

Plan:

- ① Linear regresjon
- ② Statistiske egenskaper
- ③ Tolkning av R^2

Pensum:

[L] 7.3
(ilke regel 7.6, 7.7)

Musk: ① Veiledning kl. 11 i dag.

② Neste forelesning: Repetisjon

Prøve-eksamen

Ekstra forelesning: Før-eksamen.

Repetisjon:

Korrelasjonskoeffisient: (i et utvalg)

$$R = \frac{S_{xy}}{S_x \cdot S_y}$$

$$S_{xy} = \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

(utvalgs-kovarians)

$$S_x = \sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$S_y = \sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (y_i - \bar{y})^2}$$

utvalgs-
std. avvik

$$-1 \leq R \leq 1$$

$R > 0$: positiv
sammenheng

$R < 0$: negativ
sammenheng

$|R|$: nær 1 gir sterk sammenheng
nær 0 gir svak —||—

$$\begin{aligned} (n-1) S_{xy} &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \\ (n-1) S_x^2 &= \sum_{i=1}^n x_i^2 - n \bar{x}^2 \\ (n-1) S_y^2 &= \sum_{i=1}^n y_i^2 - n \bar{y}^2 \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum (x_i y_i - \bar{x} y_i - x_i \bar{y} + \bar{x} \bar{y}) = \sum x_i y_i - \bar{x} \cdot \sum y_i \\ &\quad - \bar{y} \sum x_i + \bar{x} \bar{y} \cdot n = \sum x_i y_i - \bar{x} \cdot n \bar{y} - \bar{y} \cdot n \bar{x} + n \bar{x} \bar{y} \end{aligned}$$

① Lineær regresjon

X: forklaringsvariabel (kontrollert)

Y: responsvariabel (stokastisk)

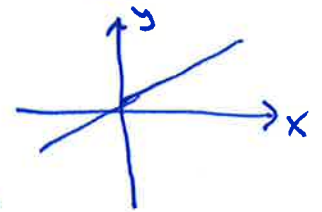
datasett

i	x_i	y_i
1	x_1	y_1
2	x_2	y_2
⋮		
n	x_n	y_n

Antar: Lineær sammenheng

$$Y = \alpha + \beta X$$

for (ukjente) parametre α, β .



α : skjæring med y-aksen

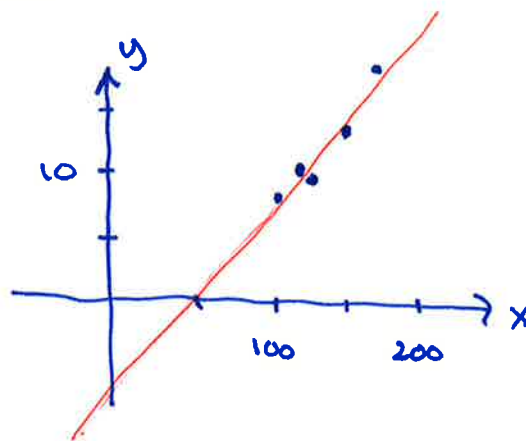
β : helling / stign.tall

Ekst: Boligpriser

x: ant. kvadr. areal

y: pris (mill. kr)

x	y
109	9.5
181	15.9
112	9.39
102	8.75
158	11.5



Multipel lineær regresjon: Mer enn én forklaringsvariabel

Tilpasser du rette linjen til datasettet?

Modell: $Y = \alpha + \beta X + \varepsilon$

$$Y_1 = \alpha + \beta X_1 + \varepsilon_1$$

$$Y_2 = \alpha + \beta X_2 + \varepsilon_2$$

\vdots

$$Y_n = \alpha + \beta X_n + \varepsilon_n$$

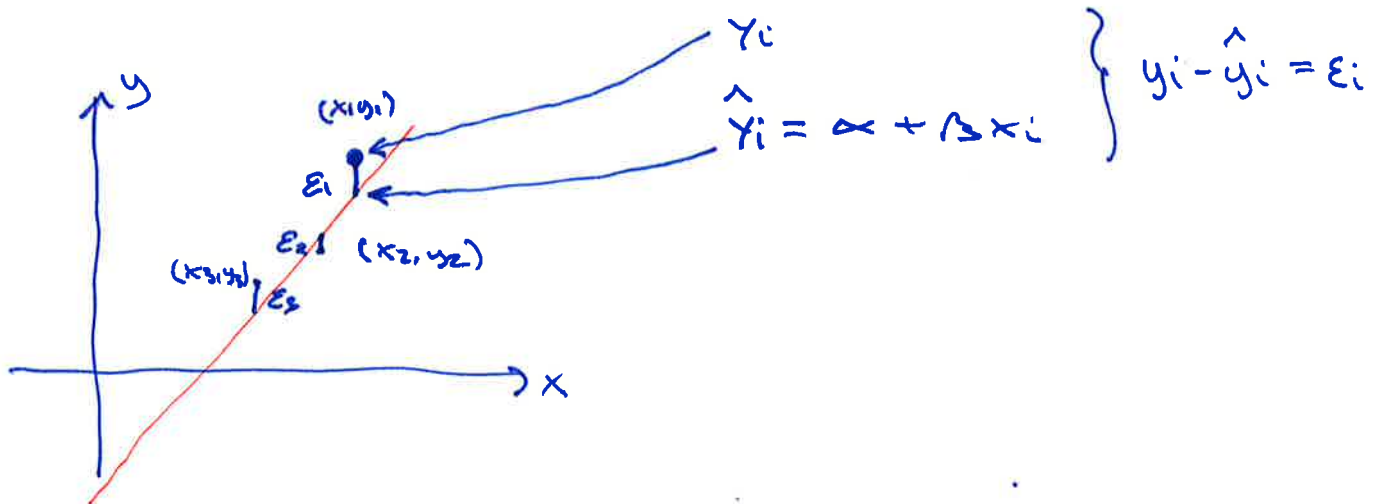
ε : feil-ledd

Antagelser:

X : gitt

$\varepsilon_i \sim N(0, \sigma^2)$
for en konstant σ

$\varepsilon_1, \dots, \varepsilon_n$: uavhengige



Beste tilpassning:
(Minste kvadraters)
metode

Velg α, β slik at

$$\varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2$$

er minst mulig.

← kvadratiske
feilfunksjon

Minste kvadraters metode :

Finn α, β slik at $E_1^2 + E_2^2 + \dots + E_n^2$ er minimal

$$\begin{aligned} \min f(\alpha, \beta) &= E_1^2 + E_2^2 + \dots + E_n^2 = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2 \\ &= (y_1 - \alpha - \beta x_1)^2 + (y_2 - \alpha - \beta x_2)^2 + \dots \\ &= \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \end{aligned}$$

Husk!
 $\hat{y}_i = \alpha + \beta x_i$

Løsning:

$$f'_\alpha = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i) \cdot (-1)$$

$$= -2 \cdot \sum_{i=1}^n (y_i - \alpha - \beta x_i) = 0 \quad | : (-2)$$

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

$$n \cdot \bar{y} - n\alpha - \beta \cdot n\bar{x} = 0$$

$$\boxed{\bar{y} = \alpha + \beta \bar{x}}$$

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

$$f'_\beta = \sum_{i=1}^n 2(y_i - \alpha - \beta x_i) \cdot (-x_i) = 0$$

$$-2 \sum_{i=1}^n (x_i y_i - \alpha x_i - \beta x_i^2) = 0 \quad | : (-2)$$

$$\sum_{i=1}^n x_i y_i - \alpha \sum_{i=1}^n x_i - \beta \sum_{i=1}^n x_i^2 = 0$$

$$\sum x_i y_i - (\bar{y} - \beta \bar{x}) \cdot \sum x_i - \beta \cdot \sum x_i^2 = 0$$

$$\sum x_i y_i - \bar{y} \sum x_i = \beta (\sum x_i^2 - \bar{x} \sum x_i)$$

$$\beta = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{(n-1) \cdot S_{xy}}{(n-1) S_x^2} = \frac{S_{xy}}{S_x^2}$$

Stasjonære plott: (α, β) sitt ved:

$$\beta = \frac{S_{xy}}{S_x^2} = \frac{S_{xy}}{S_x \cdot S_y} \cdot \frac{S_y}{S_x}$$

↑
r

$$\beta = r \cdot \frac{S_y}{S_x}$$

$$\alpha = \bar{Y} - \beta \cdot \bar{X}$$

Merk: ① Dette gir alltid et stasjonært plott, og vi kan vise at det er et globalt minimum for $f(\alpha, \beta) = \epsilon_1^2 + \dots + \epsilon_n^2$

Se detaljer neste side

→ $H(f) = \begin{pmatrix} 2n & 2n\bar{x} \\ 2n\bar{x} & 2\sum x_i^2 \end{pmatrix}$ det > 0
A > 0

② Vi kaller på $\hat{\alpha}, \hat{\beta}$ som estimer for α, β .

Beste estimat for $(\alpha, \beta) =$ estimatorene for (α, β)

→
$$\hat{\beta} = r \cdot \frac{S_y}{S_x} = \frac{S_{xy}}{S_x^2}$$

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

Eks:

x	y
109	9.5
181	15.9
112	9.39
102	8.75
158	11.5

$\hat{\beta} = 0.0790$ (79.000 kr/m²) ← \bar{y}_{um} SWAP

$\hat{\alpha} = 0.555$ ← $\bar{x}_{w,b}$ SWAP

Beste tilpasning = regressionslinjen

$Y = 0.555 + 0.079 X$

Forklaring: (α, β) slik ved $\beta = \frac{S_{xy}}{S_x^2}$, $\alpha = \bar{y} - \beta \bar{x}$

(det stasjonære pkt.) gir minimum:

$$f''_{\alpha\alpha} = -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i)'_{\alpha} = -2 \cdot \sum_{i=1}^n (-1) = -2 \cdot (-1) \cdot n = \underline{2n}$$

$$f''_{\alpha\beta} = -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i)'_{\beta} = -2 \sum_{i=1}^n (-x_i) = 2 \sum x_i = \underline{2n \bar{x}}$$

$$f''_{\beta\beta} = \sum_{i=1}^n \left[2(y_i - \alpha - \beta x_i)(-x_i) \right]_{\beta} = 2 \sum_{i=1}^n (-x_i y_i + \alpha x_i + \beta x_i^2)_{\beta}$$

$$= \underline{2 \sum_{i=1}^n x_i^2}$$

↓

$$H(f) = \begin{pmatrix} 2n & 2n \bar{x} \\ 2n \bar{x} & 2 \sum_{i=1}^n x_i^2 \end{pmatrix}$$

konstant matrise siden
f har grad 2 i (α, β)

$$\det = (2n \cdot 2 \sum x_i^2) - (2n \bar{x})^2$$

$$= 4n \cdot (\sum x_i^2) - 4n^2 \bar{x}$$

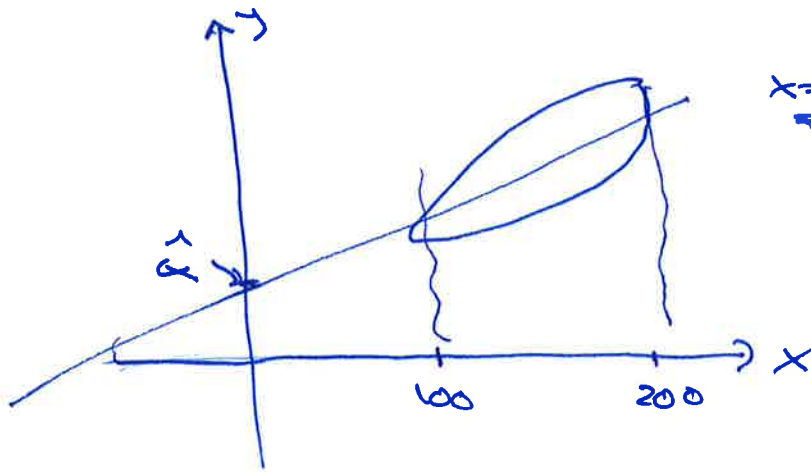
$$= 4n (\sum x_i^2 - n \bar{x})$$

$$= 4n ((n-1) S_x^2) = 4n(n-1) S_x^2 > 0$$

$$A = 2n > 0$$

↓

(α, β) er minimumsplet



$x=100: \hat{y} = 0.555 + 7.9 = 8.455$
 (forventet pris for leil. på 100 m²)

$Y = \alpha + \beta X + \epsilon$
 $E(Y) = E(\alpha + \beta X + \epsilon)$
 $= \alpha + \beta X + 0$
 $= \alpha + \beta X = \hat{y}$
 $Var(Y) = Var(\epsilon) = \sigma^2$

$Y = 0.555 + 0.079X$
 ↑ ↑
 $\hat{\alpha}$ $\hat{\beta}$

$SS_E = \sum_{i=1}^n \epsilon_i^2$
 $= \sum_{i=1}^n (y_i - \hat{y}_i)^2$

For en gitt x :
 $Y \sim N(\hat{y}, \sigma^2)$

$x=112: \hat{y} = 9.397$
 $112 \begin{matrix} \hat{y} \\ y, m \end{matrix} \rightarrow 9.397$

2) Statistiske egenskaper

Resultat: $E(\hat{\alpha}) = \alpha, E(\hat{\beta}) = \beta$ ← forventningsrette estimatorer

$SE(\hat{\alpha}) = \sqrt{Var(\hat{\alpha})}$ $SE(\hat{\alpha})^2 = \frac{\sigma^2 \cdot \sum x_i^2}{n(\sum x_i^2 - n\bar{x}^2)}$
 $SE(\hat{\beta}) = \sqrt{Var(\hat{\beta})}$ $SE(\hat{\beta})^2 = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2}$

Praktisk bruk: Må estimere σ^2

$s^2 = \frac{1}{n-2} \sum_{i=1}^n \epsilon_i^2$ ← estimat for σ^2
 $= \frac{1}{n-2} \cdot SSE$

Hvordan finne std. feilen i $\hat{\alpha}$, $\hat{\beta}$:

- Finn regresjonslinjen $(\hat{\alpha}, \hat{\beta})$

- Regn ut $SSE = \sum \varepsilon_i^2 = \sum (y_i - \hat{y}_i)^2$

og sett $s^2 = \frac{1}{n-2} \cdot SSE$.

- Bruk $\sigma^2 = s^2$ i formelene for $\text{Var}(\hat{\alpha})$, $\text{Var}(\hat{\beta})$:

$$\text{Var}(\hat{\alpha}) = \frac{\sigma^2 \sum x_i^2}{n(\sum x_i^2 - n\bar{x}^2)} \Rightarrow \text{SE}(\hat{\alpha}) = \sqrt{\text{Var}(\hat{\alpha})}$$

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2} = \frac{\sigma^2}{(n-1) \cdot s_x^2} \Rightarrow \text{SE}(\hat{\beta}) = \sqrt{\text{Var}(\hat{\beta})}$$

Bruk:

Hypotesetest for β : $\left. \begin{array}{l} H_0: \beta = 0 = \beta_0 \\ H_1: \beta \neq 0 = \beta_0 \end{array} \right\}$

Testobservator: $T = \frac{\hat{\beta} - \beta_0}{\text{SE}(\hat{\beta})} = \frac{\hat{\beta}}{\text{SE}(\hat{\beta})}$

Forkastningsområde: $|T| > t_{\alpha/2}$ $\alpha = \text{sign. nivå}$

\uparrow
t-fordelt med $n-2$ d.f.

Konfidensintervall for β :

$$\hat{\beta} \pm t_{\alpha/2}^{n-2} \cdot \text{SE}(\hat{\beta})$$

\uparrow
 $n-2$ d.f.
(frihetsgrader)

③ Tolkning av r^2 Oppdeling
av avvik:

$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i) \quad \dots$$

totalt avvik	avvik forklart av modellen	feil-ledd $\varepsilon_i =$ tilfeldig avvik
-----------------	-------------------------------------	--

Man kan vise:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Defin:

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2 = (n-1) s_y^2$$

total sum of squares

$$SS_E = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2$$

Summen av kvadratene
av feil-ledd (residuals)

$$SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Summen av kvadratene
av forklarte feil

||

$$SS_T = SS_R + SS_E$$

Merke: Jeg følger her skrivemåten i boken. I litteraturen er det vanlig å bruke TSS for SS_T (total sum of squares), ~~RSS~~ for SS_E (residual sum of squares) og ESS for SS_R (explained sum of squares). Altså:

$$SS_E = RSS$$

$$SS_R = ESS$$

Man kan vise:

$$r^2 = \frac{SSR}{SST}$$

↑
kvadraten av
korrelasjons-
koeffisienten r

↑
andelen av total variasjon
som forklares av modellen

Dette betyr:

- i) r^2 kan tolkes som andelen av variasjonen i Y som forklares av modellen (regresjonslinjen)
- ii) Når vi skal regne ut SS_E for å finne $SE(\hat{\beta})$ kan vi bruke følgende formel:

$$\begin{aligned} SS_E &= SS_T - SSR \\ &= SS_T - r^2 \cdot SS_T \\ &= SS_T \cdot (1 - r^2) \\ SS_E &= \underline{\underline{(n-1) s_y^2 \cdot (1 - r^2)}} \end{aligned}$$