

Statistic	Formula	Extreme?	R
Leverage	$h_i = \frac{(X_i - \bar{X})^2}{\sum_{j=1}^n (X_j - \bar{X})^2} + \frac{1}{n} = \mathbf{X}_i^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}_i$	$> \frac{2p}{n}$ or $.2-.5 =$ moderate, $> .5$ high	<code>hatvalues</code>
DFFITs	$\frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{MSE_{(i)} h_{ii}}}$	$> 1$ for med-sized data sets, $> 2\sqrt{\frac{p}{n}}$ for large data sets	<code>dffits</code>
Cook's Distance	$D_i = \frac{\sum_{j=1}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{pMSE}$	$\geq 1$	<code>cooks.distance</code>
DFBETAS	$\frac{b_k - b_{k(i)}}{\sqrt{MSE_{(i)} c_{kk}}}$ $c_{kk} = (\mathbf{X}^t \mathbf{X})_{kk}^{-1}$	$> 1$ for med-sized data sets, $> 2/\sqrt{n}$ for large data sets	<code>dfbetas</code>
Resid	$e_i = (Y_i - \hat{Y}_i)$		<code>resid</code>
Semi-studentized Resid	$\frac{e_i}{\sqrt{MSE}}$	outside $(-2,2)$	<code>rstandard</code>
(Internally) Studentized Resid	$\frac{e_i}{\sqrt{MSE\sqrt{1-h_{ii}}}}$	outside $(-2,2)$	
Deleted Studentized Resid	$\frac{e_i}{\sqrt{MSE_{(i)}\sqrt{1-h_{ii}}}}$	outside $(-2,2)$	<code>rstudent</code>
VIF	$(1 - R_k^2)^{-1}$ $R_k^2$ from $X_k$ regressed on $(p - 2)$ vars	$\max(VIF) > 10$ $\text{mean}(VIF) \gg 1$	<code>vif</code> package car

Notes:

- The first four statistics are measures of how **influential** the value is. Leverage measures the distance of the explanatory variables from the average. Cook's distances, and the derivatives, are a measure of how much the predicted values change when the point is removed from the model.
- The residual statistics are measures of how well the regression line fits the value. A residual is the distance from the point to the line. We standardize the residual in different ways. The studentized residuals contain the more accurate measure of standard error.
- The VIF measures the degree of collinearity between the explanatory variables. Collinear variables indicates that we should be cautious interpreting any coefficients.  $\text{mean}(VIF) \gg 1$  is meant to indicate that the average VIF is considerably larger than 1.
- Any value containing a "(*i*)" indicates that the  $i^{\text{th}}$  point was removed before calculating the value. For example,  $MSE_{(i)}$  is the  $MSE$  for the full model containing all the data **except** the  $i^{\text{th}}$  point.
- Most of the functions are in R under a general heading of `influence.measures`. The `vif` function is in the `car` package.