

**regress postestimation** — Postestimation tools for regress

<a href="#">Postestimation commands</a>	<a href="#">Predictions</a>	<a href="#">margins</a>
<a href="#">DFBETA influence statistics</a>	<a href="#">Tests for violation of assumptions</a>	<a href="#">Variance inflation factors</a>
<a href="#">Measures of effect size</a>	<a href="#">Methods and formulas</a>	<a href="#">Acknowledgments</a>
<a href="#">References</a>	<a href="#">Also see</a>	

## Postestimation commands

The following postestimation commands are of special interest after `regress`:

Command	Description
<code>dfbeta</code>	DFBETA influence statistics
<code>estat hettest</code>	tests for heteroskedasticity
<code>estat imtest</code>	information matrix test
<code>estat ovtest</code>	Ramsey regression specification-error test for omitted variables
<code>estat szroeter</code>	Szroeter's rank test for heteroskedasticity
<code>estat vif</code>	variance inflation factors for the independent variables
<code>estat esize</code>	$\eta^2$ , $\varepsilon^2$ , and $\omega^2$ effect sizes
<code>estat moran</code>	Moran's test of residual correlation with nearby residuals
<code>lassogof</code>	calculate goodness-of-fit predictions

These commands are not appropriate with `svy` estimation results.

The following standard postestimation commands are also available:

Command	Description
<code>contrast</code>	contrasts and ANOVA-style joint tests of estimates
<code>estat ic</code>	Akaike's, consistent Akaike's, corrected Akaike's, and Schwarz's Bayesian information criteria (AIC, CAIC, AICC, and BIC)
<code>estat summarize</code>	summary statistics for the estimation sample
<code>estat vce</code>	variance-covariance matrix of the estimators (VCE)
<code>estat (svy)</code>	postestimation statistics for survey data
<code>estimates</code>	cataloging estimation results
<code>etable</code>	table of estimation results
* <code>forecast</code>	dynamic forecasts and simulations
* <code>hausman</code>	Hausman's specification test
<code>lincom</code>	point estimates, standard errors, testing, and inference for linear combinations of coefficients
<code>linktest</code>	link test for model specification
* <code>lrtest</code>	likelihood-ratio test
<code>margins</code>	marginal means, predictive margins, marginal effects, and average marginal effects
<code>marginsplot</code>	graph the results from margins (profile plots, interaction plots, etc.)
<code>nlcom</code>	point estimates, standard errors, testing, and inference for nonlinear combinations of coefficients
<code>predict</code>	predictions and their SEs, leverage statistics, distance statistics, etc.
<code>predictnl</code>	point estimates, standard errors, testing, and inference for generalized predictions
<code>pwcompare</code>	pairwise comparisons of estimates
<code>suest</code>	seemingly unrelated estimation
<code>test</code>	Wald tests of simple and composite linear hypotheses
<code>testnl</code>	Wald tests of nonlinear hypotheses

---

\* `forecast`, `hausman`, and `lrtest` are not appropriate with `svy` estimation results. `forecast` is also not appropriate with `mi` estimation results.

# Predictions

## Description for predict

`predict` creates a new variable containing predictions such as linear predictions, residuals, standardized residuals, Studentized residuals, Cook's distance, leverage, probabilities, expected values, DFBETAs for *varname*, standard errors, COVRATIOs, DFITS, and Welsch distances.

## Menu for predict

Statistics > Postestimation

## Syntax for predict

```
predict [type] newvar [if] [in] [, statistic]
```

<i>statistic</i>	Description
Main	
<code>xb</code>	linear prediction; the default
<code>residuals</code>	residuals
<code>score</code>	score; equivalent to <code>residuals</code>
<code>rstandard</code>	standardized residuals
<code>rstudent</code>	Studentized (jackknifed) residuals
<code>cooksd</code>	Cook's distance
<code>leverage   hat</code>	leverage (diagonal elements of hat matrix)
<code>pr(<i>a</i>,<i>b</i>)</code>	$\Pr(y_j   a < y_j < b)$
<code>e(<i>a</i>,<i>b</i>)</code>	$E(y_j   a < y_j < b)$
<code>ystar(<i>a</i>,<i>b</i>)</code>	$E(y_j^*), y_j^* = \max\{a, \min(y_j, b)\}$
* <code>dfbeta(<i>varname</i>)</code>	DFBETA for <i>varname</i>
<code>stdp</code>	standard error of the linear prediction
<code>stdf</code>	standard error of the forecast
<code>stdr</code>	standard error of the residual
* <code>covratio</code>	COVRATIO
* <code>dfits</code>	DFITS
* <code>welsch</code>	Welsch distance

Unstarred statistics are available both in and out of sample; type `predict ... if e(sample) ...` if wanted only for the estimation sample. Starred statistics are calculated only for the estimation sample, even when `if e(sample)` is not specified.

`rstandard`, `rstudent`, `cooksd`, `leverage`, `dfbeta()`, `stdf`, `stdr`, `covratio`, `dfits`, and `welsch` are not available if any `vce()` other than `vce(ols)` was specified with `regress`.

`xb`, `residuals`, `score`, and `stdp` are the only options allowed with `svy` estimation results.

where *a* and *b* may be numbers or variables; *a* missing ( $a \geq .$ ) means  $-\infty$ , and *b* missing ( $b \geq .$ ) means  $+\infty$ ; see [U] 12.2.1 Missing values.

## Options for predict

Main

`xb`, the default, calculates the linear prediction.

`residuals` calculates the residuals.

`score` is equivalent to `residuals` in linear regression.

`rstandard` calculates the standardized residuals.

`rstudent` calculates the Studentized (jackknifed) residuals.

`cooksdi` calculates the Cook's  $D$  influence statistic (Cook 1977).

`leverage` or `hat` calculates the diagonal elements of the projection (“hat”) matrix.

`pr(a,b)` calculates  $\Pr(a < \mathbf{x}_j \mathbf{b} + u_j < b)$ , the probability that  $y_j | \mathbf{x}_j$  would be observed in the interval  $(a, b)$ .

*a* and *b* may be specified as numbers or variable names; *lb* and *ub* are variable names;

`pr(20,30)` calculates  $\Pr(20 < \mathbf{x}_j \mathbf{b} + u_j < 30)$ ;

`pr(lb,ub)` calculates  $\Pr(*lb* < \mathbf{x}_j \mathbf{b} + u_j < *ub*)$ ; and

`pr(20,ub)` calculates  $\Pr(20 < \mathbf{x}_j \mathbf{b} + u_j < *ub*)$ .

*a* missing ( $a \geq .$ ) means  $-\infty$ ; `pr(. ,30)` calculates  $\Pr(-\infty < \mathbf{x}_j \mathbf{b} + u_j < 30)$ ;

`pr(lb,30)` calculates  $\Pr(-\infty < \mathbf{x}_j \mathbf{b} + u_j < 30)$  in observations for which  $lb \geq .$

and calculates  $\Pr(*lb* < \mathbf{x}_j \mathbf{b} + u_j < 30)$  elsewhere.

*b* missing ( $b \geq .$ ) means  $+\infty$ ; `pr(20, .)` calculates  $\Pr(+\infty > \mathbf{x}_j \mathbf{b} + u_j > 20)$ ;

`pr(20,ub)` calculates  $\Pr(+\infty > \mathbf{x}_j \mathbf{b} + u_j > 20)$  in observations for which  $ub \geq .$

and calculates  $\Pr(20 < \mathbf{x}_j \mathbf{b} + u_j < *ub*)$  elsewhere.

`e(a,b)` calculates  $E(\mathbf{x}_j \mathbf{b} + u_j \mid a < \mathbf{x}_j \mathbf{b} + u_j < b)$ , the expected value of  $y_j | \mathbf{x}_j$  conditional on  $y_j | \mathbf{x}_j$  being in the interval  $(a, b)$ , meaning that  $y_j | \mathbf{x}_j$  is truncated. *a* and *b* are specified as they are for `pr()`.

`ystar(a,b)` calculates  $E(y_j^*)$ , where  $y_j^* = a$  if  $\mathbf{x}_j \mathbf{b} + u_j \leq a$ ,  $y_j^* = b$  if  $\mathbf{x}_j \mathbf{b} + u_j \geq b$ , and  $y_j^* = \mathbf{x}_j \mathbf{b} + u_j$  otherwise, meaning that  $y_j^*$  is censored. *a* and *b* are specified as they are for `pr()`.

`dfbeta(varname)` calculates the DFBETA for *varname*, the difference between the regression coefficient when the *j*th observation is included and excluded, said difference being scaled by the estimated standard error of the coefficient. *varname* must have been included among the regressors in the previously fitted model. The calculation is automatically restricted to the estimation subsample.

`stdp` calculates the standard error of the prediction, which can be thought of as the standard error of the predicted expected value or mean for the observation's covariate pattern. The standard error of the prediction is also referred to as the standard error of the fitted value.

`stdf` calculates the standard error of the forecast, which is the standard error of the point prediction for 1 observation. It is commonly referred to as the standard error of the future or forecast value. By construction, the standard errors produced by `stdf` are always larger than those produced by `stdp`; see [Methods and formulas](#).

`stdr` calculates the standard error of the residuals.

`covratio` calculates COVRATIO (Belsley, Kuh, and Welsch 1980), a measure of the influence of the *j*th observation based on considering the effect on the variance–covariance matrix of the estimates. The calculation is automatically restricted to the estimation subsample.

`dfits` calculates DFITS (Welsch and Kuh 1977) and attempts to summarize the information in the leverage versus residual-squared plot into one statistic. The calculation is automatically restricted to the estimation subsample.

`welsch` calculates Welsch distance (Welsch 1982) and is a variation on `dfits`. The calculation is automatically restricted to the estimation subsample.

## Remarks and examples for predict

Remarks are presented under the following headings:

*Terminology*  
*Fitted values and residuals*  
*Prediction standard errors*  
*Prediction with weighted data*  
*Leverage statistics*  
*Standardized and Studentized residuals*  
*DFITS, Cook's Distance, and Welsch Distance*  
*COVRATIO*

## Terminology

Many of these commands concern identifying influential data in linear regression. This is, unfortunately, a field that is dominated by jargon, codified and partially begun by Belsley, Kuh, and Welsch (1980). In the words of Chatterjee and Hadi (1986, 416), “Belsley, Kuh, and Welsch’s book, *Regression Diagnostics*, was a very valuable contribution to the statistical literature, but it unleashed on an unsuspecting statistical community a computer speak (à la Orwell), the likes of which we have never seen.” Things have only gotten worse since then. Chatterjee and Hadi’s (1986, 1988) own attempts to clean up the jargon did not improve matters (see Hoaglin and Kempthorne [1986], Velleman [1986], and Welsch [1986]). We apologize for the jargon, and for our contribution to the jargon in the form of inelegant command names, we apologize most of all.

Model *sensitivity* refers to how estimates are affected by subsets of our data. Imagine data on  $y$  and  $x$ , and assume that the data are to be fit by the regression  $y_i = \alpha + \beta x_i + \epsilon_i$ . The regression estimates of  $\alpha$  and  $\beta$  are  $a$  and  $b$ , respectively. Now imagine that the estimated  $a$  and  $b$  would be different if a small portion of the dataset, perhaps even one observation, were deleted. As a data analyst, you would like to think that you are summarizing tendencies that apply to all the data, but you have just been told that the model you fit is unduly influenced by one point or just a few points and that, as a matter of fact, there is another model that applies to the rest of the data—a model that you have ignored. The search for subsets of the data that, if deleted, would change the results markedly is a predominant theme of this entry.

There are three key issues in identifying model sensitivity to individual observations, which go by the names *residuals*, *leverage*, and *influence*. In our  $y_i = a + bx_i + e_i$  regression, the residuals are, of course,  $e_i$ —they reveal how much our fitted value  $\hat{y}_i = a + bx_i$  differs from the observed  $y_i$ . A point  $(x_i, y_i)$  with a corresponding large residual is called an outlier. Say that you are interested in outliers because you somehow think that such points will exert undue influence on your estimates. Your feelings are generally right, but there are exceptions. A point might have a huge residual and yet not affect the estimated  $b$  at all. Nevertheless, studying observations with large residuals almost always pays off.

$(x_i, y_i)$  can be an outlier in another way—just as  $y_i$  can be far from  $\hat{y}_i$ ,  $x_i$  can be far from the center of mass of the other  $x$ ’s. Such an “outlier” should interest you just as much as the more traditional outliers. Picture a scatterplot of  $y$  against  $x$  with thousands of points in some sort of mass

at the lower left of the graph and one point at the upper right of the graph. Now, run a regression line through the points—the regression line will come close to the point at the upper right of the graph and may in fact, go through it. That is, this isolated point will not appear as an outlier as measured by residuals because its residual will be small. Yet this point might have a dramatic effect on our resulting estimates in the sense that, were you to delete the point, the estimates would change markedly. Such a point is said to have high leverage. Just as with traditional outliers, a high leverage point does not necessarily have an undue effect on regression estimates, but if it does not, it is more the exception than the rule.

Now, all of this is a most unsatisfactory state of affairs. Points with large residuals may, but need not, have a large effect on our results, and points with small residuals may still have a large effect. Points with high leverage may, but need not, have a large effect on our results, and points with low leverage may still have a large effect. Can you not identify the influential points and simply have the computer list them for you? You can, but you will have to define what you mean by “influential”.

“Influential” is defined with respect to some statistic. For instance, you might ask which points in your data have a large effect on your estimated  $a$ , which points have a large effect on your estimated  $b$ , which points have a large effect on your estimated standard error of  $b$ , and so on, but do not be surprised when the answers to these questions are different. In any case, obtaining such measures is not difficult—all you have to do is fit the regression excluding each observation one at a time and record the statistic of interest which, in the day of the modern computer, is not too onerous. Moreover, you can save considerable computer time by doing algebra ahead of time and working out formulas that will calculate the same answers as if you ran each of the regressions. (Ignore the question of pairs of observations that, together, exert undue influence, and triples, and so on, which remains largely unsolved and for which the brute force fit-every-possible-regression procedure is not a viable alternative.)

## Fitted values and residuals

Typing `predict newvar` with no options creates `newvar` containing the fitted values. Typing `predict newvar, resid` creates `newvar` containing the residuals.

### ▷ Example 1

Continuing with [example 1](#) from [\[R\] regress](#), we wish to fit the following model:

$$\text{mpg} = \beta_0 + \beta_1 \text{weight} + \beta_2 \text{foreign} + \epsilon$$

```
. use https://www.stata-press.com/data/r18/auto
(1978 automobile data)
. regress mpg weight foreign
```

Source	SS	df	MS	Number of obs	=	74
Model	1619.2877	2	809.643849	F(2, 71)	=	69.75
Residual	824.171761	71	11.608053	Prob > F	=	0.0000
				R-squared	=	0.6627
				Adj R-squared	=	0.6532
Total	2443.45946	73	33.4720474	Root MSE	=	3.4071

  

mpg	Coefficient	Std. err.	t	P> t	[95% conf. interval]
weight	-.0065879	.0006371	-10.34	0.000	-.0078583    -.0053175
foreign	-1.650029	1.075994	-1.53	0.130	-3.7955    .4954422
_cons	41.6797	2.165547	19.25	0.000	37.36172    45.99768

That done, we can now obtain the predicted values from the regression. We will store them in a new variable called `pmpg` by typing `predict pmpg`. Because `predict` produces no output, we will follow that by summarizing our predicted and observed values.

```
. predict pmpg
(option xb assumed; fitted values)
. summarize pmpg mpg
```

Variable	Obs	Mean	Std. dev.	Min	Max
pmpg	74	21.2973	4.709779	9.794333	29.82151
mpg	74	21.2973	5.785503	12	41

◀

## ▶ Example 2: Out-of-sample predictions

We can just as easily obtain predicted values from the model by using a wholly different dataset from the one on which the model was fit. The only requirement is that the data have the necessary variables, which here are `weight` and `foreign`.

Using the data on two new cars (the Pontiac Sunbird and the Volvo 260) from `newautos.dta`, we can obtain out-of-sample predictions (or forecasts) by typing

```
. use https://www.stata-press.com/data/r18/newautos, clear
(New automobile models)
. predict pmpg
(option xb assumed; fitted values)
. list, divider
```

	make	weight	foreign	pmpg
1.	Pont. Sunbird	2690	Domestic	23.95829
2.	Volvo 260	3170	Foreign	19.14607

The Pontiac Sunbird has a predicted mileage rating of 23.96 mpg, whereas the Volvo 260 has a predicted rating of 19.15 mpg. In comparison, the actual mileage ratings are 24 for the Pontiac and 17 for the Volvo.

◀

## Prediction standard errors

`predict` can calculate the standard error of the forecast (`stdf` option), the standard error of the prediction (`stdp` option), and the standard error of the residual (`stdr` option). It is easy to confuse `stdf` and `stdp` because both are often called the prediction error. Consider the prediction  $\hat{y}_j = \mathbf{x}_j \mathbf{b}$ , where  $\mathbf{b}$  is the estimated coefficient (column) vector and  $\mathbf{x}_j$  is a (row) vector of independent variables for which you want the prediction. First,  $\hat{y}_j$  has a variance due to the variance of the estimated coefficient vector  $\mathbf{b}$ ,

$$\text{Var}(\hat{y}_j) = \text{Var}(\mathbf{x}_j \mathbf{b}) = s^2 h_j$$

where  $h_j = \mathbf{x}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_j'$  and  $s^2$  is the mean squared error of the regression. Do not panic over the algebra—just remember that  $\text{Var}(\hat{y}_j) = s^2 h_j$ , whatever  $s^2$  and  $h_j$  are. `stdp` calculates this quantity. This is the error in the prediction due to the uncertainty about  $\mathbf{b}$ .

If you are about to hand this number out as your forecast, however, there is another error. According to your model, the true value of  $y_j$  is given by

$$y_j = \mathbf{x}_j \mathbf{b} + \epsilon_j = \hat{y}_j + \epsilon_j$$

and thus the  $\text{Var}(y_j) = \text{Var}(\hat{y}_j) + \text{Var}(\epsilon_j) = s^2 h_j + s^2$ , which is the square of `stdf`. `stdf`, then, is the sum of the error in the prediction plus the residual error.

`stdr` has to do with an analysis-of-variance decomposition of  $s^2$ , the estimated variance of  $y$ . The standard error of the prediction is  $s^2 h_j$ , and therefore  $s^2 h_j + s^2(1 - h_j) = s^2$  decomposes  $s^2$  into the prediction and residual variances.

### ► Example 3: Standard error of the forecast

Returning to our model of `mpg` on `weight` and `foreign`, we previously predicted the mileage rating for the Pontiac Sunbird and Volvo 260 as 23.96 and 19.15 mpg, respectively. We now want to put a standard error around our forecast. Remember, the data for these two cars were in `newautos.dta`:

```
. use https://www.stata-press.com/data/r18/newautos, clear
(New automobile models)
. predict pmpg
(option xb assumed; fitted values)
. predict se_pmpg, stdf
. list, divider
```

	make	weight	foreign	pmpg	se_pmpg
1.	Pont. Sunbird	2690	Domestic	23.95829	3.462791
2.	Volvo 260	3170	Foreign	19.14607	3.525875

Thus, an approximate 95% confidence interval for the mileage rating of the Volvo 260 is  $19.15 \pm 2 \cdot 3.53 = [12.09, 26.21]$ .



### ▣ Prediction with weighted data

`predict` can be used after frequency-weighted (`fweight`) estimation, just as it is used after unweighted estimation. The technical note below concerns the use of `predict` after analytically weighted (`aweight`) estimation.

#### ▣ Technical note

After analytically weighted estimation, `predict` is willing to calculate only the prediction (no options), residual (`residual` option), standard error of the prediction (`stdp` option), and diagonal elements of the projection matrix (`hat` option). For analytically weighted estimation, the standard error of the forecast and residuals, the standardized and Studentized residuals, and Cook's  $D$  are not statistically well-defined concepts.



### Leverage statistics

In addition to providing fitted values and the associated standard errors, the `predict` command can also be used to generate various statistics used to detect the influence of individual observations. This section provides a brief introduction to leverage (`hat`) statistics, and some of the following subsections discuss other influence statistics produced by `predict`.



### ► Example 4: Diagonal elements of projection matrix

The diagonal elements of the projection matrix, obtained by the `hat` option, are a measure of distance in explanatory variable space. `leverage` is a synonym for `hat`.

```
. use https://www.stata-press.com/data/r18/auto, clear
(1978 automobile data)

. regress mpg weight foreign
(output omitted)

. predict xdist, hat

. summarize xdist, detail
```

Leverage			
	Percentiles	Smallest	
1%	.0192325	.0192325	
5%	.0192686	.0192366	
10%	.0193448	.019241	Obs 74
25%	.0220291	.0192686	Sum of wgt. 74
50%	.0383797		Mean .0405405
		Largest	Std. dev. .0207624
75%	.0494002	.0880814	
90%	.0693432	.099715	Variance .0004311
95%	.0880814	.099715	Skewness 1.159745
99%	.1003283	.1003283	Kurtosis 4.083313

Some 5% of our sample has an `xdist` measure in excess of 0.08. Let's force them to reveal their identities:

```
. list foreign make mpg if xdist>.08, divider
```

	foreign	make	mpg
24.	Domestic	Ford Fiesta	28
26.	Domestic	Linc. Continental	12
27.	Domestic	Linc. Mark V	12
43.	Domestic	Plym. Champ	34
64.	Foreign	Peugeot 604	14

To understand why these cars are on this list, we must remember that the explanatory variables in our model are `weight` and `foreign` and that `xdist` measures distance in this metric. The Ford Fiesta and the Plymouth Champ are the two lightest domestic cars in our data. The Lincolns are the two heaviest domestic cars, and the Peugeot is the heaviest foreign car. ◀

See `lvr2plot` in [R] [regress postestimation diagnostic plots](#) for information on a leverage-versus-squared-residual plot.

### Standardized and Studentized residuals

The terms standardized and Studentized residuals have meant different things to different authors. In Stata, `predict` defines the standardized residual as  $\hat{e}_{s_i} = e_i / (s\sqrt{1 - h_i})$  and the Studentized residual as  $r_i = e_i / (s_{(i)}\sqrt{1 - h_i})$ , where  $s_{(i)}$  is the root mean squared error of a regression with the  $i$ th observation removed. Stata's definition of the Studentized residual is the same as the one given in [Bollen and Jackman \(1990, 264\)](#) and is what [Chatterjee and Hadi \(1988, 74\)](#) call the "externally Studentized" residual. Stata's "standardized" residual is the same as what [Chatterjee and Hadi \(1988, 74\)](#) call the "internally Studentized" residual.

Standardized and Studentized residuals are attempts to adjust residuals for their standard errors. Although the  $\epsilon_i$  theoretical residuals are homoskedastic by assumption (that is, they all have the same variance), the calculated  $e_i$  are not. In fact,

$$\text{Var}(e_i) = \sigma^2(1 - h_i)$$

where  $h_i$  are the leverage measures obtained from the diagonal elements of hat matrix. Thus, observations with the greatest leverage have corresponding residuals with the smallest variance.

Standardized residuals use the root mean squared error of the regression for  $\sigma$ . Studentized residuals use the root mean squared error of a regression omitting the observation in question for  $\sigma$ . In general, Studentized residuals are preferable to standardized residuals for purposes of outlier identification. Studentized residuals can be interpreted as the  $t$  statistic for testing the significance of a dummy variable equal to 1 in the observation in question and 0 elsewhere (Belsley, Kuh, and Welsch 1980). Such a dummy variable would effectively absorb the observation and so remove its influence in determining the other coefficients in the model. Caution must be exercised here, however, because of the simultaneous testing problem. You cannot simply list the residuals that would be individually significant at the 5% level—their joint significance would be far less (their joint significance level would be far greater).

### ► Example 5: Standardized and Studentized residuals

In the *Terminology* section of *Remarks and examples for predict*, we distinguished residuals from leverage and speculated on the impact of an observation with a small residual but large leverage. If we adjust the residuals for their standard errors, however, the adjusted residual would be (relatively) larger and perhaps large enough so that we could simply examine the adjusted residuals. Taking our price on weight and foreign##c.mpg model from [example 1](#) of [\[R\] regress postestimation diagnostic plots](#), we can obtain the in-sample standardized and Studentized residuals by typing

```
. use https://www.stata-press.com/data/r18/auto, clear
(1978 automobile data)
. regress price weight foreign##c.mpg
(output omitted)
. predict esta if e(sample), rstandard
. predict estu if e(sample), rstudent
```

In the *lvr2plot* section of [\[R\] regress postestimation diagnostic plots](#), we discovered that the VW Diesel has the highest leverage in our data, but a corresponding small residual. The standardized and Studentized residuals for the VW Diesel are

```
. list make price esta estu if make=="VW Diesel"
```

	make	price	esta	estu
71.	VW Diesel	5,397	.6142691	.6114758

The Studentized residual of 0.611 can be interpreted as the  $t$  statistic for including a dummy variable for VW Diesel in our regression. Such a variable would not be significant.

## DFITS, Cook's Distance, and Welsch Distance

DFITS (Welsch and Kuh 1977), Cook's Distance (Cook 1977), and Welsch Distance (Welsch 1982) are three attempts to summarize the information in the leverage versus residual-squared plot into one statistic. That is, the goal is to create an index that is affected by the size of the residuals—outliers—and the size of  $h_i$ —leverage. Viewed mechanically, one way to write DFITS (Bollen and Jackman 1990, 265) is

$$\text{DFITS}_i = r_i \sqrt{\frac{h_i}{1 - h_i}}$$

where  $r_i$  are the Studentized residuals. Thus, large residuals increase the value of DFITS, as do large values of  $h_i$ . Viewed more traditionally, DFITS is a scaled difference between predicted values for the  $i$ th case when the regression is fit with and without the  $i$ th observation, hence the name.

The mechanical relationship between DFITS and Cook's Distance,  $D_i$  (Bollen and Jackman 1990, 266), is

$$D_i = \frac{1}{k} \frac{s_{(i)}^2}{s^2} \text{DFITS}_i^2$$

where  $k$  is the number of variables (including the constant) in the regression,  $s$  is the root mean squared error of the regression, and  $s_{(i)}$  is the root mean squared error when the  $i$ th observation is omitted. Viewed more traditionally,  $D_i$  is a scaled measure of the distance between the coefficient vectors when the  $i$ th observation is omitted.

The mechanical relationship between DFITS and Welsch's Distance,  $W_i$  (Chatterjee and Hadi 1988, 123), is

$$W_i = \text{DFITS}_i \sqrt{\frac{n-1}{1-h_i}}$$

The interpretation of  $W_i$  is more difficult because it is based on the empirical influence curve. Although DFITS and Cook's distance are similar, the Welsch distance measure includes another normalization by leverage.

Belsley, Kuh, and Welsch (1980, 28) suggest that DFITS values greater than  $2\sqrt{k/n}$  deserve more investigation, and so values of Cook's distance greater than  $4/n$  should also be examined (Bollen and Jackman 1990, 265–266). Through similar logic, the cutoff for Welsch distance is approximately  $3\sqrt{k}$  (Chatterjee and Hadi 1988, 124).

### ► Example 6: DFITS influence measure

Continuing with our model of price on weight and foreign##c.mpg, we can obtain the DFITS influence measure:

```
. predict e if e(sample), resid
. predict dfits, dfits
```

We did not specify `if e(sample)` in computing the DFITS statistic. DFITS is available only over the estimation sample, so specifying `if e(sample)` would have been redundant. It would have done no harm, but it would not have changed the results.

Our model has  $k = 5$  independent variables ( $k$  includes the constant) and  $n = 74$  observations; following the  $2\sqrt{k/n}$  cutoff advice, we type

```
. list make price e dfits if abs(dfits) > 2*sqrt(5/74), divider
```

	make	price	e	dfits
12.	Cad. Eldorado	14,500	7271.96	.9564455
13.	Cad. Seville	15,906	5036.348	1.356619
24.	Ford Fiesta	4,389	3164.872	.5724172
27.	Linc. Mark V	13,594	3109.193	.5200413
28.	Linc. Versailles	13,466	6560.912	.8760136
42.	Plym. Arrow	4,647	-3312.968	-.9384231

We calculate Cook's distance and list the observations greater than the suggested  $4/n$  cutoff:

```
. predict cooks if e(sample), cooks
. list make price e cooks if cooks > 4/74, divider
```

	make	price	e	cooks
12.	Cad. Eldorado	14,500	7271.96	.1492676
13.	Cad. Seville	15,906	5036.348	.3328515
24.	Ford Fiesta	4,389	3164.872	.0638815
28.	Linc. Versailles	13,466	6560.912	.1308004
42.	Plym. Arrow	4,647	-3312.968	.1700736

Here we used `if e(sample)` because Cook's distance is not restricted to the estimation sample by default. It is worth comparing this list with the preceding one.

Finally, we use Welsch distance and the suggested  $3\sqrt{k}$  cutoff:

```
. predict wd, welsch
. list make price e wd if abs(wd) > 3*sqrt(5), divider
```

	make	price	e	wd
12.	Cad. Eldorado	14,500	7271.96	8.394372
13.	Cad. Seville	15,906	5036.348	12.81125
28.	Linc. Versailles	13,466	6560.912	7.703005
42.	Plym. Arrow	4,647	-3312.968	-8.981481

Here we did not need to specify `if e(sample)` because `welsch` automatically restricts the prediction to the estimation sample.

◀

## COVRATIO

COVRATIO (Belsley, Kuh, and Welsch 1980) measures the influence of the  $i$ th observation by considering the effect on the variance–covariance matrix of the estimates. The measure is the ratio of the determinants of the covariances matrix, with and without the  $i$ th observation. The resulting formula is

$$\text{COVRATIO}_i = \frac{1}{1 - h_i} \left( \frac{n - k - \hat{e}_{s_i}^2}{n - k - 1} \right)^k$$

where  $\hat{e}_{s_i}$  is the standardized residual.

For noninfluential observations, the value of COVRATIO is approximately 1. Large values of the residuals or large values of leverage will cause deviations from 1, although if both are large, COVRATIO may tend back toward 1 and therefore not identify such observations (Chatterjee and Hadi 1988, 139).

Belsley, Kuh, and Welsch (1980) suggest that observations for which

$$|\text{COVRATIO}_i - 1| \geq \frac{3k}{n}$$

are worthy of further examination.

### ► Example 7: COVRATIO influence measure

Using our model of price on weight and foreign##c.mpg, we can obtain the COVRATIO measure and list the observations outside the suggested cutoff by typing

```
. predict covr, covratio
. list make price e covr if abs(covr-1) >= 3*5/74, divider
```

	make	price	e	covr
12.	Cad. Eldorado	14,500	7271.96	.3814242
13.	Cad. Seville	15,906	5036.348	.7386969
28.	Linc. Versailles	13,466	6560.912	.4761695
43.	Plym. Champ	4,425	1621.747	1.27782
53.	Audi 5000	9,690	591.2883	1.206842
57.	Datsun 210	4,589	19.81829	1.284801
64.	Peugeot 604	12,990	1037.184	1.348219
66.	Subaru	3,798	-909.5894	1.264677
71.	VW Diesel	5,397	999.7209	1.630653
74.	Volvo 260	11,995	1327.668	1.211888

The covratio option automatically restricts the prediction to the estimation sample.

## margins

### Description for margins

`margins` estimates margins of response for linear predictions.

### Menu for margins

Statistics > Postestimation

### Syntax for margins

```
margins [marginlist] [, options]  
margins [marginlist] , predict(statistic ...) [options]
```

<i>statistic</i>	Description
<code>xb</code>	linear prediction; the default
<code>pr(<i>a,b</i>)</code>	not allowed with margins
<code>e(<i>a,b</i>)</code>	not allowed with margins
<code><u>ystar</u>(<i>a,b</i>)</code>	not allowed with margins
<code><u>residuals</u></code>	not allowed with margins
<code><u>score</u></code>	not allowed with margins
<code><u>rstandard</u></code>	not allowed with margins
<code><u>rstudent</u></code>	not allowed with margins
<code><u>cooksd</u></code>	not allowed with margins
<code><u>leverage</u>   <u>hat</u></code>	not allowed with margins
<code><u>dfbeta</u>(<i>varname</i>)</code>	not allowed with margins
<code><u>stdp</u></code>	not allowed with margins
<code><u>stdf</u></code>	not allowed with margins
<code><u>stdr</u></code>	not allowed with margins
<code><u>covratio</u></code>	not allowed with margins
<code><u>dfits</u></code>	not allowed with margins
<code><u>welsch</u></code>	not allowed with margins

Statistics not allowed with margins are functions of stochastic quantities other than  $e(b)$ .

For the full syntax, see [R] [margins](#).

## DFBETA influence statistics

### Description for dfbeta

`dfbeta` will calculate one, more than one, or all the DFBETAs after `regress`. Although `predict` will also calculate DFBETAs, `predict` can do this for only one variable at a time. `dfbeta` is a convenience tool for those who want to calculate DFBETAs for multiple variables. The names for the new variables created are chosen automatically and begin with the letters `_dfbeta_`.

### Menu for dfbeta

Statistics > Linear models and related > Regression diagnostics > DFBETAs

### Syntax for dfbeta

```
dfbeta [indepvar [indepvar [...]]] [, stub(name) ]
```

### Option for dfbeta

`stub(name)` specifies the leading characters `dfbeta` uses to name the new variables to be generated. The default is `stub(_dfbeta_)`.

### Remarks and examples for dfbeta

DFBETAs are perhaps the most direct influence measure of interest to model builders. DFBETAs focus on one coefficient and measure the difference between the regression coefficient when the  $i$ th observation is included and excluded, the difference being scaled by the estimated standard error of the coefficient. [Belsley, Kuh, and Welsch \(1980, 28\)](#) suggest observations with  $|DFBETA_i| > 2/\sqrt{n}$  as deserving special attention, but it is also common practice to use 1 ([Bollen and Jackman 1990, 267](#)), meaning that the observation shifted the estimate at least one standard error.

#### ► Example 8: DFBETAs influence measure; the `dfbeta()` option

Using our model of `price` on `weight` and `foreign##c.mpg`, let's first ask which observations have the greatest impact on the determination of the coefficient on `1.foreign`. We will use the suggested  $2/\sqrt{n}$  cutoff:

```
. use https://www.stata-press.com/data/r18/auto, clear
(1978 automobile data)
. regress price weight foreign##c.mpg
(output omitted)
```

```
. sort foreign make
. predict dfor, dfbeta(1.foreign)
. list make price foreign dfor if abs(dfor) > 2/sqrt(74), divider
```

	make	price	foreign	dfor
12.	Cad. Eldorado	14,500	Domestic	-.5290519
13.	Cad. Seville	15,906	Domestic	.8243419
28.	Linc. Versailles	13,466	Domestic	-.5283729
42.	Plym. Arrow	4,647	Domestic	-.6622424
43.	Plym. Champ	4,425	Domestic	.2371104
64.	Peugeot 604	12,990	Foreign	.2552032
69.	Toyota Corona	5,719	Foreign	-.256431

The Cadillac Seville shifted the coefficient on `1.foreign` 0.82 standard deviations!

Now let us ask which observations have the greatest effect on the mpg coefficient:

```
. predict dmpg, dfbeta(mpg)
. list make price mpg dmpg if abs(dmpg) > 2/sqrt(74), divider
```

	make	price	mpg	dmpg
12.	Cad. Eldorado	14,500	14	-.5970351
13.	Cad. Seville	15,906	21	1.134269
28.	Linc. Versailles	13,466	14	-.6069287
42.	Plym. Arrow	4,647	28	-.8925859
43.	Plym. Champ	4,425	34	.3186909

Once again, we see the Cadillac Seville heading the list, indicating that our regression results may be dominated by this one car.

◀

### ► Example 9: DFBETAs influence measure; the `dfbeta` command

We can use `predict, dfbeta()` or the `dfbeta` command to generate the DFBETAs. `dfbeta` makes up names for the new variables automatically and, without arguments, generates the DFBETAs for all the variables in the regression:

```
. dfbeta
Generating DFBETA variables ...
  _dfbeta_1: DFBETA weight
  _dfbeta_2: DFBETA 1.foreign
  _dfbeta_3: DFBETA mpg
  _dfbeta_4: DFBETA 1.foreign#c.mpg
```

`dfbeta` created four new variables in our dataset: `_dfbeta_1`, containing the DFBETAs for `weight`; `_dfbeta_2`, containing the DFBETAs for `mpg`; and so on. Had we wanted only the DFBETAs for `mpg` and `weight`, we might have typed

```
. dfbeta mpg weight
Generating DFBETA variables ...
  _dfbeta_5: DFBETA weight
  _dfbeta_6: DFBETA mpg
```



In the example above, we typed `dfbeta mpg weight` instead of `dfbeta`; if we had typed `dfbeta` followed by `dfbeta mpg weight`, here is what would have happened:

```
. dfbeta
Generating DFBETA variables ...
  _dfbeta_7: DFBETA weight
  _dfbeta_8: DFBETA 1.foreign
  _dfbeta_9: DFBETA mpg
  _dfbeta_10: DFBETA 1.foreign#c.mpg
. dfbeta mpg weight
Generating DFBETA variables ...
  _dfbeta_11: DFBETA weight
  _dfbeta_12: DFBETA mpg
```

`dfbeta` would have made up different names for the new variables. `dfbeta` never replaces existing variables—it instead makes up a different name, so we need to pay attention to `dfbeta`'s output.

◀

## Tests for violation of assumptions

### Description for estat hettest

`estat hettest` performs three versions of the Breusch–Pagan (1979) and Cook–Weisberg (1983) test for heteroskedasticity. All three versions of this test present evidence against the null hypothesis that  $t = \mathbf{0}$  in  $\text{Var}(e) = \sigma^2 \exp(\mathbf{z}t)$ . In the `normal` version, performed by default, the null hypothesis also includes the assumption that the regression disturbances are independent-normal draws with variance  $\sigma^2$ . The normality assumption is dropped from the null hypothesis in the `iid` and `fstat` versions, which respectively produce the score and  $F$  tests discussed in *Methods and formulas*. If `varlist` is not specified, the fitted values are used for  $\mathbf{z}$ . If `varlist` or the `rhs` option is specified, the variables specified are used for  $\mathbf{z}$ .

### Menu for estat

Statistics > Postestimation

### Syntax for estat hettest

```
estat hettest [varlist] [, rhs [normal | iid | fstat] mtest [(spec)]]
```

`collect` is allowed with `estat hettest`; see [U] 11.1.10 Prefix commands.

### Options for estat hettest

`rhs` specifies that tests for heteroskedasticity be performed for the right-hand-side (explanatory) variables of the fitted regression model. The `rhs` option may be combined with a `varlist`.

`normal`, the default, causes `estat hettest` to compute the original Breusch–Pagan/Cook–Weisberg test, which assumes that the regression disturbances are normally distributed.

`iid` causes `estat hettest` to compute the  $N * R^2$  version of the score test that drops the normality assumption.

`fstat` causes `estat hettest` to compute the  $F$ -statistic version that drops the normality assumption.

`mtest` [*spec*] specifies that multiple testing be performed. The argument specifies how *p*-values are adjusted. The following specifications, *spec*, are supported:

<code>bonferroni</code>	Bonferroni's multiple testing adjustment
<code>holm</code>	Holm's multiple testing adjustment
<code>sidak</code>	Šidák's multiple testing adjustment
<code>noadjust</code>	no adjustment is made for multiple testing

`mtest` may be specified without an argument. This is equivalent to specifying `mtest(noadjust)`; that is, tests for the individual variables should be performed with unadjusted *p*-values. By default, `estat hestest` does not perform multiple testing. `mtest` may not be specified with `iid` or `fstat`.

## Description for estat imtest

`estat imtest` performs an information matrix test for the regression model and an orthogonal decomposition into tests for heteroskedasticity, skewness, and kurtosis due to [Cameron and Trivedi \(1990\)](#); White's test for homoskedasticity against unrestricted forms of heteroskedasticity ([1980](#)) is available as an option. White's test is usually similar to the first term of the Cameron–Trivedi decomposition.

## Menu for estat

Statistics > Postestimation

## Syntax for estat imtest

```
estat imtest [ , preserve white ]
```

`collect` is allowed with `estat imtest`; see [\[U\] 11.1.10 Prefix commands](#).

## Options for estat imtest

`preserve` specifies that the data in memory be preserved, all variables and cases that are not needed in the calculations be dropped, and at the conclusion the original data be restored. This option is costly for large datasets. However, because `estat imtest` has to perform an auxiliary regression on  $k(k+1)/2$  temporary variables, where  $k$  is the number of regressors, it may not be able to perform the test otherwise.

`white` specifies that White's original heteroskedasticity test also be performed.

## Description for estat ovtest

`estat ovtest` performs two versions of the [Ramsey \(1969\)](#) regression specification-error test (RESET) for omitted variables. This test amounts to fitting  $y = \mathbf{x}\mathbf{b} + \mathbf{z}\mathbf{t} + u$  and then testing  $\mathbf{t} = \mathbf{0}$ . If the `rhs` option is not specified, powers of the fitted values are used for  $\mathbf{z}$ . If `rhs` is specified, powers of the individual elements of  $\mathbf{x}$  are used.

## Menu for estat

Statistics > Postestimation

## Syntax for estat ovtest

```
estat ovtest [ , rhs ]
```

`collect` is allowed with `estat ovtest`; see [U] 11.1.10 Prefix commands.

## Option for estat ovtest

`rhs` specifies that powers of the right-hand-side (explanatory) variables be used in the test rather than powers of the fitted values.

## Description for estat szroeter

`estat szroeter` performs Szroeter's rank test for heteroskedasticity for each of the variables in *varlist* or for the explanatory variables of the regression if `rhs` is specified.

## Menu for estat

Statistics > Postestimation

## Syntax for estat szroeter

```
estat szroeter [varlist] [ , rhs mtest(spec) ]
```

Either *varlist* or `rhs` must be specified.

## Options for estat szroeter

`rhs` specifies that tests for heteroskedasticity be performed for the right-hand-side (explanatory) variables of the fitted regression model. The `rhs` option may be combined with a *varlist*.

`mtest(spec)` specifies that multiple testing be performed. The argument specifies how *p*-values are adjusted. The following specifications, *spec*, are supported:

<code><u>bonferroni</u></code>	Bonferroni's multiple testing adjustment
<code><u>holm</u></code>	Holm's multiple testing adjustment
<code><u>sidak</u></code>	Šidák's multiple testing adjustment
<code><u>noadjust</u></code>	no adjustment is made for multiple testing

`estat szroeter` always performs multiple testing. By default, it does not adjust the *p*-values.

## Remarks and examples for estat hettest, estat imtest, estat ovtest, and estat szroeter

We introduce some regression diagnostic commands that are designed to test for certain violations that `rvfplot` (see [R] [regress postestimation diagnostic plots](#)) less formally attempts to detect. `estat ovtest` provides Ramsey’s test for omitted variables—a pattern in the residuals. `estat hettest` provides a test for heteroskedasticity—the increasing or decreasing variation in the residuals with fitted values, with respect to the explanatory variables, or with respect to yet other variables. The score test implemented in `estat hettest` (Breusch and Pagan 1979; Cook and Weisberg 1983) performs a score test of the null hypothesis that  $b = 0$  against the alternative hypothesis of multiplicative heteroskedasticity. `estat szroeter` provides a rank test for heteroskedasticity, which is an alternative to the score test computed by `estat hettest`. Finally, `estat imtest` computes an information matrix test, including an orthogonal decomposition into tests for heteroskedasticity, skewness, and kurtosis (Cameron and Trivedi 1990). The heteroskedasticity test computed by `estat imtest` is similar to the general test for heteroskedasticity that was proposed by White (1980). Cameron and Trivedi (2022, chap. 3) discuss most of these tests and provides more examples.

### ► Example 10: estat ovtest, estat hettest, estat szroeter, and estat imtest

We use our model of price on weight and foreign##c.mpg.

```
. use https://www.stata-press.com/data/r18/auto, clear
(1978 automobile data)

. regress price weight foreign##c.mpg
(output omitted)

. estat ovtest

Ramsey RESET test for omitted variables
Omitted: Powers of fitted values of price
H0: Model has no omitted variables
F(3, 66) = 7.77
Prob > F = 0.0002

. estat hettest

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity
Assumption: Normal error terms
Variable: Fitted values of price
H0: Constant variance
      chi2(1) = 6.50
Prob > chi2 = 0.0108
```

Testing for heteroskedasticity in the right-hand-side variables is requested by specifying the `rhs` option. By specifying the `mtest(bonferroni)` option, we request that tests be conducted for each of the variables, with a Bonferroni adjustment for the  $p$ -values to accommodate our testing multiple hypotheses.

```
. estat hettest, rhs mtest(bonf)
Breusch-Pagan/Cook-Weisberg test for heteroskedasticity
Assumption: Normal error terms
H0: Constant variance
```

Variable	chi2	df	p
weight	15.24	1	0.0004*
foreign			
Foreign	6.15	1	0.0525*
mpg	9.04	1	0.0106*
foreign#			
c.mpg			
Foreign	6.02	1	0.0566*
Simultaneous	15.60	4	0.0036

\* Bonferroni-adjusted *p*-values

```
. estat szroeter, rhs mtest(holm)
Szroeter's test for homoskedasticity
H0: Variance constant
Ha: Variance monotonic in variables
```

Variable	chi2	df	p
weight	17.07	1	0.0001*
foreign			
Foreign	6.15	1	0.0131*
mpg	11.45	1	0.0021*
foreign#			
c.mpg			
Foreign	6.17	1	0.0260*

\* Holm-adjusted *p*-values

Finally, we request the information matrix test, which is a conditional moments test with second-, third-, and fourth-order moment conditions.

```
. estat imtest
Cameron & Trivedi's decomposition of IM-test
```

Source	chi2	df	p
Heteroskedasticity	18.86	10	0.0420
Skewness	11.69	4	0.0198
Kurtosis	2.33	1	0.1273
Total	32.87	15	0.0049

We find evidence for omitted variables, heteroskedasticity, and nonnormal skewness.

So, why bother with the various graphical commands when the tests seem so much easier to interpret? In part, it is a matter of taste: both are designed to uncover the same problem, and both are, in fact, going about it in similar ways. One is based on a formal calculation, whereas the other is based on personal judgment in evaluating a graph. On the other hand, the tests are seeking evidence of specific problems, whereas judgment is more general. The careful analyst will use both.

We performed the omitted-variable test first. Omitted variables are a more serious problem than heteroskedasticity or the violations of higher moment conditions tested by `estat imtest`. If this

were not a manual, having found evidence of omitted variables, we would never have run the `estat hettest`, `estat szroeter`, and `estat imtest` commands, at least not until we solved the omitted-variable problem.

◀

## □ Technical note

`estat ovtest` and `estat hettest` both perform two flavors of their respective tests. By default, `estat ovtest` looks for evidence of omitted variables by fitting the original model augmented by  $\hat{y}^2$ ,  $\hat{y}^3$ , and  $\hat{y}^4$ , which are the fitted values from the original model. Under the assumption of no misspecification, the coefficients on the powers of the fitted values will be zero. With the `rhs` option, `estat ovtest` instead augments the original model with powers (second through fourth) of the explanatory variables (except for dummy variables).

`estat hettest`, by default, looks for heteroskedasticity by modeling the variance as a function of the fitted values. If, however, we specify a variable or variables, the variance will be modeled as a function of the specified variables. In our example, if we had, a priori, some reason to suspect heteroskedasticity and that the heteroskedasticity is a function of a car's weight, then using a test that focuses on weight would be more powerful than the more general tests such as White's test or the first term in the Cameron–Trivedi decomposition test.

`estat hettest`, by default, computes the original Breusch–Pagan/Cook–Weisberg test, which includes the assumption of normally distributed errors. [Koenker \(1981\)](#) derived an  $N * R^2$  version of this test that drops the normality assumption. [Wooldridge \(2020, 270\)](#) gives an  $F$ -statistic version that does not require the normality assumption.

□

## Stored results for `estat hettest`, `estat imtest`, and `estat ovtest`

`estat hettest` stores the following results for the (multivariate) score test in `r()`:

Scalars

<code>r(chi2)</code>	$\chi^2$ test statistic
<code>r(df)</code>	#df for the asymptotic $\chi^2$ distribution under $H_0$
<code>r(p)</code>	$p$ -value

`estat hettest`, `fstat` stores results for the (multivariate) score test in `r()`:

Scalars

<code>r(F)</code>	test statistic
<code>r(df_m)</code>	#df of the test for the $F$ distribution under $H_0$
<code>r(df_r)</code>	#df of the residuals for the $F$ distribution under $H_0$
<code>r(p)</code>	$p$ -value

`estat hettest` (if `mtest` is specified) and `estat szroeter` store the following in `r()`:

Matrices

<code>r(mtest)</code>	a matrix of test results, with rows corresponding to the univariate tests
	<code>mtest[.,1]</code> $\chi^2$ test statistic
	<code>mtest[.,2]</code> #df
	<code>mtest[.,3]</code> unadjusted $p$ -value
	<code>mtest[.,4]</code> adjusted $p$ -value (if an <code>mtest()</code> adjustment method is specified)

Macros

<code>r(mtmethode)</code>	adjustment method for $p$ -value
---------------------------	----------------------------------

`estat imtest` stores the following in `r()`:

Scalars

<code>r(chi2_t)</code>	IM-test statistic ( $= r(chi2_h) + r(chi2_s) + r(chi2_k)$ )
<code>r(df_t)</code>	df for limiting $\chi^2$ distribution under $H_0$ ( $= r(df_h) + r(df_s) + r(df_k)$ )
<code>r(chi2_h)</code>	heteroskedasticity test statistic
<code>r(df_h)</code>	df for limiting $\chi^2$ distribution under $H_0$
<code>r(chi2_s)</code>	skewness test statistic
<code>r(df_s)</code>	df for limiting $\chi^2$ distribution under $H_0$
<code>r(chi2_k)</code>	kurtosis test statistic
<code>r(df_k)</code>	df for limiting $\chi^2$ distribution under $H_0$
<code>r(chi2_w)</code>	White's heteroskedasticity test (if <code>white</code> specified)
<code>r(df_w)</code>	df for limiting $\chi^2$ distribution under $H_0$

`estat ovtest` stores the following in `r()`:

Scalars

<code>r(p)</code>	two-sided $p$ -value
<code>r(F)</code>	$F$ statistic
<code>r(df)</code>	degrees of freedom
<code>r(df_r)</code>	residual degrees of freedom

## Variance inflation factors

### Description for `estat vif`

`estat vif` calculates the centered or uncentered variance inflation factors (VIFs) for the independent variables specified in a linear regression model.

### Menu for `estat`

Statistics > Postestimation

### Syntax for `estat vif`

```
estat vif [ , uncentered ]
```

### Option for `estat vif`

`uncentered` requests the computation of the uncentered variance inflation factors. The uncentered VIFs are often used to detect the collinearity of the regressors with the constant. `uncentered` must be specified if the regression model did not include a constant term because centered VIFs are not appropriate for these models.

### Remarks and examples for `estat vif`

Problems arise in regression when the predictors are highly correlated. In this situation, there may be a significant change in the regression coefficients if you add or delete an independent variable. The estimated standard errors of the fitted coefficients are inflated, or the estimated coefficients may not be statistically significant even though a statistical relation exists between the dependent and independent variables.

Data analysts rely on these facts to check informally for the presence of multicollinearity. `estat vif`, another command for use after `regress`, calculates the variance inflation factors and tolerances for each of the independent variables.

The output shows the variance inflation factors together with their reciprocals. Some analysts compare the reciprocals with a predetermined tolerance. In the comparison, if the reciprocal of the VIF is smaller than the tolerance, the associated predictor variable is removed from the regression model. However, most analysts rely on informal rules of thumb applied to the VIF; see [Chatterjee and Hadi \(2012\)](#). According to these rules, there is evidence of multicollinearity if

1. The largest VIF is greater than 10 (some choose a more conservative threshold value of 30).
2. The mean of all the VIFs is considerably larger than 1.

### ► Example 11: `estat vif`

We examine a regression model fit using the ubiquitous automobile dataset:

```
. use https://www.stata-press.com/data/r18/auto
(1978 automobile data)
```

```
. regress price mpg rep78 trunk headroom length turn displ gear_ratio
```

Source	SS	df	MS	Number of obs	=	69
Model	264102049	8	33012756.2	F(8, 60)	=	6.33
Residual	312694909	60	5211581.82	Prob > F	=	0.0000
				R-squared	=	0.4579
				Adj R-squared	=	0.3856
Total	576796959	68	8482308.22	Root MSE	=	2282.9

price	Coefficient	Std. err.	t	P> t	[95% conf. interval]
mpg	-144.84	82.12751	-1.76	0.083	-309.1195 19.43948
rep78	727.5783	337.6107	2.16	0.035	52.25638 1402.9
trunk	44.02061	108.141	0.41	0.685	-172.2935 260.3347
headroom	-807.0996	435.5802	-1.85	0.069	-1678.39 64.19062
length	-8.688914	34.89848	-0.25	0.804	-78.49626 61.11843
turn	-177.9064	137.3455	-1.30	0.200	-452.6383 96.82551
displacement	30.73146	7.576952	4.06	0.000	15.5753 45.88762
gear_ratio	1500.119	1110.959	1.35	0.182	-722.1303 3722.368
_cons	6691.976	7457.906	0.90	0.373	-8226.058 21610.01

```
. estat vif
```

Variable	VIF	1/VIF
length	8.22	0.121614
displacement	6.50	0.153860
turn	4.85	0.205997
gear_ratio	3.45	0.290068
mpg	3.03	0.330171
trunk	2.88	0.347444
headroom	1.80	0.554917
rep78	1.46	0.686147
Mean VIF	4.02	

The results are mixed. Although we have no VIFs greater than 10, the mean VIF is greater than 1, though not considerably so. We could continue the investigation of collinearity, but given that other authors advise that collinearity is a problem only when VIFs exist that are greater than 30 (contradicting our rule above), we will not do so here.



### ► Example 12: estat vif, with strong evidence of multicollinearity

This example comes from a dataset described in [Kutner, Nachtsheim, and Neter \(2004, 257\)](#) that examines body fat as modeled by caliper measurements on the triceps, midarm, and thigh.

```
. use https://www.stata-press.com/data/r18/bodyfat
(Body fat)
```

```
. regress bodyfat tricep thigh midarm
```

Source	SS	df	MS	Number of obs	=	20
Model	396.984607	3	132.328202	F(3, 16)	=	21.52
Residual	98.4049068	16	6.15030667	Prob > F	=	0.0000
				R-squared	=	0.8014
				Adj R-squared	=	0.7641
Total	495.389513	19	26.0731323	Root MSE	=	2.48

bodyfat	Coefficient	Std. err.	t	P> t	[95% conf. interval]
triceps	4.334085	3.015511	1.44	0.170	-2.058512 10.72668
thigh	-2.856842	2.582015	-1.11	0.285	-8.330468 2.616785
midarm	-2.186056	1.595499	-1.37	0.190	-5.568362 1.19625
_cons	117.0844	99.78238	1.17	0.258	-94.44474 328.6136

```
. estat vif
```

Variable	VIF	1/VIF
triceps	708.84	0.001411
thigh	564.34	0.001772
midarm	104.61	0.009560
Mean VIF	459.26	

Here we see strong evidence of multicollinearity in our model. More investigation reveals that the measurements on the thigh and the triceps are highly correlated:

```
. correlate triceps thigh midarm
(obs=20)
```

	triceps	thigh	midarm
triceps	1.0000		
thigh	0.9238	1.0000	
midarm	0.4578	0.0847	1.0000

If we remove the predictor `tricep` from the model (because it had the highest VIF), we get

```
. regress bodyfat thigh midarm
```

Source	SS	df	MS	Number of obs	=	20
Model	384.279748	2	192.139874	F(2, 17)	=	29.40
Residual	111.109765	17	6.53586854	Prob > F	=	0.0000
				R-squared	=	0.7757
				Adj R-squared	=	0.7493
Total	495.389513	19	26.0731323	Root MSE	=	2.5565

bodyfat	Coefficient	Std. err.	t	P> t	[95% conf. interval]
thigh	.8508818	.1124482	7.57	0.000	.6136367 1.088127
midarm	.0960295	.1613927	0.60	0.560	-.2444792 .4365383
_cons	-25.99696	6.99732	-3.72	0.002	-40.76001 -11.2339

```
. estat vif
```

Variable	VIF	1/VIF
midarm	1.01	0.992831
thigh	1.01	0.992831
Mean VIF	1.01	

Note how the coefficients change and how the estimated standard errors for each of the regression coefficients become much smaller. The calculated value of  $R^2$  for the overall regression for the subset model does not appreciably decline when we remove the correlated predictor. Removing an independent variable from the model is one way to deal with multicollinearity. Other methods include ridge regression, weighted least squares, and restricting the use of the fitted model to data that follow the same pattern of multicollinearity. In economic studies, it is sometimes possible to estimate the regression coefficients from different subsets of the data by using cross-section and time series.

◀

All examples above demonstrated the use of centered VIFs. As pointed out by [Belsley \(1991\)](#), the centered VIFs may fail to discover collinearity involving the constant term. One solution is to use the uncentered VIFs instead. According to the definition of the uncentered VIFs, the constant is viewed as a legitimate explanatory variable in a regression model, which allows one to obtain the VIF value for the constant term.

### ► Example 13: estat vif, with strong evidence of collinearity with the constant term

Consider the extreme example in which one of the regressors is highly correlated with the constant. We simulate the data and examine both centered and uncentered VIF diagnostics after fitted regression model as follows.

```
. use https://www.stata-press.com/data/r18/extreme_collin
. regress y one x z
```

Source	SS	df	MS	Number of obs	=	100
Model	223801.985	3	74600.6617	F(3, 96)	=	2710.27
Residual	2642.42124	96	27.5252213	Prob > F	=	0.0000
Total	226444.406	99	2287.31723	R-squared	=	0.9883
				Adj R-squared	=	0.9880
				Root MSE	=	5.2464

  

y	Coefficient	Std. err.	t	P> t	[95% conf. interval]
one	-3.278582	10.5621	-0.31	0.757	-24.24419 17.68702
x	2.038696	.0242673	84.01	0.000	1.990526 2.086866
z	4.863137	.2681036	18.14	0.000	4.330956 5.395319
_cons	9.760075	10.50935	0.93	0.355	-11.10082 30.62097

  

```
. estat vif
```

Variable	VIF	1/VIF
z	1.03	0.968488
x	1.03	0.971307
one	1.00	0.995425
Mean VIF	1.02	

```
. estat vif, uncentered
```

Variable	VIF	1/VIF
one	402.94	0.002482
_cons	401.26	0.002492
z	2.93	0.341609
x	1.13	0.888705
Mean VIF	202.06	

According to the values of the centered VIFs (1.03, 1.03, 1.00), no harmful collinearity is detected in the model. However, by the construction of these simulated data, we know that `one` is highly collinear with the constant term. As such, the large values of uncentered VIFs for `one` (402.94) and `_cons` (401.26) reveal high collinearity of the variable `one` with the constant term.

◀

## Measures of effect size

### Description for estat esize

`estat esize` calculates effect sizes for linear models after `regress` or `anova`. By default, `estat esize` reports  $\eta^2$  estimates (Kerlinger and Lee 2000), which are equivalent to  $R^2$  estimates. If the option `epsilon` is specified, `estat esize` reports  $\epsilon^2$  estimates (Grissom and Kim 2012). If the option `omega` is specified, `estat esize` reports  $\omega^2$  estimates (Grissom and Kim 2012). Both  $\epsilon^2$  and  $\omega^2$  are adjusted  $R^2$  estimates. Confidence intervals for  $\eta^2$  estimates are estimated by using the noncentral  $F$  distribution (Smithson 2001). See Kline (2013) or Thompson (2006) for further information.

### Menu for estat

Statistics > Postestimation

### Syntax for estat esize

```
estat esize [ , epsilon omega level(#) ]
```

`collect` is allowed with `estat esize`; see [U] 11.1.10 Prefix commands.

### Options for estat esize

`epsilon` specifies that the  $\epsilon^2$  estimates of effect size be reported. The default is  $\eta^2$  estimates.

`omega` specifies that the  $\omega^2$  estimates of effect size be reported. The default is  $\eta^2$  estimates.

`level`(#) specifies the confidence level, as a percentage, for confidence intervals. The default is `level(95)` or as set by `set level`; see [U] 20.8 Specifying the width of confidence intervals.

## Remarks and examples for estat esize

Whereas  $p$ -values are used to assess the statistical significance of a result, measures of effect size are used to assess the practical significance of a result. Effect sizes can be broadly categorized as “measures of group differences” (the  $d$  family) and “measures of association” (the  $r$  family); see Ellis (2010, table 1.1). The  $d$  family includes estimators such as Cohen’s  $D$ , Hedges’s  $G$ , and Glass’s  $\Delta$  (also see [R] esize). The  $r$  family includes estimators such as the point-biserial correlation coefficient,  $\eta^2$ ,  $\epsilon^2$ , and  $\omega^2$ . For an introduction to the concepts and calculation of effect sizes, see Kline (2013) or Thompson (2006). For a more detailed discussion, see Kirk (1996), Ellis (2010), Cumming (2012), Grissom and Kim (2012), and Kelley and Preacher (2012).

### ► Example 14: Calculating effect sizes for a linear regression model

Suppose we fit a linear regression model for low-birthweight infants.

```
. use https://www.stata-press.com/data/r18/lbw
(Hosmer & Lemeshow data)
. regress bwt smoke i.race
```

Source	SS	df	MS	Number of obs	=	189
Model	12346897.6	3	4115632.54	F(3, 185)	=	8.69
Residual	87568400.9	185	473342.708	Prob > F	=	0.0000
				R-squared	=	0.1236
				Adj R-squared	=	0.1094
Total	99915298.6	188	531464.354	Root MSE	=	688

  

bwt	Coefficient	Std. err.	t	P> t	[95% conf. interval]
smoke	-428.0254	109.0033	-3.93	0.000	-643.0746 -212.9761
race					
Black	-450.54	153.066	-2.94	0.004	-752.5194 -148.5607
Other	-454.1813	116.436	-3.90	0.000	-683.8944 -224.4683
_cons	3334.858	91.74301	36.35	0.000	3153.86 3515.855

We can use the `estat esize` command to calculate  $\eta^2$  for the entire model and a partial  $\eta^2$  for each term in the model.

```
. estat esize
Effect sizes for linear models
```

Source	Eta-squared	df	[95% conf. interval]
Model	.1235736	3	.0399862 .2041365
smoke	.0769345	1	.0193577 .1579213
race	.0908394	2	.0233037 .1700334

Note: Eta-squared values for individual model terms are partial.

The overall model effect size is 0.124. This means that roughly 12.4% of the variation in `bwt` is explained by the model. The partial effect size for `smoke` is 0.077. This means that roughly 7.7% of the variation in `bwt` is explained by `smoke` after you remove the variation explained by all other terms.

The omega option causes estat esize to report  $\omega^2$  and partial  $\omega^2$ .

```
. estat esize, omega
Effect sizes for linear models
```

Source	Omega-squared	df
Model	.1088457	3
smoke	.0715877	1
race	.0806144	2

Note: Omega-squared values for individual model terms are partial.



### ▷ Example 15: Calculating effect size for an ANOVA model

We can use estat esize after ANOVA models as well.

```
. anova bwt smoke race
```

	Number of obs =	189	R-squared =	0.1236	
	Root MSE =	687.999	Adj R-squared =	0.1094	
Source	Partial SS	df	MS	F	Prob>F
Model	12346898	3	4115632.5	8.69	0.0000
smoke	7298536.6	1	7298536.6	15.42	0.0001
race	8749453.3	2	4374726.6	9.24	0.0001
Residual	87568401	185	473342.71		
Total	99915299	188	531464.35		

```
. estat esize
Effect sizes for linear models
```

Source	Eta-squared	df	[95% conf. interval]	
Model	.1235736	3	.0399862	.2041365
smoke	.0769345	1	.0193577	.1579213
race	.0908394	2	.0233037	.1700334

Note: Eta-squared values for individual model terms are partial.



### □ Technical note

$\eta^2$  was developed in the context of analysis of variance. Thus, the published research on the calculation of confidence intervals focuses on cases where the numerator degrees of freedom are relatively small (for example,  $df < 20$ ).

Some combinations of the  $F$  statistic, numerator degrees of freedom, and denominator degrees of freedom yield confidence limits that do not contain the corresponding estimated value for an  $\eta^2$ . This problem is most commonly observed for larger numerator degrees of freedom.

Nothing in the literature suggests alternative methods for constructing confidence intervals in such cases; therefore, we recommend cautious interpretation of confidence intervals for  $\eta^2$  when the numerator degrees of freedom are greater than 20.



## Stored results for estat esize

`estat esize` stores the following results in `r()`:

### Scalars

`r(level)` confidence level

### Matrices

`r(eseize)` a matrix of effect sizes, confidence intervals, degrees of freedom, and  $F$  statistics with rows corresponding to each term in the model

<code>eseize[.,1]</code>	$\eta^2$
<code>eseize[.,2]</code>	lower confidence bound for $\eta^2$
<code>eseize[.,3]</code>	upper confidence bound for $\eta^2$
<code>eseize[.,4]</code>	$\varepsilon^2$
<code>eseize[.,5]</code>	$\omega^2$
<code>eseize[.,6]</code>	numerator degrees of freedom
<code>eseize[.,7]</code>	denominator degrees of freedom
<code>eseize[.,8]</code>	$F$ statistic

## Methods and formulas

See [Hamilton \(2013, chap. 7\)](#), [Kohler and Kreuter \(2012, sec. 9.3\)](#), or [Baum \(2006, chap. 5\)](#) for an overview of using Stata to perform regression diagnostics. See [Peracchi \(2001, chap. 8\)](#) for a mathematically rigorous discussion of diagnostics.

Methods and formulas are presented under the following headings:

*predict*  
*Special-interest postestimation commands*

## predict

Assume that you have already fit the regression model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$$

where  $\mathbf{X}$  is  $n \times k$ .

Denote the previously estimated coefficient vector by  $\mathbf{b}$  and its estimated variance matrix by  $\mathbf{V}$ . `predict` works by recalling various aspects of the model, such as  $\mathbf{b}$ , and combining that information with the data currently in memory. Let  $\mathbf{x}_j$  be the  $j$ th observation currently in memory, and let  $s^2$  be the mean squared error of the regression.

If the user specified weights in `regress`, then  $\mathbf{X}'\mathbf{X}$  in the following formulas is replaced by  $\mathbf{X}'\mathbf{D}\mathbf{X}$ , where  $\mathbf{D}$  is defined in [Weighted regression](#) under *Methods and formulas* in [\[R\] regress](#).

Let  $\mathbf{V} = s^2(\mathbf{X}'\mathbf{X})^{-1}$ . Let  $k$  be the number of independent variables including the intercept, if any, and let  $y_j$  be the observed value of the dependent variable.

The *predicted value* (`xb` option) is defined as  $\hat{y}_j = \mathbf{x}_j\mathbf{b}$ .

Let  $\ell_j$  represent a lower bound for an observation  $j$  and  $u_j$  represent an upper bound. The probability that  $y_j|\mathbf{x}_j$  would be observed in the interval  $(\ell_j, u_j)$ —the `pr( $\ell, u$ )` option—is

$$P(\ell_j, u_j) = \Pr(\ell_j < \mathbf{x}_j\mathbf{b} + e_j < u_j) = \Phi\left(\frac{u_j - \hat{y}_j}{s}\right) - \Phi\left(\frac{\ell_j - \hat{y}_j}{s}\right)$$

where for the `pr( $\ell, u$ )`, `e( $\ell, u$ )`, and `ystar( $\ell, u$ )` options,  $\ell_j$  and  $u_j$  can be anywhere in the range  $(-\infty, +\infty)$ .

The option `e(l, u)` computes the expected value of  $y_j|\mathbf{x}_j$  conditional on  $y_j|\mathbf{x}_j$  being in the interval  $(\ell_j, u_j)$ , that is, when  $y_j|\mathbf{x}_j$  is truncated. It can be expressed as

$$E(\ell_j, u_j) = E(\mathbf{x}_j\mathbf{b} + e_j \mid \ell_j < \mathbf{x}_j\mathbf{b} + e_j < u_j) = \hat{y}_j - s \frac{\phi\left(\frac{u_j - \hat{y}_j}{s}\right) - \phi\left(\frac{\ell_j - \hat{y}_j}{s}\right)}{\Phi\left(\frac{u_j - \hat{y}_j}{s}\right) - \Phi\left(\frac{\ell_j - \hat{y}_j}{s}\right)}$$

where  $\phi$  is the normal density and  $\Phi$  is the cumulative normal.

You can also compute `ystar(l, u)`—the expected value of  $y_j|\mathbf{x}_j$ , where  $y_j$  is assumed censored at  $\ell_j$  and  $u_j$ :

$$y_j^* = \begin{cases} \ell_j & \text{if } \mathbf{x}_j\mathbf{b} + e_j \leq \ell_j \\ \mathbf{x}_j\mathbf{b} + e_j & \text{if } \ell_j < \mathbf{x}_j\mathbf{b} + e_j < u_j \\ u_j & \text{if } \mathbf{x}_j\mathbf{b} + e_j \geq u_j \end{cases}$$

This computation can be expressed in several ways, but the most intuitive formulation involves a combination of the two statistics just defined:

$$y_j^* = P(-\infty, \ell_j)\ell_j + P(\ell_j, u_j)E(\ell_j, u_j) + P(u_j, +\infty)u_j$$

A diagonal element of the projection matrix (`hat`) or (`leverage`) is given by

$$h_j = \mathbf{x}_j(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_j'$$

The *standard error of the prediction* (the `stdp` option) is defined as  $s_{p_j} = \sqrt{\mathbf{x}_j\mathbf{V}\mathbf{x}_j'}$

and can also be written as  $s_{p_j} = s\sqrt{h_j}$ .

The *standard error of the forecast* (`stdf`) is defined as  $s_{f_j} = s\sqrt{1 + h_j}$ .

The *standard error of the residual* (`stdr`) is defined as  $s_{r_j} = s\sqrt{1 - h_j}$ .

The *residuals* (`residuals`) are defined as  $\hat{e}_j = y_j - \hat{y}_j$ .

The *standardized residuals* (`rstandard`) are defined as  $\hat{e}_{s_j} = \hat{e}_j/s_{r_j}$ .

The *Studentized residuals* (`rstudent`) are defined as

$$r_j = \frac{\hat{e}_j}{s_{(j)}\sqrt{1 - h_j}}$$

where  $s_{(j)}$  represents the root mean squared error with the  $j$ th observation removed, which is given by

$$s_{(j)}^2 = \frac{s^2(n - k)}{n - k - 1} - \frac{\hat{e}_j^2}{(n - k - 1)(1 - h_j)}$$

where  $n$  is the number of observations and  $k$  is the number of right-hand-side variables (including the constant).

Cook's  $D$  (`cooksD`) is given by

$$D_j = \frac{\hat{e}_{s_j}^2 (s_{p_j}/s_{r_j})^2}{k} = \frac{h_j \hat{e}_j^2}{ks^2(1 - h_j)^2}$$

DFITS (`dfits`) is given by

$$\text{DFITS}_j = r_j \sqrt{\frac{h_j}{1 - h_j}}$$

Welsch distance (`welsch`) is given by

$$W_j = \frac{r_j \sqrt{h_j(n-1)}}{1 - h_j}$$

COVRATIO (`covratio`) is given by

$$\text{COVRATIO}_j = \frac{1}{1 - h_j} \left( \frac{n - k - \hat{e}_{s_j}^2}{n - k - 1} \right)^k$$

The DFBETAS (`dfbeta`) for a particular regressor  $x_i$  are given by

$$\text{DFBETA}_j = \frac{r_j u_j}{\sqrt{U^2(1 - h_j)}}$$

where  $u_j$  are the residuals obtained from a regression of  $x_i$  on the remaining  $x$ 's and  $U^2 = \sum_j u_j^2$ .

## Special-interest postestimation commands

The omitted-variable test (Ramsey 1969) reported by `estat ovtest` fits the regression  $y_i = \mathbf{x}_i \mathbf{b} + \mathbf{z}_i \mathbf{t} + u_i$  and then performs a standard  $F$  test of  $\mathbf{t} = \mathbf{0}$ . The default test uses  $\mathbf{z}_i = (\hat{y}_i^2, \hat{y}_i^3, \hat{y}_i^4)$ . If `rhs` is specified,  $\mathbf{z}_i = (x_{1i}^2, x_{1i}^3, x_{1i}^4, x_{2i}^2, \dots, x_{mi}^4)$ . In either case, the variables are normalized to have minimum 0 and maximum 1 before powers are calculated.

The test for heteroskedasticity (Breusch and Pagan 1979; Cook and Weisberg 1983) models  $\text{Var}(e_i) = \sigma^2 \exp(\mathbf{z} \mathbf{t})$ , where  $\mathbf{z}$  is a variable list specified by the user, the list of right-hand-side variables, or the fitted values  $\mathbf{x} \hat{\boldsymbol{\beta}}$ . The test is of  $\mathbf{t} = \mathbf{0}$ . Mechanically, `estat hettest` fits the augmented regression  $\hat{e}_i^2 / \hat{\sigma}^2 = a + \mathbf{z}_i \mathbf{t} + v_i$ .

The original Breusch–Pagan/Cook–Weisberg version of the test assumes that the  $e_i$  are normally distributed under the null hypothesis which implies that the score test statistic  $S$  is equal to the model sum of squares from the augmented regression divided by 2. Under the null hypothesis,  $S$  has the  $\chi^2$  distribution with  $m$  degrees of freedom, where  $m$  is the number of columns of  $\mathbf{z}$ .

Koenker (1981) derived a score test of the null hypothesis that  $\mathbf{t} = \mathbf{0}$  under the assumption that the  $e_i$  are independent and identically distributed (i.i.d.). Koenker showed that  $S = N * R^2$  has a large-sample  $\chi^2$  distribution with  $m$  degrees of freedom, where  $N$  is the number of observations and  $R^2$  is from the augmented regression and  $m$  is the number of columns of  $\mathbf{z}$ . `estat hettest, iid` produces this version of the test.

Wooldridge (2020, 270) showed that an  $F$  test of  $\mathbf{t} = \mathbf{0}$  in the augmented regression can also be used under the assumption that the  $e_i$  are i.i.d. `estat hettest, fstat` produces this version of the test.

Szroeter's class of tests for homoskedasticity against the alternative that the residual variance increases in some variable  $x$  is defined in terms of

$$H = \frac{\sum_{i=1}^n h(x_i) e_i^2}{\sum_{i=1}^n e_i^2}$$



where  $h(x)$  is some weight function that increases in  $x$  (Szroeter 1978).  $H$  is a weighted average of the  $h(x)$ , with the squared residuals serving as weights. Under homoskedasticity,  $H$  should be approximately equal to the unweighted average of  $h(x)$ . Large values of  $H$  suggest that  $e_i^2$  tends to be large where  $h(x)$  is large; that is, the variance indeed increases in  $x$ , whereas small values of  $H$  suggest that the variance actually decreases in  $x$ . `estat szroeter` uses  $h(x_i) = \text{rank}(x_i \text{ in } x_1 \dots x_n)$ ; see Judge et al. [1985, 452] for details. `estat szroeter` displays a normalized version of  $H$ ,

$$Q = \sqrt{\frac{6n}{n^2 - 1}} H$$

which is approximately  $N(0, 1)$  distributed under the null (homoskedasticity).

`estat hetttest` and `estat szroeter` provide adjustments of  $p$ -values for multiple testing. The supported methods are described in [R] `test`.

`estat imtest` performs the information matrix test for the regression model, as well as an orthogonal decomposition into tests for heteroskedasticity  $\delta_1$ , nonnormal skewness  $\delta_2$ , and nonnormal kurtosis  $\delta_3$  (Cameron and Trivedi 1990; Long and Trivedi 1993). The decomposition is obtained via three auxiliary regressions. Let  $e$  be the regression residuals,  $\hat{\sigma}^2$  be the maximum likelihood estimate of  $\sigma^2$  in the regression,  $n$  be the number of observations,  $X$  be the set of  $k$  variables specified with `estat imtest`, and  $R_{\text{un}}^2$  be the uncentered  $R^2$  from a regression.  $\delta_1$  is obtained as  $nR_{\text{un}}^2$  from a regression of  $e^2 - \hat{\sigma}^2$  on the cross products of the variables in  $X$ .  $\delta_2$  is computed as  $nR_{\text{un}}^2$  from a regression of  $e^3 - 3\hat{\sigma}^2 e$  on  $X$ . Finally,  $\delta_3$  is obtained as  $nR_{\text{un}}^2$  from a regression of  $e^4 - 6\hat{\sigma}^2 e^2 - 3\hat{\sigma}^4$  on  $X$ .  $\delta_1$ ,  $\delta_2$ , and  $\delta_3$  are asymptotically  $\chi^2$  distributed with  $1/2k(k+1)$ ,  $K$ , and 1 degree of freedom. The information test statistic  $\delta = \delta_1 + \delta_2 + \delta_3$  is asymptotically  $\chi^2$  distributed with  $1/2k(k+3)$  degrees of freedom. White's test for heteroskedasticity is computed as  $nR^2$  from a regression of  $\hat{u}^2$  on  $X$  and the cross products of the variables in  $X$ . This test statistic is usually close to  $\delta_1$ .

`estat vif` calculates the centered variance inflation factor ( $\text{VIF}_c$ ) (Chatterjee and Hadi 2012, 248–251) for  $x_j$ , given by

$$\text{VIF}_c(x_j) = \frac{1}{1 - \hat{R}_j^2}$$

where  $\hat{R}_j^2$  is the square of the centered multiple correlation coefficient that results when  $x_j$  is regressed against all other explanatory variables, including the constant.

The uncentered variance inflation factor ( $\text{VIF}_{uc}$ ) (Belsley 1991, 28–29) for  $x_j$  is given by

$$\text{VIF}_{uc}(x_j) = \frac{1}{1 - \tilde{R}_j^2}$$

where  $\tilde{R}_j^2$  is the square of the uncentered multiple correlation coefficient that results when  $x_j$  is regressed against all other explanatory variables and a constant of 1. If the original regression model was fit without a constant, the constant would also be omitted from the regression of  $x_j$ .

The methods and formulas for `estat esize` are described in *Methods and formulas* of [R] `esize`.

## Acknowledgments

`estat ovtest` and `estat hetttest` are based on programs originally written by Richard Goldstein. `estat imtest`, `estat szroeter`, and the current version of `estat hetttest` were written by Jeroen Weesie of the Department of Sociology at Utrecht University, The Netherlands. `estat imtest` is based in part on code written by J. Scott Long of the Department of Sociology at Indiana University, coauthor of the Stata Press book *Regression Models for Categorical and Limited Dependent Variables*, and author of the Stata Press book *The Workflow of Data Analysis Using Stata*.

## References

- Adkins, L. C., and R. C. Hill. 2011. *Using Stata for Principles of Econometrics*. 4th ed. Hoboken, NJ: Wiley.
- Baum, C. F. 2006. *An Introduction to Modern Econometrics Using Stata*. College Station, TX: Stata Press.
- Belsley, D. A. 1991. *Conditional Diagnostics: Collinearity and Weak Data in Regression*. New York: Wiley.
- Belsley, D. A., E. Kuh, and R. E. Welsch. 1980. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. New York: Wiley.
- Bollen, K. A., and R. W. Jackman. 1990. Regression diagnostics: An expository treatment of outliers and influential cases. In *Modern Methods of Data Analysis*, ed. J. Fox and J. S. Long, 257–291. Newbury Park, CA: Sage.
- Breusch, T. S., and A. R. Pagan. 1979. A simple test for heteroscedasticity and random coefficient variation. *Econometrica* 47: 1287–1294. <https://doi.org/10.2307/1911963>.
- Cameron, A. C., and P. K. Trivedi. 1990. The information matrix test and its applied alternative hypotheses. Working paper 372, University of California–Davis, Institute of Governmental Affairs.
- . 2022. *Microeconometrics Using Stata*. 2nd ed. College Station, TX: Stata Press.
- Chatterjee, S., and A. S. Hadi. 1986. Influential observations, high leverage points, and outliers in linear regression. *Statistical Science* 1: 379–393. <https://doi.org/10.1214/ss/1177013622>.
- . 1988. *Sensitivity Analysis in Linear Regression*. New York: Wiley.
- . 2012. *Regression Analysis by Example*. 5th ed. New York: Wiley.
- Cook, R. D. 1977. Detection of influential observation in linear regression. *Technometrics* 19: 15–18. <https://doi.org/10.1080/00401706.1977.10489493>.
- Cook, R. D., and S. Weisberg. 1982. *Residuals and Influence in Regression*. New York: Chapman and Hall/CRC.
- . 1983. Diagnostics for heteroscedasticity in regression. *Biometrika* 70: 1–10. <https://doi.org/10.2307/2335938>.
- Cox, N. J. 2004. *Speaking Stata: Graphing model diagnostics*. *Stata Journal* 4: 449–475.
- Cumming, G. 2012. *Understanding the New Statistics: Effect Sizes, Confidence Intervals, and Meta-Analysis*. New York: Routledge.
- DeMaris, A. 2004. *Regression with Social Data: Modeling Continuous and Limited Response Variables*. Hoboken, NJ: Wiley.
- Ellis, P. D. 2010. *The Essential Guide to Effect Sizes: Statistical Power, Meta-Analysis, and the Interpretation of Research Results*. Cambridge: Cambridge University Press.
- Grissom, R. J., and J. J. Kim. 2012. *Effect Sizes for Research: Univariate and Multivariate Applications*. 2nd ed. New York: Routledge.
- Hamilton, L. C. 2013. *Statistics with Stata: Updated for Version 12*. 8th ed. Boston: Brooks/Cole.
- Hill, R. C., W. E. Griffiths, and G. C. Lim. 2018. *Principles of Econometrics*. 5th ed. Hoboken, NJ: Wiley.
- Hoaglin, D. C., and P. J. Kempthorne. 1986. Comment [on Chatterjee and Hadi 1986]. *Statistical Science* 1: 408–412. <https://doi.org/10.1214/ss/1177013627>.
- Hoaglin, D. C., and R. E. Welsch. 1978. The hat matrix in regression and ANOVA. *American Statistician* 32: 17–22. <https://doi.org/10.1080/00031305.1978.10479237>.
- Huber, C. 2013. Measures of effect size in Stata 13. *The Stata Blog: Not Elsewhere Classified*. <http://blog.stata.com/2013/09/05/measures-of-effect-size-in-stata-13/>.
- Judge, G. G., W. E. Griffiths, R. C. Hill, H. Lütkepohl, and T.-C. Lee. 1985. *The Theory and Practice of Econometrics*. 2nd ed. New York: Wiley.
- Kelley, K., and K. J. Preacher. 2012. On effect size. *Psychological Methods* 17: 137–152. <https://doi.org/10.1037/a0028086>.
- Kerlinger, F. N., and H. B. Lee. 2000. *Foundations of Behavioral Research*. 4th ed. Belmont, CA: Wadsworth.
- Kirk, R. E. 1996. Practical significance: A concept whose time has come. *Educational and Psychological Measurement* 56: 746–759. <https://doi.org/10.1177/0013164496056005002>.
- Kline, R. B. 2013. *Beyond Significance Testing: Statistics Reform in the Behavioral Sciences*. 2nd ed. Washington, DC: American Psychological Association.

- Koenker, R. 1981. A note on studentizing a test for heteroskedasticity. *Journal of Econometrics* 17: 107–112. [https://doi.org/10.1016/0304-4076\(81\)90062-2](https://doi.org/10.1016/0304-4076(81)90062-2).
- Kohler, U., and F. Kreuter. 2012. *Data Analysis Using Stata*. 3rd ed. College Station, TX: Stata Press.
- Kutner, M. H., C. J. Nachtsheim, and J. Neter. 2004. *Applied Linear Regression Models*. 4th ed. New York: McGraw–Hill/Irwin.
- Lindsey, C., and S. J. Sheather. 2010a. Optimal power transformation via inverse response plots. *Stata Journal* 10: 200–214.
- . 2010b. Model fit assessment via marginal model plots. *Stata Journal* 10: 215–225.
- Long, J. S., and P. K. Trivedi. 1993. Some specification tests for the linear regression model. *Sociological Methods and Research* 21: 161–204. Reprinted in *Testing Structural Equation Models*, ed. K. A. Bollen and J. S. Long, pp. 66–110. Newbury Park, CA: Sage.
- Paracchi, F. 2001. *Econometrics*. Chichester, UK: Wiley.
- Ramsey, J. B. 1969. Tests for specification errors in classical linear least-squares regression analysis. *Journal of the Royal Statistical Society, Series B* 31: 350–371. <https://doi.org/10.1111/j.2517-6161.1969.tb00796.x>.
- Ramsey, J. B., and P. Schmidt. 1976. Some further results on the use of OLS and BLUS residuals in specification error tests. *Journal of the American Statistical Association* 71: 389–390. <https://doi.org/10.1080/01621459.1976.10480355>.
- Rousseeuw, P. J., and A. M. Leroy. 1987. *Robust Regression and Outlier Detection*. New York: Wiley.
- Smithson, M. 2001. Correct confidence intervals for various regression effect sizes and parameters: The importance of noncentral distributions in computing intervals. *Educational and Psychological Measurement* 61: 605–632. <https://doi.org/10.1177/00131640121971392>.
- Studenmund, A. H. 2017. *Using Econometrics: A Practical Guide*. 7th ed. Boston: Pearson.
- Szroeter, J. 1978. A class of parametric tests for heteroscedasticity in linear econometric models. *Econometrica* 46: 1311–1327. <https://doi.org/10.2307/1913831>.
- Thompson, B. 2006. *Foundations of Behavioral Statistics: An Insight-Based Approach*. New York: Guilford Press.
- Velleman, P. F. 1986. Comment [on Chatterjee and Hadi 1986]. *Statistical Science* 1: 412–413. <https://doi.org/10.1214/ss/1177013628>.
- Velleman, P. F., and R. E. Welsch. 1981. Efficient computing of regression diagnostics. *American Statistician* 35: 234–242. <https://doi.org/10.2307/2683296>.
- Weisberg, S. 2014. *Applied Linear Regression*. 4th ed. Hoboken, NJ: Wiley.
- Welsch, R. E. 1982. Influence functions and regression diagnostics. In *Modern Data Analysis*, ed. R. L. Launer and A. F. Siegel, 149–169. New York: Academic Press.
- . 1986. Comment [on Chatterjee and Hadi 1986]. *Statistical Science* 1: 403–405. <https://doi.org/10.1214/ss/1177013625>.
- Welsch, R. E., and E. Kuh. 1977. Linear Regression Diagnostics. Technical Report 923–977, Massachusetts Institute of Technology, Cambridge, MA.
- White, H. L., Jr. 1980. A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica* 48: 817–838. <https://doi.org/10.2307/1912934>.
- Wooldridge, J. M. 2020. *Introductory Econometrics: A Modern Approach*. 7th ed. Boston: Cengage.

## Also see

- [R] **regress** — Linear regression
- [R] **regress postestimation diagnostic plots** — Postestimation plots for regress
- [R] **regress postestimation time series** — Postestimation tools for regress with time series
- [SP] **estat moran** — Moran’s test of residual correlation with nearby residuals
- [LASSO] **lassogof** — Goodness of fit after lasso for prediction
- [U] **20 Estimation and postestimation commands**

Stata, Stata Press, and Mata are registered trademarks of StataCorp LLC. Stata and Stata Press are registered trademarks with the World Intellectual Property Organization of the United Nations. Other brand and product names are registered trademarks or trademarks of their respective companies. Copyright © 1985–2023 StataCorp LLC, College Station, TX, USA. All rights reserved.

