

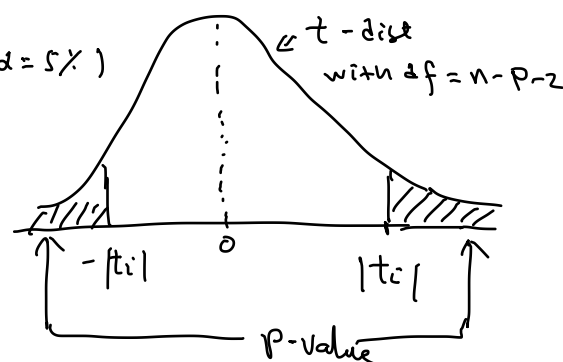
§1. Outliers. (cont'd from lecture 11)

* Suppose we have some prior information and aim to test the i^{th} case (pre-specified).

H_0 : i^{th} case is not outlier; H_A : i^{th} case is outlier.

We can use the p-value calculated from the t-test in lecture 11.

Declare i^{th} case is an
 { outlier if p-value $< \alpha$ (e.g. $\alpha = 5\%$)
 not outlier otherwise.



* If we don't have prior information and want to find/test an outlier then we usually choose the critical value (for the p-value)

to be $\frac{\alpha}{n} \times 100\%$ (usually $\alpha = 5\%$)

• This is called Bonferroni Correction (BC), popularly used in Multiple Testing Problems

• In this case,

H_0 : No outlier in the data, i.e. ~~case~~ 1, ..., case n all not outliers

H_A : at least one outlier.

• Our BC test procedure:
 (based on the t-test)

to control the Type I error.

For $i=1, \dots, n$

Step 1 we compute t -statistics values t_1, \dots, t_n (from lecture 11)

Step 2 from $t_i \stackrel{H_0}{\sim} t\text{-dist with } df = n-p-2$, we obtain their p -values
 $i=1, \dots, n$
 p_1, \dots, p_n

Step 3 If all p -values $\geq \frac{\alpha}{n}$ \Rightarrow no outliers

$$\Leftrightarrow \min_{i=1, \dots, n} p_i \geq \frac{\alpha}{n}$$

otherwise ($\min p_i < \frac{\alpha}{n}$) \Rightarrow exist outlier.

(Moreover, for $i=1, \dots, n$, if $p_i < \frac{\alpha}{n} \Rightarrow$ i th case is an outlier.

- We next show that why the critical value is taken as $\frac{\alpha}{n}$.

We want to control the Type I error, s.t.

$$P\left(\underset{\substack{\uparrow \\ \min p_i \\ i=1, \dots, n}}{\text{Reject } H_0} \mid \underset{\substack{\uparrow \\ \text{no outlier}}}{H_0 \text{ is true}} \right) \leq \alpha$$

$\min p_i < \alpha^*$ \leftarrow critical value

and we want to find α^* s.t. the above inequality holds.

- BC takes $\alpha^* = \frac{\alpha}{n}$, and this ensures the above inequality.

Why?

$$P\left(\min_i p_i < \frac{\alpha}{n} \mid H_0 \right)$$

$$= 1 - P\left(\min_i p_i \geq \frac{\alpha}{n} \mid H_0 \right)$$

$$= 1 - P\left(p_i \geq \frac{\alpha}{n} \text{ for all } i=1, \dots, n \mid H_0 \right)$$

If all p -values are independent \Rightarrow

$$= 1 - \prod_{i=1}^n P\left(p_i \geq \frac{\alpha}{n} \right)$$

Fact: under H_0 , the p -value follows a uniform distribution on $[0, 1]$ (practice)

$$= 1 - \prod_{i=1}^n \left(1 - \frac{\alpha}{n} \right)$$

$$\begin{aligned}
 & \text{If } P_i \sim \text{Unif}[0,1] \\
 & P(P_i < \alpha) = \alpha \\
 & \alpha \in [0,1]
 \end{aligned}
 \quad
 \begin{aligned}
 &= 1 - \left(1 - \frac{\alpha}{n}\right)^n \\
 &\leq 1 - \left(1 - \frac{\alpha}{n} \cdot n\right) \\
 &= \alpha.
 \end{aligned}$$

Example: If $n=65$, $p=3$. then $t_i \stackrel{H_0}{\sim} t\text{-dist df} = 65-3-2 = 60$

$$P(P_i < \alpha = 5\% \mid H_0) \approx 5\%.$$

$\Leftrightarrow \approx |t_i| > 2$
approx.

However,

$$P(\min_{i=1, \dots, n} P_i < 5\% \mid H_0) = 96.4\%$$

$$\text{If we use BC. } P(\min_i P_i < \frac{5\%}{65} \mid H_0) \leq 5\%$$

* A more efficient way to compute t_i 's.

$$\begin{aligned}
 t_i & \stackrel{\text{from lecture 11}}{=} \frac{\hat{e}_i}{\hat{\sigma}_{\hat{\beta}_i} \sqrt{1-h_{ii}}} \quad \leftarrow \text{from reg with all } n \text{ obs.} \\
 & \quad \uparrow \\
 & \quad \text{from reg with all but the } i\text{th obs.} \\
 & = r_i \underbrace{\left(\frac{n-p-2}{n-p-1-r_i^2} \right)^{\frac{1}{2}}}_{(**)}
 \end{aligned}$$

where r_i is the standardized residual for the i th case

$$r_i = \frac{\hat{e}_i}{\hat{\sigma} \sqrt{1-h_{ii}}}$$

where $\hat{\sigma}, \hat{e}_i$ are from regression with all data.

If we use the first equation (above), we need to do $(n+1)$ regressions
(from lecture 11) to get t_1, \dots, t_n .

On the other hand, using the second equation (**), we only need to do
one regression

§ 2. Influential cases. (§ 9. Applied linear Reg).

- A single case or small groups can be strongly influence the fit of a regression model.

Cases whose removal will cause major changes in the regression analysis are called influential cases.

- How to detect ?

We will use Cook's Distance
D

- For i th case, its Cook's distance is

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})^T (X^T X) (\hat{\beta}_{(i)} - \hat{\beta})}{(p+1) \hat{\sigma}^2}$$

$$\left[\text{Recall that } \text{Var}(\hat{\beta} | X) = \hat{\sigma}^2 (X^T X)^{-1} \right]$$

where $\hat{\beta}, \hat{\sigma}^2$ are OLS estimates from the regression with all observations

$\hat{\beta}_{(i)}$ is the OLS estimates with all observation but the i th obs.

$$y_j \sim x_j \text{ for } j=1, \dots, i-1, i+1, \dots, n$$

- In practice, usually 1 is used as a critical value for potential influential cases.

- A more efficient way to compute D_i is

$$D_i = \frac{1}{p+1} r_i^2 \frac{h_{ii}}{1-h_{ii}}$$

where r_i is the standardized residual and h_{ii} is the leverage