

Chapter 8 · Outliers and Influential Observations

Sometimes, while most of the observations fit the model and meet G-M conditions at least approximately, some of the observations do not. This occurs when there is something wrong with the observations or if the model is faulty.

As far as the observations are concerned, there could have been a mistake in inputting or recording data. A few observations might reflect conditions or situations different from those under which other observations were obtained.

The leverage.

The leverages h_{ii} are the diagonal elements of $H = X(X^T X)^{-1} X^T$. $h_{ii} = x_i^T (X^T X)^{-1} x_i$

The key features of a leverage h_{ii} is that it describes how far away the individual data point is from the centroid of all data points in the space of independent variables, i.e., how far removed x_i is from $\bar{x} = \frac{\sum x_i}{n}$.

$$\hat{y} = Hy$$

$$\hat{y}_i = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j.$$

When h_{ii} is large, y_i influence \hat{y}_i

$$\begin{aligned} h_{ii} &= x_i' (X' X)^{-1} x_i \\ &= \frac{1}{n} + (x_i - \bar{x})' \left[\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})' \right]^{-1} (x_i - \bar{x}) \end{aligned}$$

$\sum_{i=1}^n h_{ii} = k+1$, if we had the option of choosing
(k : # of variables)

independent variable values we would choose them

So as to make each $h_{ii} = \frac{k+1}{n}$.

In practice, a point i is a leverage point if

(1) $h_{ii} > \frac{2(k+1)}{n}$ criterion

(2) $h_{ii} > \frac{3(k+1)}{n}$ for $k+1 > 6$ and $n-k-1 > 12$.

(3) $h_{ii} > 0.5$

The Residuals

For the purpose of detecting observations that do not belong to the model, more valuable than the residuals are the Studentized residuals, often called RSTUDENT and defined by

$$e_i^* = \frac{e_i}{s(i) \sqrt{1-h_{ii}}}$$

$$e_i = y_i - \hat{x}_i \hat{\beta}$$

$$\text{var}(e_i) = \sigma^2 (1-h_{ii})$$

$$e_i^{(10)} = \frac{e_i}{\sigma \sqrt{1-h_{ii}}}$$

where e_i is a residual.

$s(i)$ is equivalent to s if least squares is

Run after deleting the i th case.

$y_{(i)}$ and $X_{(i)}$ the results of subtracting the i th row from y and X

dim: $y_{(i)} : (n-1) \times 1$

$X_{(i)} : (n-1) \times (k+1)$

$$\hat{\beta}_{(i)} = (X'_{(i)} X_{(i)})^{-1} X'_{(i)} y_{(i)}$$

$$(n-k-2) S^2_{(i)} = \sum_{\substack{l=1 \\ l \neq i}}^n [y_l - X'_l \hat{\beta}_{(i)}]^2$$

Theorem:

(1) Relations between $\hat{\beta}$ et $\hat{\beta}_{(i)}$:

$$\hat{\beta} - \hat{\beta}_{(i)} = \frac{(X'X)^{-1} x_{il}}{1-h_{ii}}$$

If $\hat{\beta} - \hat{\beta}_{(i)}$ is large,
 i is an outlier.

(2) Under G-M conditions,

$S_{(i)}$ is an unbiased and consistent
estimate of σ^2

(3) Relations between $s_{(i)}^2$ and s^2 :

$$(n-k-2) s_{(i)}^2 = (n-k-1) s^2 - \frac{e_i^2}{1-h_{ii}}$$

Detecting Outliers and Points That do not belong to the Model.

An outlier is an observation that is poorly explained by the model and has a high residual $|e_i^*| > 2$.

For a search for points that might not fit the model, the studentized residuals are much more useful. One could examine a listing or a plot of these against case numbers. One can easily identify from these which e_i^* are significant at some given level. e.g 5%.

(1) An outlier if $|e_i^*| > t_{n-k-2}; \frac{\alpha}{2}$

(2) Normal probability plot.

(3) **box-plot**
to detect outliers.

Influential Observations.

An influential point is an observation that has a major impact on the parameter estimation process. It should be noted that not all the influential points have high values of studentized residuals.

Measures of influence:

$$(1) DFBETA_i = \hat{\beta} - \hat{\beta}_{(i)} = \frac{(X'X)^{-1} X_i e_i}{1-h_{ii}}$$

$$(2) DFBETA_{ij} = \hat{\beta}_j - \hat{\beta}_{j(i)} = \frac{a_{ji} e_i}{1-h_{ii}}$$

$$a_{ji} = [(X'X)^{-1} X_i]_j \quad (X'X)^{-1} X_i = (a_{0i}, \dots, a_{ki})'$$

$$(3) DFFIT_j = \hat{y}_i - \hat{y}_{i(i)} = X_i' \hat{\beta} - X_i' \hat{\beta}_{(i)} = \frac{h_{ii} e_i}{1-h_{ii}}$$

(4) DFBETAS_{ij} (standard DFBETA_{ij}) =

$$\frac{\hat{\beta}_j - \beta_j e_i}{S_{ii} \sqrt{g_{jj}}} = \frac{a_{ji} e_i}{S_{ii} (1-h_{ii}) \sqrt{g_{jj}}} = \frac{a_{ji} e_i^*}{\sqrt{g_{jj} (1-h_{ii})}}$$

g_{jj} : j^{th} diagonal element of $(X'X)^{-1}$.

DFBETAS_{ij} to examine the i^{th} observation's influence on the estimation.

criterion: $\frac{2}{\sqrt{n}}$ if $> \frac{2}{\sqrt{n}}$, requires investigation

(5) DFFITS_i (Standard DFFIT_i)

$$\frac{\sqrt{h_{ii}} \cdot e_i}{S_{ii}(1-h_{ii})} = \sqrt{\frac{h_{ii}}{1-h_{ii}}} \cdot e_i^* \quad e_i^* = \frac{e_i}{S_{ii} \sqrt{1-h_{ii}}}$$

$$\text{var}(\hat{y}_i(i)) = S_{ii} h_{ii} \quad h_{ii} = \frac{(k+1)}{n}, |e_i^*|^2 = 2.$$

to inform the change in the predicted value with and without the observation.

criterion level: if $> 2 \sqrt{\frac{k+1}{n-k-1}}$ \Rightarrow requires investigation.

(6) A very frequently used measure of influence can be defined as the distance between $\hat{\beta}$ and $\hat{\beta}_{(i)}$ standardized by the estimated covariance matrix $S^2(X'X)^{-1}$ of $\hat{\beta}$.

Cook's distance:
$$\frac{(\hat{\beta} - \hat{\beta}_{(i)})'(X'X)(\hat{\beta} - \hat{\beta}_{(i)})}{(k+1)S^2}$$

$\sim F_{k+1; n-k-1}$.

$$\because \hat{\beta} - \hat{\beta}_{(i)} = \frac{(X'X)^{-1}x_i e_i}{1-h_{ii}}$$

$$\Rightarrow \text{Cook's distance} = \frac{h_{ii} (e_i^{(s)})^2}{(k+1)(1-h_{ii})}$$

influential observation when Cook's distance is large.

$$F_{k+1, n-k-1; \alpha=0.5} \approx 1. \text{ or } \frac{4}{n} \text{ (criterion)}$$

(7) Covariance Ratio measures the effect on the estimated covariance matrix of $\hat{\beta}$.

$$\frac{\det [S^2(i) (X'_{(i)} X_{(i)})^{-1}]}{\det [S^2 (X'X)^{-1}]}$$

A value of this ratio close to 1 would indicate lack of influence of the i^{th} data point.

$$\det [X'_{(i)} X_{(i)}] = (1-h_{ii}) \det (X'X)$$

$$\Rightarrow \text{Covariance Ratio} = \frac{(n-k-1)^{k-1}}{(1-h_{ii})(n-k-2+\ell_i^{*2})^{k+1}}$$

when $|\text{Covariance Ratio} - 1| > \frac{3(k+1)}{n} \Rightarrow \text{investigation}$

- Outliers that deviate from the model are detected by analysis of the studentized residuals ℓ_i^* . They do not necessarily influence the plot of regression line.

- Leverage points can have an impact on the estimation of the model. They have high value of h_{ii} . Their distance from the center of gravity is large.
- Influential observations have high DFBETAS and/or DFFITS. and/or Cook's distance. They play a large role in the regression line.