

## 9 Diagnostics and model building

### 9.1 Model validation and diagnostics

The model is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

#### 1. Residuals

$$\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}} = \mathbf{y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{y} = (\mathbf{I} - \mathbf{H})\boldsymbol{\epsilon}$$

where  $\mathbf{x}$  is  $n \times (k + 1)$  and the hat matrix (projection matrix) is

$$\mathbf{H} = \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}'$$

In terms of the elements  $h_{ij}$  of  $\mathbf{H}$ ,

$$\hat{\epsilon}_i = \epsilon_i - \sum_{j=1}^n h_{ij}\epsilon_j, \quad i = 1, 2, \dots, n.$$

#### Properties

- (a)  $E(\hat{\boldsymbol{\epsilon}}) = \mathbf{0}$  (residual mean is the same as the error mean)
- (b)  $\text{Cov}(\hat{\boldsymbol{\epsilon}}) = \sigma^2(\mathbf{I} - \mathbf{H})$  (residuals are NOT independent)
- (c)  $\text{Cov}(\hat{\boldsymbol{\epsilon}}, \mathbf{y}) = \sigma^2(\mathbf{I} - \mathbf{H})$  (residuals correlated with the observations)
- (d)  $\text{Cov}(\hat{\boldsymbol{\epsilon}}, \hat{\mathbf{Y}}) = \mathbf{0}$  (residuals uncorrelated with the predicted values)
- (e)  $\bar{\hat{\epsilon}} = \sum_{i=1}^n \hat{\epsilon}_i/n = \hat{\boldsymbol{\epsilon}}' \mathbf{1}/n = 0$
- (f)  $\hat{\boldsymbol{\epsilon}}' \mathbf{y} = SSE$
- (g)  $\hat{\boldsymbol{\epsilon}}' \hat{\mathbf{Y}} = 0$
- (h)  $\hat{\boldsymbol{\epsilon}}' \mathbf{X} = \mathbf{0}'$  ( $\hat{\boldsymbol{\epsilon}}'$  is orthogonal to each column of  $\mathbf{X}$ )

prove by yourself!

## 2. Model in centered form

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \epsilon_i \\ &= \alpha + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + \cdots + \beta_k(x_{ik} - \bar{x}_k) + \epsilon_i \end{aligned}$$

In matrix form, the centered model is

$$\mathbf{y} = (\mathbf{1}, \mathbf{X}_c) \begin{pmatrix} \alpha \\ \boldsymbol{\beta}_1 \end{pmatrix} + \boldsymbol{\epsilon},$$

where

$$\boldsymbol{\beta}_1 = (\beta_1, \beta_2, \dots, \beta_k)',$$

and

$$\mathbf{X}_c = \left( \mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{X}_1$$

It can be shown from previous notes that the least squares estimators of the parameters are

$$\begin{aligned} \hat{\alpha} &= \bar{y} \\ \hat{\boldsymbol{\beta}}_1 &= (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} \end{aligned}$$

Hence the predicted value is

$$\hat{\mathbf{y}} = \bar{y} \mathbf{1} + \mathbf{X}_c (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \mathbf{y} = \left( \frac{1}{n} \mathbf{1}' \mathbf{y} \right) \mathbf{1} + \mathbf{H}_c \mathbf{y} = \left( \frac{1}{n} \mathbf{J} + \mathbf{H}_c \right) \mathbf{y}$$

Since  $\hat{\mathbf{y}} = \mathbf{H} \mathbf{y}$ , hence

$$\mathbf{H} = \frac{1}{n} \mathbf{J} + \mathbf{H}_c = \frac{1}{n} \mathbf{J} + \mathbf{X}_c (\mathbf{X}'_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \quad (*)$$

$$\mathbf{H} = (h_{ij})_{i,j=1 \dots n}$$

$$\begin{aligned} h_{ii} &\geq \frac{1}{n} \\ \text{positive definite} & \quad \hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}' \mathbf{X})) \\ \mathbf{H} = \mathbf{H}^T & \Rightarrow h_{ii} = \sum_{j=1}^n h_{ij}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2 \\ & \quad \hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} \\ & \quad = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \\ & \quad = \mathbf{H} \mathbf{y} \end{aligned}$$

$$\begin{aligned} h_{ii} &\geq \frac{1}{n} \\ \text{ith row of } \mathbf{H} & \quad \text{if } h_{ii} = 0 \Rightarrow h_{ij} = 0 \quad \Rightarrow \quad h_{ii} \leq 1 \\ \text{if } h_{ii} &\geq 1 \quad \Rightarrow \quad h_{ij} = 0 \quad \Rightarrow \quad h_{ii} \leq 1 \\ \Rightarrow & \quad \frac{1}{n} \leq h_{ii} \leq 1 \end{aligned}$$

$\Rightarrow$

3. Hat matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \{h_{ij}\}$  (Let  $\mathbf{x}$  be a matrix with full column rank and with  $\mathbf{1}$  as its first column)

Properties

(a)  $1/n \leq h_{ii} \leq 1$  for  $i = 1, 2, \dots, n.$

(b)  $-.5 \leq h_{ij} \leq .5$  for all  $j \neq i.$

(c)  $h_{ii} = 1/n + (\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)'(\mathbf{x}_c' \mathbf{x}_c)^{-1}(\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)$ , where  $\mathbf{x}_{1i}' = (x_{i1}, x_{i2}, \dots, x_{ik})$ ,  $\bar{\mathbf{x}}_1' = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k)$ , and  $(\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)'$  is the  $i$ th row of the centered matrix  $\mathbf{x}_c$ .

(d)  $\text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii} = k + 1.$

$$h_{ii} = \underline{h_{ii}}^2 + \underline{h_{ij}}^2 + \sum_{r \neq i, j} h_{ir}^2$$

$$h_{ii} - h_{ii}^2 = h_{ij}^2 + \sum_{r \neq i, j} h_{ir}^2 \geq h_{ij}^2$$

$$h_{ij}^2 \leq h_{ii} - h_{ij}^2 = \left(\frac{1}{2}\right)^2 - (h_{ii} - \frac{1}{2})^2 \leq \frac{1}{4}$$

Residuals

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i \quad i=1, \dots, n$$

#### 4. Outliers

(a) Variance of the residuals is not constant

i. Studentized residual

ii. Studentized deleted residual (externally studentized residual)

(b) Deleted residuals

(c) Press (prediction sum of squares)

$$\varepsilon = y - \hat{X}\hat{\beta}, \text{ residual } \hat{\varepsilon} = y - \hat{X}\hat{\beta} = (\hat{I} - \hat{H})\hat{y}$$

$$\text{Var}(\hat{\varepsilon}) = (\hat{I} - \hat{H})\sigma^2 = (\hat{I} - \hat{H})\hat{\varepsilon}$$

$$\text{Var}(\hat{\varepsilon}_i) = \text{Var}(y_i - \hat{y}_i) = \sigma^2(1 - h_{ii})$$

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma} \sqrt{1-h_{ii}}} = \frac{\hat{\varepsilon}_i}{s \sqrt{1-h_{ii}}}, \quad s^2 = \hat{\sigma}^2 = \frac{\text{SSE}}{n-k-1}$$

$$t_i = \frac{\hat{\varepsilon}_i}{s_{(i)} \sqrt{1-h_{ii}}}, \quad \underline{s_{(i)}}^2 = \frac{\text{SSE}_{(i)}}{(n-1)-k-1} = \frac{\sum_{j \neq i} (y_j - \hat{y}_j)^2}{n-k-2}$$

(b) deleted residuals

$$\hat{\varepsilon}_{(i)} = y_i - \hat{y}_{(i)} = y_i - \hat{x}_{(i)} \hat{\beta}_{(i)}$$

estimate of  $\beta$

(c) PRESS (prediction sum of squares)

using data set

$$\{y_j, x_{ij}, j \neq i, i=1, \dots, n\}$$

$$= \sum_{i=1}^n \hat{\varepsilon}_{(i)}^2 = \sum_{i=1}^n \left( \frac{\hat{\varepsilon}_i}{1-h_{ii}} \right)^2$$

CROSS Validation

leave one out CV

$$CV = \frac{1}{n} \sum \hat{\varepsilon}_{(i)}^2 \rightarrow \text{MSE} = E(y_i - \hat{y}_{(i)})^2$$

Remark 9.1

if it is an influential observation, if

$\leftarrow$  This is correct!

$$D_i \gg r_{(k+1)}^2$$

### 5. Influential Observations

(a) Leverage  $h_{ii}$

(b) Cook's distance

make a large difference  
of  $\hat{\beta}$  with or without  
outliers using the observation /

$$\hat{\beta}_1, \hat{\beta}_{(ci)}$$

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

$$(A) \quad \hat{y} = \mathbf{X} \hat{\beta}, \quad \hat{y}_{(i)} = \sum_{j=1}^n h_{ij} y_j$$

$$= h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$$

$$1 = h_{ii} + \frac{\sum_{j \neq i} h_{ij}^2}{h_{ii}} \Rightarrow \begin{cases} h_{ii} \rightarrow 1, \text{ other } h_{ij}^2 \\ \text{will be very small.} \end{cases}$$

If  $h_{ii}$  is large,  $\Rightarrow$  the  $i$ -th observation is a high leverage point

$$\frac{1}{n} \sum h_{ii} = \frac{k+1}{n}$$

if  $h_{ii} \geq 2 \cdot \frac{k+1}{n} \rightarrow$  high leverage point

guideline

$$Dffits_i = \frac{\hat{y}_i - \hat{y}_{(i)}}{\hat{\sigma}_{(i)} \sqrt{h_{ii}}}$$

$$\hat{y}_{(i)} = \mathbf{x}' \hat{\beta}_{(i)}$$

guideline

the  $i$ -th observation is influential if

$$|Dffits_i| \geq 2 \cdot \sqrt{\frac{k+1}{68n}}$$

standardized difference of fitted value with or without using the  $i$ -th observation

**9.2 More complex models**

**9.3 Multicollinearity**

**9.4 Variable selection**