

Enkel lineær, forts.

$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma^2)$ uavh.

MKE $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$, $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

Form. uttrykk med varians $\frac{\sigma^2}{S_{xx}}$ og $\sigma^2 (\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})$

Spesielt \bar{y} og $\hat{\beta}_1$ ukorreleret

Dersom vi er i at $t_1 = \frac{\hat{\beta}_1}{\sqrt{\frac{\sigma^2}{S_{xx}}}} \sim t_{n-2}$ finner vi at $H_0: \beta_1 = 0$, da $t^2 = \frac{SSE}{n-2} = \frac{1}{n-2} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

$\Rightarrow F = t_1^2 \sim F_{1, n-2}$

Rammen for ANOVA-tabellen for 1-vekt. analyse.

Kvadrattotsum oppspaltning, n^2 minst

$SST = \sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2$

da $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1 (x_i - \bar{x})$

Dersom, kan vi at $\hat{\beta}_1^2 S_{xx} = \sum (\bar{y} - \hat{y}_i)^2 = SSR$

hva det om normal til seg.

Aktuelt $SST = SSR + SSE$

Til SST svarer $n-1$ frihetsgrader

SSE	$n-2$	$n-1$
SSR	1	1

okt 24-12:16

Dersom $\beta_1 = 0$, $\hat{\beta}_1 \sim N(0, \sigma^2/S_{xx})$

$\Rightarrow E[\hat{\beta}_1^2 S_{xx}] = \text{Var}(\hat{\beta}_1 \sqrt{S_{xx}}) = \sigma^2$

$E[SSR]$

Anova tabell

Kilde	SS	df	MS	F	P-vert
x	SSR	1	$SSR/1 = SSR$	SSR/MSE	
Residual	SSE	$n-2$	$MSE = \frac{SSE}{n-1}$		
Total	SST	$n-1$			

okt 24-12:35

Prediksjon + KI / PI

Ex) $x_i =$ barnets alder

$Y_i =$ lange funksjon (FEV_1, PEF_1, \dots)

Ng verdi x^* for forklaringsvariabel

Predikert verdi $\hat{y}_{x^*} = \hat{\beta}_0 + \hat{\beta}_1 x^* = E(\hat{y}_{x^*})$ $\left. \begin{matrix} \text{Estimert } E(Y_{x^*}) \\ \text{Predikert } Y_{x^*} \end{matrix} \right\}$ Punktverdi Y_{x^*}

Sen fremt på minimumspråk

$V_{\hat{y}_{x^*}} = V_{\hat{\beta}_0 + \hat{\beta}_1 x^*} = V_{\bar{y} + \hat{\beta}_1 (x^* - \bar{x})}$

$= V_{\bar{y}} + (x^* - \bar{x})^2 V_{\hat{\beta}_1}$

$= \sigma^2 (\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}})$

Har videre at $\frac{\hat{y}_{x^*} - E(\hat{y}_{x^*})}{\sqrt{\frac{\sigma^2}{S_{xx}} (1 + \frac{(x^* - \bar{x})^2}{S_{xx}})}} \sim t_{n-2}$

hvis $\varepsilon_i \sim N(0, \sigma^2)$ uavh

Rekke liden til at $\hat{y}_{x^*} \pm t_{\alpha/2} \sqrt{\dots}$

at $(1-2)100\%$ KI \pm $\beta_0 + \beta_1 x^*$

hvis $t_{\alpha} = (1-\alpha)100$ persentil i t_{n-2}

okt 24-12:43

For prediksjonsintervall nær vi på

$Y_{x^*} = \hat{y}_{x^*} + \varepsilon^*$, $\varepsilon^* \sim N(0, \sigma^2)$ uavh $\varepsilon_1, \dots, \varepsilon_n$

$E[Y_{x^*}] = E[\hat{y}_{x^*}] + E(\varepsilon^*) = \beta_0 + \beta_1 x^*$

$V_{Y_{x^*}} = V_{\hat{y}_{x^*}} + V_{\varepsilon^*}$

$= \sigma^2 (\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}) + \sigma^2$

$= \sigma^2 (1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}})$

Tilsv. $\frac{Y_{x^*} - (\beta_0 + \beta_1 x^*)}{\sqrt{\sigma^2 (1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}})}} \sim t_{n-2}$

$\Rightarrow \hat{y}_{x^*} \pm t_{\alpha/2} \sqrt{\sigma^2 (1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}})}$ \pm ut $(1-\alpha)100\%$ prediksjonsintervall for Y_{x^*}

okt 24-13:17

14.6 Modellrisikale

Y_i skiver den lineære regresjon modellen som

$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma^2)$ uavh.

Modellen kan bygges opp i 4 elementer

- (1) Linearitet $E(Y_i) = \beta_0 + \beta_1 x_i$
- (2) Konstant varians $\text{Var}(Y_i) = \text{Var}(\varepsilon_i) = \sigma^2$
- (3) Uavhengighet
- (4) Normalitet, $\varepsilon_i \sim N(0, \sigma^2)$, $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

Elementene a rett opp etter viktighet

Hvis (1) svikter - alvorlig - en hele modellen uvidelig.

Hvis (1) holder, men (2) eller (3) svikter så en parameter est $\hat{\beta}_0$ og $\hat{\beta}_1$ fortsatt foru. uttrykk, men standardfeil kan være gitt estimert. \Rightarrow Problema ved inferens.

Hvis (1), (2) og (3) holder, (4) svikter, og datasettet er stort (n stor), så vil testen KI likevel være tilnærmet OK. Likevel: Pass på utligner.

okt 24-13:31

Modell- og funksjonsrisikale ved å kontrollere

- (1) Predikerte verdi $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- (2) Residualer $e_i = Y_i - \hat{y}_i$

Norm gjennomsnittet i stedet for e_i

Standard: samle residualer $e_i = \frac{e_i}{\sqrt{1 - \frac{(x_i - \bar{x})^2}{S_{xx}}}}$

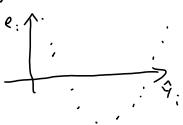
Men finne de standard: samle residualer ved e_i

kontrollere $\text{Var}(e_i) = \text{Var}(\hat{y}_i - Y_i) = \sigma^2 (1 - \frac{(x_i - \bar{x})^2}{S_{xx}})$

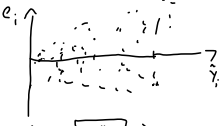
Shed i ven.

okt 24-13:40

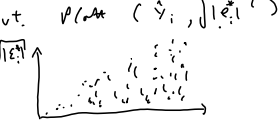
Sjekk (1) linearitet: Plot residuala e_i mot predikerte verdier \hat{y}_i .
 Kurvatur tyder på ikke linearitet.



Sjekk (2) Konstant varians
 Same plott e_i vs \hat{y}_i .
 Hvis større/småere varians i e_i med større \hat{y}_i indikerer dette at variansen øker med $E(Y_i) = \beta_0 + \beta_1 X_i$.



evt. plott $(\hat{y}_i, \sqrt{|e_i|})$
 Tyder også på ikke-konstant varians.



okt 24-13:50

Sjekk (3) uavhengighet
 Tidsserie $Y_i = \text{obs. dag } i$
 (a) $P(Y_i = e_i^*)$ ikke automatisk
 (b) $P(Y_{i-1} = e_{i-1}^*, Y_i = e_i^*)$; R

Sjekk (4) normalitet
 QQ plott, box plott, histogram
 over $e_i = (e_i^*)$.

okt 24-13:58