



Norges teknisk-naturvitenskapelige universitet  
Institutt for matematiske fag

TMA4240 Statistikk  
Høst 2020

Anbefalte oppgaver 12  
Løsningskisse

### Oppgave 1

- a) Minste kvadraters metode tilpasser en linje til punktene ved å velge den linja som minimerer kvadratsummen

$$\text{SSE} = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

av avstanden fra hvert punkt til linja. Derivasjon av SSE med hensyn på parametrene  $\alpha$  og  $\beta$  gir

$$\frac{d\text{SSE}}{d\alpha} = -2 \sum_i (y_i - \alpha - \beta x_i) \quad \text{og} \quad \frac{d\text{SSE}}{d\beta} = -2 \sum_i x_i (y_i - \alpha - \beta x_i).$$

Setter vi de deriverte lik null, får vi

$$\sum_{i=1}^n (y_i - \alpha - \beta x_i) = 0 \quad \text{og} \quad \sum_{i=1}^n x_i (y_i - \alpha - \beta x_i) = 0,$$

og, når vi deler på  $n$ ,

$$\bar{y} - \alpha - \beta \bar{x} = 0 \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n x_i y_i - \alpha \bar{x} - \beta \cdot \frac{1}{n} \sum_{i=1}^n x_i^2 = 0.$$

Løser den første likningen for  $\alpha$ , og får

$$\alpha = \bar{y} - \beta \bar{x},$$

som innsatt i den andre likningen gir

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n x_i y_i - (\bar{y} - \beta \bar{x}) \bar{x} - \beta \cdot \frac{1}{n} \sum_{i=1}^n x_i^2 &= 0, \\ \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{y} \bar{x} + \beta \left( \bar{x}^2 - \frac{1}{n} \sum_{i=1}^n x_i^2 \right) &= 0 \Rightarrow \beta = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2}. \end{aligned}$$

Ganger vi med  $n$  i teller og nevner i det siste uttrykket, får vi

$$\beta = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}.$$

For få de oppgitte estimatorene bytter vi ut  $y_i$  med den tilsvarende tilfeldige variabelen  $Y_i$ , altså

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i Y_i - n\bar{x}\bar{Y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \quad \text{og} \quad \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x}.$$

b) Utgangspunktet er

$$P\left(-t_{n-2,0.025} < \frac{(\hat{\beta} - \beta)}{s/\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} < t_{n-2,0.025}\right) = 0.95$$

Løser hver av ulikhetene for  $\beta$  og får

$$\begin{aligned} -t_{n-2,0.025} < \frac{(\hat{\beta} - \beta)}{s/\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} &\Rightarrow -\frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} < \hat{\beta} - \beta \\ &\Rightarrow \hat{\beta} + \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} > \beta \end{aligned}$$

og

$$\begin{aligned} \frac{(\hat{\beta} - \beta)}{s/\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} < t_{n-2,0.025} &\Rightarrow \hat{\beta} - \beta < \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \\ &\Rightarrow \beta > \hat{\beta} - \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}. \end{aligned}$$

Dermed har vi

$$P\left(\hat{\beta} - \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} < \beta < \hat{\beta} + \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}\right) = 0.95,$$

Og konfidensintervallet blir altså

$$\left(\hat{\beta} - \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, \hat{\beta} + \frac{st_{n-2,0.025}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}\right).$$

Vi har  $n = 28$  og tabelloppslag gir kvantilen  $t_{n-2,0.025} = t_{26,0.025} = 2.056$ . Innsetting av tallverdier gir estimatet  $\hat{\beta} = -5942/36517 = -0.16$ . Vinnertiden forventes å forkortes med  $4 \cdot 0.16 = 0.64$  sekunder mellom etterfølgende olympiske leker. Videre er 95%-konfidensintervallet for stigningstallet lik  $(-0.199, -0.126)$ .

c) Vi lar  $x_0 = 2016$ , og vi har  $\hat{\alpha} = 109.26 + 0.1627 \cdot 1954.5 = 427.3$ . Den predikerte tiden er  $\hat{Y}_0 = \hat{\alpha} + 2016\hat{\beta} = 427.3 - 2016 \cdot 0.1627 = 99.3$ , altså ca. 1 minutt og 39 sekunder. Vinnertiden i 2016 har 95%-prediksjonsintervall

$$\hat{Y}_0 \pm t_{n-2,0.025} s \sqrt{1 + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} + 1/n}.$$

Med tallverdier innsatt blir det  $(91.8, 106.7)$ .

d) Vi har  $\hat{Y}_0 = \hat{\alpha} + x_0\hat{\beta} = 90$ , som betyr at

$$x_0 = \frac{90 - \hat{\alpha}}{\hat{\beta}} = \frac{90 - 427.3}{-0.1627} = 2073.$$

De første olympiske sommerlekene etter 2073 finner sted i 2076.

**Modellantakelser:** Det ser ut til at vinnertidene følger en ikkelinear trend i tid. Om vi bruker den tilpassede modellen til å ekstrapolere bakover i tid, ser vi at den tilsier at vinnertiden i år 0 ville vært 427 sekunder, hvilket er urimelig. Ekstrapolerer vi tilstrekkelig langt framover i tid, predikerer modellen dessuten negative vinnertider, hvilket er umulig.

Modellantakelsene kan kontrolleres ved hjelp av residualplott. Ser residualene  $e_i = Y_i - \hat{Y}_i$  ut til å ha en trend? Ifølge modellen bør de være nærmest uavhengige og identisk normalfordelt.

## Oppgave 2

a)  $Y \sim n(y; 500, 80)$ . Transformerer  $Y$  til standard  $N(0, 1)$ -normalfordeling.

$$\begin{aligned} \text{Prob}(Y > 550) &= \text{Prob}\left(\frac{Y - 500}{80} > \frac{550 - 500}{80}\right) = \text{Prob}\left(Z > \frac{5}{8}\right) \\ &= 1 - \text{Prob}\left(Z \leq \frac{5}{8}\right) = 1 - \Phi(0.625) = 1 - 0.734 = 0.266. \end{aligned}$$

$Y_1 - Y_2 \sim n(y; 0, \sqrt{2} \cdot 80)$ . (Lineærkombinasjonen av to uavhengige normalfordelinger er normalfordelt, sjekk forventningsverdi og varians ved de vanlige regnereglene.)

Da kan vi regne ut sannsynligheten for at målingene avviker med mer enn 80 g/tonn.

$$\begin{aligned} \text{Prob}(|Y_1 - Y_2| > 80) &= 1 - \text{Prob}(-80 < Y_1 - Y_2 < 80) \\ &= 1 - \text{Prob}\left(\frac{-80}{80\sqrt{2}} < \frac{Y_1 - Y_2}{80\sqrt{2}} < \frac{80}{80\sqrt{2}}\right) \\ &= 1 - \text{Prob}\left(-\frac{\sqrt{2}}{2} < Z < \frac{\sqrt{2}}{2}\right) = 2\text{Prob}\left(Z \leq \frac{-\sqrt{2}}{2}\right) = 2\Phi(-0.707) \\ &= 2 \cdot 0.24 = 0.48. \end{aligned}$$

b) Setter inn  $\bar{x} = 20$ ,  $x_1 = \dots = x_5 = 0$  og  $x_6 = \dots = x_{10} = 40$  i uttrykket for  $B$ .

$$\begin{aligned} B &= \frac{\sum_{j=1}^5 -20Y_j + \sum_{j=6}^{10} 20Y_j}{\sum_{j=1}^{10} 20^2} = \frac{20 \left( \sum_{j=6}^{10} Y_j - \sum_{j=1}^5 Y_j \right)}{10 \cdot 20^2} \\ &= \frac{\sum_{j=6}^{10} Y_j - \sum_{j=1}^5 Y_j}{200}, \text{ som skulle vises.} \end{aligned}$$

$$A = \bar{Y} - B\bar{x} = \frac{1}{10} \sum_{j=1}^{10} Y_j - \frac{20}{200} \left( \sum_{j=6}^{10} Y_j - \sum_{j=1}^5 Y_j \right) = \frac{1}{5} \sum_{j=1}^5 Y_j.$$

$A$  er skjæringspunktet regresjonslinja har med  $y$ -aksen. Det er kanskje ikke så rart at gjennomsnittet av målingene ved  $x = 0$  er et estimat for denne verdien? (I hvert fall når målingene bare er gjort for to  $x$ -verdier.)

$$\text{Var}(B) = \frac{1}{200^2} \left( \sum_{j=6}^{10} \text{Var}(Y_j) + \sum_{j=1}^5 \text{Var}(Y_j) \right) = \frac{10\sigma^2}{200^2} = \frac{\sigma^2}{4000}.$$

- c) Med bare to målepunkter, kan vi estimere variansen i hver ende for seg, dvs at vi beregner  $s_V^2$  og  $s_E^2$ . (Husk at målingene ikke har samme forventningsverdi i de to endene av gruva, så vi kan ikke se på alle som ett datasett.) Ettersom vi antar samme varians i begge ender, er gjennomsnittet av  $s_V^2$  og  $s_E^2$  et godt estimat for  $\sigma^2$ .

Mer formelt, vi har en to-utvalgssituasjon, og kan da bruke  $s_p^2$  fra pensum. Denne sikrer  $\chi^2$ -fordeling og T-fordeling. Brukes estimatoren for variansen fra regresjonsanalysen, får en også samme resultat.

$$\begin{aligned} s^2 &= \frac{1}{2} (s_V^2 + s_E^2) = \frac{1}{2} \left( \frac{\sum_{j=1}^5 (y_j - \bar{y}_V)^2}{5-1} + \frac{\sum_{j=6}^{10} (y_j - \bar{y}_E)^2}{5-1} \right) \\ &= \frac{1}{8} \left( \sum_{j=1}^5 (y_j - \bar{y}_V)^2 + \sum_{j=6}^{10} (y_j - \bar{y}_E)^2 \right) = \frac{26064 + 22720}{8} = 6098. \end{aligned}$$

Hypotesene blir:  $H_0: \beta = 12$  mot  $H_1: \beta > 12$ .

Vi baserer testen på estimatoren  $B$ . Siden variansen til  $B$  er ukjent, bruker vi estimatet  $S_B^2 = \frac{s^2}{4000} = 1.525$  i stedet for  $\frac{\sigma^2}{4000}$ .

Testobservatoren,  $\frac{B-12}{S_B}$ , er T-fordelt med 8 frihetsgrader. Det er  $n - 2$  frihetsgrader denne gangen, fordi vi bruker "pooled" varians, eller, som sagt, variansestimatoren fra regresjonsanalysen. (Estimert varians er basert på to gjennomsnitt,  $\bar{y}_V$  og  $\bar{y}_E$ . Da er det ikke så urimelig at vi mister to frihetsgrader?) Med oppgitte data blir stigningstallet

$$b = \frac{\sum_{j=6}^{10} y_j - \sum_{j=1}^5 y_j}{200} = \frac{\bar{y}_E - \bar{y}_V}{40} = 17.$$

Gjennomfører hypotesetesten.

$$\frac{b - 12}{s_B} = \frac{17 - 12}{\sqrt{1.525}} = 4.05 > t_{0.05,8} = 1.86,$$

som betyr at vi forkaster nullhypotesen på signifikansnivå 5%.

- d) Fra det første uttrykket for  $B$  får vi

$$\text{Var}(B) = \frac{\sigma^2}{\sum_{j=1}^n (x_j - \bar{x})^2}.$$

Variansen er liten for  $\sum_{j=1}^n (x_j - \bar{x})^2$  stor. Altså vil vi ha alle  $|x_j - \bar{x}|$  så store som mulig. Når  $\bar{x}$  er fast, bør  $x_j$ -ene legges til endene, som i denne oppgaven. (Det kan være andre grunner til å spre målepunktene, f.eks. for å vurdere om dataene tilnærmet følger en rett linje, her var det antatt kjent.)

$$\text{Var}(Y_0 - \hat{Y}_0) = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \right) = \frac{11}{10} \cdot \sigma^2$$

når  $x_0 = \bar{x}$ . Punkttestimatet blir  $\hat{y}_0 = a + bx_0 = \bar{y}_V + 17 \cdot 20 = 470$ .

Vi benytter fortsatt estimatet  $S^2$  for  $\sigma^2$ , derfor fortsatt T-fordeling med  $n - 2$  frihetsgrader. Prediksjonsintervallet blir derfor

$$(\hat{y}_0 \pm t_{0.025,8} \cdot s \sqrt{\frac{11}{10}}) = (470 \pm 2.306 \cdot \sqrt{6098} \cdot \sqrt{1.1}) = (281.1, 658.9).$$

Den nye målingen, 600 g/tonn, ligger innenfor prediksjonsintervallet, så vi kan ikke konkludere med at den eller modellen er urimelig.