

## PART 2: LINEAR REGRESSION

THM4267 L10

17.02.2017

[F3.2]

### Distribution of $\hat{\epsilon}$ (residuals)

$$\text{Residuals: } \hat{\epsilon} = Y - \hat{Y} = Y - X\hat{\beta} = Y - \underbrace{X(X^T X)^{-1} X^T Y}_{\hat{\beta}}$$

$$= Y - HY = (I - H)Y$$

$$E(\hat{\epsilon}) = E((I - H)Y) = (I - H) \underbrace{E(Y)}_{X\beta} = 0$$

$$= (I - H)X\beta = 0$$

project onto  
the space orth. to column  
space of  $X$

$$E(\hat{\epsilon}) = 0$$

$$Y = X\beta + \epsilon$$

$$E(Y) = X\beta$$

$$\text{Cov}(Y) = \text{Cov}(\epsilon) = \sigma^2 I$$

$$\text{or } (I - H)X\beta = (X - HX)\beta = 0$$

$$HX = \underbrace{X(X^T X)^{-1} X^T}_{\text{X}} X = X$$

$$\begin{aligned} \text{Cov}(\hat{\epsilon}) &= \text{Cov}((I - H)Y) = (I - H) \underbrace{\text{Cov}(Y)}_{\sigma^2 I} (I - H)^T \\ &= \sigma^2 (I - H) I (I - H) = \underline{\sigma^2 (I - H)} \end{aligned}$$

$$\text{Assume } \epsilon \sim N_n(0, \sigma^2 I) \Rightarrow \hat{\epsilon} \sim N_n(0, \sigma^2 (I - H))$$

NB:  $\text{rank}(H) = p$ ,  $\text{rank}(I - H) = n - p$ , which means  $(I - H)^{-1}$  does not exist and we use the singular version of the normal pdf.

## Distribution of SSE and $\hat{\sigma}^2$

$$\text{SSE} = \hat{\Sigma}^T \hat{\Sigma} = Y^T (I - H) (I - H) Y$$

$\uparrow$

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \hat{\Sigma} = (I - H) Y$$

$$\underline{\text{SSE} = Y^T (I - H) Y}$$

remember  $\text{rank}(I - H) = n - p$ .

RecEx3.p3 looks at the distribution of  $\frac{1}{\sigma^2} Y^T (I - H) Y = \frac{\text{SSE}}{\sigma^2}$   
by using the result on quadratic forms from Part 1  
(see slide)

$$Y \sim N_n(\bar{X}\beta, \sigma^2 I)$$

$$Y^* = \frac{1}{\sigma} (Y - \bar{X}\beta) \sim N_n(0, I)$$

$$\underbrace{Y^{*T} (I - H) Y^*}_{\begin{array}{l} (I - H)(Y - \bar{X}\beta) \\ (I - H)Y - \underbrace{(I - H)\bar{X}\beta}_0 \end{array}} = \frac{\frac{1}{\sigma^2} Y^T (I - H) Y}{\parallel} \underbrace{\frac{\text{SSE}}{\sigma^2} \sim \chi^2_{n-p}}$$

$\hat{\Sigma}^T \hat{\Sigma}$

$$\text{Now: } \hat{\sigma}^2 = \frac{1}{n-p} \text{SSE} \Leftrightarrow \text{SSE} = (n-p) \hat{\sigma}^2$$

$$V = \frac{\text{SSE}}{\sigma^2} = \frac{(n-p) \hat{\sigma}^2}{\sigma^2} \sim \underline{\chi^2_{n-p}}$$

$$E(V) = n-p \quad \text{Var}(V) = 2(n-p)$$

Is  $\hat{\sigma}^2$  an unbiased estimator?

$$\begin{aligned} E(\hat{\sigma}^2) &= E\left(\frac{1}{n-p} \underbrace{\text{SSE}}_{\sigma^2 \cdot V}\right) = \frac{1}{n-p} E(\sigma^2 V) \\ &= \frac{\sigma^2}{n-p} \underbrace{E(V)}_{n-p} = \underline{\underline{\sigma^2}} \quad \text{unbiased.} \end{aligned}$$

This is true when we assume  $\varepsilon \sim N$ .

If we do not assume  $\varepsilon \sim N$ , then we can use the trace-formula

$$\begin{aligned} E(\text{SSE}) &= E(Y^\top (I-H)Y) \quad E(Y) = \bar{X}\beta \\ &\quad \text{Cov}(Y) = \sigma^2 I \\ &= \text{tr}((I-H)\sigma^2 I) + (\bar{X}\beta)^\top \underbrace{(I-H)\bar{X}\beta}_0 \\ &= (n-p)\sigma^2 \rightarrow 0 \\ E(\hat{\sigma}^2) &= E\left(\frac{\text{SSE}}{n-p}\right) = \underline{\underline{\sigma^2}} \end{aligned}$$

## Inference about one $\beta_j$

Ex: Acid rain       $\beta_1 = \text{effect of SO}_4 \text{ on pH of lake}$

$$\hat{\beta}_1 = -0.315$$
$$SD(\hat{\beta}_1) = \sqrt{\sigma^2 \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right]_{1,1}} \quad \begin{matrix} \text{[diagonal elem]} \\ \text{corresp. to } \frac{SO_4}{x_1} \end{matrix}$$
$$\hat{SD}(\hat{\beta}_1) = \sqrt{\sigma^2 \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right]_{1,1}}$$

||  
St. Error  $\Rightarrow 0.0587$  in printout

$\hat{\sigma}$ : "Residual standard error" = 0.1165

$$\hat{\sigma} = \sqrt{\frac{SSE}{n-p}} \quad \begin{matrix} n=26 \\ p=8 \end{matrix} \quad n-p = 18$$

"on 18 degrees of freedom".

To find a confidence interval for  $\beta_j$  - or to test hypotheses about  $\beta_j$  we need to know the distribution of a statistic involving  $\hat{\beta}_j$  and  $\beta_j$  - with no other unknown parameters.

$$\hat{\beta}_j \sim N_1 \left( \beta_j, \sigma^2 \underbrace{\left[ (\mathbf{X}^T \mathbf{X})^{-1} \right]_{jj}}_{C_{jj}} \right)$$

and  $\hat{\sigma}^2 = \frac{SSE}{n-p}$  where  $\frac{SSE}{\sigma^2} \sim \chi^2_{n-p}$ .

Then :

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{C_{jj}\hat{\sigma}^2}} \xrightarrow{E(\hat{\beta}_j)} \sim N(0, 1)$$

$\downarrow SD(\hat{\beta}_j)$

$$SD(\hat{\beta}_j) = \sqrt{C_{jj}} \cdot \hat{\sigma} \quad \text{and } \hat{\beta}_j \text{ and } \hat{\sigma}^2 \text{ are independent} \\ (\text{to be shown})$$

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

General result :

$$\frac{N(0, 1)}{\sqrt{\frac{\chi^2_q}{q}}} \sim t_q$$

We have:

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

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

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

$\hat{\beta}_j$  and  $\hat{\sigma}^2$  need to be independent in order that this holds.  $\Rightarrow$  RecEx3.p3 + slides



Use  $T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{C_{jj}} \hat{\sigma}} \sim t_{n-p}$  for inference.

- a) Find a 95% confidence interval (CI) for  $\beta_j$

b) Test:

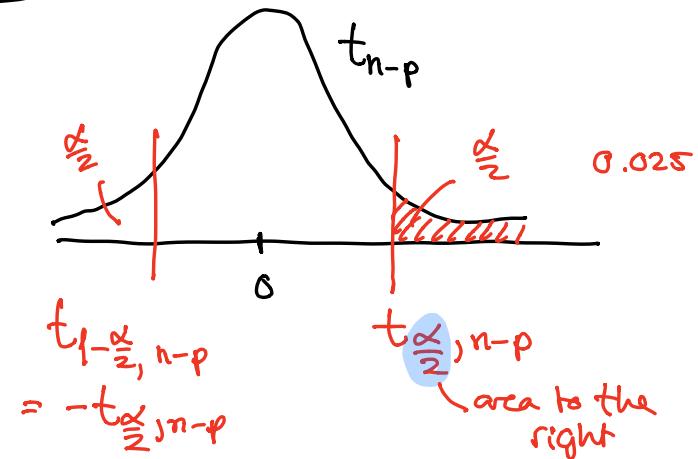
$$H_0: \beta_j = 0 \quad \text{vs} \quad H_1: \beta_j \neq 0$$

at significance level 0.05.

Ex: Acid rain  $\underline{\beta_1}$  = effect of SO<sub>4</sub> on pH.

a) 95% CI:

$$(1-\alpha) \cdot 100$$



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

↑

0.025                                    0.025                                    0.95

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{C_{jj}} \cdot \hat{\sigma}}$$

$$P(\underbrace{\hat{\beta}_j - t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{C_{jj}} \cdot \hat{\sigma}}_{\hat{\beta}_L} < \beta_j < \underbrace{\hat{\beta}_j + t_{\frac{\alpha}{2}, n-p} \sqrt{C_{jj}} \cdot \hat{\sigma}}_{\hat{\beta}_U}) = 1-\alpha$$

7

$L$  = lower  
 $U$  = upper

$$\text{Ex: } \begin{array}{l} \hat{\beta}_1 = -0.315 \\ \sqrt{C_{jj}} \sigma = 0.058 \\ n=26, r=8 \\ t_{0.025, 18} = 2.1 \end{array} \left. \right\} \begin{array}{l} -0.315 \pm 2.1 \cdot 0.058 \\ = \underline{\underline{[-0.44, -0.19]}} \end{array}$$

How do you interpret this interval?

→ Strong belief (95%) that  $\beta_j$  is in interval:

We see that 0 is not in the interval - what does this mean?  $\Rightarrow$  Reject  $H_0: \beta_j = 0$  vs  $H_1: \beta_j \neq 0$   
at sign. level 5%