

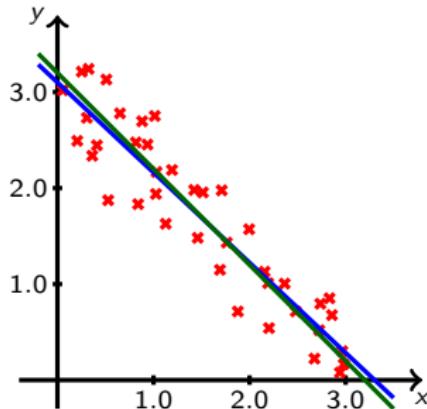
# **Enkel lineær regresjon og minste kvadraters metode**

TMA4240/TMA4245 Statistikk

Håkon Tjelmeland  
Institutt for matematiske fag  
Norges teknisk-naturvitenskapelige universitet

## Enkel lineær regresjon

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

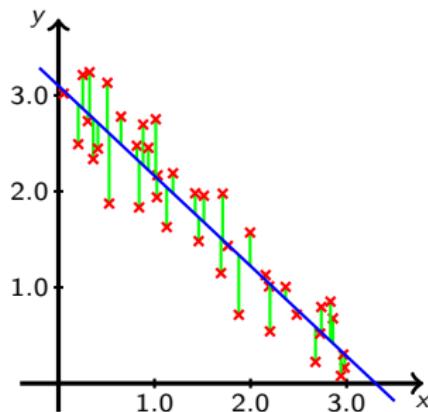


- ★ Ønsker å tilpasse en rett linje til de observerte parene
- ★ Sannsynlighetsmodell:  $Y_i = \alpha + \beta x_i + \varepsilon_i$  der  $E[\varepsilon_i] = 0$  og  $\text{Var}[\varepsilon_i] = \sigma^2$ 
  - $\alpha$  og  $\beta$  er ukjente parametere, vi ønsker å estimere disse
  - vi betrakter  $y_i$ 'er som realisasjoner av stokastiske variable  $Y_1, Y_2, \dots, Y_n$
  - vi betrakter  $x_i$  som tall (altså ikke stokastiske variable)
- ★ Hvordan finne gode estimatorer for  $\alpha$  og  $\beta$ ?
  - minste kvadraters metode
- ★ Merk: Vi har to linjer
  - den ukjente sanne linja  $y = \alpha + \beta x$
  - den estimerte linja  $y = \hat{\alpha} + \hat{\beta}x$

## Minste kvadraters metode

- \* Idé: Måler «avstanden» mellom observasjonene og estimert linje ved

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



- velger verdier for  $\hat{\alpha}$  og  $\hat{\beta}$  slik at SSE minimeres
- \* Matematisk: Finner  $\hat{\alpha}$  og  $\hat{\beta}$  ved å minimere SSE med hensyn på  $\hat{\alpha}$  og  $\hat{\beta}$ .
  - må løse ligningssystemet

$$\frac{\partial SSE}{\partial \hat{\alpha}} = 0, \quad \frac{\partial SSE}{\partial \hat{\beta}} = 0$$

med hensyn på  $\hat{\alpha}$  og  $\hat{\beta}$

# Utregning av estimatorer

- ★ Husk: Vi minimere

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

- ★ Partiellderiverte

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

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

- ★ Ligningssystem (normalligningene)

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

- ★ Løsning

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i Y_i - \frac{1}{n} (\sum_{i=1}^n x_i) (\sum_{i=1}^n Y_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad \text{og} \quad \hat{\alpha} = \frac{1}{n} \sum_{i=1}^n Y_i - \hat{\beta} \cdot \frac{1}{n} \sum_{i=1}^n x_i$$

- ★ Omskriving (bruker  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  og  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ )

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

# 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  $E[\varepsilon_i] = 0$  og  $\text{Var}[\varepsilon_i] = \sigma^2$
- ★ Minste kvadraters metode
  - bestemmer estimatorer  $\hat{\alpha}$  og  $\hat{\beta}$  ved å minimere

$$\text{SSE} = \sum_{i=1}^n \left( Y_i - (\hat{\alpha} + \hat{\beta} x_i) \right)^2$$

- estimatorer

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

- ★ Merk: Vi har enda ikke diskutert om det er fornuftig å tilpasse ei rett linje
- ★ Hvorfor kalles det **enkel lineær regresjon**?
  - regresjon: ønsker å finne hvordan  $y$  varierer som funksjon av  $x$
  - lineær:  $Y_i$  er en lineær funksjon av parametrene  $\alpha$  og  $\beta$
  - enkel: modellen har kun en forklaringsvariabel,  $x$