

Chapter 8: Multiple Regression: Model Validation and Diagnostics

1 Residuals

Consider the linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ again. The residual is defined as

$$\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} := \mathbf{y} - \hat{\mathbf{y}},$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. The fitted value is given by

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y},$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is called the hat matrix. The hat matrix has the following properties:

$$\mathbf{H}\mathbf{X} = \mathbf{X},$$

$$\mathbf{j} = \mathbf{H}\mathbf{j}, \quad \mathbf{x}_i = \mathbf{H}\mathbf{x}_i, \quad i = 1, 2, \dots, k,$$

and

$$\hat{\boldsymbol{\epsilon}} = (\mathbf{I} - \mathbf{H})\mathbf{y}.$$

In addition,

$$\hat{\boldsymbol{\epsilon}} = (\mathbf{I} - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = (\mathbf{I} - \mathbf{H})\boldsymbol{\epsilon}.$$

The following are some of the properties of $\hat{\boldsymbol{\epsilon}}$:

$$E(\hat{\boldsymbol{\epsilon}}) = \mathbf{0},$$

$$\text{Cov}(\hat{\boldsymbol{\epsilon}}) = \sigma^2(\mathbf{I} - \mathbf{H}),$$

$$\text{Cov}(\hat{\boldsymbol{\epsilon}}, \mathbf{y}) = \sigma^2(\mathbf{I} - \mathbf{H}),$$

$$\text{Cov}(\hat{\boldsymbol{\epsilon}}, \hat{\mathbf{y}}) = \mathbf{0},$$

$$\bar{\hat{\boldsymbol{\epsilon}}} = \sum_{i=1}^n \hat{\epsilon}_i / n = \hat{\boldsymbol{\epsilon}}^T \mathbf{j} / n = 0,$$

$$\hat{\boldsymbol{\epsilon}}^T \mathbf{y} = SSE = \mathbf{y}^T (\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{y} = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y},$$

$$\hat{\boldsymbol{\epsilon}}^T \hat{\mathbf{y}} = 0,$$

$$\hat{\boldsymbol{\epsilon}}^T \mathbf{X} = \mathbf{0}.$$

The above formulas imply that $\hat{\boldsymbol{\epsilon}}$ is orthogonal to $\hat{\boldsymbol{y}}$ and each column of \mathbf{X} , while it is correlated with \mathbf{y} . Therefore, we use the scatter plots $\hat{\boldsymbol{\epsilon}}$ versus $\hat{\boldsymbol{y}}$ to examine the pattern of the residual.

If the model is incorrect, various plots involving residuals may show departures from the fitted model such as outliers, curvature, or nonconstant variance.

2 The Hat Matrix

For the centered model,

$$\mathbf{y} = \alpha \mathbf{j} + \mathbf{X}_c \boldsymbol{\beta}_1 + \boldsymbol{\epsilon},$$

where $\mathbf{X}_c = (\mathbf{I} - \frac{1}{n} \mathbf{J}) \mathbf{X}_1$, and $\hat{\boldsymbol{y}}$ becomes

$$\hat{\boldsymbol{y}} = \hat{\alpha} \mathbf{j} + \mathbf{X}_c \hat{\boldsymbol{\beta}}_1,$$

and we can define

$$\mathbf{H}_c = \mathbf{X}_c (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T.$$

Therefore,

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

By arbitrariness of \mathbf{y} , we have

$$\mathbf{H} = \frac{1}{n} \mathbf{J} + \mathbf{H}_c = \frac{1}{n} \mathbf{J} + \mathbf{X}_c (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T.$$

Theorem 2.1. *If \mathbf{X} is $n \times (k+1)$ of rank $k+1 < n$, and if the first column of \mathbf{X} is \mathbf{j} , then the element h_{ij} of $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ have the following properties:*

- $1/n \leq h_{ii} \leq 1$ for $i = 1, 2, \dots, n$.
- $-0.5 \leq h_{ij} \leq 0.5$ for all $j \neq i$.
- $h_{ii} = 1/n + (\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)^T (\mathbf{X}_c \mathbf{X}_c)^{-1} (\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)$, where $\mathbf{x}_{1i}^T = (x_{i1}, x_{i2}, \dots, x_{ik})$, $\bar{\mathbf{x}}_1^T = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k)$, and $(\mathbf{x}_{1i} - \bar{\mathbf{x}}_1)^T$ is the i th row of the centered matrix \mathbf{X}_c .
- $\text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii} = k+1$.

Proof. (i) The lower bound follows from the relationship $\mathbf{H} = \frac{1}{n} \mathbf{J} + \mathbf{H}_c = \frac{1}{n} \mathbf{J} + \mathbf{X}_c (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T$, where $\mathbf{X}_c^T \mathbf{X}_c$ is positive definite. The upper bound follows from the property $H = H^2$, which implies

$$h_{ii} = \mathbf{h}_i^T \mathbf{h}_i = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2,$$

or, equivalently,

$$1 = h_{ii} + \sum_{j \neq i} h_{ij}^2 / h_{ii},$$

which implies $h_{ii} \leq 1$.

(ii) Since $h_{ii} = \mathbf{h}_i^T \mathbf{h}_i = h_{ii}^2 + h_{ij}^2 + \sum_{r \neq i,j} h_{ir}^2$, we have

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

and thus $h_{ij}^2 \leq h_{ii} - h_{ii}^2$. Since the maximum value of $h_{ii} - h_{ii}^2$ is $1/4$, we have $-0.5 \leq h_{ij} \leq 0.5$. \square

3 Outliers

For outlier analysis, we need to keep in mind that the variance of the residuals is not constant:

$$\text{Var}(\hat{\epsilon}) = \sigma^2(1 - h_{ii}).$$

An additional verification that large values of h_{ii} are accompanied by small residuals is provided by the property:

$$\frac{1}{n} \leq h_{ii} + \frac{\hat{\epsilon}_i^2}{\hat{\boldsymbol{\epsilon}}^T \hat{\boldsymbol{\epsilon}}} \leq 1.$$

There are two common methods of scaling of the residuals:

- Standardized residual:

$$r_i = \frac{\hat{\epsilon}_i}{s\sqrt{1 - h_{ii}}},$$

where $s = \sqrt{SSE/(n - k - 1)}$.

- Studentized residual:

$$t_i = \frac{\hat{\epsilon}_i}{s_{(i)}\sqrt{1 - h_{ii}}},$$

where $s_{(i)}$ is the standard error computed with the $n - 1$ samples remaining after omitting the i th sample. Such a residual is also called a studentized deleted residual or externally studentized residual. Alternatively, t_i can be calculated as

$$t_i = r_i \left(\frac{n - p - 1}{n - p - r_i^2} \right)^{1/2},$$

where $p = k + 1$ denotes the number of columns of \mathbf{X} .

For deleted residuals, we have the following relationships:

$$\hat{\epsilon}_{(i)} = y_i - \hat{y}_{(i)} = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{(i)},$$

where $\hat{\boldsymbol{\beta}}_{(i)} = (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{y}_{(i)}$, $\mathbf{X}_{(i)}$ is an $(n - 1) \times (k + 1)$ matrix. In addition, we have

$$\hat{\boldsymbol{\beta}}_{(i)} = \hat{\boldsymbol{\beta}} - \frac{\hat{\epsilon}_i}{1 - h_{ii}} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i.$$

The deleted residual can also be expressed as

$$\hat{\epsilon}_{(i)} = \frac{\hat{\epsilon}_i}{1 - h_{ii}},$$

and

$$t_i = \frac{\hat{\epsilon}_{(i)}}{s_{(i)}}.$$

The deleted sample variance $s_{(i)}^2$ can be expressed as

$$s_{(i)}^2 = \frac{SSE_{(i)}}{n - k - 2} = \frac{SSE - \hat{\epsilon}_i^2 / (1 - h_{ii})}{n - k - 2}.$$

The n deleted residuals can be used for model validation or selection by defining the prediction sum of squares (PRESS):

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

To use PRESS to compare alternative models when the objective is prediction, preference would be shown to models with small values of PRESS.

4 Influential Observations and Leverage

4.1 Leverage

To investigate the influence of each observation, we begin with $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$, the element of which are

$$\hat{y}_i = h_{ii}y_i + \sum_{j \neq i} h_{ij}y_j.$$

Therefore, if h_{ii} is large (close to 1), then h_{ij} 's, $j \neq i$, are small, and y_i contributes much more than others to \hat{y}_i . Hence, h_{ii} is called the leverage of y_i .

By Theorem 2.1, the average value of h_{ii} 's is $(k + 1)/n$. Hoaglin and Welsch (1978) suggest that a point with $h_{ii} > 2(k + 1)/n$ is a high leverage point. Alternatively, we can simply examine any observations whose value of h_{ii} is unusually large relative to the other values of h_{ii} .

4.2 Cook's Distance

To formalize the influence of an observation, we consider the effect of its deletion on $\boldsymbol{\beta}$ and $\hat{\mathbf{y}}$. This is measured by Cook's distance

$$D_i = \frac{(\hat{\boldsymbol{\beta}}_{(i)} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\hat{\boldsymbol{\beta}}_{(i)} - \hat{\boldsymbol{\beta}})}{(k + 1)s^2} = \frac{(\hat{\mathbf{y}}_{(i)} - \hat{\mathbf{y}})^T (\hat{\mathbf{y}}_{(i)} - \hat{\mathbf{y}})}{(k + 1)s^2}.$$

Therefore, if D_i is large, the observation i has substantial influence on both $\hat{\boldsymbol{\beta}}$ and $\hat{\mathbf{y}}$. A more computationally convenient form of D_i is given by

$$D_i = \frac{r_i^2}{k + 1} \left(\frac{h_{ii}}{1 - h_{ii}} \right).$$