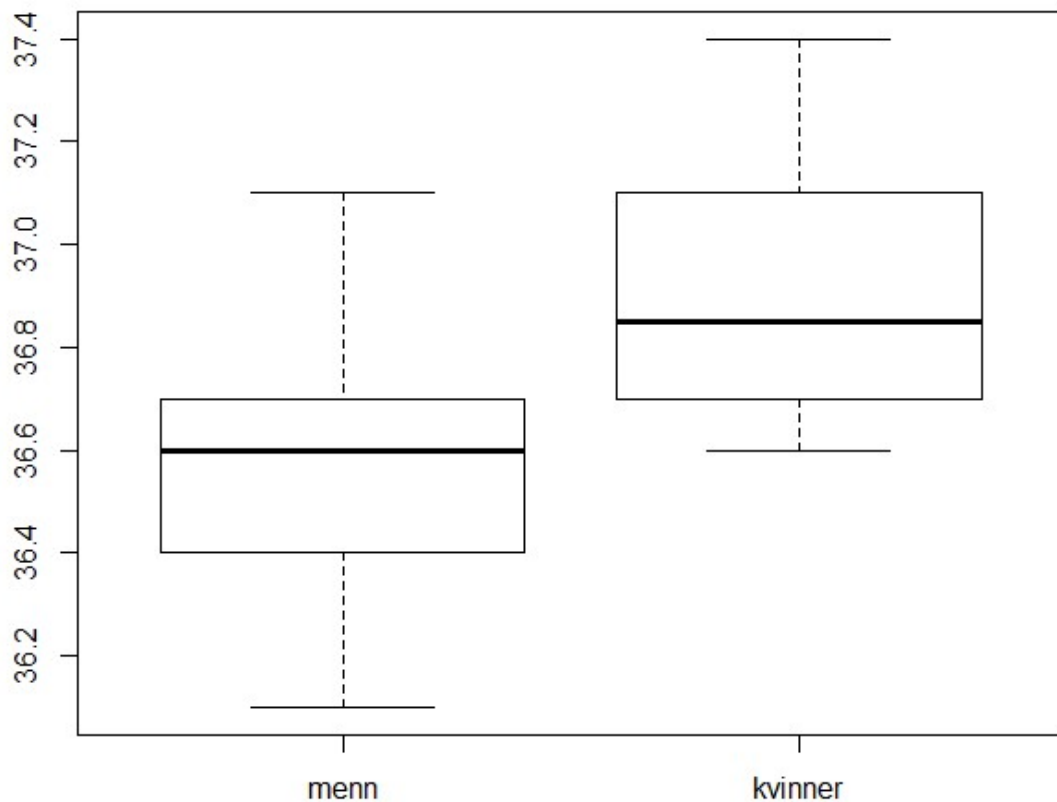


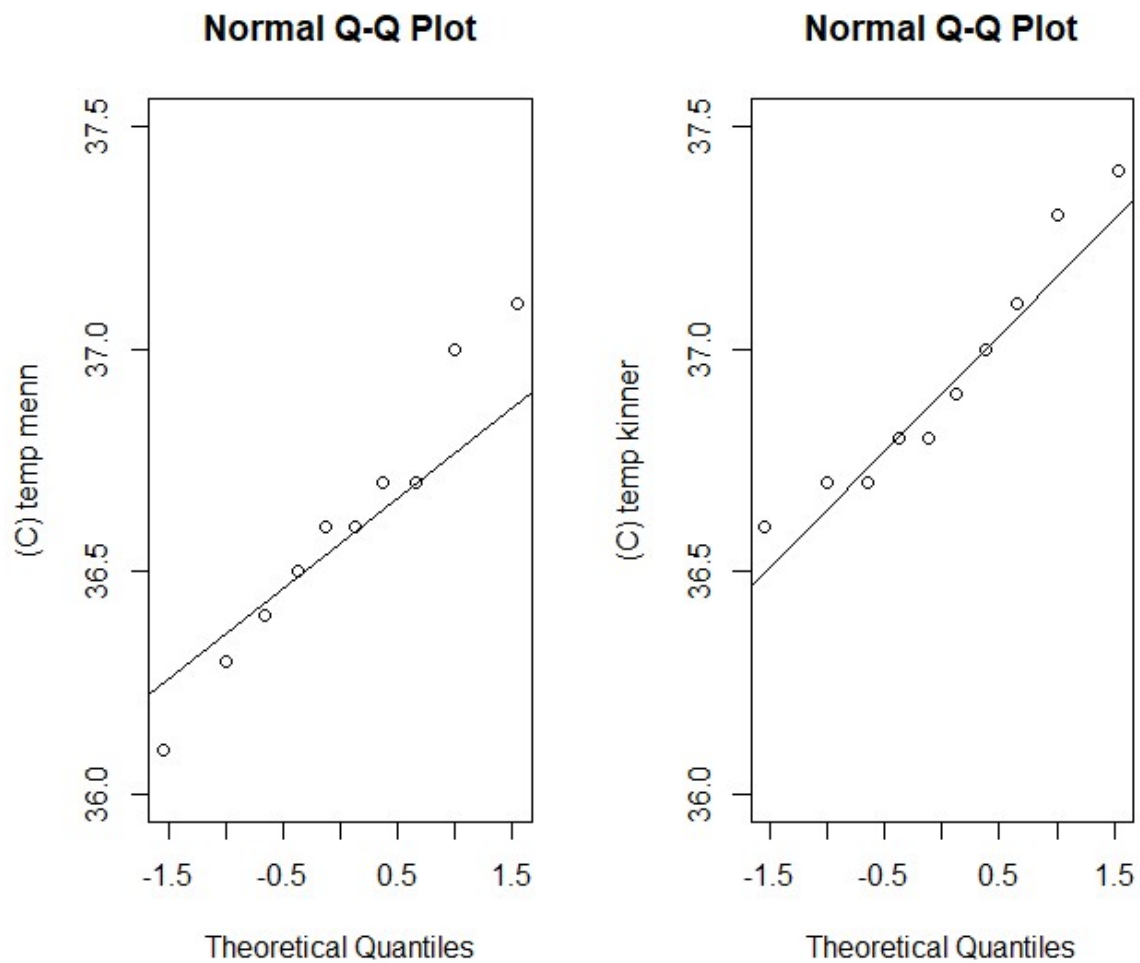
Oppg1

- a) Lag boksplott som viser fordelingen av observasjonene. Kommenter hva du finner.



Ser at kvinner sin snitt temperatur er høyere en menns. Og at kvinner har et høymaksverdi, mens min verdien er nærmere snittet. Dermed er kvinners fordeling skjev. Menn fordelingen er mindre skjev med maks, og min verdier ganske jevnt rundt snittet. Dog snittet er litt positivt skjevt.

- b) Lag normalfordelingsplott for de to observasjonssettene. Kommenter hva du ser.



Ser at både kvinner og menn fordelingene er ganske normalt fordelte med få punkter lengre unna enn resten.

- c) Anta at variansen er den samme for de to utvalgene, og test med niv^o 5% om det er noen forskjell i forventet kroppstemperatur. Beregn ogs^a P-verdien, og lag et 95% konfidensintervall for forventet forskjell.

Vi sette $H_0: \mu_{\text{menn}} = \mu_{\text{kvinner}}$, $H_a: \mu_{\text{menn}} \neq \mu_{\text{kvinner}}$

For å gjøre testen trenger jeg t_{obs} , og t_{lim} :

b1CI

$$t_{obs} = \frac{|\bar{x}_1 - \bar{x}_2| - \overbrace{(\mu_1 - \mu_2)}^{=0, \text{ p.g.a } H_0}}{\sqrt{S_p \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

setter vi inn tall får vi (Regnet i "r")
2,59

$$t_{lin} = t_{(1-\alpha/2, n-1)}$$

gir oss

$$t_{lin} = 2,26$$

Dvs $t_{lin} < t_{obs}$, som vil si at det er en signifikant forskjell. P.G.A vi forkaster H_0

For å finne P bruker jeg r koden $P = 2 * pt(t_{obs}, nx+ny-2)$.

Dvs. $P = 1.981519$

Vi kan finne CI ved
setter for ryddighetskyld

$$t_{\alpha/2, n_x + n_y - 2} = t$$

$$1 - \alpha = P \left(-t < \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)}} < t \right)$$

Lösen p.h.m $\mu_x - \mu_y$

setter $\mu_x - \mu_y = \mu_D$,

$$1 - \alpha = P \left(-t < \frac{\bar{X} - \bar{Y} - \mu_D}{\left(Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right) \right)^{1/2}} < t \right)$$

$$= P \left(-t \sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)} < \bar{X} - \bar{Y} - \mu_D < t \sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)} \right)$$

$$= P \left(\bar{X} - \bar{Y} - t \sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)} < \mu_D < \bar{X} - \bar{Y} + t \sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)} \right)$$

Dermed har vi et CI for $\mu_D = \mu_x - \mu_y$

$$\bar{X} - \bar{Y} \pm t_{\alpha/2, n_x + n_y - 2} \sqrt{Sp \left(\frac{1}{n_x} + \frac{1}{n_y} \right)}$$

Setter vi inn verdiene får vi ("r")

$$[0.062, 0.598]$$

- d) Gjennomfør testen og beregn P-verdien også i det tilfellet der man ikke antar felles varians. Diskuter og forklar resultatene.

Hvis det ikke er samme varians bruker vi

$$t_{obs} = \frac{|\bar{x}_1 - \bar{x}_2| - |\mu_1 - \mu_2|}{\sqrt{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}}}$$

0Det gir oss en t_{obs} verdi på 3.36483. som er en noe større verdi enn det vi hadde når vi antok lik varians. De er fortsatt ganske lik noe vi ville antatt da variansen for de to fordelingene er rimelig like.

Vi estimerer frihetsgraden med:

$$n - 1 = \frac{(s_x^2/n_x + s_y^2/n_y)^2}{\frac{(s_x^2/n_x)^2}{n_x - 1} + \frac{(s_y^2/n_y)^2}{n_y - 1}}$$

Så runder vi $n-1$ til nærmeste hele tall. Da har vi estimerte frihetsgrader.

Estimatet er $v = 18$ som stemmer overens med den tidligere utregningen da vi regnet med at Variansene var det samme.

Oppg2

b)) Kall forventet forskjell mellom Twin A og Twin B for μ_D . Sett opp nullhypotese og alternativ hypotese for å besvare spørsmålet om forskjell i IQ. Finn en egnet testobservator, og beregn dennes numeriske verdi. Beregn så tilhørende P-verdi. Spesifiser antall frihetsgrader i fordelingen du bruker. Formuler din konklusjon på testen

	N	Mean	StDev	SE Mean
Twin A	31	93.32	15.41	2.77
Twin B	31	96.58	13.84	2.49
Difference	31	-3.26	8.81	1.58

$$t_{obs} = \frac{\mu_D - A_D}{\sigma_D / \sqrt{n}}$$

Setter inn verdiene

$$t_{obs} = \frac{-3,26 - 0}{8,81 / \sqrt{31}} = -2,06$$

Dermed trenger vi P til å sammenligne med $\alpha = 0,05$

$$P = 2 P(t_{n-1} \geq t_{\alpha/2})$$

Setter inn verdiene

$$P = 2 P(t_{30} \geq -2,06)$$

$$\text{"r"} = 0,046$$

Dermed har vi $P < \alpha$ så vi forkaster H_0

II

Vi har $n-1$ frihetsgrader.

c) Finn et 95% konfidensintervall for μ_D . Hva betyr det at dette intervallet dekker kun negative verdier? Forklar kort om sammenhengen mellom tosidig testing og konfidensintervaller.

For å finne 95% CI for μ_D bruker jeg

$$\bar{d} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}$$

Setter inn tall og får

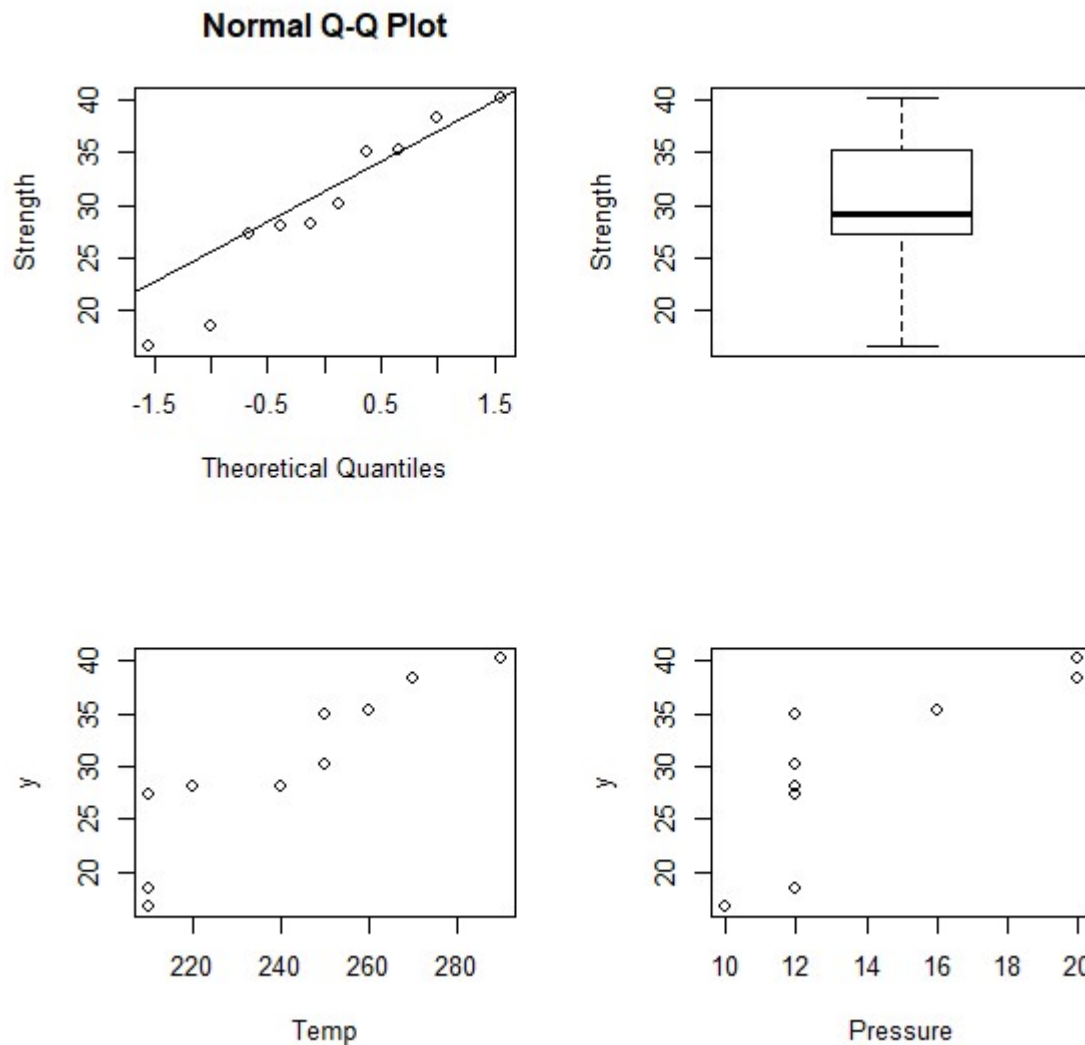
$$-3,26 \pm -2,04 \cdot \frac{8,81}{5,57}$$

"r"
= $[-6,49, -0,028]$

Vår forventningsverdi er negativ da får vi et konfidensintervall som er sentrert rundt forventningsverdien.

Oppg3

- a) Utfør ulike typer eksplorativ data analyse (e.g. histogram, qqplot (probability plot), spredningsplott, dvs. scatterplot) p°a datasettet og beskriv hva du ser.

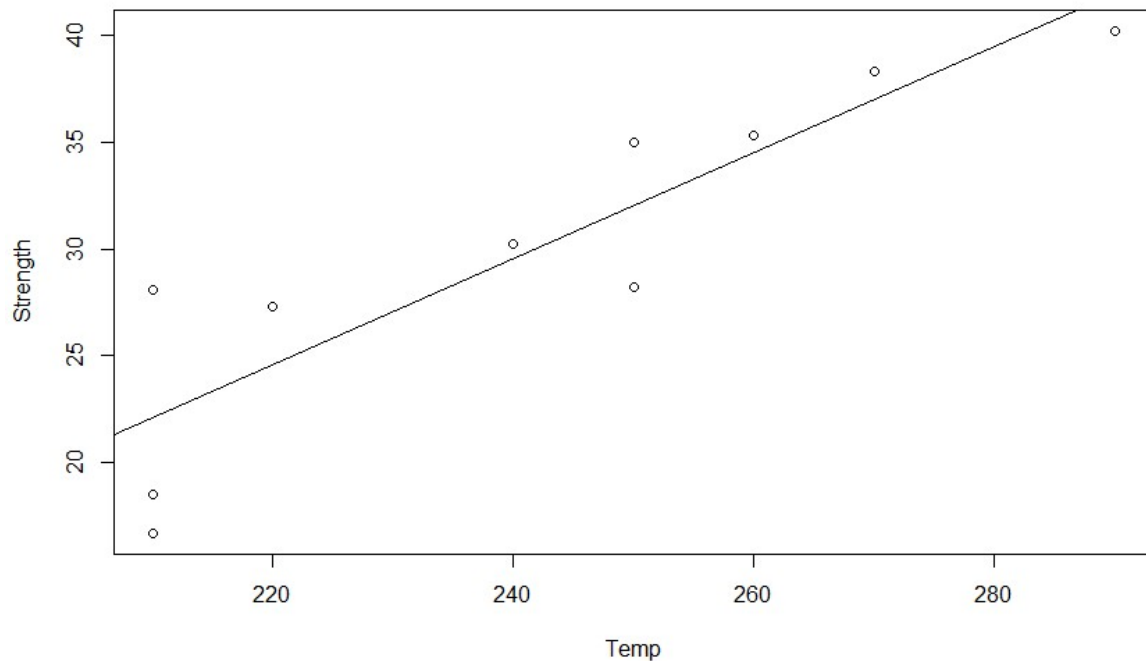


Fra Normal qq Plot ser vi at Strength er tilnærmet normalfordelt. Siden det ser ut som punktene følger normalen.

Fra boxplottet ser vi at Strength er positivt skjevt fordelt. Og vi har noen veldig lave verdier sammenlignet med forventingsverdien.

Fra de to siste plottene, er Strength(y) plottet mot Temp(x1) og Pressure(x2). Nå kan vi se hvilke som ser ut til å være mest linjer. Vi ser at y(temp) ser mer linjer ut, enn y(Pressure), som har et diskontinuerlig hop rundt $x_2 = 12$.

- b) Vi vil i første omgang konsentrere oss om forklaringsvariablen Temperature. Utfør en enkel lineær regresjon i dette tilfellet. Plott data sammen med den tilpassende regresjonslinjen og kommentér resultatene.



Ved å bruke minste kvadraters problem kan vi finne regresjonslinjen.

Vi skal finne $\hat{\beta}_0 + \hat{\beta}_1 x$

Vi finner β_1 ved

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}} \quad (12.2)$$

$$\hat{\beta}_0 = \frac{\sum y_i - \hat{\beta}_1 \sum x_i}{n} = \bar{y} - \hat{\beta}_1 \bar{x} \quad (12.3)$$

Ved r kan vi bruke `lm(x1, y, xlab='Temp', ylab='Strength')` for å finne β_{1_hat} og β_{0_hat}

(kalt $\hat{\beta}_1$ og $\hat{\alpha}_1$ i koden). Får da verdiene $\hat{\beta}_1=0.2474189$, $\hat{\beta}_0=-29.8479549$. Fra plottet ser regresjonen rimelig ut, og det virker som regresjonen er en ok estimator for $y(x_1)$.

C) Lag konfidensintervall for regresjonskoeffisienten for Temperature. Gir intervallet indikasjon på om temperatur er en viktig forklaringsvariabel?

CI for regresjonsintervallet finner vi ved:

$$\begin{aligned} CI &= \hat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot s_b \\ \text{Setter vi i inn tall får vi} \\ &= 0,247 \pm 2,306 \cdot 0,098 \\ &= [0,020, 0,4743] \end{aligned}$$

Vi ser at min tallet 0.02 ikke er veldig nærme -1 og maks tallet er 0.47 som da er nærmere 1, dette avkrefter ikke at vi kan bruke dette som en modell, men det er ikke veldig bra.