

19 Lecture 19: March 12

Last time

- Dummy-Variable regression
- Interaction

Today

- Midterm exam starts next Friday
- Unusual and influential data

Unusual and influential data

Linear models make strong assumptions about the structure of data, assumptions that often do not hold in applications. The method of least squares can be very sensitive to the structure of the data and may be markedly influenced by one or a few unusual observations.

Outliers

In simple regression analysis, an outlier is an observation whose response-variable value is *conditionally* unusual *given* the value of the explanatory variable: see Figure 19.1.

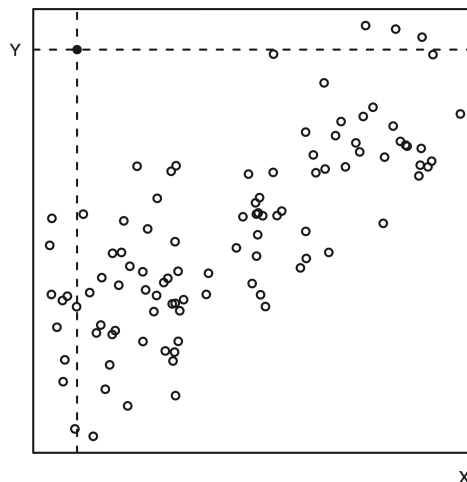


Figure 19.1: The black point is a regression outlier because it combines a relatively large value of Y with a relatively small value of X , even though neither its X -value nor its Y -value is unusual individually. Because of the positive relationship between Y and X , points with small X -values also tend to have small Y -values, and thus the black point is far from other points with similar X -values. JF Figure 11.1.

Unusual data are problematic in linear models fit by least squares because they can unduly

influence the results of the analysis. Their presence may be a signal that the model fails to capture important characteristics of the data.

Figure 19.2 illustrates some distinctions for the simple-regression model $Y = \beta_0 + \beta_1 X + \epsilon$.

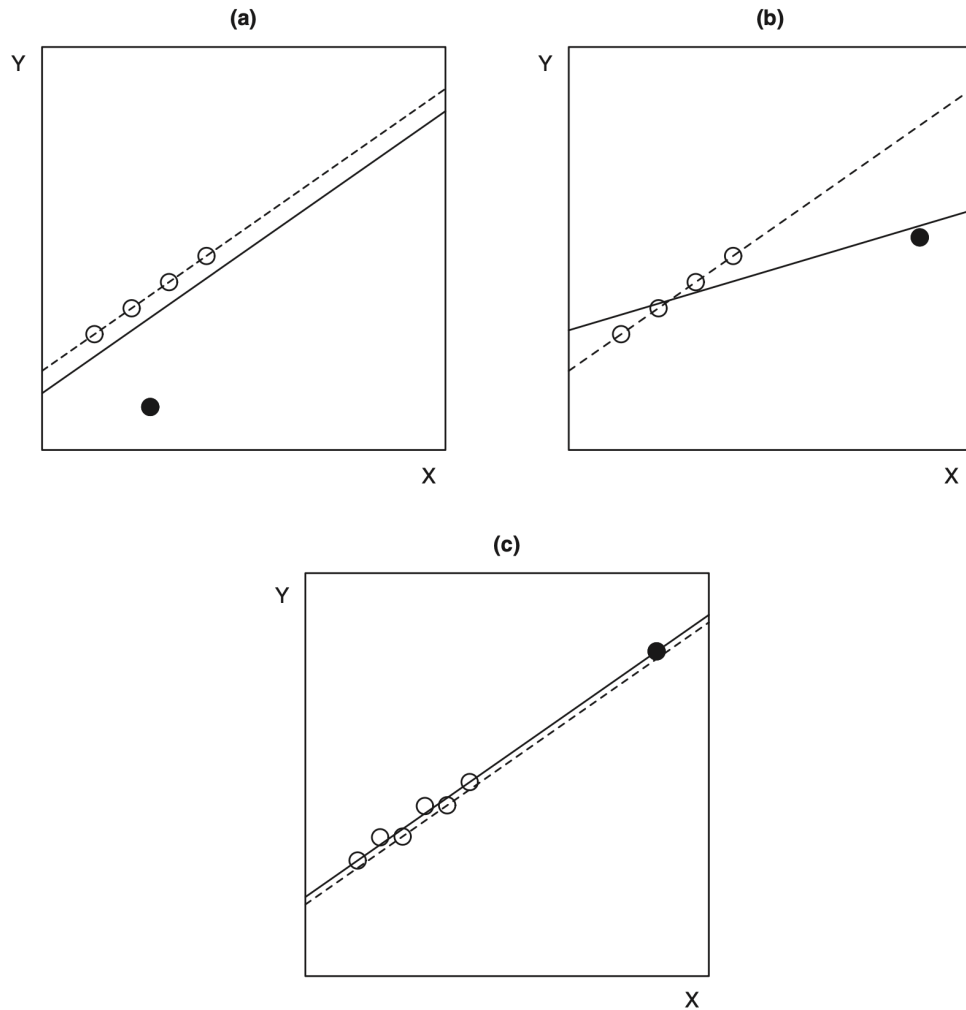


Figure 19.2: Leverage and influence in simple regression. In each graph, the solid line gives the least-squares regression for all the data, while the broken line gives the least-squares regression with the unusual data point (the black circle) omitted. (a) An outlier near the mean of X has low leverage and little influence on the regression coefficients. (b) An outlier far from the mean of X has high leverage and substantial influence on the regression coefficients. (c) A high-leverage observation in line with the rest of the data does not influence the regression coefficients. In panel (c), the two regression lines are separated slightly for visual effect but are, in fact, coincident JF Figure 11.2.

Some qualitative distinctions between outliers and high leverage observations:

- An outlier is a data point whose response Y does not follow the general trend of the rest of the data.

- A data point has high leverage if it has “extreme” predictor X values:
 - With a single predictor, an extreme X value is simply one that is particularly high or low.
 - With multiple predictors, extreme X values may be particularly high or low for one or more predictors, or may be “unusual” combinations of predictor values.

And the influence of a data point is the combination of leverage and discrepancy (“outlyingness”) through the following heuristic formula:

$$\text{Influence on coefficients} = \text{Leverage} \times \text{Discrepancy}.$$

Assessing leverage: hat-values

The hat-value h_i is a common measure of leverage in regression. They are named because it is possible to express the fitted values \hat{Y}_j (“Y-hat”) in terms of the observed values Y_i :

$$\hat{Y}_j = h_{1j}Y_1 + h_{2j}Y_2 + \cdots + h_{jj}Y_j + \cdots + h_{nj}Y_n = \sum_{i=1}^n h_{ij}Y_i.$$

The weight h_{ij} captures the contribution of observation Y_i to the fitted value \hat{Y}_j : If h_{ij} is large, then the i th observation can have a considerable impact on the j th fitted value. With the least square solutions, for the fitted values:

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

we (already) get the hat matrix:

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

Properties:

- (idempotent) $\mathbf{H} = \mathbf{H}\mathbf{H}$
- $h_i \equiv h_{ii} = \sum_{j=1}^n h_{ij}^2$
- $\frac{1}{n} \leq h_i \leq 1$ (a [proof](#) by Mohammad Mohammadi)
- $\bar{h} = (p + 1)/n$

In the case of SLR, the hat-values are:

$$h_i = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{j=1}^n (X_j - \bar{X})^2}$$

Detecting outliers: studentized residuals

The variance of the residuals ($\hat{\epsilon}_i = Y_i - \hat{Y}_i$) do not have equal variances (even if the errors ϵ_i have equal variances):

$$\text{Var}(\hat{\epsilon}) = \text{Var}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \text{Var}[(\mathbf{I} - \mathbf{H})\mathbf{Y}] = (\mathbf{I} - \mathbf{H})\text{Var}(\mathbf{Y})(\mathbf{I} - \mathbf{H}) = \sigma^2(\mathbf{I} - \mathbf{H})$$

so that for $\hat{\epsilon}_i$,

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

High-leverage observations tend to have small residuals (in other words, these observations can pull the regression surface toward them).

The standardized residual (sometimes called internally studentized residual)

$$\hat{\epsilon}'_i \equiv \frac{\hat{\epsilon}_i}{\hat{\sigma}\sqrt{1 - h_i}},$$

however, does not follow a t -distribution, because the numerator and denominator are not independent.

Suppose, we refit the model deleting the i th observation, obtaining an estimate $\hat{\sigma}_{(-i)}$ of σ that is based on the remaining $n - 1$ observations. Then the studentized residual (sometimes called externally studentized residual)

$$\hat{\epsilon}^*_i \equiv \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(-i)}\sqrt{1 - h_i}}$$

has an independent numerator and denominator and follows a t -distribution with $n - p - 2$ degrees of freedom.

The studentized and the standardized residuals have the following relationship (Beckman and Trussell, 1974):

$$\hat{\epsilon}^*_i = \hat{\epsilon}'_i \sqrt{\frac{n - p - 2}{n - p - 1 - \hat{\epsilon}'_i{}^2}}$$

For large n ,

$$\hat{\epsilon}^*_i \approx \hat{\epsilon}'_i \approx \frac{\hat{\epsilon}_i}{\hat{\sigma}}$$

Test for outlier

It is of our interest to pick the studentized residual $\hat{\epsilon}^*_{max}$ with the largest absolute value among $\hat{\epsilon}^*_1, \hat{\epsilon}^*_2, \dots, \hat{\epsilon}^*_n$ to test for outlier. However, by doing so, we are effectively picking the biggest of n test statistics such that it is not legitimate simply to use t_{n-p-2} to find a p -value. We need a correction on the p -value because of multiple-comparisons.

Suppose that we have $p' = \Pr(t_{n-p-2} > |\hat{\epsilon}^*_{max}|)$, the p -value before correction. Then the Bonferroni adjusted p -value is $p = np'$.

Measuring influence

Influence on the regression coefficients combines leverage and discrepancy. The most direct measure of influence simply expresses the impact on each coefficient of deleting each observation in turn:

$$D_{ij} = \hat{\beta}_j - \tilde{\beta}_{j(-i)} \quad \text{for } i = 1, \dots, n \text{ and } j = 0, 1, \dots, p$$

where $\hat{\beta}_j$ are the least-squares coefficients calculated for all the data, and the $\tilde{\beta}_{j(-i)}$ are the least-squares coefficients calculated with the i th observation omitted. To assist in interpretation, it is useful to scale the D_{ij} by (deleted) coefficient standard errors:

$$D_{ij}^* = \frac{D_{ij}}{\widehat{SE}_{(-i)}(\tilde{\beta}_{j(-i)})}$$

Following Belsley, Kuh, and Welsh (1980), the D_{ij} are often termed $DFBETA_{ij}$, and D_{ij}^* are called $DFBETAS_{ij}$. One problem associated with using D_{ij} or D_{ij}^* is their large number: $n(p+1)$ of each.

Cook's distance is another popular measure, calculated as

$$D_i = \frac{\sum_{j=1}^n (\tilde{y}_{j(-i)} - \hat{y}_j)^2}{(p+1)\hat{\sigma}^2} = \frac{\hat{\epsilon}_i'^2}{p+1} \times \frac{h_i}{1-h_i}$$

In effect, the first term in the formula for Cook's D is a measure of discrepancy, and the second is a measure of leverage. We look for values of D_i that stand out from the rest.

A similar measure suggested by Belsley et al. (1980)

$$DFFITS_i = \hat{\epsilon}_i^* \frac{h_i}{1-h_i}$$

Except for unusual data configurations, Cook's $D_i \approx DFFITS_i^2/(p+1)$.

Numerical cutoffs (suggested)

Diagnostic statistic	Cutoff value
h_i	$2\bar{h} = \frac{2(p+1)}{n}$, ($3\bar{h}$ for small sample)
D_{ij}^*	$ D_{ij}^* > 1$ or 2 ($2/\sqrt{n}$ for large samples)
Cook's D_i	$D_i > \frac{4}{n-p-1}$
DFFITS	$ DFFITS_i > 2\sqrt{\frac{p+1}{n-p-1}}$