

Residuals and Model Diagnosis (§ 9. Applied Linear Regression book)

① Linear regression (p predictors)

Review

- For $i = 1, \dots, n$.

$$Y_i = \underbrace{\beta_0 + \beta_1 X_{i,1} + \dots + \beta_p X_{i,p}}_{\text{linear combination}} + e_i$$

$$= \underbrace{X_i^T \beta}_{\text{linear combination}} + e_i$$

where $X_i = \begin{pmatrix} 1 \\ X_{i,1} \\ \vdots \\ X_{i,p} \end{pmatrix}$, $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$

$$\Rightarrow Y = X\beta + e$$

where $Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$, $X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix}_{n \times (p+1)}$, $e = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}_{n \times 1}$

- Assumptions on statistical error (e)

(1) $E[e | X] = 0$ (Equivalently, for $i = 1, \dots, n$
 $E[e_i | X] = 0, \text{Var}(e_i | X) = \sigma^2$)
 (2) $\text{Var}(e | X) = \sigma^2 I_n$ and e_i 's are uncorrelated
(independent)

Note that Assumption (1) implies that for $i = 1, \dots, n$

$$E[e_i X_i] = 0 \quad \text{and} \quad \text{Cov}(e_i, X_i) = 0$$

This is because

$$\begin{aligned} \text{(1)} \Rightarrow E[e_i | X] = 0 &\Rightarrow E[E[e_i X_i | X]] = 0 \\ &\Rightarrow \underbrace{E[e_i X_i]}_{E[E[e_i X_i | X]]} = 0 \end{aligned}$$

$$\text{Cov}(e_i, X_i) = E[\underbrace{e_i X_i}_{\text{II}}] - \underbrace{E[e_i] \cdot E[X_i]}_{\text{I}} = 0.$$

- For OLS estimator $\hat{\beta} = (X^T X)^{-1} X^T Y$. (Assuming $(X^T X)^{-1}$ exists)

From previous lecture, we have

$$E(\hat{\beta} | X) = \beta, \quad \text{Var}(\hat{\beta} | X) = \sigma^2 (X^T X)^{-1}.$$

$\hat{\beta}$ is an unbiased estimator of β

$$(2) \text{ Properties of Residuals} \quad \hat{\mathbf{e}} = \begin{pmatrix} \hat{e}_1 \\ \vdots \\ \hat{e}_n \end{pmatrix}_{n \times 1} = \mathbf{y} - \hat{\mathbf{y}} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \hat{\beta} \\ \vdots \\ y_n - \mathbf{x}_n^T \hat{\beta} \end{pmatrix}_{n \times 1}$$

Prop:

(2.1) sample version. Sample mean of residuals = 0 (from earlier lecture)

$$\bar{\hat{e}} \triangleq \frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0 \quad (\text{can prove using normal equations})$$

v Normal equations:

$$\Leftrightarrow \underbrace{\begin{pmatrix} \mathbf{X}^T & \hat{\mathbf{e}} \end{pmatrix}}_{(p+1) \times 1} = \mathbf{0}$$

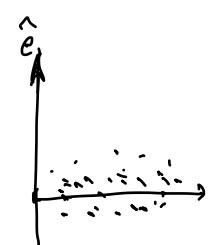
$$\Leftrightarrow \underbrace{\begin{pmatrix} x_1 & x_2 & \dots & x_n \end{pmatrix}}_{\text{red}} \begin{pmatrix} \hat{e}_1 \\ \hat{e}_2 \\ \vdots \\ \hat{e}_n \end{pmatrix} = 0$$

$$\Leftrightarrow \sum_{i=1}^n \hat{e}_i x_i = 0 \quad (*)$$

- Prop: The Sample Covariance between predictors and residuals = 0.
 x_1, \dots, x_n $\hat{e}_1, \dots, \hat{e}_n$

Proof: Sample covariance $= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(\hat{e}_i - \bar{\hat{e}})$

$\uparrow \text{sample mean of } x$ $\uparrow \text{sample mean of } \hat{e}$
 $\frac{1}{n} \sum_{i=1}^n x_i$ $\frac{1}{n} \sum_{i=1}^n \hat{e}_i$

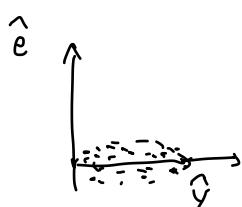


if linear reg. Assumptions hold

$$= \frac{1}{n-1} \left[\underbrace{\sum_{i=1}^n x_i \hat{e}_i}_{0 \text{ by } (*)} - \underbrace{\sum_{i=1}^n \bar{x} \hat{e}_i}_{\bar{x} \sum_{i=1}^n \hat{e}_i = 0} \right] = 0$$

$$= 0.$$

- Prop: The sample covariance between fitted values and residuals = 0
 $\hat{y}_1, \dots, \hat{y}_n$ $\hat{e}_1, \dots, \hat{e}_n$



if linear reg. Assumptions hold.

Proof: Sample Cov of \hat{Y} and \hat{e}

$$= \frac{1}{n-1} \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}}) (\hat{e}_i - \bar{\hat{e}})$$

Note that $\bar{\hat{y}} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ "Sample mean of $\hat{y}_1 \dots \hat{y}_n$

because $\bar{\hat{y}} - \bar{y} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) - \bar{\hat{e}}_i = 0$ $\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$

$$= \frac{1}{n-1} \sum_{i=1}^n (\hat{y}_i - \bar{y}) \hat{e}_i$$

$$= \frac{1}{n-1} \left[\underbrace{\sum_{i=1}^n \hat{y}_i \hat{e}_i}_{\hat{y}^\top \hat{e}} - \underbrace{\sum_{i=1}^n \bar{y} \hat{e}_i}_{\bar{y} \sum_{i=1}^n \hat{e}_i = 0} \right] = 0.$$

$$\hat{y} = X\hat{\beta}$$

$$\hat{\beta}^\top X^\top \hat{e} = 0$$

by normal equations.

(2.2) population version Next we will study the properties of \hat{e} by treating it as random variables.

$$\hat{e}_{n \times 1} = Y_{n \times 1} - \underbrace{\hat{y}_{n \times 1}}_{\text{population mean}} = Y - \underbrace{X(X^\top X)^{-1} X^\top Y}_{\text{population regression function}} = (I_n - \underbrace{X(X^\top X)^{-1} X^\top}_{\text{H}}) Y.$$

$$= X\beta + e - X\hat{\beta}$$

$$= \underbrace{X\beta + e}_{\text{population error}} - \underbrace{X\hat{\beta}}_{I_n e} - \underbrace{X(X^\top X)^{-1} X^\top e}_{\beta + (X^\top X)^{-1} X^\top e} = (X^\top X)^{-1} X^\top (X\beta + e) = (X^\top X)^{-1} X^\top e.$$

$$= \underbrace{(I_n - X(X^\top X)^{-1} X^\top)}_{\text{H}} e$$

$$= (I_n - \underbrace{X(X^\top X)^{-1} X^\top}_{\text{H}}) e$$

$H \triangleq X(X^\top X)^{-1} X^\top$ is often called the "Hat Matrix".

Note from the above derivation $\hat{y} = H Y$

Some properties of H : $H X = X$

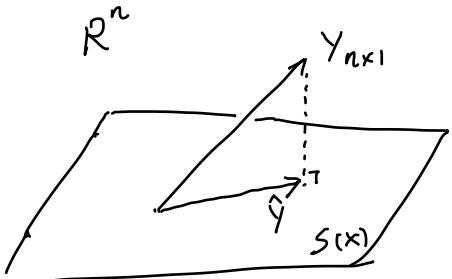
because $Hx = \underbrace{x(x^T x)^{-1}}_{\text{because } Hx = x(x^T x)^{-1} x^T x} x^T x = x$.

$$\cdot HH = H \text{ because } HH = \underbrace{x(x^T x)^{-1}}_{\text{because } HH = x(x^T x)^{-1} x^T x(x^T x)^{-1} x^T} x^T x(x^T x)^{-1} x^T = H.$$

$$\cdot H^T = H$$

- From the geometry interpretation of OLS, ($n > p+1$)

$\hat{y}_{n \times 1}$ is the closest point in the subspace $S(x)$, expanded by the column vectors of X matrix, to our response vector $y_{n \times 1}$, and the projection matrix is H . ($\hat{y} = Hy$).



✓ In summary, $\hat{e} = (I_n - H)y = (I_n - H)e$.

Next, we will calculate $E[\hat{e}|x]$ and $\text{Var}(\hat{e}|x)$

based on $E(e|x) = \underbrace{\text{Var}(e|x)}_{\text{linear reg. model assumptions.}} = \sigma^2 I$

Prop 1: $E(\hat{e}|x) = 0$.

$$\Rightarrow \text{Cov}(\hat{e}_i, x_i|x) = E[\hat{e}_i x_i|x] - E[\hat{e}_i|x] \cdot E[x_i|x]$$

Similarly $\text{Cov}(\hat{e}_i, x_j|x) = 0$.

Proof: $E[\hat{e}|x] = E[(I_n - H)e|x] = 0$.

Prop 2: $\text{Var}(\hat{e}|x) = \sigma^2 (I_n - H)$.

Proof: $\text{Var}(\hat{e}|x) = \text{Var}((I_n - H)e|x)$

For A only depending on x.

$$\text{Var}(Ae|x) = A \text{Var}(e|x) A^T \stackrel{?}{=} (I_n - H) \underbrace{\text{Var}(e|x)}_{\sigma^2 I_n} (I_n - H)^T$$

$$= \sigma^2 (I_n - H)(I_n - H)$$

$$= \sigma^2 (I_n - H - H + \underbrace{H^T H}_{H})$$

$$= \sigma^2 (I_n - H)$$

- Let $H = (h_{ij})_{n \times n}$ with h_{ij} being the (i, j) element.

Then h_{ii} corresponds to the i^{th} diagonal element of H .

By prop 2. $\text{Var}(\hat{e}_i | x) = \sigma^2(1 - h_{ii})$ for $i = 1, \dots, n$.

- The h_{ii} is called the leverage of the i^{th} observation.

from $H = X(X^T X)^{-1} X^T$, we know $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$

$$(1) \quad h_{ii} = x_i^T (X^T X)^{-1} x_i \quad \begin{array}{l} \text{recall } x_i = \begin{pmatrix} x_{i,1} \\ \vdots \\ x_{i,p} \end{pmatrix} \in \mathbb{R}^{(p+1) \times 1} \\ \text{here } x_i = (1, x_{i,1}, \dots, x_{i,p}) \in \mathbb{R}^{(p+1) \times 1} \end{array}$$

more generally

$$h_{ij} = x_i^T (X^T X)^{-1} x_j.$$

- Some properties of h_{ii}

$$\textcircled{1} \quad 0 \leq h_{ii} \leq 1$$

$$\textcircled{2} \quad \text{observations with } h_{ii} \text{ will have smaller } \text{Var}(\hat{e}_i | x) \quad (\text{larger } h_{ii} \text{ is smaller } \text{Var}(\hat{e}_i | x))$$

If $h_{ii}=1$, then $\hat{e}_i=0$ and $\hat{y}_i=y_i$

that is, (x_i, y_i) is on the fitted OLS regression line/plane.

$$\textcircled{3} \quad \text{Consider the Simple Linear Regression } (Y_i, x_i \in \mathbb{R}), i=1, \dots, n$$

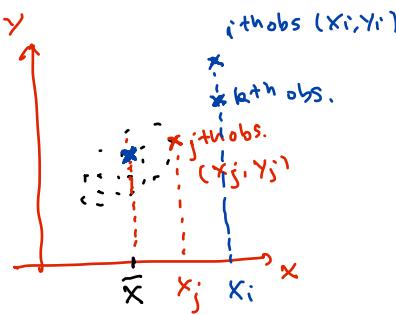
from equation (1) above, we have

$$h_{ii} = (1, x_i) (X^T X)^{-1} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \text{ where } X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2}$$

$$\left[\begin{array}{l} \text{practice} \\ \text{problem} \end{array} \right] = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \rightarrow \text{fixed for all } i = 1, \dots, n$$

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

- This implies $1 \geq h_{ii} \geq \frac{1}{n}$ if $x_i = \bar{x}$ then $h_{ii} = \frac{1}{n}$.



$$h_{jj} < h_{ii} = h_{kk}.$$

Prop 3: $\text{Cov}(\hat{\epsilon}_{n \times 1}, \hat{y}_{n \times 1} | x) = 0.$

Proof: $\text{Cov}(\hat{\epsilon}, \hat{y} | x)$

$$= \text{Cov}((I_n - H)y, Hy | x)$$

For matrix A, B only depend on x , we have

$$\text{Cov}(AY, BY | x)$$

$$= A \cdot \text{Var}(Y | x) \cdot B^T$$

$$= (I_n - H) \underbrace{\text{Var}(Y | x)}_{\sigma^2 I_n} H^T$$

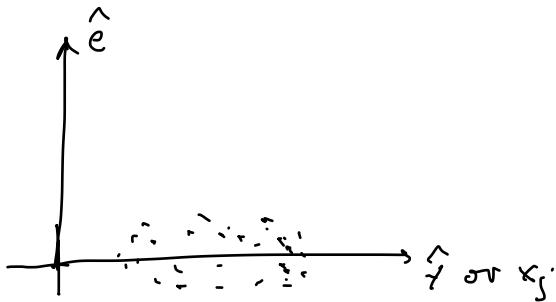
$$= \sigma^2 (I_n - H) H$$

$$= \sigma^2 (H - H^T H) = 0 . \quad \text{#}$$

* Application of Residuals (§ 9.1 Applied Linear Regression)

- Residuals are generally used in scatterplots of $\hat{\epsilon}$ against the fitted value \hat{y} or some predictors x_j , $j=1, \dots, p$.

[Residual Plot]

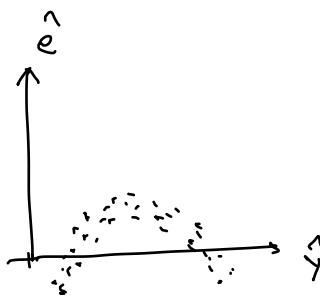


- If our regression model is true (meaning our mean and covariance assumptions both hold)

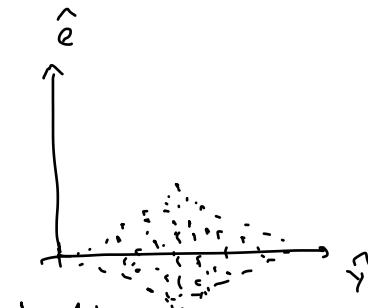
Prop 1 & 3 \Rightarrow the Scatterplot of $\hat{\epsilon}$ vs any x or \hat{y} should have a constant mean function = 0.

Prop 2 \Rightarrow Residual should appear to have
"Approximately" constant variance

E.g:



indicating the failure of our
linear mean function assumption



indicating
that Constant Variance assumption
is problematic.

Prop 4 (Note that all props assume linear regression model assumptions hold)
p predictors.

$$E[\hat{\sigma}^2 | x] = \sigma^2$$

where $\hat{\sigma}^2$ is the OLS estimator $\hat{\sigma}^2 = \frac{RSS}{n-(p+1)} = \frac{\sum_{i=1}^n \hat{e}_i^2}{n-(p+1)}$.

That is, $\hat{\sigma}^2$ is an unbiased estimator of σ^2 . ($\sigma^2 = \text{Var}(e_i | x)$)

Proof: We know $\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{e}_i^2}{n-(p+1)} = \frac{\hat{e}^\top \hat{e}}{n-(p+1)}$, $\hat{e} = \begin{pmatrix} \hat{e}_1 \\ \vdots \\ \hat{e}_n \end{pmatrix}$

$$\begin{aligned} \hat{e} &= (I_n - H)e \quad \checkmark \\ &= \frac{[(I_n - H)e]^\top (I_n - H)e}{n-(p+1)} \\ &= \frac{e^\top (I_n - H) (I_n - H) e}{n-(p+1)} \end{aligned}$$

Note that $(I_n - H)(I_n - H) = I_n - 2H + \underbrace{HH^\top}_{H^2} = I_n - H$

$$\begin{aligned} &\downarrow \\ &= \frac{e^\top (I_n - H) e}{n-(p+1)} \end{aligned}$$

Then

$$\begin{aligned} E[\hat{\sigma}^2 | x] &= \frac{1}{n-(p+1)} E\left[e^\top (I_n - H) e | x\right] \\ &= \frac{1}{n-(p+1)} E\left[\text{Tr}\left(\underbrace{e^\top (I_n - H) e}_{\text{red}}\right) | x\right] \end{aligned}$$

For A, B.

$$\text{Tr}(AB) = \text{Tr}(BA)$$

$$= \frac{1}{n-(p+1)} E \left[\text{Tr}((I_{n-H})ee^T) | x \right]$$

$$= \frac{1}{n-(p+1)} \text{Tr} \left[E \left[(I_{n-H})ee^T | x \right] \right]$$

Note that $E[ee^T | x] = \text{Var}(e | x) = \sigma^2 I_n$

$$(\text{Var}(e | x) = E[(e - E(e|x))(e^T - E(e^T|x)) | x]) \quad \text{from our variance assumption}$$

$$= \frac{1}{n-(p+1)} \text{Tr} \left[(I_{n-H}) \sigma^2 I_n \right]$$

$$= \sigma^2 \cdot \frac{\text{Tr}(I_{n-H})}{n-(p+1)}$$

From above, To show $E(\hat{\sigma}^2 | x) = \sigma^2$, we only need to show

$$(*) \quad \cdots \quad \text{Tr}(I_{n-H}) = n-(p+1).$$

To show (*), Note that

$$\text{Tr}(I_{n-H}) = \text{Tr}(I_n) - \text{Tr}(H)$$

$$= n - \text{Tr} \left(\underset{n \times (p+1)}{X} (X^T X)^{-1} \underset{(p+1) \times n}{X^T} \right)$$

$$= n - \text{Tr} \left((X^T X)^{-1} \underset{(p+1) \times (p+1)}{X^T X} \right)$$

$$= n - (p+1).$$

Note that from the above proof, we have

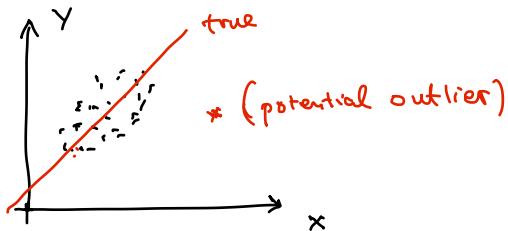
$$\text{Tr}(H) = p+1, \quad E[\text{RSS} | x] = \sigma^2(n-(p+1)).$$

$$\Rightarrow \sum_{i=1}^n h_{ii} = p+1.$$

③ Detect Outliers (§ 9.4 Applied Linear Reg.)

(Linear Reg. model)

- Outliers: observations that do not follow the same model as the majority and often called "Outliers"



- Cases with large residuals are candidates for outliers
 - Not all large residuals are outliers
 - Outliers may not be bad.
 - Some outliers will have longer influence on regression estimates than others

↑
(will be discussed in a later lecture)

- We use the Studentized Residual to detect outliers.

For the i^{th} observation, is

$$t_i = \frac{\hat{e}_i}{\hat{\sigma}_{(i)} \sqrt{1-h_{ii}}} \quad \begin{matrix} \leftarrow \text{Recall } \sqrt{\text{Var}(\hat{e}_i | X)} \\ = \sqrt{\hat{\sigma}^2 (1-h_{ii})} \end{matrix}$$

where $\hat{\sigma}_{(i)}$ is the estimated σ without using the i^{th} observation.
 Fit regression $(Y_j, X_j)_{j=1, \dots, n}, i, \dots, n \Rightarrow \hat{\sigma}_{(i)}^2 = \frac{\text{RSS}_{(i)}}{(n-1)-(p+1)}$

If the i^{th} observation is not an outlier (the rest also are not outliers)
 and we assume normal statistical error,

then

$$t_i \sim t\text{-distribution}$$

$$\text{with d.f.} = (n-1)-(p+1)$$

$$= n-p-2.$$