

Enkel lineær regresjon og SME

TMA4240/TMA4245 Statistikk

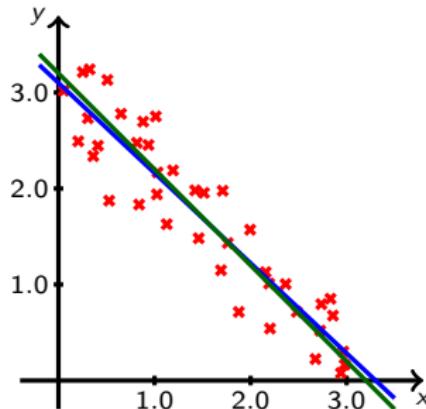
Håkon Tjelmeland

Institutt for matematiske fag

Norges teknisk-naturvitenskapelige universitet

Enkel lineær regresjon

- * Situasjon: Har observasjonspar $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$



- * Ønsker å tilpasse en rett linje til de observerte parene
- * Sannsynlighetsmodell: $Y_i = \alpha + \beta x_i + \varepsilon_i$ der $E[\varepsilon_i] = 0$ og $\text{Var}[\varepsilon_i] = \sigma^2$
 - α og β er ukjente parametere, vi ønsker å estimere disse
 - vi betrakter y_i 'er som realisasjoner av stokastiske variabler Y_1, Y_2, \dots, Y_n
 - vi betrakter x_i som tall (altså ikke stokastiske variabler)
- * For å bruke sannsynlighetsmaksimeringsprinsippet må vi også anta hvilken sannsynlighetsfordeling ε_i 'ene har
 - anta $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ uavhengige og $\varepsilon_i \sim N(0, \sigma^2)$
- * Dermed: Y_1, Y_2, \dots, Y_n er uavhengige og
$$Y_i \sim N(\alpha + \beta x_i, \sigma^2)$$
- * Merk: Vi har tre parametre vi trenger å estimere, α , β og σ^2

SME i enkel lineær regresjon

$$\text{Normalfordeling: } f(x) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp \left\{ -\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right\}$$

- Husk: Y_1, Y_2, \dots, Y_n er uavhengige og $Y_i \sim N(\alpha + \beta x_i, \sigma^2)$
- Rimelighetsfunksjon

$$L(\alpha, \beta, \sigma^2) = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - (\alpha + \beta x_i))^2 \right\} \right]$$

- Log-rimelighetsfunksjonen

$$\begin{aligned} \ell(\alpha, \beta, \sigma^2) &= \sum_{i=1}^n \left[-\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (y_i - (\alpha + \beta x_i))^2 \right] \\ &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 \end{aligned}$$

- Partielle deriverte

$$\frac{\partial \ell}{\partial \alpha} = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(Y_i - (\alpha + \beta x_i)) \cdot (-1) = \dots = \frac{1}{\sigma^2} \left[\sum_{i=1}^n nY_i - n\alpha - \beta \sum_{i=1}^n x_i \right]$$

$$\frac{\partial \ell}{\partial \beta} = \frac{1}{2\sigma^2} \sum_{i=1}^n 2(Y_i - (\alpha + \beta x_i)) \cdot (-x_i) = \dots = \frac{1}{\sigma^2} \left[\sum_{i=1}^n x_i Y_i - \alpha \sum_{i=1}^n x_i - \beta \sum_{i=1}^n x_i^2 \right]$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 = \frac{1}{2\sigma^2} \left[-n + \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 \right]$$

SME i enkel lineær regresjon

- ★ Ligningssystemet vi må løse blir

$$n\alpha + \beta \sum_{i=1}^n x_i = \sum_{i=1}^n Y_i$$

$$\alpha \sum_{i=1}^n x_i + \beta \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i$$

$$\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 = n$$

– merk: De to første er normalligningene fra minste kvadraters metode

- ★ Log-rimelighetsfunksjonen har sitt maksimum for

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

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

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

- ★ SME er

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{\alpha} + \hat{\beta} x_i))^2$$

Oppsummering

★ Enkel lineær regresjon

- har observasjonspar $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- ønsker å finne en rett linje som passer med dataene
- betrakter y_1, y_2, \dots, y_n som realisasjoner av stokastiske variabler Y_1, Y_2, \dots, Y_n
- betrakter x_1, x_2, \dots, x_n som tall
- sannsynlighetsmodell: $Y_i = \alpha + \beta x_i + \varepsilon_i$ der $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ uavhengige og $\varepsilon_i \sim N(0, \sigma^2)$

★ SME

- bestemmer estimatorer $\hat{\alpha}$, $\hat{\beta}$ og $\hat{\sigma}^2$ ved å maksimere log-rimelighetsfunksjonen
- estimatorer

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{\alpha} + \hat{\beta} x_i))^2$$

★ Merk:

- vi har enda ikke diskutert om det er fornuftig å tilpasse ei rett linje
- vi skal senere se at $\hat{\sigma}^2$ er forventningsskjev