

## How to detect an outlier (whether or not it exists)

by Bo E. Honoré<sup>1</sup>

Consider the model

$$y_i = x_i' \beta + \gamma d_i + \varepsilon_i$$

where  $(x_i, \varepsilon_i)$  is i.i.d.  $E[\varepsilon_i | x_i] = 0$  and

$$d_i = \begin{cases} 1 & \text{if } i = 1 \\ 0 & \text{if } i \neq 1 \end{cases}$$

It is tempting to test whether the first observation is an outlier by testing whether  $\gamma = 0$ . As we will see, this can be a good idea if you *want* to conclude that it is, but not necessarily if you actually want to know “the truth.”

Let  $z_i = (x_i, d_i)$ . The OLS estimator of  $(\beta, \gamma)$  minimizes

$$\sum_{i=1}^n (y_i - x_i' b - \gamma d_i)^2$$

It is clear that the solution to this would be

$$\begin{aligned} \hat{\gamma} &= y_1 - x_1' \hat{\beta} \\ \hat{\beta} &= \left( \sum_{i=2}^n x_i x_i' \right)^{-1} \sum_{i=2}^n x_i y_i \end{aligned}$$

hence  $\hat{\beta}$  is consistent etc.

The heteroskedasticity-consistent standard errors of the OLS estimator of  $\beta$  and  $\gamma$  are the square roots of the diagonal of

$$\begin{aligned} \hat{V} &= \left( \sum_{i=1}^n z_i z_i' \right)^{-1} \left( \sum_{i=1}^n e_i^2 z_i z_i' \right) \left( \sum_{i=1}^n z_i z_i' \right)^{-1} \\ &= \begin{pmatrix} \sum_{i=1}^n x_i x_i' & \sum_{i=1}^n x_i d_i \\ \sum_{i=1}^n d_i x_i' & \sum_{i=1}^n d_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n e_i^2 x_i x_i' & \sum_{i=1}^n e_i^2 x_i d_i \\ \sum_{i=1}^n e_i^2 d_i x_i' & \sum_{i=1}^n e_i^2 d_i^2 \end{pmatrix} \\ &\quad \begin{pmatrix} \sum_{i=1}^n x_i x_i' & \sum_{i=1}^n x_i d_i \\ \sum_{i=1}^n d_i x_i' & \sum_{i=1}^n d_i^2 \end{pmatrix}^{-1} \\ &= \begin{pmatrix} \sum_{i=1}^n x_i x_i' & x_1 \\ x_1' & 1 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=2}^n e_i^2 x_i x_i' & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \sum_{i=1}^n x_i x_i' & x_1 \\ x_1' & 1 \end{pmatrix}^{-1} \end{aligned}$$

---

<sup>1</sup>Comments should be emailed to bo@honore.com

because  $e_1 = 0$  ( $e_i$  is the residual for the  $i$ 'th observation)

To simplify the notation, let

$$\begin{pmatrix} \sum_{i=1}^n x_i x_i' & x_1 \\ x_1' & 1 \end{pmatrix}^{-1} = \begin{pmatrix} a & b \\ b' & c \end{pmatrix}$$

then

$$\begin{aligned} \hat{V} &= \begin{pmatrix} a & d \\ d' & c \end{pmatrix} \begin{pmatrix} \sum_{i=2}^n e_i^2 x_i x_i' & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} a & d \\ d' & c \end{pmatrix} \\ &= \begin{pmatrix} a (\sum_{i=2}^n e_i^2 x_i x_i') a & a (\sum_{i=2}^n e_i^2 x_i x_i') d \\ d' (\sum_{i=2}^n e_i^2 x_i x_i') a & d' (\sum_{i=2}^n e_i^2 x_i x_i') d \end{pmatrix} \end{aligned}$$

The estimated variance of  $\hat{\gamma}$  is  $d' (\sum_{i=2}^n e_i^2 x_i x_i') d$ . Clearly the term in the middle is of order  $n$ . Moreover

$$d = x_1' \left( \sum_{i=1}^n x_i x_i' - x_1 x_1' \right)^{-1}$$

which is of order  $\frac{1}{n}$ . It therefore follows that the estimated variance of  $\hat{\gamma}$  is of order  $\frac{1}{n}$  and since  $\hat{\gamma}$  will converge to  $\gamma + \varepsilon_1$  this means that (unless  $\gamma + \varepsilon_1$  happens to equal 0) the calculated t-statistic will converge to  $\pm\infty$  with probability 1 whether or not  $\gamma = 0$ .