

3.3.2 → Leverage + Residuals

Overview → to better understand what makes an observation unusual or influential, lets discuss these concepts ideas in relation to leverage, residuals + Cook's distance

→ leverage
 → The leverage of an observation measures its influence in predicting the response. We denote the i th leverage as h_i , which is the i th diagonal entry of the hat matrix H . Recall the hat matrix

$$H = \begin{bmatrix} h_1 & h_{12} & \dots & h_{1n} \\ h_{21} & h_{22} & \dots & h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \dots & h_{nn} \end{bmatrix}$$

→ In the case of SLR, $h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$
 \downarrow
 $\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}}$

→ other facts about leverage

- S_{xx} is a function of the x_i 's but not y
- The larger h_i is, the more unusual the set of x_1, \dots, x_n values is relative to other observations
- $\sum_{i=1}^n h_i = p+1$

→ There is nothing inherently bad about an observation with a large leverage. The concern is whether \hat{y} is mainly driven by a few observations that dominate a vast amount of the "total leverage" ($\sum h_i$), to the point that \hat{y} would alter drastically without these observations. Our rule of thumb is that the i th observation is a high leverage point if

$$h_i > 3 \left(\frac{p+1}{n} \right)$$

→ Residuals

→ Recall that a residual is $y - \hat{y}$ (actual - predicted). This means it has the same units as y . Consequently, the value of an extreme residual depends on the units. For this reason, it is common to analyze a version of the residuals that is unitless. There are two types to discuss:

→ Standardized residuals + studentized residuals

→ Standardized residuals are the residuals divided by an estimated standard error. The i th standardized residual is

$$e_{std,i} = \frac{e_i}{\sqrt{MSE(1-h_i)}}$$

→ Since the standard error is estimated, the standardized residuals are approximate realizations of the standard normal distribution, provided the model is correct

→ Studentized residuals are similar to standardized residuals, except they are divided by a different estimated standard error.

The i th (external) studentized residual is

$$e_{stud,i} = \frac{e_i}{\sqrt{MSE_{(i)}(1-h_i)}}$$

→ While $MSE_{(i)}$ is the MSE of the regression that excludes the i th observation from the training data.

→ The studentized residuals are realizations of a t -stat, provided the model is correct. Recall $t \rightarrow z$ as $n \rightarrow \infty$

$$\Rightarrow z \sim N(0,1) \Rightarrow$$

$$\Rightarrow 1.96 \in [-2, 2]$$

→ This allows us to choose a specific cutoff in defining an outlier, regardless of the unit of the response. It is rather common to find the rule of thumb of $> |2|$ to identify outliers. Moreover, both types of residuals are equally effective in determining outliers.

→ Cook's distance

→ If we wish to combine leverage + residuals into a single measure, we can compute either DFITS or Cook's distance for each observation.

$$DFITS_i = e_{std,i} \sqrt{\frac{h_i}{1-h_i}}$$

The i th Cook's distance is

$$D_i = \frac{DFITS_i^2}{p+1}$$

$$= \frac{e_{std,i}^2 h_i}{(p+1)(1-h_i)}$$

$$= \frac{e_i^2 h_i}{MSE(p+1)(1-h_i)^2}$$

$$e_{std,i} = \frac{e_i}{\sqrt{MSE(1-h_i)}}$$

→ one rule of thumb is that the i th observation is an influential point if D_i exceeds unity, i.e. $D_i > 1$

→ Assignment 1

→ Q1) Given $n=15$, $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)$, $b_0 = 0.0715$, $e_{std,0} = 0.0032$, $b_1 = 3^2$

$$D_0 = \frac{DFITS_0^2}{p+1}$$

$$\downarrow$$

$$= \frac{e_{std,0}^2 h_0}{(p+1)(1-h_0)}$$

$$0.0715 = \frac{0.0032^2 h_0}{1(1-h_0)}$$

$$\frac{0.0715}{0.0032^2} = \frac{h_0}{1-h_0}$$

$$L = L$$

$$C - Ch_0 = Ch_0$$

$$C = h_0(1+L)$$

$$\Rightarrow h_0 = \frac{C}{1+C}$$

$$\downarrow$$

$$= 0.17379$$

→ Q2) Given $p=12$ (excluding intercept)

$$SSE = 618$$

$$n_{(p)} = 0.35$$

$$b_0 = 9.5$$

$$d_{10} = 5$$

→ $n=?$

$$\rightarrow d_{10} = \frac{e_{10}^2 h_{10}}{\frac{SSE}{n-p} (p+1)(1-h_{10})}$$

$$= \frac{9.5^2 (0.35)}{\frac{618}{n-13} (13+1)(1-0.35)^2}$$

$$\downarrow$$

$$= \frac{(n-13) 9.5^2 (0.35)}{618 (13)(1-0.35)^2}$$

$$\Rightarrow n = 550.74 \Rightarrow n \geq 551$$

3.5.4 → VIF

→ Overview → one way to measure multicollinearity is the variance inflation factor (VIF). To calculate the VIF for the j th predictor, first run a SLR w/ x_j acting as the response variable, predicted by the rest of the $p-1$ predictors. Denote the coefficient of determination of that regression as R_j^2 . Then

$$VIF_j = \frac{1}{1-R_j^2}$$

→ Intuitively, R_j^2 is like a suitable measure for whether x_j is almost a linear combination of the other predictors. A reason to use VIF instead comes from how $SSE(p)$ can be written as

$$SSE(p) = \sqrt{VIF_j} \sqrt{\frac{MSE}{1-R_j^2}}$$

→ See the regression of y on remaining $p-1$ vars

→ Since $R_j^2 = 0.8 \Rightarrow VIF_j = 2.5$, it is suggested that we should be concerned w/ multicollinearity if any of the p VIFs is 5 or greater

→ Another inspection way to detect multicollinearity

→ calculate all β_j 's + s_{xy} (simple correlation between the response & each predictor)

→ If signs differ, then the contradiction is possibly the result of multicollinearity (i.e. some information from multiple predictors)