

Multiple regression model:

- ① Estimation
- ② Prediction
- ③ Model selection  $\rightarrow$  need transformation
- ④ Diagnostics

## 9 Diagnostics and model building

### 9.1 Model validation and diagnostics

The model is 确认

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

$$\begin{aligned}\hat{\mathbf{y}} &= (\mathbf{I} - \mathbf{H})\mathbf{y} = (\mathbf{I} - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= (\mathbf{X}\boldsymbol{\beta} - \mathbf{H}\mathbf{X}\boldsymbol{\beta}) + (\mathbf{I} - \mathbf{H})\boldsymbol{\epsilon} \\ (\mathbf{H}\mathbf{X} = \mathbf{X}) &= (\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta}) + (\mathbf{I} - \mathbf{H})\boldsymbol{\epsilon} \\ &= (\mathbf{I} - \mathbf{H})\boldsymbol{\epsilon}\end{aligned}$$

#### 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  $\mathbf{H}$ ,  $n \times n$

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

#### Properties

- proof by yourself {
- (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  $\mathbf{X}$ )

Properties:

(a)  $E(\hat{\underline{\epsilon}}) = E[(\underline{I} - \underline{H})\underline{\epsilon}] = 0$  ( $\underline{I} - \underline{H}$ ) is symmetric & idempotent.

(b)  $\text{Var}(\hat{\underline{\epsilon}}) = \text{Var}[(\underline{I} - \underline{H})\underline{\epsilon}] = (\underline{I} - \underline{H})\sigma^2 \underline{I} (\underline{I} - \underline{H})' = (\underline{I} - \underline{H})\sigma^2$

(c).

(d)

(e)

(f)  $\hat{\underline{\epsilon}}'\underline{y} = (\underline{y} - \underline{X}\hat{\underline{\beta}})'\underline{y} = \underline{y}'\underline{y} - \underline{\hat{\beta}}'\underline{X}'\underline{y} = \text{SSE}$

(g)

(h)

## 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

$$\underline{y} = (\mathbf{1}, \mathbf{X}_c) \begin{pmatrix} \alpha \\ \beta_1 \end{pmatrix} + \epsilon, \quad \mathbf{X}_c = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots \\ \vdots & \vdots & \ddots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots \end{pmatrix}$$

where

$$\beta_1 = (\beta_1, \beta_2, \dots, \beta_k)'$$

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

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

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

nonnegative  $\rightarrow$  elements in diagonal  $> 0$

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

Positive definite

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

$$h_{ii} \geq \frac{1}{n}$$

$$\mathbf{H} = \begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \cdots & h_{nn} \end{pmatrix}$$

$$\mathbf{H} = \mathbf{H}^2 \Rightarrow h_{ii} = \sum_{j=1}^n h_{ij}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$$

i-th row of H

$$\Rightarrow 1 = h_{ii} + \frac{\sum_{j \neq i} h_{ij}^2}{h_{ii}} \Rightarrow h_{ii} \leq 1$$

$$\Rightarrow \frac{1}{n} \leq h_{ii} \leq 1$$

$$\hat{\beta} \sim N(\hat{\beta}, \sigma^2 (\mathbf{X}' \mathbf{X})^{-1})$$

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta} = \underbrace{\mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'}_{\mathbf{H}} \mathbf{y} = \mathbf{H} \mathbf{y}$$

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  $i = 1, 2, \dots, n$ . leverage
- (b)  $-0.5 \leq h_{ij} \leq 0.5$  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)  $\text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii} = k + 1$ . (d).  $\text{tr}(\mathbf{H}) = \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = k+1$

$$h_{ii} = \sum_{r=1}^n h_{ir}^2$$

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

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

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

$$\text{i.e. } -\frac{1}{2} \leq h_{ij} \leq \frac{1}{2}$$

Prediction:  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} h_{11} & \dots & h_{1n} \\ h_{21} & \dots & h_{2n} \\ \vdots & \ddots & \vdots \\ h_{n1} & \dots & h_{nn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}, \quad i=1, \dots, n$$

$$\varepsilon_i \sim N(0, \sigma^2)$$

$$\text{or } \varepsilon_i \sim N(0, \sigma^2)$$

(i), (ii) are two ways to scale the residuals  $\Rightarrow$  same variance

#### 4. Outliers

(a) Variance of the residuals is not constant  $\text{Cov}(\hat{\varepsilon}) = \sigma^2(I - H) \Rightarrow \text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{ii})$

i. Studentized residual

ii. Studentized deleted residual (externally studentized residual)

(b) Deleted residuals

(c) Press (prediction sum of squares)

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

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

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

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

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

(当一个女孩被分到男孩子中时, 会被分离的更远一些)

(b) deleted residuals:

$$\hat{\varepsilon}_{(i)} = y_i - \hat{y}_{(i)} = y_i - x_i' \hat{\beta}_{(i)}, \quad \text{estimate of } \beta \text{ using data set } \{y_j, x_j, j \neq i, i=1, \dots, n\}$$

(c) Press:

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

Cross-Validation

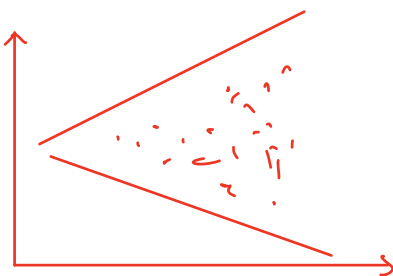
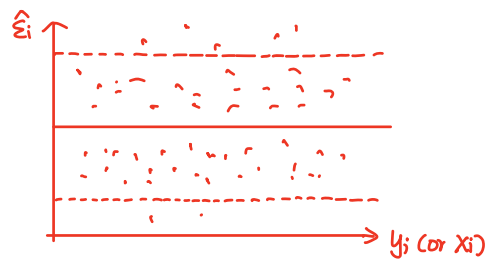
leave-one-out CV

$$CV = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_{(i)}^2 \rightarrow MSE = E[(y_i - \hat{y}_i)^2]$$

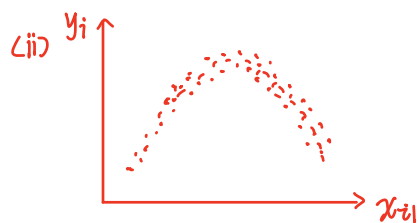
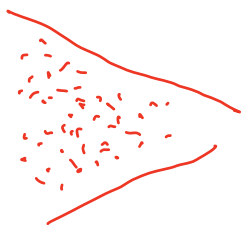
$$CV = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_{(i)}^2 \rightarrow MSE = E(y_i - \hat{y}_i)^2$$

Remark 9.1:

(i) check constant variance:



$$\text{Var}(e_i) \propto \hat{y}_i$$



(iii) Q-Q plot p-value.

Residual analysis.

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$\varepsilon_i = y_i - f(x_{i1}, \dots, x_{ik}, \beta)$$

$$\hat{\varepsilon}_i = y_i - \hat{y}_i = y_i - f(x_{i1}, \dots, x_{ik}, \hat{\beta}) \quad \beta = (\beta_0, \dots, \beta_k, \sigma^2)$$

Linear model:

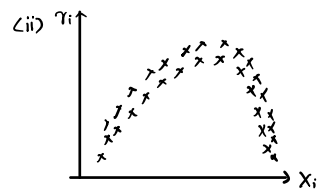
i) constant variance  $\text{Var}(y_i) = \text{Var}(\varepsilon_i) = \sigma^2$

ii) linearity

iii) Normality

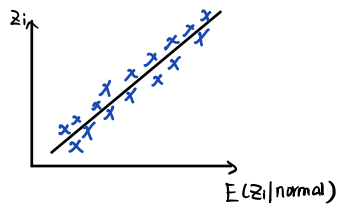
Residual plots.

i).



not linear  $\rightarrow$  maybe need  $X_i^2$  or ...

iii) Q-Q plot.  $|\varepsilon_i| \rightarrow$  ordered  $z_i$



outliers: make a large difference

## 5. Influential Observations

**Guideline**

(a) Leverage  $h_{ii}$

(b) Cook's distance

it is an influential observation, if  $D_i \geq F_{0.05, k+1, n-k-1}$

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

$$(a) \hat{y} = H y \quad \hat{y}_i = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$$

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

if  $h_{ii}$  is large  $\Rightarrow$  the  $i$ th observation is a high leverage point.

Q: How to judge whether  $h_{ii}$  is large

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

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

$$\text{Diffits} = \frac{\hat{y}_i - \hat{y}_{(i)}}{\hat{\sigma}_{(i)} \sqrt{h_{ii}}} \quad \text{where } \hat{y}_{(i)} = \tilde{x}_i' \tilde{\beta}_{(i)}$$

standardized deleted residual  
(standardized difference of fitted value with or without using the  $i$ th observation)

**Guideline:** the  $i$ th observation is influential if  $|\text{diffits}_i| > 2\sqrt{\frac{k+1}{n}}$

$$S_{(i)}^2 = \frac{SSE_{(i)}}{n-k-2}, \quad SSE_{(i)} = y_{(i)}' y_{(i)} - \tilde{\beta}_{(i)}' \tilde{x}_{(i)}' y_{(i)}$$

or  $= SSE - \hat{\epsilon}_i^2 / h_{ii}$