

12: Multiple linear regression (2nd week)

04.09.2017

So far: what to remember

Model:
$$\underset{n \times 1}{\hat{Y}} = \underset{n \times p}{\sum} \underset{p \times 1}{\beta} + \underset{n \times 1}{\epsilon}$$

where X (design matrix) has full rank (rank p) $n \gg p$

Normal model: $\epsilon \sim N_n(0, \sigma^2 I) \Rightarrow Y \sim N_n(X\beta, \sigma^2 I)$
 $\epsilon_1, \dots, \epsilon_n$ i.i.d $N(0, \sigma^2)$

ML estimation:
$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

 (LS)

$$\sim N_p(\beta, \sigma^2 (X^T X)^{-1})$$

REML estimation:
$$\hat{\sigma}^2 = \frac{1}{n-p} SSE$$

$$\hat{\epsilon} = Y - \hat{Y} = Y - X\hat{\beta}$$

$$\frac{\hat{\sigma}^2 (n-p)}{\sigma^2} \sim \chi^2_{n-p}$$

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

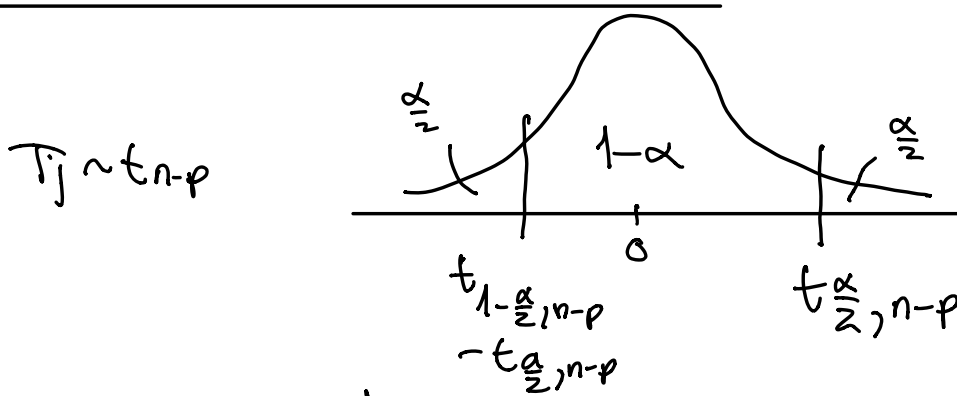
Elements of $\hat{\beta}$:
$$\hat{\beta}_j \sim N_1(\beta_j, \sigma^2 \underbrace{(X^T X)^{-1}}_{C_{jj}})$$

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{C_{jj}} \cdot \sigma} \sim N(0, 1)$$

$$T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{\frac{\hat{\sigma}^2 (n-p)}{\sigma^2} \cdot \frac{1}{n-p}}} = \frac{\hat{\beta}_j - \beta_j}{\sqrt{C_{jj}} \cdot \hat{\sigma}} \sim t_{n-p}$$

$$\frac{N(0,1)}{\sqrt{\frac{\chi^2_{n-p}}{n-p}}} \sim t_{n-p} \leftarrow \text{general rule}$$

Confidence interval: β_j (CI)



$$P(-t_{\frac{\alpha}{2}, n-p} \leq T_j \leq t_{\frac{\alpha}{2}, n-p}) = 1-\alpha$$

and solve for β_j
in the middle

Interpretation: 95% CI

fix X

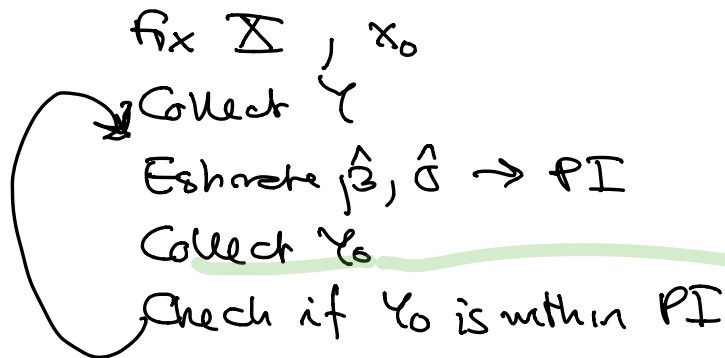
- Collect Y
- Construct CI from $\hat{\beta}_j$ and $\hat{\sigma}^2$
- Check if true β_j in CI

On average 95% of the CI will contain the true β_j .

Prediction interval (PI)

THA4315
H2w2

Interpretation : 95% PI



On average 95% of the PI's will contain the collected Y_0 .

Single hypothesis testing

$H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$ (think area)

	H_0 true $\beta_j = 0$	H_0 false $\beta_j \neq 0$
Reject H_0 say " $\beta_j \neq 0$ "	Type I error	Correct
Not reject H_0 say " $\beta_j = 0$ "	Correct	Type II error

Linear hypotheses

$$H_0: C\beta = d \quad \text{vs} \quad C\beta \neq d$$

$r \times p$ $p \times 1$
 $r \times 1$

A: large model

B: small model (rt0 true)

SSE

↑
A

SSE_{rt0}

↑
B

$$F_{obs} = \frac{\frac{1}{r}(SSE_{rt0} - SSE)}{\frac{SSE}{n-p}} \sim F_{r, n-p}$$

Anova tables: sequential change in SSE
(β_0 always in here)

$SSE(\beta_1, \beta_2, \dots, \beta_k) \leftarrow$ SSE full model

$SSE(\beta_1) =$ SSE when only β_1 (with x_1) in

$SSE(\beta_1, \beta_2) =$ SSE when both β_1 and β_2 (x_1 and x_2) in

$$SSE(\beta_2 | \beta_1) = SSE(\beta_1, \beta_2) - SSE(\beta_1)$$

\uparrow
added effect of β_2 compared to model
with β_1

$$F = \frac{\frac{1}{r} SSE(\beta_2 | \beta_1)}{\frac{SSE(\beta_1, \beta_2, \dots, \beta_k)}{n-p}}$$

\leftarrow added effect

full model \nearrow

$r =$ difference in
number of parameters
between model with (β_1, β_2)
and (β_1)

This is used in 'anova(tst)' in R.

Called type 1.

(Can not - as far as I understand - be written as
an hypothesis test with $Q_p = d$.)