

*Wooldridge, Introductory Econometrics, 4th ed.*

## **Chapter 8: Heteroskedasticity**

In laying out the standard regression model, we made the assumption of **homoskedasticity** of the regression error term: that its variance is assumed to be constant in the population, conditional on the explanatory variables. The assumption of homoskedasticity fails when the variance changes in different segments of the population: for instance, if the variance of the unobserved factors influencing individuals' saving increases with their level of income. In such a case, we say that the error process is **heteroskedastic**. This does not affect the optimality of ordinary least squares for the computation of point estimates—and the assumption of homoskedasticity did not underly our derivation of the OLS formulas. But if this assumption is not tenable, we may not be able to rely

on the interval estimates of the parameters—on their confidence intervals, and  $t$ -statistics derived from their estimated standard errors. Indeed, the Gauss-Markov theorem, proving the optimality of least squares among linear unbiased estimators of the regression equation, does not hold in the presence of heteroskedasticity. If the error variance is not constant, then OLS estimators are no longer BLUE.

How, then, should we proceed? The classical approach is to test for heteroskedasticity, and if it is evident, try to model it. We can derive modified least squares estimators (known as **weighted least squares**) which will regain some of the desirable properties enjoyed by OLS in a homoskedastic setting. But this approach is sometimes problematic, since there are many plausible ways in which the error variance may differ in segments of the population—depending on some of the explanatory variables

in our model, or perhaps on some variables that are not even in the model. We can use weighted least squares effectively if we can derive the correct weights, but may not be much better off if we cannot convince ourselves that our application of weighted least squares is valid.

Fortunately, fairly recent developments in econometric theory have made it possible to avoid these quandaries. Methods have been developed to adjust the estimated standard errors in an OLS context for **heteroskedasticity of unknown form**—to develop what are known as **robust** standard errors. Most statistical packages now support the calculation of these robust standard errors when a regression is estimated. If heteroskedasticity is a problem, the robust standard errors will differ from those calculated by OLS, and we should take the former as more appropriate. How can you compute these robust standard errors? In Stata,

one merely adds the option `,robust` to the `regress` command. The ANOVA F-table will be suppressed (as will the adjusted  $R^2$  measure), since neither is valid when robust standard errors are being computed, and the term “robust” will be displayed above the standard errors of the coefficients to remind you that robust errors are in use.

How are robust standard errors calculated? Consider a model with a single explanatory variable. The OLS estimator can be written as:

$$b_1 = \beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2}$$

This gives rise to an estimated variance of the slope parameter:

$$Var(b_1) = \frac{\sum (x_i - \bar{x})^2 \sigma_i^2}{\left(\sum (x_i - \bar{x})^2\right)^2} \quad (1)$$

This expression reduces to the standard expression from Chapter 2 if  $\sigma_i^2 = \sigma^2$  for all observations:

$$\text{Var}(b_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

But if  $\sigma_i^2 \neq \sigma^2$  this simplification cannot be performed on (1). How can we proceed? Halbert White showed (in a famous article in *Econometrica*, 1980) that the unknown error variance of the  $i^{\text{th}}$  observation,  $\sigma_i^2$ , can be consistently estimated by  $e_i^2$ —that is, by the square of the OLS residual from the original equation. This enables us to compute robust variances of the parameters—for instance, (1) can now be computed from OLS residuals, and its square root will be the robust standard error of  $b_1$ . This carries over to multiple regression; in the general case of  $k$  explanatory variables,

$$\text{Var}(b_j) = \frac{\sum r_{ij}^2 e_i^2}{\left(\sum (x_{ij} - \bar{x}_j)^2\right)^2} \quad (2)$$

where  $e_i^2$  is the square of the  $i^{\text{th}}$  OLS residual, and  $r_{ij}$  is the  $i^{\text{th}}$  residual from regressing variable  $j$  on all other explanatory variables. The square root of this quantity is the **heteroskedasticity-robust standard error**, or the “White” standard error, of the  $j^{\text{th}}$  estimated coefficient. It may be used to compute the **heteroskedasticity-robust  $t$ -statistic**, which then will be valid for tests of the coefficient even in the presence of heteroskedasticity of unknown form. Likewise,  $F$ -statistics, which would also be biased in the presence of heteroskedasticity, may be consistently computed from the regression in which the robust standard errors of the coefficients are available.

If we have this better mousetrap, why would we want to report OLS standard errors—which would be subject to bias, and thus unreliable, if there is a problem of heteroskedasticity? If

(and only if) the assumption of homoskedasticity is valid, the OLS standard errors are preferred, since they will have an exact  $t$ -distribution at any sample size. The application of robust standard errors is justified as the sample size becomes large. If we are working with a sample of modest size, and the assumption of homoskedasticity is tenable, we should rely on OLS standard errors. But since robust standard errors are very easily calculated in most statistical packages, it is a simple task to estimate both sets of standard errors for a particular equation, and consider whether inference based on the OLS standard errors is fragile. In large data sets, it has become increasingly common practice to report the robust standard errors.

## **Testing for heteroskedasticity**

We may want to demonstrate that the model we have estimated does not suffer from heteroskedasticity, and justify reliance on OLS and

OLS standard errors in this context. How might we evaluate whether homoskedasticity is a reasonable assumption? If we estimate the model via standard OLS, we may then base a test for heteroskedasticity on the OLS residuals. If the assumption of homoskedasticity, conditional on the explanatory variables, holds, it may be written as:

$$H_0 : \text{Var} (u|x_1, x_2, \dots, x_k) = \sigma^2$$

And a test of this null hypothesis can evaluate whether the variance of the error process appears to be independent of the explanatory variables. We cannot observe the variances of each observation, of course, but as above we can rely on the squared OLS residual,  $e_i^2$ , to be a consistent estimator of  $\sigma_i^2$ . One of the most common tests for heteroskedasticity is derived from this line of reasoning: the



**Breusch–Pagan** test. The BP test involves regressing the squares of the OLS residuals on a set of variables—such as the original explanatory variables—in an auxiliary regression:

$$e_i^2 = d_0 + d_1x_1 + d_2x_2 + \dots d_kx_k + v \quad (3)$$

If the magnitude of the squared residual—a consistent estimator of the error variance of that observation—is not related to any of the explanatory variables, then this regression will have no explanatory power: its  $R^2$  will be small, and its ANOVA  $F$ -statistic will indicate that it does not explain any meaningful fraction of the variation of  $e_i^2$  around its own mean. (Note that although the OLS residuals have mean zero, and are in fact uncorrelated by construction with each of the explanatory variables, that does not apply to their squares). The

Breusch–Pagan test can be conducted by either the ANOVA  $F$ –statistic from (3), or by a large-sample form known as the Lagrange multiplier statistic:  $LM = n \times R^2$  from the auxiliary regression. Under  $H_0$  of homoskedasticity,  $LM \sim \chi_k^2$ .

The Breusch–Pagan test can be computed with the `estat hetttest` command after `regress`.

```
regress price mpg weight length
estat hetttest
```

which would evaluate the residuals from the regression for heteroskedasticity, with respect to the original explanatory variables. The null hypothesis is that of homoskedasticity; if a small  $p$ –value is received, the null is rejected in favor of heteroskedasticity (that is, the auxiliary regression (which is not shown) had a meaningful amount of explanatory power). The routine displays the  $LM$  statistic and its  $p$ –value

versus the  $\chi_k^2$  distribution. If a rejection is received, one should rely on robust standard errors for the original regression. Although we have demonstrated the Breusch–Pagan test by employing the original explanatory variables, the test may be used with any set of variables—including those not in the regression, but suspected of being systematically related to the error variance, such as the size of a firm, or the wealth of an individual.

The Breusch-Pagan test is a special case of **White's general test for heteroskedasticity**. The sort of heteroskedasticity that will damage OLS standard errors is that which involves correlations between squared errors and explanatory variables. White's test takes the list of explanatory variables  $\{x_1, x_2, \dots, x_k\}$  and augments it with squares and cross products of each of these variables. The White test then runs an auxiliary regression of  $e_i^2$  on the

explanatory variables, their squares, and their cross products. Under the null hypothesis, none of these variables should have any explanatory power, if the error variances are not systematically varying. The White test is another *LM* test, of the  $n \times R^2$  form, but involves a much larger number of regressors in the auxiliary regression. In the example above, rather than just including `mpg weight length`, we would also include `mpg2, weight2, length2, mpg×weight, mpg×length, and weight×length`: 9 regressors in all, giving rise to a test statistic with a  $\chi^2_{(9)}$  distribution.

How can you perform White's test? Give the command `ssc install whitetst` (you only need do this once) and it will install this routine in Stata. The `whitetst` command will automatically generate these additional variables and perform the test after a `regress` command. Since Stata knows what explanatory variables

were used in the regression, you need not specify them; just give the command `whitetst` after `regress`. You may also use the `fitted` option to base the test on powers of the predicted values of the regression rather than the full list of regressors, squares and cross products.

## **Weighted least squares estimation**

As an alternative to using heteroskedasticity-robust standard errors, we could transform the regression equation if we had knowledge of the form taken by heteroskedasticity. For instance, if we had reason to believe that:

$$\text{Var}(u|x) = \sigma^2 h(x)$$

where  $h(x)$  is some function of the explanatory variables that could be made explicit (e.g.

$h(x) = \text{income}$ ), we could use that information to properly specify the correction for heteroskedasticity. What would this entail? Since in this case we are saying that  $\text{Var}(u|x) \propto \text{income}$ , then the standard deviation of  $u_i$ , conditional on  $\text{income}_i$ , is  $\sqrt{\text{income}_i}$ . Thus could be used to perform **weighted least squares**: a technique in which we transform the variables in the regression, and then run OLS on the transformed equation. For instance, if we were estimating a simple savings function from the dataset `saving.dta`, in which `sav` is regressed on `inc`, and believed that there might be heteroskedasticity of the form above, we would perform the following transformations:

```
gen sd=sqrt(inc)
gen wsav=sav/sd
gen kon=1/sd
gen winc=inc/sd
regress wsav kon winc,noc
```

Note that there is no constant term in the weighted least squares (WLS) equation, and that the coefficient on `winc` still has the same connotation: that of the marginal propensity to save. In this case, though, we might be thankful that Stata (and most modern packages) have a method for estimating WLS models by merely specifying the form of the weights:

```
regress sav inc [aw=1/inc]
```

In this case, the “aw” indicates that we are using “analytical weights”—Stata’s term for this sort of weighting—and the analytical weight is specified to be the inverse of the observation variance (not its standard error). If you run this regression, you will find that its coefficient estimates and their standard errors are identical to those of the transformed equation—with less hassle than the latter, in which the summary statistics (F-statistic,  $R^2$ , predicted

values, residuals, etc.) pertain to the transformed dependent variable ( $w_{sav}$ ) rather than the original variable.

The use of this sort of WLS estimation is less popular than it was before the invention of “White” standard errors; in theory, the transformation to homoskedastic errors will yield more attractive properties than even the use of “White” standard errors, conditional on our proper specification of the form of the heteroskedasticity. But of course we are not sure about that, and imprecise treatment of the errors may not be as attractive as the less informed technique of using the robust estimates.

One case in which we do know the form of the heteroskedasticity is that of *grouped data*, in which the data we are using has been aggregated from microdata into groups of different sizes. For instance, a dataset with 50



states' average values of income, family size, etc. calculated from a random sample of the U.S. population will have widely varying precision in those average values. The mean values for a small state will be computed from relatively few observations, whereas the counterpart values for a large state will be more precisely estimated. Since we know that the standard error of the mean is  $\sigma/\sqrt{n}$ , we recognize how this effect will influence the precision of the estimates. How, then, can we use this dataset of 50 observations while dealing with the known heteroskedasticity of the states' errors? This too is weighted least squares, where the weight on the individual state should be its population. This can be achieved in Stata by specifying "frequency weights"—a variable containing the number of observations from which each sample observation represents. If we had state-level data on saving, income and population, we might regress `saving income [fw=pop]` to achieve this weighting.

One additional observation regarding heteroskedasticity. We often see, in empirical studies, that an equation has been specified in some ratio form—for instance, with per capita dependent and independent variables for data on states or countries, or in terms of financial ratios for firm- or industry-level data. Although there may be no mention of heteroskedasticity in the study, it is very likely that these ratio forms have been chosen to limit the potential damage of heteroskedasticity in the estimated model. There can certainly be heteroskedasticity in a per-capita form regression on country-level data, but it is much less likely to be a problem than it would be if, say, the levels of GDP were used in that model. Likewise, scaling firms' values by total assets, or total revenues, or the number of employees will tend to mitigate the difficulties caused by extremes in scale between large corporations and corner stores. Such models should still be examined for their errors' behavior, but the popularity of the ratio form in these instances is an implicit consideration of potential heteroskedasticity.