

# XI. Unusual Observations

0

Whole chapter  $\rightarrow M: Y|X \sim (X\beta, \sigma^2 I_n)$   
 $\text{rank}(X_{n \times k}) = k$

Interest: EXPLORATORY identification  
of unusual observations

• summary of (standard) notation

1

$\rightarrow$  all probabilistic statements will be  
conditional given  $X$

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \equiv \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$$

(both capital  
and small  
letters used)

as usual  $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \equiv$  vector of outcomes.

1

## 11.1 Leave-one-out and outlier model

2

NOTATION: for chosen  $t \in \{1, \dots, n\}$

$$Y_{(-t)} = Y \text{ without } Y_t$$

$$x_t^T = \text{row } t \text{ of the model matrix } X$$

$$X_{(-t)} = \text{matrix } X \text{ without row } t$$

$$j_t = (0, \dots, \underset{\substack{\uparrow \\ \text{place } t}}{1}, \dots, 0)^T \in \mathbb{R}^n$$

### Def 11.1 Leave-one-out model

3

The  $t^{\text{th}}$  leave-one-out model is a linear model

$$M_{(-t)}: Y_{(-t)} | X_{(-t)} \sim (X_{(-t)} \beta, \sigma^2 I_{n-1})$$

### Def 11.2 Outlier model

The  $t^{\text{th}}$  outlier model is a linear model

$$M_t^{\text{out}}: Y | X \sim (X\beta + j_t j_t^{\text{out}}, \sigma^2 I_n)$$

Meaning of the outlier model

$$E(Y_s | X_s = x_s) = x_s^T \beta, \quad s \neq t$$

$$E(Y_t | X_t = x_t) = \underbrace{x_t^T \beta}_{\text{can be whatever}} + j_t^{\text{out}}$$

$\equiv t^{\text{th}}$  observation does not follow model  $M$

21

## Lemma 11.1 Three equivalent statements

4

While assuming  $\text{rank}(X_{(-t)}) = k$ , the following three statements are equivalent:

(i)  $\text{rank}(X) = \text{rank}(X_{(-t)})$ , i.e.  $X_t \in \mathcal{N}(X_{(-t)}^T)$

(ii)  $m_{tt} > 0$ ;

(iii)  $\text{rank}(X_{j_t}) = k+1$ .

Further remarks:

(ii)  $\Leftrightarrow h_{tt} < 1$

(iii)  $\Leftrightarrow$  The  $t^{\text{th}}$  outlier model is of full-rank.

Proof: see the notes (exercise in linear algebra),  
not requested for exam

NOTATION (once conditions of Lemma 11.1 satisfied):

5

$$M_{(-t)}: Y_{(-t)} | X_{(-t)} \sim (X_{(-t)} \beta, \sigma^2 I_{m-1})$$

$$\rightarrow \hat{\beta}_{(-t)}, \hat{Y}_{(-t)}, SSE_{(-t)}, MSE_{(-t)}, \dots$$

$$M_t^{\text{out}}: Y | X \sim (X \beta + j_t \beta_t^{\text{out}}, \sigma^2 I_m)$$

$$\rightarrow \hat{\beta}_t^{\text{out}}, \hat{Y}_t^{\text{out}}, SSE_t^{\text{out}}, MSE_t^{\text{out}}, \dots$$

$$((\hat{\beta}_t^{\text{out}})^T, \hat{\beta}_t^{\text{out}})^T = \text{LSE of } ((\beta_t^{\text{out}})^T, \beta_t^{\text{out}})$$

~~etc~~

3

# Lemma 1.2 Equivalence of the outlier model and the leave-one-out model 6

1. The residual sums of squares in models  $M_{(t)}$  and  $M_t^{\text{out}}$  are the same, i.e.,  $SS_{e, (t)} = SS_{e, t}^{\text{out}}$ .
2. Vector  $\hat{\beta}_{(t)}$  solves the normal equations of model  $M_{(t)}$  if and only if a vector  $(\hat{\beta}_t^{\text{out}}, \hat{\eta}_t^{\text{out}})^T$  solves the normal equations of model  $M_t^{\text{out}}$ , where  $\hat{\beta}_t^{\text{out}} = \hat{\beta}_{(t)}$ ,  $\hat{\eta}_t^{\text{out}} = Y_t - X_t^T \hat{\beta}_{(t)}$ .

Proof: Remember: solution to normal equations must minimize the corresponding sum of squares.

$$\begin{aligned}
 M_t^{\text{out}}: SS_t^{\text{out}}(\beta, \eta_t^{\text{out}}) &= \|Y - X\beta - \mathbf{j}_t \eta_t^{\text{out}}\|^2 = \\
 &= \underbrace{\sum_{s \neq t} (Y_s - X_s^T \beta)^2}_{SS_{(t)}(\beta)} + (Y_t - X_t^T \beta - \eta_t^{\text{out}})^2 \\
 \|Y_{(t)} - X_{(t)}\beta\|^2 &= SS_{(t)}(\beta) = \text{sum of squares to be minimized in } M_{(t)}
 \end{aligned}$$

That is,  $SS_t^{\text{out}}(\beta, \eta_t^{\text{out}}) = SS_{(t)}(\beta) + \underbrace{(Y_t - X_t^T \beta - \eta_t^{\text{out}})^2}_{\geq 0}$

• For arbitrary  $\beta \in \mathbb{R}^k$ ,  $(Y_t - X_t^T \beta - \eta_t^{\text{out}})^2 = 0$  if  $\eta_t^{\text{out}} = Y_t - X_t^T \beta$ .

• Hence  $\min_{\beta, \eta_t^{\text{out}}} SS_t^{\text{out}}(\beta, \eta_t^{\text{out}}) = \min_{\beta} SS_{(t)}(\beta)$ .

and  $\hat{\beta}_{(t)}$  minimizes  $SS_{(t)}(\beta)$

$$\Rightarrow ((\hat{\beta}_{(t)}), \underbrace{Y_t - X_t^T \hat{\beta}_{(t)}}_{\hat{\eta}_t^{\text{out}}})^T \text{ minimizes } SS_t^{\text{out}}(\beta, \eta_t^{\text{out}})$$

NOTATION: leave-one-out least squares estimators of the response expectations

7

of  $m_{t,t} > 0$  for all  $t=1, \dots, n$ :

•  $\hat{Y}_{[t]} = X_t^T \hat{\beta}_{(-t)} = \text{LSE of } E(Y_t | X_t = x_t) \text{ based on a model } M_{(-t)}$

•  $\hat{Y}_{[0]} = (\hat{Y}_{[1]}, \dots, \hat{Y}_{[n]})^T = \text{LSE of } \mu = X\beta = E(Y|X), \text{ where each element } \mu_t \text{ of } \mu \text{ is estimated using model } M_{(-t)} \text{ (= data without obs. } t \text{)}$

$M_t^{\text{out}}: Y|X \sim (X\beta + j_t j_t^{\text{out}}, \sigma^2 I_n)$

8

= model with added regression to model

$M: Y|X \sim (X\beta, \sigma^2 I_n)$

Lemma 9.1 (under the condition  $m_{t,t} > 0$  = everything of full rank)

$$\hat{\beta}_t^{\text{out}} \stackrel{9.1}{=} \underbrace{(j_t^T M j_t)^{-1}}_{m_{t,t}} \underbrace{j_t^T U}_{U_t} = \frac{U_t}{m_{t,t}} \stackrel{11.2}{=} Y_t - X_t^T \hat{\beta}_{(-t)} \stackrel{\text{notation}}{=} Y_t - \hat{Y}_{[t]}$$

$$\hat{\beta}_{(-t)} \stackrel{11.2}{=} \hat{\beta}_t^{\text{out}} \stackrel{9.1}{=} \hat{\beta} - \frac{U_t}{m_{t,t}} (X^T X)^{-1} X_t$$

$$S_{Se, (t)} \stackrel{M.2}{=} S_{Se, t}^{out} \stackrel{9.1}{=} S_{Se} - \frac{U_t^2}{m_{tt}}$$

$$\stackrel{\text{algebra}}{=} S_{Se} - M_{Se} (U_t^{std})^2$$

$$\text{remember } U_t^{std} = \frac{U_t}{\sqrt{M_{Se} m_{tt}}}$$

$$\text{Further } M_{Se, (t)} = \frac{S_{Se, (t)}}{n-k-1}, \quad M_{Se, t}^{out} = \frac{S_{Se, t}^{out}}{n-(k+1)}$$

$$\Rightarrow \frac{M_{Se, (t)}}{M_{Se}} \stackrel{M.2}{=} \frac{M_{Se, t}^{out}}{M_{Se}} \stackrel{\text{algebra}}{=} \frac{n-k - (U_t^{std})^2}{n-k-1}$$

$$\text{Finally: } \hat{y}_t^{out} \stackrel{9.1}{=} \hat{y}_t + \frac{U_t}{m_{tt}} m_{tt}$$

↑  
row (or column)  
of matrix  $M$ .

Lemma 11.3 Quantities of the outlier and leave-one-out model expressed using quantities of the original model 9

Suppose that for given  $t \in \{1, \dots, n\}$ ,  $m_{tt} > 0$ . The following quantities of the outlier model  $M_t^{\text{out}}$  and the leave-one-out model  $M_{(-t)}$  are expressible using the quantities of the original model  $M$  as follows.

$$\hat{y}_t^{\text{out}} = y_t - x_t^T \hat{\beta}_{(-t)} = y_t - \hat{y}_{[t]} = \frac{U_t}{m_{tt}}$$

$$\hat{\beta}_{(-t)} = \hat{\beta}_t^{\text{out}} = \hat{\beta} - \frac{U_t}{m_{tt}} (X^T X)^{-1} x_t$$

$$SSE_{(-t)} = SSE_t^{\text{out}} = SSE - \frac{U_t^2}{m_{tt}} = SSE - MSE (U_t^{\text{std}})^2$$

$$\frac{MSE_{(-t)}}{MSE} = \frac{MSE_t^{\text{out}}}{MSE} = \frac{n-k - (U_t^{\text{std}})^2}{n-k-1}$$

Proof: Previous calculations while using Lemmas 9.1 and 11.2. □

TO REMEMBER:

- Lemma 11.3 provides quantification of influence of obs.  $t$  on  $\hat{\beta}$ ,  $MSE$ , ...
- all quantities can be calculated from the original fitted model (no need to estimate models  $M_{(-t)}$  or  $M_t^{\text{out}}$  for each  $t$ )  
= speciality of linear model

Def 11.3 Deleted residual

If  $m_{tt} > 0$ , then the quantity

$$j_{t,t}^{out} = y_t - \underbrace{\hat{y}_{[t]}}_{x_t^T \hat{\beta}_{(t-1)}} = \frac{u_t}{m_{tt}}$$

=  $\hat{E}(y_t | x_t = x_t)$  based on model without observation  $t$

is called the  $t^{th}$  deleted residual of model  $M$ .

## 11.2 Outliers

11

Observation  $t$  is outlier of model  $M$

if  $E(Y_t | X_t = x_t) \neq x_t^T \beta$  (for any  $\beta$ )

! outlier is a relative notion  
with respect to some model !

That is,

outlier

$$\Leftrightarrow E(Y_t | X_t = x_t) = x_t^T \beta + \hat{\beta}_t^{\text{out}}$$

for some  $\beta$ ,  $\hat{\beta}_t^{\text{out}}$

= correction of the model  
based expectation

↓ ~~other~~ observation with unusual  $Y$  value

Possible reasons for being outlier

- data error, really unusual, ...

### SOME CALCULATIONS:

•  $M_{tt} > 0 \Leftrightarrow M_t^{\text{out}}$  is of full rank and LSE  
of  $\hat{\beta}_t^{\text{out}}$  can be calculated

•  $t^{\text{th}}$  obs. is not outlier  $\Leftrightarrow H_0: \hat{\beta}_t^{\text{out}} = 0$

• under normality (but...)

$$T_t = \frac{\hat{\beta}_t^{\text{out}}}{\sqrt{\widehat{\text{var}}(\hat{\beta}_t^{\text{out}} | X_t)}}$$

$\sim t_{n-(k+1)}$

(LSE properties in  
model  $M_t^{\text{out}}$ )

$$\begin{aligned} \bullet \text{var}(\hat{\beta}_t^{\text{out}} | X) &= \text{var}\left(\frac{U_t}{m_{tt}} | X\right) = \\ &= \frac{1}{m_{tt}^2} \underbrace{\text{var}(U_t | X)}_{\sigma^2 m_{tt}} = \frac{\sigma^2}{m_{tt}} \end{aligned}$$

$\text{var}(U_t | X) = \sigma^2 m_{tt}$  in any model  $(M_t, M_t^{\text{out}})$  where we assume  $\text{var}(Y | X) = \sigma^2 I_n$  since  $U = MY$ ,  $\text{var}(U | X) = M \underbrace{\text{var}(Y | X)}_{\sigma^2 I_n} M^T = \sigma^2 M$ .

Hence  $\text{var}(\hat{\beta}_t^{\text{out}} | X) = \frac{\sigma^2}{m_{tt}} \xrightarrow{\text{from model } M_t^{\text{out}}} = \frac{\text{MSE}_{t,t}^{\text{out}}}{m_{tt}} = \frac{\text{MSE}_{t,t-t}}{m_{tt}}$

$$\Rightarrow T_t = \frac{\hat{\beta}_t^{\text{out}}}{\sqrt{\frac{\text{MSE}_{t,t-t}}{m_{tt}}}} = \frac{U_t}{m_{tt}} \cdot \sqrt{\frac{m_{tt}}{\text{MSE}_{t,t-t}}} =$$

$$= \frac{U_t}{\sqrt{\text{MSE}_{t,t-t} m_{tt}}}$$

= almost standardized residual

where is the difference?

12

## Def M.4 Studentized residual

If  $m_{tt} > 0$  then the quantity

$$T_t = \frac{y_t - \hat{y}_{[t]}}{\sqrt{MSE_{(-t)}}} \sqrt{m_{tt}} = \frac{U_t}{\sqrt{MSE_{(-t)} m_{tt}}}$$

is called the studentized residual of model  $M$ .

Remark: Previous calculations have shown:

$$\frac{MSE_{(-t)}}{MSE} = \frac{n-k - (U_t^{std})^2}{n-k-1}$$

algebra  $\Rightarrow$  
$$T_t = \sqrt{\frac{n-k-1}{n-k - (U_t^{std})^2}} \cdot U_t^{std}$$

$\rightarrow$   $T_t$  can be calculated without fitting  $M_t^{out}$  or  $M_{(-t)}$  ◻

## Lemma M.4 On studentized residuals

Let  $Y|X \sim N_n(X\beta, \sigma^2 I_n)$ ,  $\text{rank}(X_{n \times k}) = k < n$ .

Let further  $n > k+1$ . Let for given  $t \in \{1, \dots, n\}$   $m_{tt} > 0$ .

Then

1. The  $t^{\text{th}}$  studentized residuals  $T_t$  follows the Student  $t$ -distribution with  $n-k-1$  degrees of freedom.
2. If additionally  $n > k+2$  then  $E(T_t) = 0$ .
3. If additionally  $n > k+3$  then  $\text{Var}(T_t) = \frac{n-k-1}{n-k-3}$ .

Proof: 1. previous calculations.

2. & 3. standard properties of  $t$ -distribution.

## Test for outliers

15

$$M_t^{\text{out}}: Y|X \sim N_n(X\beta + j_t \beta_t^{\text{out}}, \sigma^2 I_n)$$

$$H_0: \beta_t^{\text{out}} = 0 \quad \equiv t^{\text{th}} \text{ not outlier of } M$$

$$H_1: \beta_t^{\text{out}} \neq 0 \quad \equiv t^{\text{th}} \text{ is outlier of } M$$

$$M: Y|X \sim N_n(X\beta, \sigma^2 I_n)$$

• Under  $H_0$ :  $T_t \sim t_{n-k-1}$ .

• Multiple testing problem (see later)

→ e.g. Bonferroni corrections of p-values

$$p_t^{\text{Bonf}} = \min\{1, p_t \cdot m\}$$

$p_t \equiv$  standard p-value for a single test

Example: Cars 2004.

16-19

→ Hummer H2 and Toyota Prius are both outliers v.r.t. considered model

some facts about outliers

20-22

Many "outliers"  $\equiv$  (usually) a wrong model

12

### 11.3 Leverage points

23

$\equiv$  observations with unusual  $x$ -values

$$H = X(X^T X)^{-1} X^T = (h_{t,s})_{t,s=1,\dots,n}$$

#### Leverage

24

a diagonal element  $h_{tt}$  ( $t=1, \dots, n$ ) of the hat matrix  $H$  is called the leverage of the  $t^{\text{th}}$  observation.

Consider now a model with intercept

25

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1,k-1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{n,k-1} \end{pmatrix}$$
$$\bar{x}^1 = \frac{1}{n} \sum_{i=1}^n x_{i1} \quad \dots \quad \bar{x}^{k-1} = \frac{1}{n} \sum_{i=1}^n x_{i,k-1}$$

$\rightarrow$  center non-intercept columns

$$\tilde{X} = (x^1 - \bar{x}^1 \mathbb{1}, \dots, x^{k-1} - \bar{x}^{k-1} \mathbb{1})$$

$$\rightarrow \mathcal{M}(X) = \mathcal{M}(\mathbb{1}, \tilde{X}) \quad \& \quad \mathbb{1}^T \tilde{X} = \mathbf{0}^T$$

(sum of values  
in each column  
is equal to 0)

$$H = X(X^T X)^{-1} X^T = (1, \tilde{X}) \begin{pmatrix} (1, \tilde{X})^T (1, \tilde{X}) & 0^T \\ \tilde{X}^T 1 & \tilde{X}^T \tilde{X} \end{pmatrix}^{-1} \begin{pmatrix} 1^T \\ \tilde{X}^T \end{pmatrix}$$

$$= (1, \tilde{X}) \begin{pmatrix} 1^T 1 & 1^T \tilde{X} \\ \tilde{X}^T 1 & \tilde{X}^T \tilde{X} \end{pmatrix}^{-1} \begin{pmatrix} 1^T \\ \tilde{X}^T \end{pmatrix} =$$

$$= (1, \tilde{X}) \begin{pmatrix} \frac{1}{n} & 0^T \\ 0 & (\tilde{X}^T \tilde{X})^{-1} \end{pmatrix} \begin{pmatrix} 1^T \\ \tilde{X}^T \end{pmatrix} =$$

$$= \frac{1}{n} 1 1^T + \tilde{X} (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T$$

$$\Rightarrow h_{tt} = \frac{1}{n} +$$

$$(x_{t,1} - \bar{x}^1, \dots, x_{t,k-1} - \bar{x}^{k-1}) (\tilde{X}^T \tilde{X})^{-1} \begin{pmatrix} x_{t,1} - \bar{x}^1 \\ \vdots \\ x_{t,k-1} - \bar{x}^{k-1} \end{pmatrix}$$

$\equiv$  square of a generalized distance  
of row  $t$   $(x_{t,1}, \dots, x_{t,k-1})^T$  from  
centroid over rows of  $X$   $(\bar{x}^1, \dots, \bar{x}^{k-1})^T$ .

Remember:  $H = Q Q^T$  ( $Q =$  orthonormal vector basis of  $\mathcal{V}(X)$ )

$$\sum_{i=1}^n h_{ii} = \text{tr}(H) = \text{tr}(Q Q^T) = \text{tr}(Q^T Q) = \text{tr}(I_k) = k.$$

$$\Rightarrow \bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{k}{n}.$$

2 software

$t^{\text{th}}$  observation is a leverage point

$$\text{if } h_{tt} > \frac{3k}{n}.$$

Leverage can be INFLUENTIAL ("dangerous"):

$$\text{var}(U_t | X) = \text{var}(Y_t - \hat{Y}_t | X) = \sigma^2 m_{tt} = \sigma^2 (1 - h_{tt})$$

$t = 1, \dots, n$

• High leverage  $\Rightarrow$  low  $\text{var}(U_t | X)$   
 $= \text{var}(Y_t - \hat{Y}_t | X)$

$\Rightarrow$  the  $t^{\text{th}}$  fitted value is forced to be close to the observed response value

Example: Cars 2004

29-32