

Forelesning 14

MET1190

Plan:

- ① Lineær regresjon
- ② Statistiske egenskaper

Repetisjon:

Forelesning 12-13 : Hypotesetester

$$\begin{cases} H_0: \text{komplement til } H_1 \\ H_1: \text{det som skal "vises"} \end{cases}$$

- Metoder:
- finne forhetsutsynsmåte $t < 1$
testobservatør (basert på α)
 \Rightarrow Forkoste H_0 hvis testobservatøren
harver i forhast.omr.
 - Behold H_0 ellers
 - bruke p-verdi
 $p =$ sannsynlighet for å få en
minst like ekspen verdi
på testobs. som den realiserte

Forklare H0: $\beta < \alpha$
Bekræfte H0: ellers

⑤ Linear regression

X: forklaringsvariabel (kontrollert)

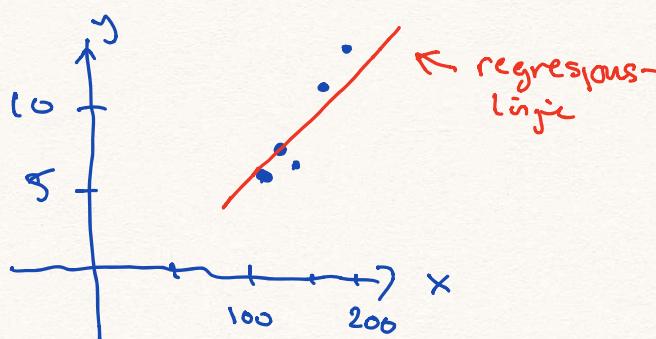
Y: responsvariabel (stokastisk)

Eks: Boligpriser

X: ant. kvar.

y: pris (mill. kr.)

x	y	\hat{y}
109	9.5	
181	15.9	
112	9.39	
102	8.75	
158	11.5	



multippel
linear regr.
→ mer enn én
forklarende var.

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

α, β : parameter med alene verdi
(α : sky. med Y-akser, β : stign. tall)

ε : feil-ledd, tilfeldig variasjon

Antagelser:

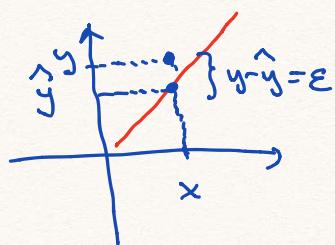
X : gitt verdi

ε : normalfordelt $N(0, \sigma^2)$ (o ulikt parameter)

$$Y = \underbrace{\alpha + \beta X}_{\text{regr. linjen } \hat{y}} + \varepsilon.$$

regr.
linjen \hat{y}

Y : normalfordelt $N(\hat{y}, \sigma^2)$



$$y = \alpha + \beta x$$

Tilpassing:

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

datasett

$$y_1 = \alpha + \beta x_1 + \varepsilon_1$$

$$y_2 = \alpha + \beta x_2 + \varepsilon_2$$

:

$$y_n = \alpha + \beta x_n + \varepsilon_n$$

Musk: α, β ulikt

Minste kvaradrekters
metode:

Ved α, β slik at

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

er mindst mulig.

Husk: α, β er de uløste (verkaledene)

$$\varepsilon_1 = y_1 - (\alpha + \beta x_1)$$

$$\varepsilon_2 = y_2 - (\alpha + \beta x_2)$$

⋮

$$\varepsilon_n = y_n - (\alpha + \beta x_n)$$

$$\varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2 = [y_1 - (\alpha + \beta x_1)]^2 + \dots + [y_n - (\alpha + \beta x_n)]^2$$

$$= y_1^2 - 2y_1 \cdot (\alpha + \beta x_1) + (\alpha + \beta x_1)^2 + \dots$$

$$\dots + y_n^2 - 2y_n \cdot (\alpha + \beta x_n) + (\alpha + \beta x_n)^2$$

$$= \underbrace{y_1^2}_{+ \dots} - \underbrace{2y_1\alpha}_{+ \dots} - \underbrace{2x_1y_1}_{+ \dots} + \underbrace{\alpha^2}_{+ \dots} + \underbrace{2\alpha\beta x_1}_{+ \dots} + \underbrace{\beta^2 x_1^2}_{+ \dots}$$

$$\begin{aligned} f(\alpha, \beta) &= n\alpha^2 + 2\alpha\beta(x_1 + x_2 + \dots + x_n) + \beta^2(x_1^2 + x_2^2 + \dots + x_n^2) \\ &\quad - 2\alpha(y_1 + y_2 + \dots + y_n) - 2\beta(x_1y_1 + x_2y_2 + \dots + x_ny_n) \\ &\quad + (y_1^2 + y_2^2 + \dots + y_n^2) \end{aligned}$$

- Mål:
- Setter inn $(x_1, y_1), \dots, (x_n, y_n)$ fra datasekket
 - løser minningsproblem
$$\min f(\alpha, \beta)$$

$$\text{Sett } f(\alpha, \beta) = \sum_{i=1}^n (\gamma_i - (\alpha + \beta x_i))^2$$

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

Vi viser optimiseringssproblemene: $\min f(\alpha, \beta)$

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

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

$$\sum y_i = \sum (\alpha + \beta x_i) = n \cdot \alpha + (\sum x_i) \beta$$

$$n \bar{y} = n \alpha + n \bar{x} \beta \quad | : n$$

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

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

$$= -2 \sum (x_i y_i - \alpha x_i - \beta x_i^2) = 0 \quad | : (-2)$$

$$\sum x_i y_i = \alpha \sum x_i + \beta \sum x_i^2$$

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

fra

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

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

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

$$\beta = \frac{s_{xy}}{s_x^2} = \frac{s_{xy}}{s_x \cdot s_y} \cdot \frac{s_y}{s_x} = r \frac{s_y}{s_x}$$

$$\beta = r \cdot \frac{s_y}{s_x}$$

Oppsummering:

$$\begin{cases} f'_{\alpha} = 0 \\ f'_{\beta} = 0 \end{cases}$$

gir

$$\begin{cases} \alpha^* = \bar{y} - \beta \bar{x} \\ \beta^* = r \cdot \frac{s_y}{s_x} \end{cases}$$

dvs
største innre
vinkel

Klassifisering:

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

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

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

$$= \sum -2x_i (-x_i) = \underline{2 \sum_{i=1}^n x_i^2}$$

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

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

$$= 4n (\sum x_i^2 - n \bar{x}^2) = 4n \cdot (n-1) s_x^2 \geq 0$$

Sånn $f(\alpha, \beta)$ er kvaravstrikk uttrykk
med $\det H(t) = 4n(n-1)s_x^2 \geq 0$
og $\text{tr } H(t) = 2n + 2\sum s_i^2 > 0$

så er det stasjonære pt.

$$\left\{ \begin{array}{l} \alpha^* = \bar{y} - \beta^* \bar{x} \\ \beta^* = r \cdot \frac{s_y}{s_x} \end{array} \right.$$

et minimum for $f(\alpha, \beta)$.

Dette ptet gir derfor den
rette linjen $y = \alpha + \beta x$ som passer
best til datasekten.

Hovedtrekk:

i) Finn stasjonære plt.: (α^*, β^*)

$$\begin{cases} f'_\alpha = 0 \\ f'_\beta = 0 \end{cases}$$

ii) Klassifisere stasjonære plt.

$$H(f)(\alpha^*, \beta^*)$$

Resultat: et stasjonært plt,
lokalt og globalt min.

Kan rive mellomregningene på notene
eller førelsesnig Resultat:

$$\beta^* = \frac{s_{xy}}{s_x^2} \quad \alpha^* = \bar{y} - \hat{\beta} \bar{x}$$

Konklusjon:

Regressjonslinjen er gitt ved $\hat{y} = \hat{\alpha} + \hat{\beta}x$
der

$$\hat{\beta} = \frac{s_{xy}}{s_x^2} = r \cdot \frac{s_y}{s_x}$$

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

$\hat{\alpha}, \hat{\beta}$: estimater for α, β

Beregning av $\hat{\alpha}$, $\hat{\beta}$:

x	y
109	9.5
108	15.9
112	9.39
102	8.75
158	11.5

$$\hat{\beta} = r \cdot \frac{s_y}{s_x}$$

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

C STAT

109 INPUT 9.5 $\Sigma +$
 :
 158 INPUT 11.5 $\Sigma +$

\bar{x}, r Swap $0.946 = r$

\bar{s}_x $38.0 = s_x$

\bar{s}_x, s_y Swap $2.92 = s_y$

\bar{x} $132.4 = \bar{x}$
 \bar{x}, \bar{y} Swap $11 = \bar{y}$

x (forsle) : forsl. var.
 y (andare) : respons var.

en kvadratnetter
 extra oljer pris
 med 79.000 kr.

$\hat{\alpha}, \hat{\beta}$ direkt:

$\hat{\beta} :$ \bar{y}, m Swap

$$\hat{\beta} = 0.079$$

$\hat{\alpha} :$ \bar{x}, b Swap

$$\hat{\alpha} = 0.555$$

Estimat för regressionslinjen:

$$y = 0.555 + 0.079x$$

Predictions!

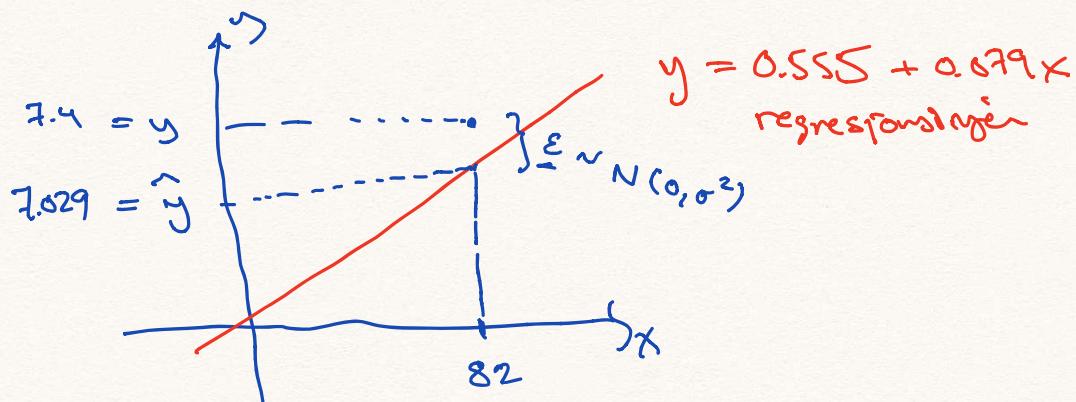
$$y = 0.555 + 0.079x$$

$$\underline{x=82:} \quad \hat{y} = 0.555 + 0.079 \cdot 82$$

$$\underline{(82 \text{ km})} \quad \underline{\text{Kalk:}} \quad 82 \quad \boxed{\hat{Y}_{\text{km}}} \rightarrow 7.029 = \hat{y}$$

$$y \sim N(\hat{y}, \sigma^2)$$

$$\underline{x=82:} \quad y \sim N(7.029, \sigma^2)$$



Husk:

Anter at $\epsilon \sim N(0, \sigma^2)$

Anter at vi har et tilfeldig utvalg

slik at $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ er uavhengige

② Statistiske egenskaper

Resultat: $E(\hat{\alpha}) = \alpha$ } forventningsrette
 $E(\hat{\beta}) = \beta$ } estimatører

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

$$\hat{\beta} = r \cdot \frac{s_y}{s_x}$$

Resultat: $\text{Var}(\hat{\alpha}) = \sigma^2 \cdot \frac{\sum_{i=1}^n x_i^2}{n \cdot (\sum_{i=1}^n x_i^2 - n \bar{x}^2)}$ } $\sigma^2 = \text{Var}(\varepsilon)$

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

$$SE(\hat{\alpha}) = \sqrt{\text{Var}(\hat{\alpha})}$$

$$SE(\hat{\beta}) = \sqrt{\text{Var}(\hat{\beta})}$$

} std. feil -
std. avvik

Praktisk bruk: Må estimer σ^2 ved å
bruke s^2

$$|| \quad s^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{SSE}{n-2}$$

$$SSE = \sum_{i=1}^n e_i^2 = \varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2$$

Eks: Hypotestest for β

(
X: förklaringsvariabel
Y: responsvariabel)

$H_0: \beta = 0 \stackrel{\beta_0 = 0}{\leftarrow}$ ingen samband
nållan X os Y
 $H_1: \beta \neq 0 \leftarrow$ det finns en samband
nållan X os Y

eller

$H_0: \beta \leq 0 \stackrel{\beta_0}{\leftarrow}$ ingen eller neg. samband
nållan X os Y
 $H_1: \beta > 0 \leftarrow$ pos. samband

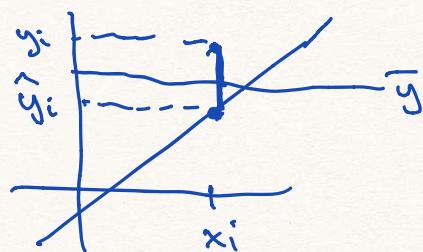
Testobserveator: T eller $\hat{\beta}$

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

- kan regna ut $\hat{\beta} = r \cdot \frac{S_y}{S_x}$ fra dataserdet
- må estinera $SE(\hat{\beta})$ för att föra T.

Tolkning av R^2 :

Oppdeling av
avvik:



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

totalt
avvik

avvik
fortalt
av
regressjønn

feilredd
 ϵ_i

Man kan vise:

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

SS_T

SS_R

SS_E

$$\epsilon_1^2 + \dots + \epsilon_n^2$$

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

$$= \frac{1}{n-1} \cdot SS_T$$

↑

$$SS_T = (n-1) S_y^2$$

↑
kan finne
dette på
kalkulator

$$SS_T = SS_R + SS_E$$

↑

$$SS_E = SS_T - SS_R$$

↑
Ønsker
å regne
ut

Man kan vise:

$$r^2 = \frac{SS_R}{SS_T}$$

↳

$$SS_R = r^2 \cdot SS_T$$

Tolkning:

r^2 er andelen av total varians forklart av regressionsmodellen

Vi kombinerer disse formlene:

$$\begin{aligned} SSE &= SS_T - SS_R \\ &= SS_T - r^2 \cdot SS_T \\ &= SS_T (1 - r^2) \end{aligned}$$

$$\boxed{\begin{aligned} SSE &= (n-1) s_y^2 \cdot (1 - r^2) \\ &= \epsilon_1^2 + \epsilon_2^2 + \dots + \epsilon_n^2 \end{aligned}}$$

Bruker dette slik:

$$\begin{aligned} \sigma^2 \text{ estimeres av } S^2 &= \frac{1}{n-2} \cdot SSE \\ &= \frac{1}{n-2} \cdot (n-1) s_y^2 \cdot (1 - r^2) \end{aligned}$$

$$\boxed{S^2 = \frac{n-1}{n-2} s_y^2 (1 - r^2)}$$

$$\underbrace{\frac{1}{n-1} SSt}_{S_y^2} = \underbrace{\frac{1}{n-1} SSR}_{\text{der durch das Modell erklärbare Varianzanteil}} + \underbrace{\frac{1}{n-1} SSE}_{\text{der durch das Modell nicht erklärbare Varianzanteil}}$$

Varianzen
(total Varianz)
 $i \gamma$

der durch das Modell
erklärbare
Varianz
erklärt
aus modellieren
(Regr. linig)

der durch das
Modell nicht
erklärbare
Varianz
erklärt
aus modellieren.