## Transformation on X

— Naive method $\quad \dfrac{\max(x_i)}{\min(x_i)} > 10 \implies \log_e(x)$

— Partial Residual ~~model~~ plot

$x_\ell$ — non-linear relationship with $y$

$$y_i \approx \beta_0 + \beta_1 x_{i_1} + \cdots + \beta_{\ell-1} x_{i,\ell-1} + \boxed{p_\ell(x_{i,\ell})} + \beta_{\ell+1} x_{i,\ell+1} + \cdots$$
$$+ \beta_p x_{i,p}$$

least squares estimates for $\beta_0, \beta_1, \cdots, \boxed{\beta_\ell}, \cdots, \beta_p$

$\implies \boxed{\hat{y}_i} \approx \hat{\beta}_0 + \hat{\beta}_1 x_{i_1} + \cdots + \hat{\beta}_{\ell-1} x_{i,\ell-1} + \boxed{p_\ell(x_{i,\ell})} + \hat{\beta}_{\ell+1} x_{i,\ell+1} + \cdots + \hat{\beta}_p x_{i,p}$

$\implies \hat{e}_i = y_i - \hat{y}_i$

$\qquad = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i_1} + \cdots + \hat{\beta}_{\ell-1} x_{i,\ell-1} + \hat{\beta}_\ell , x_{i,\ell} + \hat{\beta}_{\ell+1} x_{i,\ell+1} + \cdots + \hat{\beta}_p x_{i,p}$

$\implies \hat{e}_i = p_\ell(x_{i,\ell}) - \hat{\beta}_\ell x_{i,\ell}$

$\implies p_\ell(x_{i,\ell}) = \underbrace{\hat{e}_i + \hat{\beta}_\ell x_{i,\ell}}_{\ell^{\text{th}} \text{ partial residual}}$

Partial residual plot

$\qquad \hat{e}_i + \hat{\beta}_\ell x_{i,\ell} \quad \text{vs} \quad x_{i,\ell}$

Subjective !



— Box - tidwell transformation
$\qquad$ — power transformation

— ~~Interplot~~ Interpolation spline / smoothing spline

$\qquad$ ~~Gen~~ Generalized additive model
$\qquad \implies$ non-linear part of $x_\ell$

eg. $y$ – area $m^2$ $\sqrt{y}$ $\log(y)$ $y$

$x$ – perimeter $m$ $x$ $\log(x)$ $x^2$

Outlier

Influential Point

} $\Leftarrow$ unusual obs.

Outlier — affect intercept. $\perp$ test $H_0$ & $p$-value

Influential point — affect reg. coeff.

Fig. 6.1

**(a) Single influential observation remote from center**

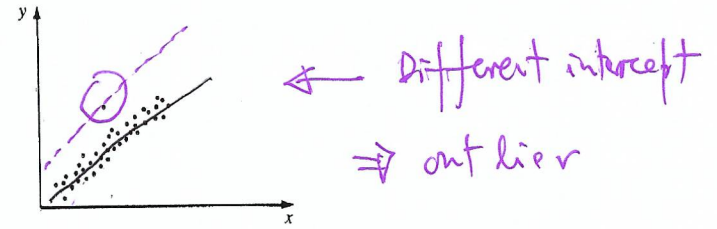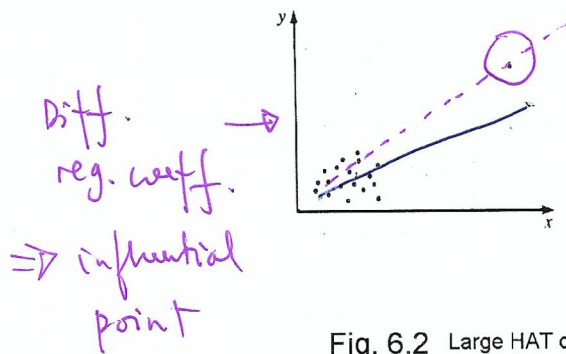**(b) Single observation with error in *y*-direction**



*Diff.*
*reg. coeff.* →
⇒ *influential point*

← *Different intercept*
⇒ *outlier*

Fig. 6.2 Large HAT diagonal but not influential observation
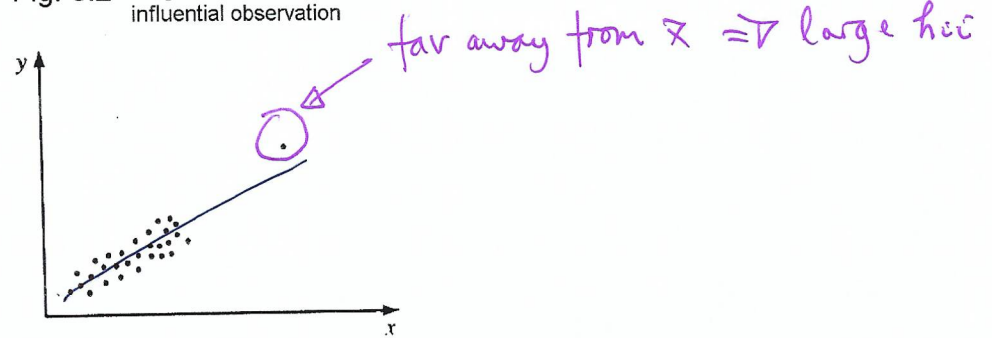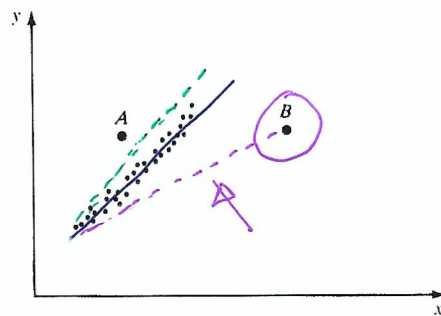


*far away from* $\bar{x}$ ⇒ *large* $h_{ii}$

Fig. 6.3 **Point *B* is clearly influential**

Is # an obs. an outlier ?

$H_0$: ith obs. is not an outlier

$$y_j = \beta_0 + \beta_1 x_{j1} + e_j \qquad j \neq i$$

$$y_i = \beta_0 + \delta + \beta_1 x_{i1} + e_i$$

$H_0$ = ith obs. is not an outlier

$$\Rightarrow H_0 : \delta = 0$$

Using results in Chapter 1

$$\Rightarrow \hat{\delta}, \ \text{s.e. of } \hat{\delta} \quad \longleftarrow \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

$$\Rightarrow t_i = \hat{\delta} / \text{s.e. of } \hat{\delta} \qquad Var(\hat{\beta}) = (X^T X)^{-1} \sigma^2$$

$$\Rightarrow t_i = \frac{\hat{e}_i}{\hat{\sigma}_{(-i)} \sqrt{1 - h_{ii}}} \quad \text{(externally studentized residual)}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \delta \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_{11} & 0 \\ \vdots & \vdots & \vdots \\ 1 & x_{i-1,1} & 0 \\ 1 & x_i & 1 \\ 1 & x_{i+1,1} & 0 \\ \vdots & \vdots & \vdots \\ 1 & x_{n,1} & 0 \end{pmatrix} \leftarrow \text{ith row}$$

$$\sim t \ (\text{\# of obs.} - \text{\# of unknown para})$$

$$\underset{n}{\uparrow} \qquad \underset{(1 + p + 1)}{\uparrow} \quad \overset{\delta}{\longleftarrow}$$

$$\frac{Res SS.}{d.f} \sim \chi^2$$

$$\Rightarrow \text{Can't reject } H_0 \Rightarrow \text{ith obs is not an outlier}$$

$$\text{Reject } H_0 \Rightarrow \text{ith obs is an outlier}$$

$$(\text{delete this obs})$$

$$t_1, \ \cdots, \ t_n \ - \ \text{multiple tests of } \underline{n \ obs.}$$

$$\Downarrow$$

$$n \ tests$$

$$H_0 : \delta_1 = 0, \ \cdots, \ H_0 : \delta_n = 0$$

Reject $H_0$ if

Bonferroni correction

$$\alpha \longrightarrow \frac{\alpha}{\text{\# of multiple tests}}$$

$$\underbrace{\phantom{\text{\# of multiple tests}}}_{n}$$

$$|t_i| \geq t_{\frac{\alpha}{2n}} (n - (p+2))$$

(4)

- If all $|t_i| < 3 \Rightarrow$ no outlier

- If $\max(t_i) >$ critical value, ~~dete~~ delete the corresponding obs.
  ~~the~~ & then check the 2nd largest of $|t_i|$
  If $\max(t_i) <$ critical value $\Rightarrow$ no outlier

## Influential point

| Name | Expression | Cutoff point |
|------|-----------|--------------|
| Student Residual | $r_i = \dfrac{\hat{e}_i}{\hat{\sigma}\sqrt{1-h_{ii}}}$ | $|r_i| > 2$ |
| Rstudent | $t_i = \dfrac{\hat{e}_i}{\hat{\sigma}_{-i}\sqrt{1-h_{ii}}}$ | $|t_i| > t_{\alpha/(2n)}$ |
| Hat Diag H | $h_{ii} = \boldsymbol{x_i}^T(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{x_i}$ | $h_{ii} > 2p'/n$ |
| Cook's D | $D_i = \left(\dfrac{t_i^2}{p'}\right)\left(\dfrac{h_{ii}}{1-h_{ii}}\right)$ | $D_i \gg 1$ |
| DFFITS | $(\text{DFFITS})_i = \dfrac{\hat{y}_i - \hat{y}_{i,-i}}{\hat{\sigma}_{-i}\sqrt{h_{ii}}} = (\text{Rstudent})_i\left(\dfrac{h_{ii}}{1-h_{ii}}\right)^{1/2}$ | $> 2\sqrt{p'/n}$ |
| $(\text{DFBETAS})_{j,i}$ | $(\text{DFBETAS})_{j,i} = \dfrac{\hat{\beta}_j - \hat{\beta}_{j,-i}}{\hat{\sigma}_{-i}\sqrt{c_{jj}}} = \dfrac{r_{j,i}}{\sqrt{\boldsymbol{r}_j'\boldsymbol{r}_j}}\dfrac{(\text{Rstudent})_i}{\sqrt{1-h_{ii}}}$ | $> 2/\sqrt{n}$ |
| Cov Ratio | $(\text{COVRATIO})_i = \dfrac{(\hat{\sigma}_{-i})^{2p'}}{\hat{\sigma}^{2p'}}\left(\dfrac{1}{1-h_{ii}}\right)$ | $> 1 + 3p'/n$ or $< 1 - 3p'/n$ |

outlier → Rstudent

Influential Point →

large $h_{ii}$
large $|t_i|$

$t_i$ (Cook's D)

$t_i$ (DFFITS)

$j = 1, \cdots, p$

$R = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T$

$p' \times n$

$r_{qs} = (q, s)^{th}$ element in R

e.g. ~~obs~~ $k^{th}$ obs — with large value of $h_{ii}$

large value of $|t_i|$

— it may be an influential point

$\hat{\beta}_j$ , $\hat{\beta}_{j,(-k)}$ changes $j = 1, \cdots, p$

If there is significant ~~diff~~. between $\hat{\beta}_j$ & $\hat{\beta}_{j(-k)}$, $j = 1, \cdots, p$

$\Rightarrow k^{th}$ obs is an influential point

1

6

Significant changes

$$\hat{\beta}_j \qquad\qquad \hat{\beta}_{j(-k)}$$

sign +ve $\Rightarrow$ sign. -ve

sign +ve/-ve $\Rightarrow$ insignificant

insignificant $\Rightarrow$ sign +ve/-ve

~~Multi multi~~

## Multicollinearity (among $X$)

Ridge regression $\qquad Var(\hat{\beta}) = \boxed{(X^T X)^{-1}} \sigma^2$

# Multicollinearity

e.g. $n = 8$

| $x_1$ | 10 | 10 | 10 | 10 | 15 | 15 | 15 | 15 |
|---|---|---|---|---|---|---|---|---|
| $x_2$ | 10 | 10 | 15 | 15 | 10 | 10 | 15 | 15 |

$\gamma_{12} = 0$ — linear independent (simple correlation coeff. between $x_1$ and $x_2$)

$$X_{i1}{}^* = \frac{X_{i1} - \bar{X}_1}{S_1}$$

$$X_{i2}{}^* = \frac{X_{i2} - \bar{X}_2}{S_2}$$

← Z-scores of $X_1$ & $X_2$

Sample mean $= 0$

Sample variance $= 1$

where $S_1^2 = \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$ and $S_2^2 = \sum_{i=1}^n (x_{i2} - \bar{x}_2)^2$

$$X^* = \begin{pmatrix} x_{11}^* & x_{12}^* \\ \vdots & \vdots \\ x_{n1}^* & x_{n2}^* \end{pmatrix}$$

$$\gamma_{12} = \sum_{i=1}^{n} (x_{i1}^* - \bar{x}_1^*)(x_{i2}^* - \bar{x}_2^*)$$

$$\sum_{i=1}^n x_{i1}^{*2} = \sum_{i=1}^n \left(\frac{x_{i1} - \bar{x}_1}{S_1}\right)^2$$

$$= \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}{S_1^2}$$

$$\sum_{i=1}^n x_{i1}^* x_{i2}^* = \sum_{i=1}^n \left(\frac{x_{i1} - \bar{x}_1}{S_1}\right)\left(\frac{x_{i2} - \bar{x}_2}{S_2}\right)$$

$$= \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{S_1 S_2}$$

$H_0 : \beta_1 = 0$

$$X^{*T} X^* = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \qquad (X^{*T} X^*)^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$t_1 = \frac{\hat{\beta}_1}{\sqrt{\text{Var}(\hat{\beta}_1)}} = \frac{\hat{\beta}_1}{\hat{\sigma}}$$

$$\text{Var}(\hat{\beta}_1) = \sigma^2 \qquad \text{Var}(\hat{\beta}_0) = \sigma^2 \qquad \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = 0$$

(first case)

e.g. $n = 8$

| $x_1$ | 10 | 11 | 11.9 | 12.7 | 13.3 | 14.2 | 14.7 | 15.0 |
|-------|----|----|------|------|------|------|------|------|
| $x_2$ | 10 | 11.4 | 12.2 | 12.5 | 13.2 | 13.9 | 14.4 | 15.0 |

$$\gamma_{12} = 0.99215 \text{ --- linear dependent}$$

$$X^{*T}X^* = \begin{pmatrix} 1 & 0.99215 \\ 0.99215 & 1 \end{pmatrix} \qquad (X^{*T}X^*)^{-1} = \begin{pmatrix} 63.94 & -63.44 \\ -63.44 & 63.94 \end{pmatrix}$$

$$\text{Var}(\hat{\beta}_1) = 63.94\sigma^2 \qquad \text{Var}(\hat{\beta}_0) = 63.94\sigma^2$$

Multicollinearity occurs when there are near linear dependenceies among the $x_j^*$ the column of $X^*$. That is, there is a set of constants (not all zero) for which $\sum_{j=1}^{p} c_j x_j^* \approx 0$ (2nd case)

Consider a regression with two perdictors:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + e_i$$
$$= \beta_0^* + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + e_i$$

$$X = \begin{pmatrix} 1 & x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 \end{pmatrix}, \qquad \beta = \begin{pmatrix} \beta_0^* \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

$$X^T X = \begin{pmatrix} n & 0 & 0 \\ 0 & \sum_{i=1}^{n}(x_{i1} - \bar{x}_1)^2 & \sum_{i=1}^{n}(x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \\ 0 & \sum_{i=1}^{n}(x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) & \sum_{i=1}^{n}(x_{i2} - \bar{x}_2)^2 \end{pmatrix} \qquad (X^T X)^{-1} = \begin{pmatrix} \frac{1}{n} & 0 & 0 \\ 0 & * & * \\ 0 & * & * \end{pmatrix}$$

$\Rightarrow X_1$ & $X_2$ strongly linear related

$\Uparrow$

problem of multicollinearity

— delete the indep. variable

— Ridge regression

**[right margin handwritten notes]**

$H_0 : \beta_1 = 0$

$$t_2 = \frac{\hat{\beta}_1 - 0}{S_e \cdot d \hat{\beta}_1}$$

$$= \frac{\hat{\beta}_1}{\sqrt{63.94\, \hat{\sigma}^2}}$$

$$= \frac{\hat{\beta}_1}{\hat{\sigma}\sqrt{63.94}}$$

$t_1 > t_2$ for the same $\hat{\beta}_1$

$\Rightarrow$ more diff. to reject $H_0$ for the 2nd case because $\gamma_{12}$ for 2nd case $= 0.99215$

$\boxed{9}$

When $p > 2$, $Var(\hat{\beta_1}) = \sigma^2 \left( \dfrac{1}{S_{x_1 x_1}} \right) \boxed{\left( \dfrac{1}{1 - R_j^2} \right)} - VIF$

$\boxed{R_j^2}$ = coeff. of determination of the $x_j$ regression on all other indep. variables $x_k$

$= 1 - \dfrac{Res\,S.S.}{total\,S.S.}$

$= \dfrac{Reg\,S.S}{total\,S.S.}$

$k \neq j$

Variance Inflation Factor

$R_j^2 \to 1 \Rightarrow Var(\hat{\beta_1}) \to$ large

$\Rightarrow t$ for testing $\beta_1 = 0 \to$ small

$- VIF > 10 \Rightarrow$ problem of multicollinearity $\to$ Can't reject $H_0$

$-$ condition index

$\underline{X^T X} \Rightarrow$ eigenvalues

$\sqrt{\dfrac{max(eigenvalue)}{each\ eigenvalue}} > 100 \Rightarrow$ problem of multicollinearity

Condition index

A condition index of 30 to 100 indicates moderate to strong collinearity

e.g.

| | VIF | |
|---|---|---|
| X1 | 9597.571 | $\leftarrow$ largest $> 10$ $\Rightarrow$ delete $X_1$ |
| X2 | 7.94 | |
| X3 | 8993.086 | |
| X4 | 23.29386 | |
| X5 | 4.27984 | |

$\Rightarrow$

| | |
|---|---|
| X2 | 7.9258 |
| X3 | 23.9268 |
| X4 | 12.7060 |
| X5 | 3.36087 |

— Descriptive stat.

~~quant~~ which one is $y$ ?   $y <$ quantitative
                                    binary   - -

quantitative variable / categorical variable

⟹ clean the data

— Full model

— Residual analysis — residual plot / Q-Q plot

ⓐ transformation on $y$ (Box - Cox transformation) ⟹ $\lambda$

   "      "   $x$  (naive, partial residual plot,
                    Box - tidewell transformation,
                    spline $\neq$ etc)

ⓑ outlier / influential point

ⓒ multicollinearity — before ~~to~~ you make transformation on $x$
                     — quantitative variable only

— Model selection

— Check residual again
  Add the originally deleted outlier / influential point
     into the best model

— hypothesis testing ?
  prediction ?