >
16 > t_p = (x.bar - y.bar) / (s_p * sqrt(1/m + 1/n))
17 > print(t_p)
18 [1] -2.444631

```

Vi får altså  $t_p = -2.44$ . Vi finner så ut om denne verdien er tilstrekkelig nok

for å forkaste  $H_0$  eller ikke ved å finne forkastningsområdet. Den kritiske  $t$ -verdien for  $\alpha/2 = 0.025$  og  $m + n - 2 = 10 + 10 - 2 = 18$  frihetsgrader finner vi fra Tabell A.6 i læreboka til være  $t_{\alpha/2, m+n-2} = t_{0.025, 18} = 2.101$ . Siden vi benytter en tosidig test, så er forkastningsområdet da lik

$$t_p \leq -2.101 \quad \text{og} \quad 2.101 \leq t_p$$

Ettersom  $t_p = -2.44 < -2.101$ , så forkaster vi da  $H_0$  med signifikansnivå  $\alpha = 0.05$ . Testen vår konkluderer da med at det er forskjell mellom temperaturen til kvinner og menn. P-verdien finner vi ved hjelp av Tabell A.7 i læreboka, hvor vi finner arealet under grafen etter 2.44 og multipliserer dette med to (siden vi har en tosidet test):

$$\begin{aligned} \text{P-verdi} &= 2[1 - \Phi_t(2.44)] \\ &= 2[0.014] \\ &= 0.028 \end{aligned}$$

Vi ser at P-verdi  $= 0.028 < 0.05 = \alpha$ , som også resulter i å forkaste  $H_0$ . Lar vi heller  $s_p = s_2$ , så får vi  $t_p = -2.76$  og P-verdi  $= 0.012$ , som også resulterer i forkastelsen av  $H_0$ . Et 95% konfidensintervall for temperaturforskjellen er:

$$\begin{aligned} &(\bar{x} - \bar{y}) \pm t_{0.025, 18} \cdot s_1 \sqrt{1/m + 1/n} \\ &(36.60 - 36.93) \pm 2.101 \cdot 0.30 \sqrt{1/10 + 1/10} \\ &-0.330 \pm 0.263 \end{aligned}$$

Dette tilsvare til intervallet  $(-0.593, -0.067)$ . Vi ser at intervallet ikke inneholder 0, som betyr at vi i dette tilfellet også konkluderer med at kroppstemperaturen mellom menn og kvinner er forskjellig ved å forkaste  $H_0$ .

**d)** I dette tilfellet antar vi forskjellige varianser, så testobservatoren blir nå:

$$t = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{s_1^2/m + s_2^2/n}} = \frac{\bar{x} - \bar{y}}{\sqrt{s_1^2/m + s_2^2/n}}$$

Dette får vi til å bli

```
1 > t = (x.bar - y.bar) / sqrt(s1^2/m + s2^2/n)
2 > print(t)
3 [1] -2.590062
```

Altså  $t = -2.59$ . Forkastningsområdet i dette tilfellet er litt annerledes ettersom  $t$  fordelingen er bestemt av Welch's formel for frihetsgradene:

$$\nu = \frac{(s_1^2/m + s_2^2/n)^2}{((s_1^2/m)^2/(m-1) + (s_2^2/n)^2/(n-1))}$$

```

1 > welch = (s1^2/m + s2^2/n)^2 / ((s1^2/m)^2 / (m - 1) + (s2^2/n
  )^2 / (n - 1))
2 > print(welch)
3 [1] 17.7338

```

For våre målinger får vi altså  $\nu = 17.7$  som vi runder ned til 17. Da er den kritiske t-verdien  $t_{0.025,17} = 2.110$ , som gir oss forkastningsområdet:

$$t \geq t_{0.025,17} = 2.110$$

$$t \leq -t_{0.025,17} = -2.110$$

Siden  $-2.59 < -2.110$  så forkaster vi  $H_0$  også i dette tilfellet. P-verdien i dette tilfellet er

$$\begin{aligned}
 \text{P-verdi} &= 2[1 - \Phi_t(2.59)] \\
 &= 2[0.009] \\
 &= 0.018
 \end{aligned}$$

Vi har altså at P-verdi = 0.018 < 0.05 =  $\alpha$ , som betyr at vi forkaster  $H_0$ . Sammenlikner dette med `t.test()` funksjonen i R:

```

1 > t.test(temp_menn, temp_kvin)
2
3      Welch Two Sample t-test
4
5 data:  temp_menn and temp_kvin
6 t = -2.5901, df = 17.734, p-value = 0.01863
7 alternative hypothesis: true difference in means is not equal
  to 0
8 95 percent confidence interval:
9  -0.59796699 -0.06203301
10 sample estimates:
11 mean of x mean of y
12   36.60    36.93

```

Vi ser at funksjonen genererer en P-verdi på 0.01863. Dette stemmer godt med det vi regnet, hvor den lille forskjellen skyldes grovere avrundingen i utregningen vår.

e) Tester  $H_0 : \sigma_1 = \sigma_2$  mot  $H_A : \sigma_1 \neq \sigma_2$  ved F-test med signifikansnivå  $\alpha = 0.02$ . Da benytter vi testobservasjonen:

$$f = s_1^2/s_2^2$$

Målingene våre får dette til å bli:

```

1 > s1^2/s2^2
2 [1] 1.279251

```

Forkastningsområdet i dette tilfellet tilsvarer de kritiske F-verdiene hvor arealet under F-grafen er  $0.02/2=0.01$  på hver side. Disse verdiene finner vi i Tabell A.8 i læreboka med frihetsgradene  $\nu_1 = 9$  og  $\nu_2 = 9$ :

$$f \geq F_{0.01,9,9} = 5.35$$

$$f \leq F_{0.99,9,9} = 0.187$$

Testobservasjonen  $f = 1.279$  ligger ikke i forkastningsområdet, så vi forkaster ikke  $H_0 : \sigma_1 = \sigma_2$ . Det er derfor rimelig å anta at variansen mellom utvalgene er like. Sjekker dette mot **var.test()** funksjonen i R:

```
1 > var.test(temp_menn, temp_kvin)
2
3      F test to compare two variances
4
5 data:  temp_menn and temp_kvin
6 F = 1.2793, num df = 9, denom df = 9, p-value = 0.7197
7 alternative hypothesis: true ratio of variances is not equal to
8 1
9 95 percent confidence interval:
10  0.3177479 5.1502577
11 sample estimates:
12 ratio of variances
1.279251
```

Vi ser at R funksjonen regnet ut en P-verdi på 0.7197. Dette ligger over signifikansnivået  $\alpha = 0.02$  som betyr at  $H_0$  ikke blir forkastet i dette tilfellet. Dette stemmer overens med konklusjonen fra vår F-test.

**f)** Et rimelig anslag for en ny måling  $X_{n+1}$  gitt et tilfeldig målingssett  $X_1, X_2, \dots, X_n$  er snittet av målingssettet,  $\bar{X}$ . Dermed er et rimelig anslag for  $X_{11} - Y_{11}$ , med gitt målingssett  $(X_1 - Y_1, X_2 - Y_2, \dots, X_{10} - Y_{10})$  lik snittet (forventningen) av målingssettet,  $E(X - Y) = E(X) - E(Y) = \bar{X} - \bar{Y}$ .

Siden  $X_{11}$ ,  $Y_{11}$ ,  $\bar{X}$  og  $\bar{Y}$  er alle normalfordelt, så er lineær kombinasjonen

$$X_{11} - Y_{11} - (\bar{X} - \bar{Y})$$

også normalfordelt. Med dette kan vi lage et prediksjonsintervall for differansen  $X_{11} - Y_{11}$ , ved å se på  $X - Y$  som ett utvalg. Da finner vi et 95% prediksjonsintervall med  $n - 1 = 9$  frihetsgrader for differansen  $X_{11} - Y_{11}$  med formelen

$$(\bar{x} - \bar{y}) \pm t_{0.025,9} \cdot s \sqrt{1 + 1/n}$$

hvor  $s$  er standardavviket til målingene  $x$  og  $y$  som vi antar er like. Den kritiske t-verdien finner vi til å være  $t_{0.025,9} = 2.262$ . Intervallet beregnes da med koden:

```

1 > x = temp_kvin
2 > y = temp_menn
3 > n = length(x)
4 > s = sd(x)
5 > t.crit = 2.262
6 > x.bar = mean(x)
7 > y.bar = mean(y)
8 > lower = (x.bar - y.bar) - t.crit * s * sqrt(1 + 1/n)
9 > upper = (x.bar - y.bar) + t.crit * s * sqrt(1 + 1/n)
10 > print(c(lower, upper))
11 [1] -0.3031356  0.9631356

```

Med dette kan vi si med 95% sikkerhet at forskjellen  $X_{11} - Y_{11}$  vil befinne seg mellom -0.30 og 0.96. I prediksjonsintervall som dette så får vi sannsynlige verdier for en eventuelt ny måling, sammenlignet med konfidensintervall hvor man får sannsynlige verdier for parametere som er av interesse.

## Oppgave 2

a) I dette tilfellet vil det passe best å se på en parret sammenligning ettersom det gir mening at tvillingene i hvert par ikke er helt uavhengige av hverandre gitt at tvillingene er biologisk tilknyttet. For videre analyse så må vi anta at selve parrene  $(A_1, B_1)$ ,  $(A_2, B_2)$ , ...,  $(A_{31}, B_{31})$  er uavhengige fra hverandre, og at forskjellen mellom parrene er tilnærmet normalfordelt.

b) Vil teste  $H_0 : \mu_D = 0$  mot  $H_A : \mu_D \neq 0$  med signifikansnivå  $\alpha = 0.05$ . Testobservatoren for denne testen er

$$t = \frac{\bar{d} - \Delta_0}{s_D / \sqrt{n}} = \frac{\bar{d}}{s_D / \sqrt{n}}$$

hvor vi har fra tabellen  $\bar{d} = -3.26$  og  $s_D = 8.81$ . Setter inn disse verdiene med  $n = 31$  og får testobservatoren:

$$t = \frac{-3.26}{8.81 / \sqrt{31}} = -2.06$$

Den kritiske t-verdien med  $\alpha/2 = 0.025$  og antall frihetsgrader  $df = n - 1 = 31 - 1 = 30$  finner vi til å være  $t_{0.025, 30} = 2.042$ . Det betyr at forkastningsområdet for denne testen er:

$$t \leq -2.042 \quad \text{og} \quad 2.042 \leq t$$

Ettersom vi fikk  $t = -2.06 < -2.042$  så forkaster vi altså  $H_0$  med signifikansnivå  $\alpha = 0.05$ . Siden testen vår er tosidet, så er den tilhørende P-verdien lik  $2[1 - \Phi_t(2.06)]$  hvor  $\Phi_t$  er arealet til venstre for t fordelingen vår med 30 frihetsgrader. Vi beregner:

```

1 > p = 2*(1-pt(2.06, 30))
2 > print(p)
3 [1] 0.04816536

```

P-verdien er altså 0.048, som er akkurat mindre enn signifikansnivået  $\alpha = 0.05$ . Vi konkluderer dermed også utifra P-verdien at  $H_0$  forkastes. Fra testen vår så betyr dette at det er en forskjell i IQ mellom tvillingene som vokser opp med biologiske foreldrene og de som vokser opp med adoptivforeldrene.

c) Lager et 95% konfidensintervall for  $\mu_D$  med formelen

$$\bar{d} \pm t_{0.025,30} \frac{s_D}{\sqrt{n}}$$

hvor disse faktorene ble bestemt i c). Beregner dette med:

```

1 > n = 31
2 > d.bar = -3.26
3 > s.D = 8.81
4 > t.crit = 2.042
5 > lower = d.bar - t.crit * s.D/sqrt(n)
6 > upper = d.bar + t.crit * s.D/sqrt(n)
7 > print(c(lower, upper))
8 [1] -6.49110298 -0.02889702

```

Med 95% sikkerhet har vi altså at  $\mu_D$  ligger i intervallet (-6.49, -0.03). Vi ser at konfidensintervallet bare inneholder negative verdier, som betyr at 0 ikke befinner seg i intervallet. Dette spiller bra med konklusjonen  $\mu_D \neq 0$  som vi fikk fra hypotesetesten. Bare negative verdier tyder også på at differansen må være negativ, som isåfall betyr at tvillingene med biologiske foreldre har lavere IQ enn de med adoptivforeldre.

Vi ser at i denne oppgaven så konkluderer både den tosidige hypotesetesten og konfidensintervallet at  $\mu_D \neq 0$ . Dette skyldes av at hypotesetesten hadde signifikantnivå 5% og konfidensintervallet var  $(100\% - 5\%) = 95\%$ . Dette gjelder generelt hvor en hypotesetest med nivå  $\alpha$  tilsvarer et  $100(1 - \alpha)\%$  konfidensintervall.

### Oppgave 3

a) Vil teste  $H_0 : p_1 - p_2 = 0$  mot  $H_A : p_1 - p_2 > 0$  med signifikansnivå  $\alpha = 0.05$ . Antar at målingene kommer fra to tilfeldig, uavhengige populasjoner. Vi ser at  $n_1\hat{p}_1 = 486$  og  $n_2\hat{p}_2 = 441$  er begge større enn 10 og at  $\hat{p}_1 = 0.162$  og  $\hat{p}_2 = 0.147$ . Dette betyr at vi kan benytte to-proporsjon z testen med testobservator:

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}(1/m + 1/n)}}$$



hvor  $\hat{p} = \frac{m}{m+n}\hat{p}_1 + \frac{n}{m+n}\hat{p}_2$  og  $\hat{q} = 1 - \hat{p}$ . Vi finner disse og beregner  $z$  med koden:

```
1 > p1 = 0.162
2 > p2 = 0.147
3 > m = 3000
4 > n = 3000
5 >
6 > p = m/(m+n) * p1 + n/(m+n) * p2
7 > q = 1 - p
8 >
9 > z = (p1 - p2) / sqrt(p*q*(1/m + 1/n))
10 > print(z)
11 [1] 1.60737
```

Vi får altså  $z = 1.607$ . Den kritiske  $z$ -verdien med  $\alpha = 0.05$  finner vi til å være  $z_{0.05} = 1.74$ . Siden testen vår er upper-tailed ( $H_A : p_1 - p_2 > 0$ ), så er forkastningsområdet  $z \geq 1.74$ . Ettersom  $1.607 < 1.74$ , forkaster vi ikke  $H_0$ . P-verdien i dette tilfellet tilsvarer  $1 - \Phi(1.607)$ , som vi finner ved:

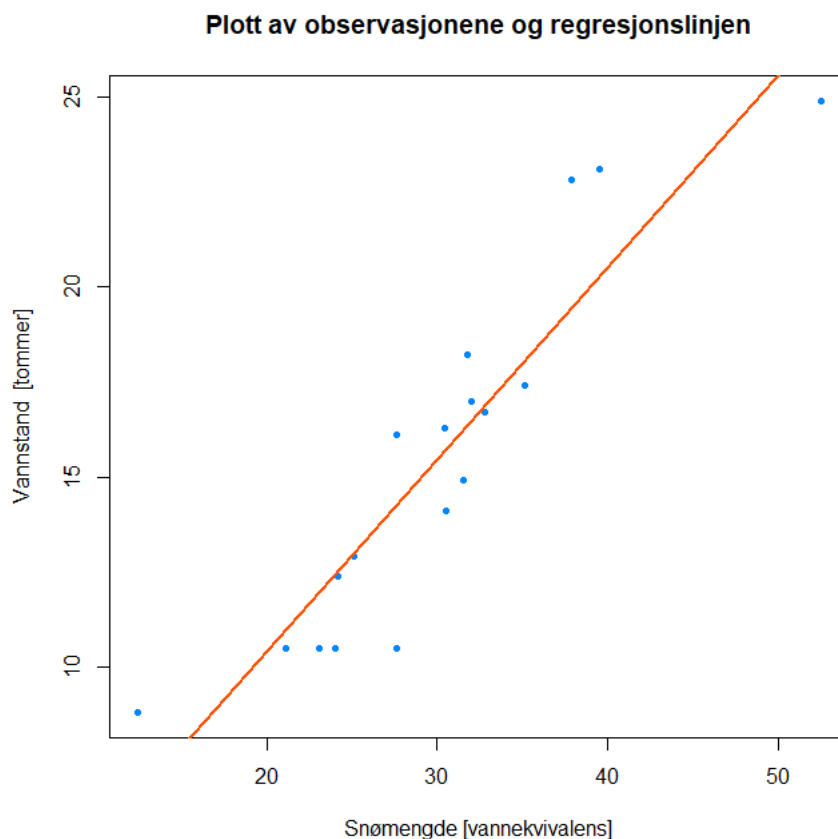
```
1 > p = 2*(1-pnorm(1.607))
2 > print(p)
3 [1] 0.1080544
```

Altså P-verdi =  $0.108 > 0.05 = \alpha$ . Siden P-verdien er større enn signifikansnivået så forkaster vi ikke  $H_0$  her heller. Testen vår konkluderer altså med at  $H_0 : p_1 - p_2 = 0$  stemmer, som betyr at det ikke er noen signifikant forskjell i tidsklemmeproblemer mellom fedre og mødre.

**b)** Kontrollerer svarene fra a) ved å benytte `prop.test()` funksjonen i R:

```
1 > prop.test(x = c(486, 441), n = c(3000, 3000))
2
3      2-sample test for equality of proportions with
      continuity correction
4
5 data:  c(486, 441) out of c(3000, 3000)
6 X-squared = 2.4701, df = 1, p-value = 0.116
7 alternative hypothesis: two.sided
8 95 percent confidence interval:
9  -0.003619808  0.033619808
10 sample estimates:
11 prop 1 prop 2
12  0.162  0.147
```

Vi ser at P-verdien her er lik 0.116, som betyr at  $H_0$  ikke forkastes. Dette samsvarer med det vi fikk i a). Legger merke til at P-verdien her er noe høyere enn den vi fikk i a), altså 0.116 mot 0.108. Denne forskjellen skyldes mest sannsynlig av avrundingen av  $z$ -verdien fra 1.60737 til 1.607.



Figur 4: *Spredningsplott av målingene (blå prikker) og regresjonslinjen vi fikk ved hjelp av `lm()` funksjonen (oransje linje).*

## Oppgave 4

a) Figur 4 viser spredningsplott av observasjonene sammen med regresjonslinjen. Vi ser at punktene former en rett linje med jevn spredning. Estimaterne vi får for koeffisientene fra `lm()`-funksjonen i R er  $\hat{\beta}_0 = 0.28$  og  $\hat{\beta}_1 = 0.5056$ . At stigningstallet er ca. 0.5 vil si at når snømengden øker én vannekvivalens så øker vannstanden med ca. 1/2 tomme. Gitt at vi ikke vet noe særlig om hvor stor elven er og hvor mye én vannekvivalens tilsvarer i forhold til elven, så er det litt vanskelig å si noe om hvor rimelig disse tallene er. Men de virker ikke helt på bærtur ihvertfall, så vi tar oss til gode med disse verdiene.

```
1 > url = "https://www.uio.no/studier/emner/matnat/math/STK1110/
  data/snoe_vann.txt"
2 > data = read.table(url, header=FALSE)
3 >
```

```

4 > snö = data$V1
5 > vann = data$V2
6 >
7 > lin.fit = lm(vann ~ snö)
8 > print(lin.fit)
9
10 Call:
11 lm(formula = vann ~ snö)
12
13 Coefficients:
14 (Intercept)          snö
15    0.2800         0.5056
16
17 >
18 > plot(snö, vann,
19 +      xlab="Snömengde [vannekvivalens]",
20 +      ylab="Vannstand [tommer]",
21 +      main="Plott av observasjonene og regresjonslinjen",
22 +      pch=20, col="#0084ff")
23 > abline(lin.fit, col="#ff5100", lwd=2)

```

b) Plotter residualene mot forklaringsvariabelen (snømengde) i Figur 5.

```

1 > # Finner standardisert residualene og plotter disse
2 > stan.res = rstandard(lin.fit)
3 > plot(snö, stan.res,
4 +      xlab="Snömengde [vannekvivalens]",
5 +      ylab="Standardisert residualer",
6 +      pch=20, col="#0084ff")

```

Et normalfordelingsplott av residualene er vist i Figur 6. Vi ser at plottet er veldig rett, som tyder på at normalitet var rett å anta.

```

1 > qqnorm(stan.res, col="#000000", pch=20)
2 > qqline(stan.res, col="#314bbe", lwd=2)

```

Den lineære regresjonsmodellen virker altså godt egnet til målingene våre.

c) Finner et estimat for variansen til feilleddene ved å benytte

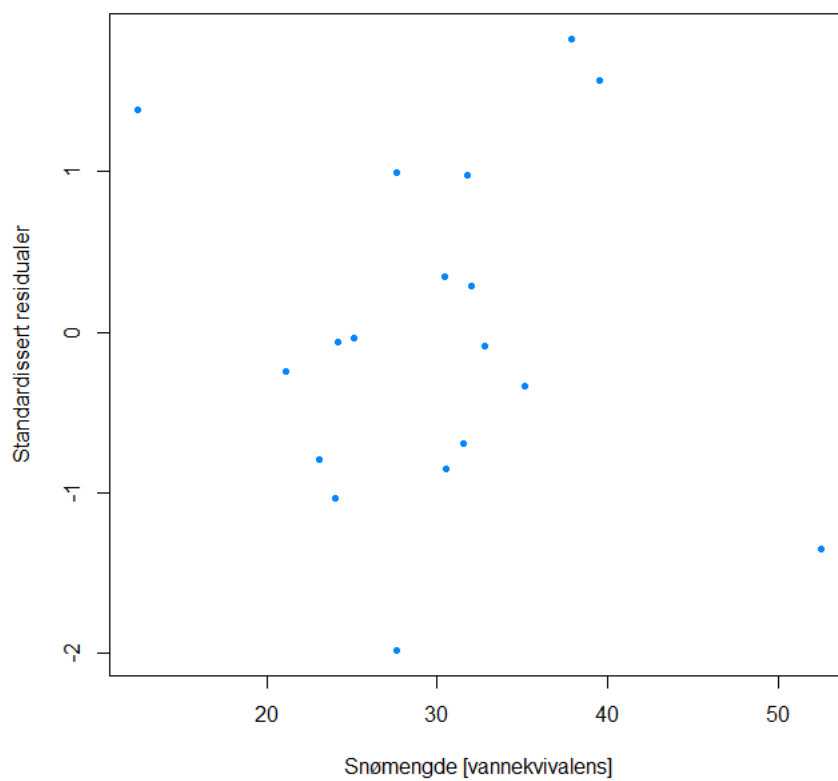
$$s_e^2 = \text{SSE}/(n - 2)$$

hvor  $n$  er antall målinger og SSE er residualkvadratsummen  $\sum_n (y_i - \hat{y}_i)^2$ .

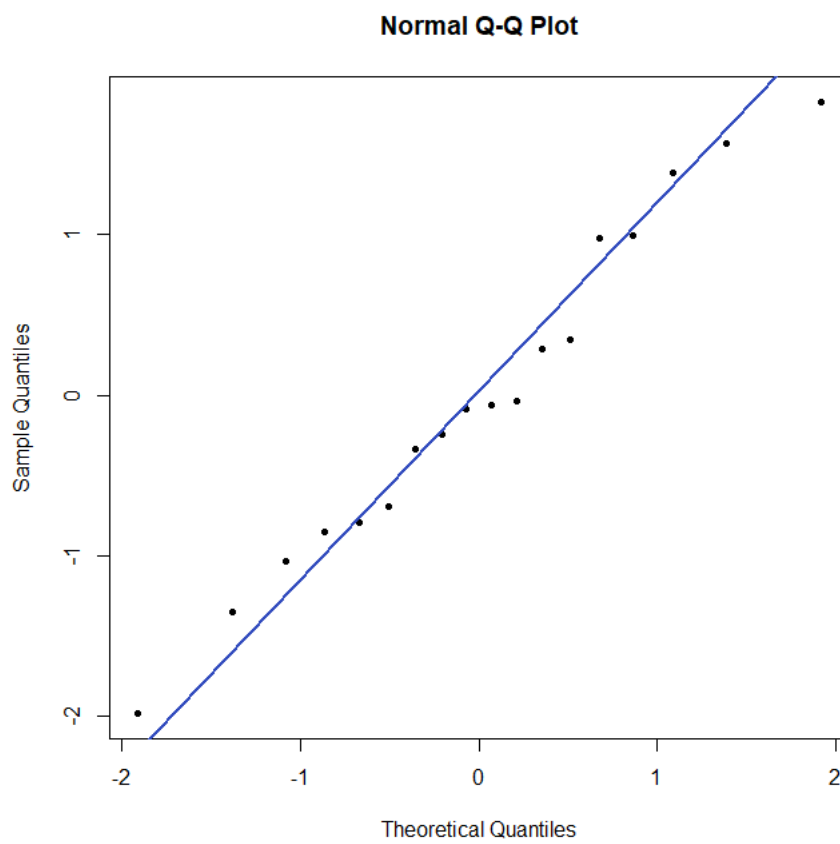
```

1 > n = length(snö)
2 > y = vann
3 > y.hat = lin.fit$fitted.values
4 > SSE = sum((y - y.hat)^2)
5 > s.e2 = SSE/(n-2)
6 > print(s.e2)
7 [1] 3.774598

```



Figur 5: *Residualplott for målingene. Vi ser at residualene ikke har noe spesielt mønster for seg og er godt fordelt over og under  $y=0$  linjen.*



Figur 6: Normalfordelingsplott av residualene. Vi ser at punktene følger en rett linje ganske godt.

Estimatet for variansen er altså  $s_e^2 = 3.77$ . Lager så et 95% konfidensintervall for stigningstallet  $\beta_1$  ved å benytte formelen

$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot \frac{s_e}{\sqrt{S_{xx}}}$$

hvor  $S_{xx} = \sum (x_i - \bar{x})^2$  og  $\hat{\beta}_1$  er estimatet for stigningstallet som vi fant i a) til å være  $\hat{\beta}_1 = 0.5056$ . Beregningene blir gjort med følgende koden

```
1 > x = snö
2 > beta = 0.5056
3 > t.crit = qt(0.025, n-2, lower.tail = FALSE)
4 > Sxx = sum((x - mean(x))^2)
5 > lower = beta - t.crit * sqrt(s.e2)/sqrt(Sxx)
6 > upper = beta + t.crit * sqrt(s.e2)/sqrt(Sxx)
7 > print(c(lower, upper))
8 [1] 0.3888443 0.6223557
```

Vi kan altså si med 95% sikkerhet at stigningstallet for den lineære relasjonen mellom vannstand og snømengde ligger innenfor intervallet (0.389, 0.622). Legger merke til at det bare inneholder positive verdier, som betyr at vannstanden øker når snømengden øker, noe som virker rimelig.

**d)** Vil teste  $H_0 : \beta_0 = 0$  mot  $H_A : \beta_0 \neq 0$  med signifikansnivå 5% ( $\alpha = 0.05$ ). Testobservasjonen for testen er

$$T = \frac{\hat{\beta}_0 - \beta_{00}}{S_{\hat{\beta}_0}} = \frac{\hat{\beta}_0}{S_{\hat{\beta}_0}}$$

hvor  $S_{\hat{\beta}_0} = S_e \sqrt{1/n + \bar{x}^2/S_{xx}}$  og  $\hat{\beta}_0$  er estimatet for skjæringspunktet som vi fant i a) til å være  $\hat{\beta}_0 = 0.28$ . Med signifikansnivå 5% så får vi en kritisk t-verdi på  $t_{\alpha/2, n-2} = t_{0.025, 16} = 2.120$ . Forkastningsområdet med  $\alpha = 0.05$  blir da:

$$t \leq -2.120 \quad \text{og} \quad t \geq 2.120$$

Testobservasjonen som kommer av våre målinger beregner vi ved:

```
1 > beta0 = 0.28
2 > s_b0 = sqrt(s.e2)*sqrt(1/n + mean(x)^2/Sxx)
3 > t = beta0 / s_b0
4 > print(t)
5 [1] 0.1635603
```

Vi får altså testobservasjonen  $t = 0.164$ . Denne verdien ligger ikke i forkastningsområdet, så vi konkluderer med at  $H_0 : \beta_0 = 0$  ikke skal forkastes. P-verdien finner vi ved

```
1 > p = 2*(1-pt(0.164, 16))
2 > print(p)
3 [1] 0.871785
```

Altså  $P\text{-verdi} = 0.87$ , som er mye høyere enn signifikansnivået  $\alpha = 0.05$ . Vi konkluderer også her at  $H_0$  ikke skal forkastes. Som sagt så fikk vi fra a) at  $\hat{\beta}_0 = 0.28$  ved hjelp av `lm()` funksjonen i R:

```
1 > lin.fit = lm(vann ~ snö)
2 > print(lin.fit)
3
4 Call:
5 lm(formula = vann ~ snö)
6
7 Coefficients:
8 (Intercept)          snö
9      0.2800         0.5056
```

Vi har altså to forskjellige antagelse for  $\hat{\beta}_0$ . Fra `lm()` funksjonen får vi  $\hat{\beta}_0 = 0.28$ , mens fra hypotesetesten får vi at  $\hat{\beta}_0 = 0$ . Dette virker kanskje litt rart, men hvis vi sammenligner dette med verdiene for vannstand (10.5, 16.7, 18.2, 17.0, ...) og snømengde (23.1, 32.8, 31.8, 32.0, ...), så ser vi at  $\hat{\beta}_0 = 0.28$  ikke er så forskjellig fra  $\hat{\beta}_0 = 0$ . Hypotesetesten vår samsvarer altså godt med resultatene fra a).