

# Estimation and forecasting: OLS, IV, IV-GMM

Christopher F Baum

*Boston College and DIW Berlin*

IMF Institute, Spring 2011

# Linear regression

A key tool in multivariate statistical inference is *linear regression*, in which we specify the conditional mean of a response variable  $y$  as a linear function of  $k$  independent variables

$$E[y|x_1, x_2, \dots, x_k] = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k \quad (1)$$

The conditional mean of  $y$  is a function of  $x_1, x_2, \dots, x_k$  with fixed parameters  $\beta_1, \beta_2, \dots, \beta_k$ . Given values for these  $\beta$ s the linear regression model predicts the average value of  $y$  in the population for different values of  $x_1, x_2, \dots, x_k$ .

This population regression function specifies that a set of  $k$  regressors in  $X$  and the stochastic disturbance  $u$  are the determinants of the response variable (or regressand)  $y$ . The model is usually assumed to contain a constant term, so that  $x_1$  is understood to equal one for each observation. We may write the linear regression model in matrix form as

$$y = X\beta + u \quad (2)$$

where  $X = \{x_1, x_2, \dots, x_k\}$ , an  $N \times k$  matrix of sample values.

The key assumption in the linear regression model involves the relationship in the population between the regressors  $X$  and  $u$ . We may rewrite Equation (2) as

$$u = y - X\beta \quad (3)$$

We assume that

$$E(u | X) = 0 \quad (4)$$

i.e., that the  $u$  process has a *zero conditional mean*. This assumption states that the unobserved factors involved in the regression function are not related in any systematic manner to the observed factors. This approach to the regression model allows us to consider both non-stochastic and stochastic regressors in  $X$  without distinction; all that matters is that they satisfy the assumption of Equation (4).

We may use the zero conditional mean assumption (Equation (4)) to define a *method of moments* estimator of the regression function. Method of moments estimators are defined by *moment conditions* that are assumed to hold on the population moments. When we replace the unobservable population moments by their sample counterparts, we derive feasible estimators of the model's parameters. The zero conditional mean assumption gives rise to a set of  $k$  moment conditions, one for each  $x$ . In the population, each regressor  $x$  is assumed to be unrelated to  $u$ , or have zero covariance with  $u$ . We may then substitute calculated moments from our sample of data into the expression to derive a method of moments estimator for  $\beta$ :

$$\begin{aligned} X'u &= 0 \\ X'(y - X\beta) &= 0 \end{aligned} \tag{5}$$

Substituting calculated moments from our sample into the expression and replacing the unknown coefficients  $\beta$  with estimated values  $b$  in Equation (5) yields the *ordinary least squares* (OLS) estimator

$$\begin{aligned} X'y - X'Xb &= 0 \\ b &= (X'X)^{-1}X'y \end{aligned} \tag{6}$$

We may use  $b$  to calculate the regression residuals:

$$e = y - Xb \tag{7}$$

Given the solution for the vector  $b$ , the additional parameter of the regression problem  $\sigma_u^2$ , the population variance of the stochastic disturbance, may be estimated as a function of the regression residuals  $e_i$ :

$$s^2 = \frac{\sum_{i=1}^N e_i^2}{N - k} = \frac{e'e}{N - k} \quad (8)$$

where  $(N - k)$  are the residual *degrees of freedom* of the regression problem. The positive square root of  $s^2$  is often termed the standard error of regression, or standard error of estimate, or root mean square error. Stata uses the last terminology and displays  $s$  as `Root MSE`.

To learn more about the sampling distribution of the OLS estimator, we must make some additional assumptions about the distribution of the stochastic disturbance  $u_i$ . In classical statistics, the  $u_i$  were assumed to be independent draws from the same normal distribution. The modern approach to econometrics drops the normality assumptions and simply assumes that the  $u_i$  are independent draws from an identical distribution (*i.i.d.*).

The normality assumption was sufficient to derive the exact finite-sample distribution of the OLS estimator. In contrast, under the *i.i.d.* assumption, one must use large-sample theory to derive the sampling distribution of the OLS estimator. The sampling distribution of the OLS estimator can be shown to be approximately normal using large-sample theory.



Specifically, when the  $u_i$  are *i.i.d.* with finite variance  $\sigma_u^2$ , the OLS estimator  $b$  has a large-sample normal distribution with mean  $\beta$  and variance  $\sigma_u^2 Q^{-1}$ , where  $Q^{-1}$  is the variance-covariance matrix of  $X$  in the population. We refer this variance-covariance matrix of the estimator as a VCE.

Because it is unknown, we need a consistent estimator of the VCE. While neither  $\sigma_u^2$  nor  $Q^{-1}$  is actually known, we can use consistent estimators of them to construct a consistent estimator of  $\sigma_u^2 Q^{-1}$ . Given that  $s^2$  consistently estimates  $\sigma_u^2$  and  $1/N(X'X)$  consistently estimates  $Q$ ,  $s^2(X'X)^{-1}$  is a VCE of the OLS estimator.

Under the assumption of *i.i.d.* errors, the celebrated Gauss–Markov theorem holds. Within the class of linear, unbiased estimators the OLS estimator has the smallest sampling variance, or the greatest precision. In that sense, it is *best*, so that “ordinary least squares is BLUE” (the *best linear unbiased estimator*) for the parameters of the regression model. If we restrict our consideration to unbiased estimators which are linear in the parameters, we cannot find a more *efficient* estimator.

The property of *efficiency* refers to the precision of the estimator. If estimator *A* has a smaller sampling variance than estimator *B*, estimator *A* is said to be *relatively efficient*. The Gauss–Markov theorem states that OLS is relatively efficient versus all other linear, unbiased estimators of the model. We must recall, though, that this statement rests upon the maintained hypotheses of an appropriately specified model and an *i.i.d.* disturbance process with a zero conditional mean, as specified in Equation (4).

As an illustration, we present regression estimates from a simple macroeconomic model, constructed with US quarterly data from the latest edition of *International Financial Statistics*. The model, of the log of real investment expenditures, should not be taken seriously. Its purpose is only to illustrate the workings of regression in Stata. In the initial form of the model, we include as regressors the log of real GDP, the log of real wages, the 10-year Treasury yield and the S&P Industrials stock index.

We present the descriptive statistics with `summarize`, then proceed to fit a regression equation.

```
.
. use usmacr1, clear
. tsset
      time variable:  yq, 1959q1 to 2010q3
      delta: 1 quarter
. summarize lrgrossinv lrgdp lrwage tr10yr S_Pindex, sep(0)
```

Variable	Obs	Mean	Std. Dev.	Min	Max
lrgrossinv	207	7.146933	.4508421	6.31017	7.874346
lrgdp	207	8.794305	.4707929	7.904815	9.50028
lrwage	207	4.476886	.1054649	4.21887	4.619725
tr10yr	207	6.680628	2.58984	2.73667	14.8467
S_Pindex	207	37.81332	40.04274	4.25073	130.258

The `regress` command, like other Stata estimation commands, requires us to specify the response variable followed by a *varlist* of the explanatory variables.

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex
```

Source	SS	df	MS	Number of obs = 207		
Model	41.3479199	4	10.33698	F( 4, 202) = 3989.87		
Residual	.523342927	202	.002590807	Prob > F = 0.0000		
Total	41.8712628	206	.203258557	R-squared = 0.9875		
				Adj R-squared = 0.9873		
				Root MSE = .0509		

lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6540464	.0414524	15.78	0.000	.5723115	.7357813
lrwage	.7017158	.1562383	4.49	0.000	.3936485	1.009783
tr10yr	.0131358	.0022588	5.82	0.000	.008682	.0175896
S_Pindex	.0020351	.0002491	8.17	0.000	.001544	.0025261
_cons	-1.911161	.399555	-4.78	0.000	-2.698994	-1.123327

The header of the regression output describes the overall model fit, while the table presents the point estimates, their precision, and interval estimates.

The regression output for this model includes the analysis of variance (ANOVA) table in the upper left, where the two sources of variation are displayed as `Model` and `Residual`. The `SS` are the Sums of Squares, with the `Residual SS` corresponding to  $e'e$  and the `Total SS` to  $\tilde{y}'\tilde{y}$  in equation (10) below.

The next column of the table reports the `df`: the degrees of freedom associated with each sum of squares. The degrees of freedom for total `SS` are  $(N - 1)$ , since the total `SS` has been computed making use of one sample statistic,  $\bar{y}$ . The degrees of freedom for the model are  $(k - 1)$ , equal to the number of slopes (or explanatory variables): one fewer than the number of estimated coefficients due to the constant term.

As discussed above, the model  $SS$  refer to the ability of the four regressors to jointly explain a fraction of the variation of  $y$  about its mean (the total  $SS$ ). The residual degrees of freedom are  $(N - k)$ , indicating that  $(N - k)$  residuals may be freely determined and still satisfy the constraint posed by the first normal equation of least squares that the regression surface passes through the multivariate point of means  $(\bar{y}, \bar{X}_2, \dots, \bar{X}_k)$ :

$$\bar{y} = b_1 + b_2\bar{X}_2 + b_3\bar{X}_3 + \dots + b_k\bar{X}_k \quad (9)$$

In the presence of the constant term  $b_1$  the first normal equation implies that  $\bar{e} = \bar{y} - \sum_i \bar{X}_i b_i$  must be identically zero. It must be stressed that this is not an assumption. This is an algebraic implication of the least squares technique which guarantees that the sum of least squares residuals (and their mean) will be very close to zero.



The last column of the ANOVA table reports the  $MS$ , the Mean Squares due to regression and error, which are merely the  $SS$  divided by the  $df$ . The ratio of the `Model`  $MS$  to `Residual`  $MS$  is reported as the ANOVA  $F$ -statistic, with numerator and denominator degrees of freedom equal to the respective  $df$  values.

This ANOVA  $F$  statistic is a test of the null hypothesis that the slope coefficients in the model are jointly zero: that is, the null model of  $y_i = \mu + u_i$  is as successful in describing  $y$  as is the regression alternative. The `Prob > F` is the tail probability or  $p$ -value of the  $F$ -statistic. In this example we may reject the null hypothesis at any conventional level of significance.

We may also note that the `Root MSE` for the regression of 0.0509, which is in the units of the response variable  $y$ , is very small relative to the mean of that variable, 7.14.

The upper right section of `regress` output contains several *goodness of fit* statistics. These statistics measure the degree to which an estimated model can explain the variation of the response variable  $y$ .

Other things equal, we should prefer a model with a better fit to the data. With the principle of parsimony in mind, we also prefer a simpler model. The mechanics of regression imply that a model with a very large number of regressors can explain  $y$  arbitrarily well.

Given the least squares residuals, the most common measure of goodness of fit, regression  $R^2$ , may be calculated (given a constant term in the regression function) as

$$R^2 = 1 - \frac{e'e}{\tilde{y}'\tilde{y}} \quad (10)$$

where  $\tilde{y} = y - \bar{y}$ : the regressand with its sample mean removed. This emphasizes that the object of regression is not the explanation of  $y'y$ , the raw sum of squares of the response variable  $y$ . That would amount to explaining why  $Ey \neq 0$ , which is often not a very interesting question. Rather, the object is to explain the variations in the response variable. That variable may be always positive—such as the level of GDP—so that it is not sensible to investigate whether the average price might be zero.

With a constant term in the model, the least squares approach seeks to explain the largest possible fraction of the sample *variation* of  $y$  about its mean (and not the associated *variance*!) The null model to which the estimated model is being contrasted is  $y = \mu + u$  where  $\mu$  is the population mean of  $y$ .

In estimating a regression, we are trying to determine whether the information in the regressors  $X$  is useful. Is the conditional expectation  $E(y|X)$  more informative than the unconditional expectation  $Ey = \mu$ ? The null model above has an  $R^2 = 0$ , while virtually *any* set of regressors will explain some fraction of the variation of  $y$  around  $\bar{y}$ , the sample estimate of  $\mu$ .  $R^2$  is that fraction in the unit interval: the proportion of the variation in  $y$  about  $\bar{y}$  explained by  $X$ .

Below the ANOVA table and summary statistics, Stata reports the coefficient estimates for each of the  $b_j$  values, along with their estimated standard errors,  $t$ -statistics, and the associated  $p$ -values labeled  $P > |t|$ : that is, the tail probability for a two-tailed test on  $b_j$  corresponding to the hypothesis  $H_0 : b_j = 0$ .

In the last two columns, a confidence interval for the coefficient estimate is displayed, with limits defined by the current setting of `level`. The `level()` option on `regress` (or other estimation commands) may be used to specify a particular level. After performing the estimation (e.g., with the default 95% level) the regression results may be redisplayed with, for instance, `regress, level(90)`. The default `level` may be either changed for the session or changed permanently with `set level n [, permanently]`.

# beta coefficients

In other social science disciplines, linear regression results are often reported in terms of estimated *beta coefficients*. This terminology is somewhat confusing for economists given their common practice of writing the regression model in terms of  $\beta$ s.

The beta coefficient is defined as  $\partial y^* / \partial X_j^*$  where the starred quantities are z-transformed or standardized variables: for instance,  $y^* = (y_i - \bar{y}) / s_y$  where  $\bar{y}$  is the sample mean and  $s_y$  is the sample standard deviation of the response variable. Thus, the beta coefficient for the  $j^{th}$  regressor tells us how many standard deviations  $y$  would change given a one standard deviation change in  $X_j$ .

This is an attractive measure in disciplines where many empirical quantities are indices lacking a natural scale. You may then rank regressors by the magnitudes of their beta coefficients because the absolute magnitude of the beta coefficient for  $X_j$  is indicative of the strength of the effect of that variable. For the regression model above, we can merely redisplay the regression using the `beta` option:

```
. regress, beta
```

Source	SS	df	MS		
Model	41.3479199	4	10.33698	Number of obs =	207
Residual	.523342927	202	.002590807	F( 4, 202) =	3989.87
Total	41.8712628	206	.203258557	Prob > F =	0.0000
				R-squared =	0.9875
				Adj R-squared =	0.9873
				Root MSE =	.0509

lrgrossinv	Coef.	Std. Err.	t	P> t	Beta
lrgdp	.6540464	.0414524	15.78	0.000	.6829896
lrwage	.7017158	.1562383	4.49	0.000	.1641515
tr10yr	.0131358	.0022588	5.82	0.000	.075458
S_Pindex	.0020351	.0002491	8.17	0.000	.1807493
_cons	-1.911161	.399555	-4.78	0.000	.

The output indicates that `lrwage` has the largest beta coefficient, in absolute terms, followed by `lrgdp`. In economic and financial applications, where most regressors have a natural scale, it is more common to compute marginal effects such as elasticities or semi-elasticities. We will discuss the `margins` command, used for those computations.



# Regression without a constant term

Stata offers the option of estimating a regression equation without a constant term with the `noconstant` option, although in general it is recommended not to use this option. Such a model makes little sense if the mean of the response variable is nonzero and all regressors' coefficients are insignificant.

Estimating a constant term in a model that does not have one causes a small loss in the efficiency of the parameter estimates. In contrast, incorrectly omitting a constant term produces inconsistent estimates. The tradeoff should be clear: include a constant term, and let the data indicate whether its estimate can be distinguished from zero.

# Recovering estimation results

The `regress` command shares the features of all estimation (e-class) commands. Saved results from `regress` can be viewed by typing `ereturn list`. All Stata estimation commands save an estimated parameter vector as matrix `e(b)` and the estimated variance-covariance matrix of the parameters as matrix `e(V)`.

One item listed in the `ereturn list` should be noted: `e(sample)`, listed as a function rather than a scalar, macro or matrix. The `e(sample)` function returns 1 if an observation was included in the estimation sample and 0 otherwise.

The `regress` command honors any *if* and *in* qualifiers and then practices case-wise deletion to remove any observations with missing values across the set  $\{y, X\}$ . Thus, the observations actually used in generating the regression estimates may be fewer than those specified in the `regress` command. A subsequent command such as `summarize regressors if (or in)` will not necessarily provide the descriptive statistics of the observations on  $X$  that entered the regression unless all regressors and the  $y$  variable are in the *varlist*.

This is particularly relevant when building models with time series data, as the use of lags, leads and differences will cause observations to be omitted from the estimation sample.

The set of observations actually used in estimation can easily be determined with the qualifier `if e(sample)`:

```
summarize regressors if e(sample)
```

will yield the appropriate summary statistics from the regression sample. It may be retained for later use by placing it in a new variable:

```
generate byte reg1sample = e(sample)
```

where we use the `byte` data type to save memory since `e(sample)` is an indicator  $\{0,1\}$  variable.

# Hypothesis testing in regression

The application of regression methods is often motivated by the need to conduct tests of hypotheses which are implied by a specific theoretical model. In this section we discuss hypothesis tests and interval estimates assuming that the model is properly specified and that the errors are independently and identically distributed (*i.i.d.*). Estimators are random variables, and their sampling distributions depend on that of the error process.

There are three types of tests commonly employed in econometrics: *Wald* tests, *Lagrange multiplier* (LM) tests, and *likelihood ratio* (LR) tests. These tests share the same large-sample distribution, so that reliance on a particular form of test is usually a matter of convenience. Any hypothesis involving the coefficients of a regression equation can be expressed as one or more restrictions on the coefficient vector, reducing the dimensionality of the estimation problem. The Wald test involves estimating the unrestricted equation and evaluating the degree to which the restricted equation would differ in terms of its explanatory power.

The LM (or *score*) test involves estimating the restricted equation and evaluating the curvature of the objective function. These tests are often used to judge whether *i.i.d.* assumptions are satisfied.

The LR test involves comparing the objective function values of the unrestricted and restricted equations. It is often employed in maximum likelihood estimation.

Consider the general form of the Wald test statistic. Given the regression equation

$$y = X\beta + u \quad (11)$$

Any set of linear restrictions on the coefficient vector may be expressed as

$$R\beta = r \quad (12)$$

where  $R$  is a  $q \times k$  matrix and  $r$  is a  $q$ -element column vector, with  $q < k$ . The  $q$  restrictions on the coefficient vector  $\beta$  imply that  $(k - q)$  parameters are to be estimated in the restricted model. Each row of  $R$  imposes one restriction on the coefficient vector; a single restriction may involve multiple coefficients.

For instance, given the regression equation

$$y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + u \quad (13)$$

We might want to test the hypothesis  $H_0 : \beta_2 = 0$ . This single restriction on the coefficient vector implies  $R\beta = r$ , where

$$\begin{aligned} R &= (0 \ 1 \ 0 \ 0) \\ r &= (0) \end{aligned} \quad (14)$$

A test of  $H_0 : \beta_2 = \beta_3$  would imply the single restriction

$$\begin{aligned} R &= (0 \ 1 \ -1 \ 0) \\ r &= (0) \end{aligned} \quad (15)$$



Given a hypothesis expressed as  $H_0 : R\beta = r$ , we may construct the Wald statistic as

$$W = \frac{1}{s^2} (Rb - r)' [R(X'X)^{-1}R']^{-1} (Rb - r) \quad (16)$$

This quadratic form makes use of the vector of estimated coefficients,  $b$ , and evaluates the degree to which the restrictions fail to hold: the magnitude of the elements of the vector  $(Rb - r)$ . The Wald statistic evaluates the sums of squares of that vector, each weighted by a measure of their precision. Its denominator is  $s^2$ , the estimated variance of the error process, replacing the unknown parameter  $\sigma_u^2$ .

Stata contains a number of commands for the construction of hypothesis tests and confidence intervals which may be applied following an estimated regression. Some Stata commands report test statistics in the normal and  $\chi^2$  forms when the estimation commands are justified by large-sample theory. More commonly, the finite-sample  $t$  and  $F$  distributions are reported.

Stata's tests do not deliver verdicts with respect to the specified hypothesis, but rather present the *p-value* (or *prob-value*) of the test. Intuitively, the *p-value* is the probability of observing the estimated coefficient(s) if the null hypothesis is true.

In `regress` output, a number of test statistics and their  $p$ -values are automatically generated: that of the ANOVA  $F$  and the  $t$ -statistics for each coefficient, with the null hypothesis that the coefficients equal zero in the population. If we want to test additional hypotheses after a regression equation, three Stata commands are particularly useful: `test`, `testparm` and `lincom`. The `test` command may be specified as

`test coeflist`

where *coeflist* contains the names of one or more variables in the regression model.

A second syntax is

`test exp = exp`

where *exp* is an algebraic expression in the names of the regressors. The arguments of `test` may be repeated in parentheses in conducting joint tests. Additional syntaxes for `test` are available for multiple-equation models.

The `testparm` command provides similar functionality, but allows wildcards in the coefficient list:

```
testparm varlist
```

where the *varlist* may contain `*` or a hyphenated expression such as `ind1-ind9`.

The `lincom` command evaluates linear combinations of coefficients:

```
lincom exp
```

where *exp* is any linear combination of coefficients that is valid in the second syntax of `test`. For `lincom`, the *exp* must *not* contain an equal sign.

If we want to test the hypothesis  $H_0 : \beta_j = 0$ , the ratio of the estimated coefficient to its estimated standard error is distributed  $t$  under the null hypothesis that the population coefficient equals zero. That ratio is displayed by `regress` as the  $t$  column of the coefficient table.

Returning to our investment equation, a test statistic for the significance of a coefficient could be produced by using the commands:

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex
```

Source	SS	df	MS	Number of obs = 207		
Model	41.3479199	4	10.33698	F( 4, 202) = 3989.87		
Residual	.523342927	202	.002590807	Prob > F = 0.0000		
Total	41.8712628	206	.203258557	R-squared = 0.9875		
				Adj R-squared = 0.9873		
				Root MSE = .0509		

lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6540464	.0414524	15.78	0.000	.5723115	.7357813
lrwage	.7017158	.1562383	4.49	0.000	.3936485	1.009783
tr10yr	.0131358	.0022588	5.82	0.000	.008682	.0175896
S_Pindex	.0020351	.0002491	8.17	0.000	.001544	.0025261
_cons	-1.911161	.399555	-4.78	0.000	-2.698994	-1.123327

```
. test lrwage
```

```
( 1) lrwage = 0
```

```
F( 1, 202) = 20.17
```

```
Prob > F = 0.0000
```

In Stata's shorthand this is equivalent to the command `test _b[lrwage] = 0` (and much easier to type). If we use the `test` command, we note that the statistic is displayed as  $F(1, N-k)$  rather than in the  $t_{N-k}$  form of the coefficient table.

As many hypotheses to which `test` may be applied involve more than one restriction on the coefficient vector—and thus more than one degree of freedom—Stata routinely displays an  $F$ -statistic.

If we cannot reject the hypothesis  $H_0 : \beta_j = 0$ , and wish to restrict the equation accordingly, we remove that variable from the list of regressors.



More generally, we may to test the hypothesis  $\beta_j = \beta_j^0 = \theta$ , where  $\theta$  is any constant value. If theory suggests that the coefficient on variable `lrgdp` should be 0.75, then we may specify that hypothesis in test:

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex
```

Source	SS	df	MS	Number of obs = 207		
Model	41.3479199	4	10.33698	F( 4, 202) = 3989.87		
Residual	.523342927	202	.002590807	Prob > F = 0.0000		
Total	41.8712628	206	.203258557	R-squared = 0.9875		
				Adj R-squared = 0.9873		
				Root MSE = .0509		

lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6540464	.0414524	15.78	0.000	.5723115	.7357813
lrwage	.7017158	.1562383	4.49	0.000	.3936485	1.009783
tr10yr	.0131358	.0022588	5.82	0.000	.008682	.0175896
S_Pindex	.0020351	.0002491	8.17	0.000	.001544	.0025261
_cons	-1.911161	.399555	-4.78	0.000	-2.698994	-1.123327

```
. test lrgdp = 0.75
```

```
( 1) lrgdp = .75
```

```
F( 1, 202) = 5.36
```

```
Prob > F = 0.0216
```

The estimated coefficient of 0.65 is distinguished from 0.75.

We might want to compute a point and interval estimate for the sum of several coefficients. We may do that with the `lincom` (linear combination) command, which allows the specification of any linear expression in the coefficients. In the context of our investment equation, let us consider an arbitrary restriction: that the coefficients on `lrdgp`, `lrwage` and `tr10yr` sum to unity, so that we may write

$$H_0 : \beta_{lrgdp} + \beta_{lrwage} + \beta_{tr10yr} = 1 \quad (17)$$

It is important to note that although this hypothesis involves *three* estimated coefficients, it only involves *one* restriction on the coefficient vector. In this case, we have unitary coefficients on each term, but that need not be so.

```
. lincom lrgdp + lrwage + tr10yr
( 1)  lrgdp + lrwage + tr10yr = 0
```

lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
(1)	1.368898	.1196203	11.44	0.000	1.133033	1.604763

The sum of the three estimated coefficients is 1.369, with an interval estimate excluding unity. The hypothesis would be rejected by a `test` command.

We may use `test` to consider equality of two of the coefficients, or to test that their ratio equals a particular value:

```
. test lrgdp = lrwage
( 1)  lrgdp - lrwage = 0
      F( 1, 202) = 0.06
      Prob > F = 0.8061

. test tr10yr = 10 * S_Pindex
( 1)  tr10yr - 10*S_Pindex = 0
      F( 1, 202) = 9.24
      Prob > F = 0.0027
```

The hypothesis that the coefficients on `lrgdp` and `lrwage` are equal cannot be rejected at the 95% level, while the test that the ratio of the `tr10yr` and `S_Pindex` coefficients equals 10 may be rejected at the 99% level. Notice that Stata rewrites both expressions into a normalized form.

# Joint hypothesis tests

All of the tests illustrated above are presented as an  $F$ -statistic with one numerator degree of freedom since they only involve one restriction on the coefficient vector. In many cases, we wish to test an hypothesis involving multiple restrictions on the coefficient vector. Although the former test could be expressed as a  $t$ -test, the latter cannot. Multiple restrictions on the coefficient vector imply a *joint test*, the result of which is not simply a box score of individual tests.

A joint test is usually constructed in Stata by listing each hypothesis to be tested in parentheses on the `test` command. As presented above, the first syntax of the `test` command, `test coeflist`, performs the joint test that two or more coefficients are jointly zero, such as  $H_0 : \beta_2 = 0$  and  $\beta_3 = 0$ .

It is important to understand that this joint hypothesis is not at all the same as  $H'_0 : \beta_2 + \beta_3 = 0$ . The latter hypothesis will be satisfied by a locus of  $\{\beta_2, \beta_3\}$  values: all pairs that sum to zero. The former hypothesis will only be satisfied at the point where *each coefficient* equals zero. The joint hypothesis may be tested for our investment equation:

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex
```

Source	SS	df	MS	Number of obs = 207		
Model	41.3479199	4	10.33698	F( 4, 202) = 3989.87		
Residual	.523342927	202	.002590807	Prob > F = 0.0000		
Total	41.8712628	206	.203258557	R-squared = 0.9875		
				Adj R-squared = 0.9873		
				Root MSE = .0509		
lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6540464	.0414524	15.78	0.000	.5723115	.7357813
lrwage	.7017158	.1562383	4.49	0.000	.3936485	1.009783
tr10yr	.0131358	.0022588	5.82	0.000	.008682	.0175896
S_Pindex	.0020351	.0002491	8.17	0.000	.001544	.0025261
_cons	-1.911161	.399555	-4.78	0.000	-2.698994	-1.123327

```
. test tr10yr S_Pindex
```

```
( 1) tr10yr = 0
```

```
( 2) S_Pindex = 0
```

```
F( 2, 202) = 35.31
```

```
Prob > F = 0.0000
```

The data overwhelmingly reject the joint hypothesis that the model excluding `tr10yr` and `S_Pindex` is correctly specified relative to the full model.

# Tests of nonlinear hypotheses

What if the hypothesis tests to be conducted cannot be written in the linear form

$$H_0 : R\beta = r \quad (18)$$

for example, if theory predicts a certain value for the product of two coefficients in the model, or for an expression such as  $(\beta_2/\beta_3 + \beta_4)$ ? Two Stata commands are analogues to those we have used above: `testnl` and `nlcom`.

The former allows specification of nonlinear hypotheses on the  $\beta$  values, but unlike `test`, the syntax `_b[ varname]` must be used to refer to each coefficient value. If a joint test is to be conducted, the equations defining each nonlinear restriction must be written in parentheses, as illustrated below.



The `nlcom` command permits us to compute nonlinear combinations of the estimated coefficients in point and interval form, similar to `lincom`. Both commands employ the *delta method*, an approximation to the distribution of a nonlinear combination of random variables appropriate for large samples which constructs Wald-type tests. Unlike tests of linear hypotheses, nonlinear Wald-type tests based on the delta method are sensitive to the scale of the  $y$  and  $X$  data.

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex
```

Source	SS	df	MS	Number of obs = 207		
Model	41.3479199	4	10.33698	F( 4, 202) = 3989.87		
Residual	.523342927	202	.002590807	Prob > F = 0.0000		
Total	41.8712628	206	.203258557	R-squared = 0.9875		
				Adj R-squared = 0.9873		
				Root MSE = .0509		

lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6540464	.0414524	15.78	0.000	.5723115	.7357813
lrwage	.7017158	.1562383	4.49	0.000	.3936485	1.009783
tr10yr	.0131358	.0022588	5.82	0.000	.008682	.0175896
S_Pindex	.0020351	.0002491	8.17	0.000	.001544	.0025261
_cons	-1.911161	.399555	-4.78	0.000	-2.698994	-1.123327

```
. testnl _b[lrgdp] * _b[lrwage] = 0.33
```

```
(1) _b[lrgdp] * _b[lrwage] = 0.33
```

```
F(1, 202) = 2.77
```

```
Prob > F = 0.0978
```

In this example, we consider a restriction on the product of the coefficients of `lrgdp` and `lrwage`. The product of these coefficients cannot be distinguished from 0.33 at the 95% level.

We may also test a joint nonlinear hypothesis:

```
. testnl (_b[lrgdp] * _b[lrwage] = 0.33) ///  
>          (_b[lrwage] / _b[tr10yr] = 100 * _b[lrgdp])  
(1)  _b[lrgdp] * _b[lrwage] = 0.33  
(2)  _b[lrwage] / _b[tr10yr] = 100 * _b[lrgdp]  
      F(2, 202) =          29.83  
      Prob > F =          0.0000
```

The joint hypothesis may be rejected at the 99% level.

# Computing residuals and predicted values

After estimating a linear regression model with `regress` we may compute the regression residuals or the predicted values.

Computation of the residuals for each observation allows us to assess how well the model has done in explaining the value of the response variable for that observation. Is the in-sample prediction  $\hat{y}_i$  much larger or smaller than the actual value  $y_i$ ?

Computation of predicted values allows us to generate in-sample predictions: the values of the response variable generated by the estimated model. We may also want to generate out-of-sample predictions: that is, apply the estimated regression function to observations that were not used to generate the estimates. This may involve hypothetical values of the regressors or actual values. In the latter case, we may want to apply the estimated regression function to a separate sample (e.g., to a different time period than that used for estimation) to evaluate its applicability beyond the regression sample.

If a regression model is well specified, it should generate reasonable predictions for any sample from the population. If out-of-sample predictions are poor, the model's specification may be too specific to the original sample.

Neither the residuals nor predicted values are calculated by Stata's `regress` command, but either may be computed immediately thereafter with the `predict` command. This command is given as

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

where *choice* specifies the quantity to be computed for each observation.

For linear regression, `predict`'s default action is the computation of predicted values. These are known as the *point predictions*, and are specified by the choice `xb`. If the residuals are required, the command

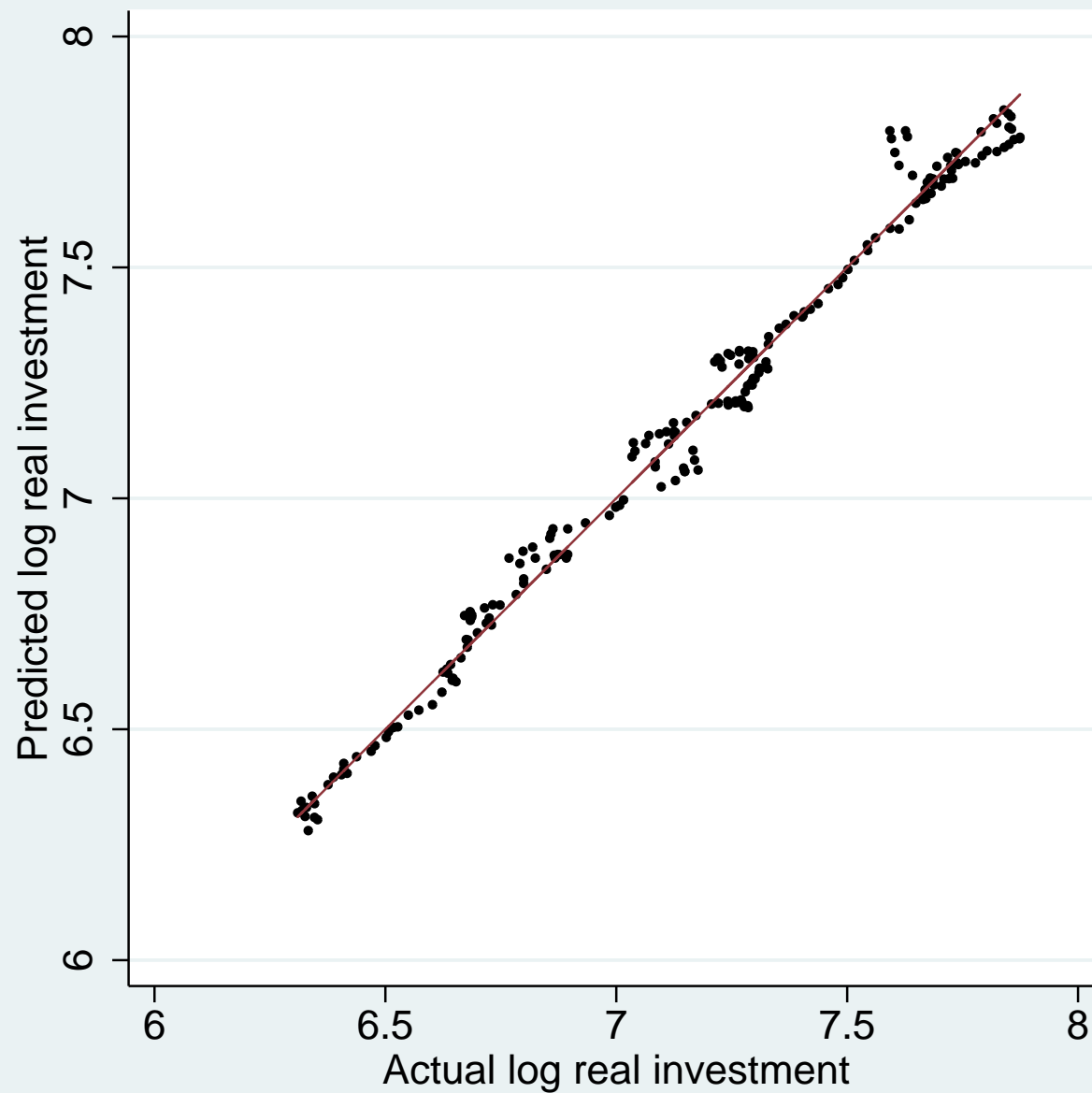
```
predict double lpriceeps, residual
```

should be used.

The regression estimates are only available to `predict` until another estimation command (e.g., `regress`) is issued. If these series are needed, they should be computed at the earliest opportunity. The use of `double` as the optional *type* in these commands ensures that the series will be generated with full numerical precision, and is strongly recommended.

We often would like to evaluate the quality of the regression fit in graphical terms. With a single regressor, a plot of actual and predicted values of  $y_i$  versus  $x_i$  will suffice. In multiple regression, the natural analogue is a plot of actual  $y_i$  versus the predicted  $\hat{y}_i$  values.

## Actual vs. predicted log real investment:

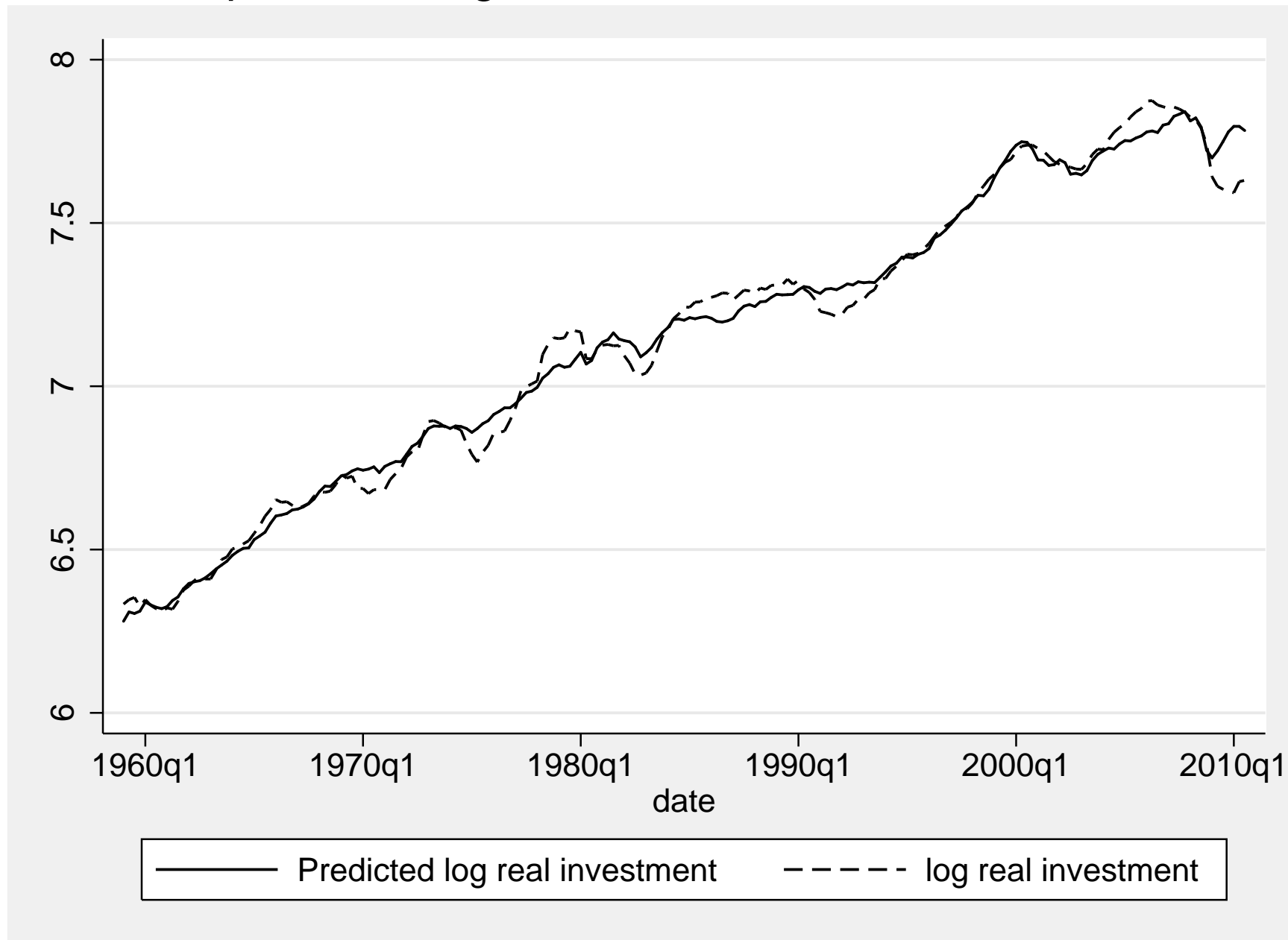




The aspect ratio has been constrained to unity so that points on the 45° line represent perfect predictions. Note that the model systematically overpredicts the log of relatively high levels of investment.

When using time series data, we may also want to examine the model's performance on a time series plot, using the `tsline` command. By using the graphics option `scheme(s2mono)` rather than the default `s2color`, we can get a graph which will reproduce well in black and white. If a graph is to be included in a document, use `graph export graphname.eps, replace`, which will be usable in high quality on any operating system. On Mac OS X systems (only), you can also export as PDF.

## Actual vs. predicted log real investment:



You might want to graph the predicted values versus NBER recession dates, using my `nbercycles` command from SSC:

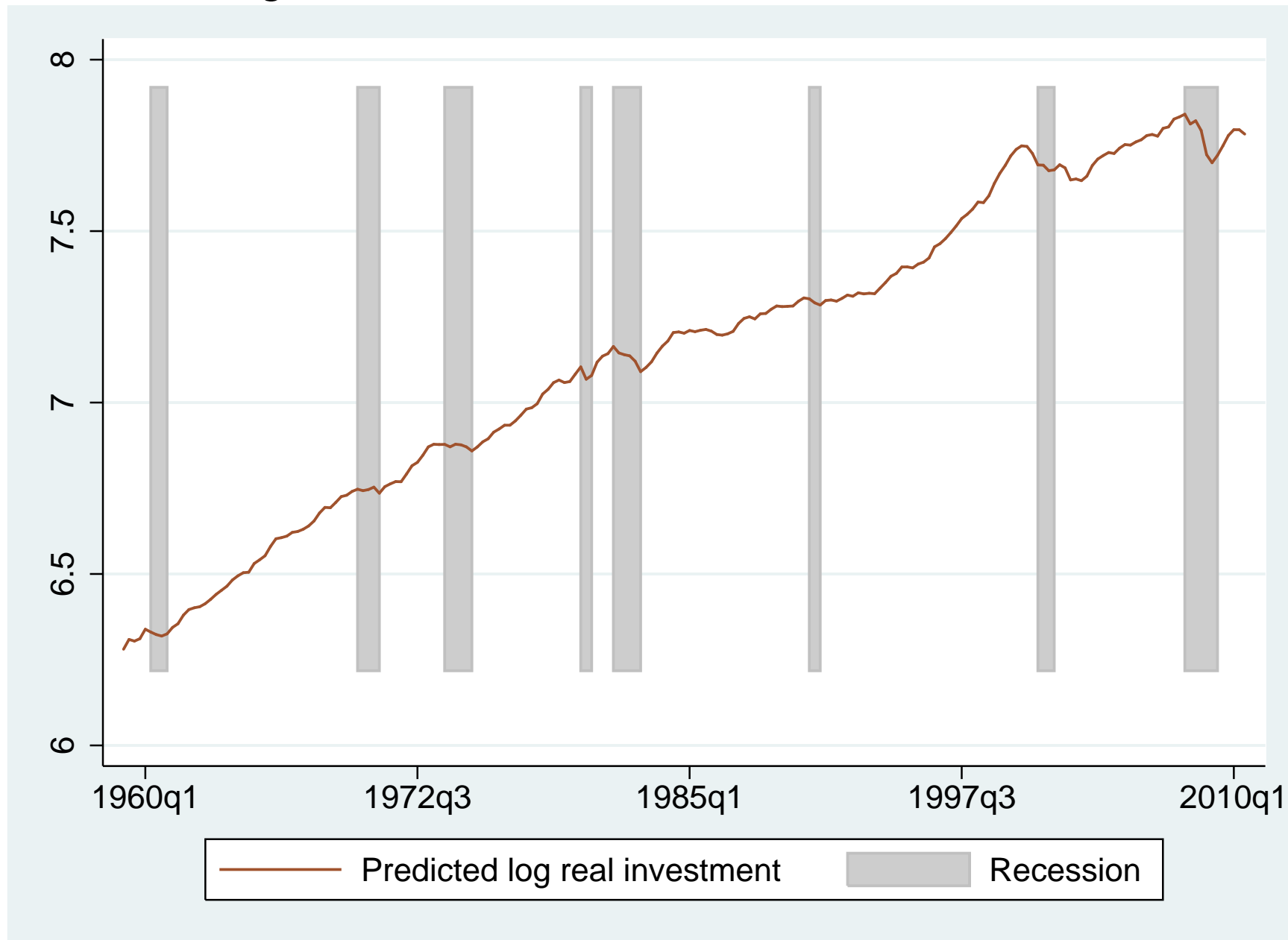
```
. // nbercycles from SSC
. nbercycles lrinvhat, file(invhatrecess.do)

Code to graph NBER recession dates written to invhatrecess.do

. * append your graph command to this file: e.g.
. * tsline timeseriesvar, xlabel(,format(%tq)) legend(order(9 1 "Recession"))
. twoway function y=7.91930087899802,range(1 4) recast(area) color(gs12) base(6
> .217729068487222) || ///
> function y=7.91930087899802,range(39 43) recast(area) color(gs12) base(6.2177
> 29068487222) || ///
> function y=7.91930087899802,range(55 60) recast(area) color(gs12) base(6.2177
> 29068487222) || ///
> function y=7.91930087899802,range(80 82) recast(area) color(gs12) base(6.2177
> 29068487222) || ///
> function y=7.91930087899802,range(86 91) recast(area) color(gs12) base(6.2177
> 29068487222) || ///
> function y=7.91930087899802,range(122 124) recast(area) color(gs12) base(6.21
> 7729068487222) || ///
> function y=7.91930087899802,range(164 167) recast(area) color(gs12) base(6.21
> 7729068487222) || ///
> function y=7.91930087899802,range(191 197) recast(area) color(gs12) base(6.21
> 7729068487222) || ///
> tsline lrinvhat , xlabel(,format(%tq)) legend(order(9 1 "Recession"))

.
end of do-file
```

## Predicted log real investment and NBER recession dates:



Like other Stata commands, `predict` will generate predictions for the entire sample. We may want to estimate a model over a subsample, and produce out-of-sample predictions, or *ex ante* forecasts. We may also want to produce interval estimates for forecasts, in- or out-of-sample. The latter may be done, after a regression, by specifying choice `stdp` for the standard error of prediction around the expected value of  $y|X_0$ .

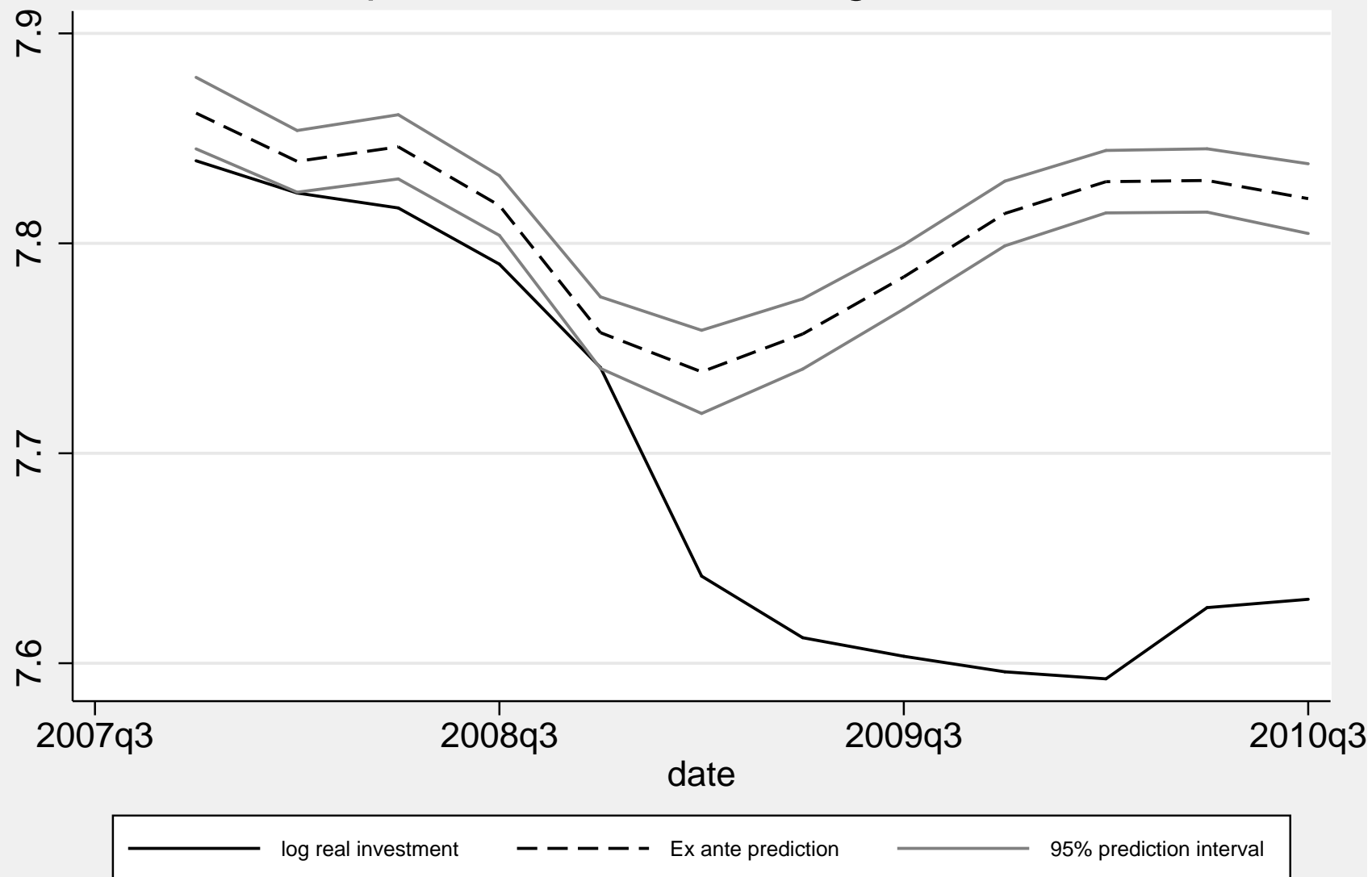
We illustrate by reestimating the investment model through 2007Q3, the calendar quarter preceding the most recent recession, and producing *ex ante* point and interval forecasts for the remaining periods. We juxtapose these point and interval estimates against the actual series during the recession and aftermath.

```
. regress lrgrossinv lrgdp lrwage tr10yr S_Pindex if tin(,2007q3)
```

Source	SS	df	MS	Number of obs = 195		
Model	37.640714	4	9.4101785	F( 4, 190) = 5512.25		
Residual	.324356548	190	.00170714	Prob > F = 0.0000		
Total	37.9650706	194	.19569624	R-squared = 0.9915		
				Adj R-squared = 0.9913		
				Root MSE = .04132		
lrgrossinv	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lrgdp	.6360608	.033894	18.77	0.000	.569204	.7029176
lrwage	.9161446	.1286431	7.12	0.000	.6623926	1.169897
tr10yr	.0074467	.0019506	3.82	0.000	.0035992	.0112942
S_Pindex	.0019152	.0002094	9.15	0.000	.0015021	.0023282
_cons	-2.663739	.3344459	-7.96	0.000	-3.323443	-2.004035

```
. predict double lrinvXA if tin(2007q4,), xb
(195 missing values generated)
. predict double lrinvSTDP if tin(2007q4,), stdp
(195 missing values generated)
. scalar tval = invttail(e(df_r), 0.025)
. generate double uplim = lrinvXA + tval * lrinvSTDP
(195 missing values generated)
. generate double lowlim = lrinvXA - tval * lrinvSTDP
(195 missing values generated)
. lab var uplim "95% prediction interval"
. lab var lowlim "95% prediction interval"
. lab var lrinvXA "Ex ante prediction"
. twoway (tsline lrgrossinv lrinvXA if tin(2007q4,)) ///
>       (rline uplim lowlim yq if tin(2007q4,)) ///
>       scheme(s2mono) legend(cols(3) size(vsmall)) ///
>       ti("Ex ante predicted vs. actual log real investment"))
```

## Ex ante predicted vs. actual log real investment





# Regression with non-i.i.d. errors

If the regression errors are independently and identically distributed (*i.i.d.*), OLS produces consistent point and interval estimates. Their sampling distribution in large samples is normal with a mean at the true coefficient values and their *VCE* is consistently estimated by the standard formula.

If the zero conditional mean assumption holds but the errors are not *i.i.d.*, OLS produces consistent estimates whose sampling distribution in large samples is still normal with a mean at the true coefficient values, but whose *VCE* cannot be consistently estimated by the standard formula.

We have two options when the errors are not *i.i.d.* First, we can use the consistent OLS point estimates with a different estimator of the *VCE* that accounts for non-*i.i.d.* errors. Alternatively, if we can specify how the errors deviate from *i.i.d.* in our regression model, we can model that process, using a different estimator that produces consistent and more efficient point estimates.

The tradeoff between these two methods is that of *robustness* versus *efficiency*. In a *robust* approach we place fewer restrictions on the estimator: the idea being that the consistent point estimates are good enough, although we must correct our estimator of their *VCE* to account for non-*i.i.d.* errors. In the *efficient* approach we incorporate an explicit specification of the non-*i.i.d.* distribution into the model. If this specification is appropriate, the additional restrictions which it implies will produce a more efficient estimator than that of the robust approach.

# Robust standard errors

We will only illustrate the robust approach. If the errors are conditionally heteroskedastic and we want to apply the robust approach, we use the Huber–White–sandwich estimator of the variance of the linear regression estimator, available in most Stata estimation commands as the `robust` option.

If the assumption of homoskedasticity is valid, the non-robust standard errors are more efficient than the robust standard errors. If we are working with a sample of modest size and the assumption of homoskedasticity is tenable, we should rely on non-robust standard errors. But since robust standard errors are very easily calculated in Stata, it is simple to estimate both sets of standard errors for a particular equation and consider whether inference based on the non-robust standard errors is fragile. In large data sets, it has become increasingly common practice to report robust (or Huber–White–sandwich) standard errors.

The alternate approach, generalized least squares (GLS), can be implemented for a model with heteroskedastic errors by specifying the form of the heteroskedasticity using Stata's weights. For this reason, GLS of this sort is often referred to as weighted least squares (WLS). To implement GLS (WLS), you must provide estimates of the error variance for each observation derived from some model of the heteroskedasticity process.

# The Newey–West estimator of the VCE

In an extension to Huber–White–sandwich robust standard errors, we may employ the *Newey–West* estimator that is appropriate in the presence of arbitrary heteroskedasticity and autocorrelation, thus known as the *HAC* estimator. Its use requires us to specify an additional parameter: the maximum order of any significant autocorrelation in the disturbance process, or the maximum lag  $L$ . One rule of thumb that has been used is to choose  $L = \sqrt[4]{N}$ . This estimator is available as the Stata command `newey`, which may be used as an alternative to `regress` for estimation of a regression with *HAC* standard errors.

Like the `robust` option, application of the *HAC* estimator does not modify the point estimates; it only affects the *VCE*. Test statistics based on the *HAC VCE* are robust to arbitrary heteroskedasticity and autocorrelation as well.

Similar to the case of pure heteroskedasticity, the GLS alternative to utilizing *HAC* standard errors is to explicitly model the nature of the serial correlation process. A common assumption is that the process is adequately represented by a first-order autoregression ( $AR(1)$ ). A regression model with  $AR(1)$  errors can be estimated by the Stata command `prais`, which implements the Prais–Winsten, Cochrane–Orcutt, Hildreth–Lu and maximum likelihood estimators. For higher-order autoregressive processes, the `arima` command may be used.

# Testing for heteroskedasticity

After estimating a regression model we may base a test for heteroskedasticity on the regression residuals. If the assumption of homoskedasticity conditional on the regressors holds, it can be expressed as:

$$H_0 : \text{Var}(u|X_2, X_3, \dots, X_k) = \sigma_u^2 \quad (19)$$

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 element of the disturbance process from samples of size one, but we can rely on the squared residual,  $e_i^2$ , to be a consistent estimator of  $\sigma_i^2$ . The logic behind any such test is that although the squared residuals will differ in magnitude across the sample, they should not be systematically related to *anything*, and a regression of squared residuals on any candidate  $Z_i$  should have no meaningful explanatory power.

One of the most common tests for heteroskedasticity is derived from this line of reasoning: the *Breusch–Pagan* test. The BP test, a Lagrange Multiplier (LM) test, involves regressing the squares of the regression residuals on a set of variables in an auxiliary regression

$$e_i^2 = d_1 + d_2 Z_{i2} + d_3 Z_{i3} + \dots d_\ell Z_{i\ell} + v_i \quad (20)$$

The Breusch–Pagan (Cook–Weisberg) test may be executed with `estat hetttest` after `regress`. If no regressor list (of  $Z$ s) is provided, `hetttest` employs the fitted values from the previous regression (the  $\hat{y}_i$  values). As mentioned above, the variables specified in the set of  $Z$ s could be chosen as measures which did not appear in the original regressor list.



We consider the potential scale-related heteroskedasticity in a cross-sectional model of median housing prices from the `hprice2a` dataset. The scale factor can be thought of as the average size of houses in each community, roughly measured by its number of rooms.

After estimating the model, we calculate three test statistics: that computed by `estat hettest` without arguments, which is the Breusch–Pagan test based on fitted values; `estat hettest` with a variable list, which uses those variables in the auxiliary regression; and White's general test statistic from `whitetst`, available from SSC.

```
. qui regress lprice rooms crime ldist
. hettest
Breusch-Pagan / Cook-Weisberg test for heteroskedasticity
    Ho: Constant variance
    Variables: fitted values of lprice
    chi2(1)          =    140.84
    Prob > chi2      =    0.0000

. hettest rooms crime ldist
Breusch-Pagan / Cook-Weisberg test for heteroskedasticity
    Ho: Constant variance
    Variables: rooms crime ldist
    chi2(3)          =    252.60
    Prob > chi2      =    0.0000

. whitetst
White's general test statistic :  144.0052  Chi-sq( 9)  P-value =  1.5e-26
```

Each of these tests indicates that there is a significant degree of heteroskedasticity related to scale in this model.

We illustrate the estimation of the model with OLS and robust standard errors.

```
. estimates table nonRobust Robust, b(%9.4f) se(%5.3f) t(%5.2f) ///
> title(Estimates of log housing price with OLS and Robust standard errors)
Estimates of log housing price with OLS and Robust standard errors
```

Variable	nonRobust	Robust
rooms	0.3072	0.3072
	0.018	0.026
	17.24	11.80
crime	-0.0174	-0.0174
	0.002	0.003
	-10.97	-6.42
ldist	0.0749	0.0749
	0.026	0.030
	2.93	2.52
_cons	7.9844	7.9844
	0.113	0.174
	70.78	45.76

legend: b/se/t

Note that the OLS standard errors are considerably smaller, biased downward, relative to the robust estimates.

# Testing for serial correlation

How might we test for the presence of serially correlated errors? Just as in the case of pure heteroskedasticity, we base tests of serial correlation on the regression residuals. In the simplest case, autocorrelated errors follow the so-called *AR(1)* model: an *autoregressive process* of order one, also known as a first-order Markov process:

$$u_t = \rho u_{t-1} + v_t, \quad |\rho| < 1 \quad (21)$$

where the  $v_t$  are uncorrelated random variables with mean zero and constant variance.

If we suspect that there might be autocorrelation in the disturbance process of our regression model, we could use the estimated residuals to diagnose it. The empirical counterpart to  $u_t$  in Equation (21) will be the  $e_t$  series produced by `predict`. We estimate the auxiliary regression of  $e_t$  on  $e_{t-1}$  without a constant term, as the residuals have mean zero.

The resulting slope estimate is a consistent estimator of the first-order autocorrelation coefficient  $\rho$  of the  $u$  process from Equation (21). Under the null hypothesis,  $\rho = 0$ , so that a rejection of this null hypothesis by this Lagrange Multiplier (*LM*) test indicates that the disturbance process exhibits *AR*(1) behavior.

A generalization of this procedure which supports testing for higher-order autoregressive disturbances is the Lagrange Multiplier (*LM*) test of Breusch and Godfrey. In this test, the regression residuals are regressed on the original  $X$  matrix augmented with  $p$  lagged residual series. The null hypothesis is that the errors are serially independent up to order  $p$ .

We illustrate the diagnosis of autocorrelation using a time series dataset `ukrates` of monthly short-term and long-term interest rates on UK government securities (Treasury bills and gilts), 1952m3–1995m12.

The model expresses the monthly change in the short rate  $r_s$ , the Bank of England's monetary policy instrument as a function of the prior month's change in the long-term rate  $r_{20}$ . The regressor and regressand are created on the fly by Stata's time series operators  $D$  and  $L$ . The model represents a monetary policy reaction function.

```
. regress D.rs LD.r20
```

Source	SS	df	MS	Number of obs = 524		
Model	13.8769739	1	13.8769739	F( 1, 522) = 52.88		
Residual	136.988471	522	.262430021	Prob > F = 0.0000		
Total	150.865445	523	.288461654	R-squared = 0.0920		
				Adj R-squared = 0.0902		
				Root MSE = .51228		
D.rs	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
r20						
LD.	.4882883	.0671484	7.27	0.000	.356374	.6202027
_cons	.0040183	.022384	0.18	0.858	-.0399555	.0479921

```
. predict double eps, residual
(2 missing values generated)
```

The Breusch–Godfrey test performed here considers the null of serial independence up to sixth order in the disturbance process, and that null is soundly rejected. We also present an unconditional test—the Ljung–Box Q test, available as command `wntestq`.

```
. estat bgodfrey, lags(6)
```

Breusch-Godfrey LM test for autocorrelation

lags( <i>p</i> )	chi2	df	Prob > chi2
6	17.237	6	0.0084

H0: no serial correlation

```
. wntestq eps
```

Portmanteau test for white noise

Portmanteau (Q) statistic =	82.3882
Prob > chi2(40) =	0.0001

Both tests decisively reject the null of no serial correlation.



Given this finding, we can generate heteroskedasticity- and autocorrelation-consistent (*HAC*) standard errors using the `newey` command, specifying 6 lags:

```
. newey D.rs LD.r20, lag(6)
```

Regression with Newey-West standard errors  
maximum lag: 6

Number of obs = 524  
F( 1, 522) = 35.74  
Prob > F = 0.0000

D.rs	Coef.	Newey-West Std. Err.	t	P> t	[95% Conf. Interval]	
r20						
LD.	.4882883	.0816725	5.98	0.000	.3278412	.6487354
_cons	.0040183	.0256542	0.16	0.876	-.0463799	.0544166

```
. estimates store NeweyWest
```

```
. estimates table nonHAC NeweyWest, b(%9.4f) se(%5.3f) t(%5.2f) ///
> title(Estimates of D.rs with OLS and Newey-West standard errors)
Estimates of D.rs with OLS and Newey-West standard errors
```

Variable	nonHAC	NeweyWest
r20		
LD.	0.4883	0.4883
	0.067	0.082
	7.27	5.98
_cons	0.0040	0.0040
	0.022	0.026
	0.18	0.16

legend: b/se/t

Note that the Newey–West standard errors are considerably larger than the OLS standard errors. OLS standard errors are biased downward in the presence of positive autocorrelation ( $\rho > 0$ ).

# Regression with indicator variables

Data come in three flavors: quantitative (or cardinal), ordinal (or ordered) and qualitative. Regression analysis handles quantitative data where both regressor and regressand may take on any real value. We also may work with *ordinal* or ordered data. They are distinguished from cardinal measurements in that an ordinal measure can only express inequality of two items, and not the magnitude of their difference.

We frequently encounter data that are purely *qualitative*, lacking any obvious ordering. If these data are coded as string variables, such as M and F for survey respondents' genders, we are not likely to mistake them for quantitative values. But in other cases, where a quality may be coded numerically, there is the potential to misuse this qualitative factor as quantitative.

In order to test the hypothesis that a qualitative factor has an effect on a response variable, we must convert the qualitative factor into a set of *indicator variables*, or dummy variables. We then conduct a *joint test* on their coefficients. If the hypothesis to be tested includes a single qualitative factor, the estimation problem may be described as a one-way analysis of variance, or *one-way ANOVA*. ANOVA models may be expressed as linear regressions on an appropriate set of indicator variables.

This notion of the equivalence of one-way ANOVA and linear regression on a set of indicator variables that correspond to a single qualitative factor generalizes to multiple qualitative factors.

If there are two qualitative factors (e.g., race and sex) that are hypothesized to affect income, a researcher would regress income on two appropriate sets of indicator variables, each representing one of the qualitative factors. This is then an example of *two-way ANOVA*.

# Using factor variables

One of the biggest innovations in Stata version 11 is the introduction of *factor variables*. Just as Stata's time series operators allow you to refer to lagged variables (`L.` or differenced variables (`D.`), the `i.` operator allows you to specify factor variables for any non-negative integer-valued variable in your dataset.

In the standard `auto` dataset, where `rep78` takes on values 1...5, you could `list rep78 i.rep78`, or `summarize i.rep78`, or `regress mpg i.rep78`. Each one of those commands produces the appropriate indicator variables 'on-the-fly': not as permanent variables in your dataset, but available for the command.

For the `list` command, the variables will be named `1b.rep78`, `2.rep78` ... `5.rep78`. The `b.` is the base level indicator, by default assigned to the smallest value. You can specify other base levels, such as the largest value, the most frequent value, or a particular value.

For the `summarize` command, only levels `2...5` will be shown; the base level is excluded from the list. Likewise, in a regression on `i.rep78`, the base level is the variable excluded from the regressor list to prevent perfect collinearity. The conditional mean of the excluded variable appears in the constant term.

# Interaction effects

If this was the only feature of factor variables (being instantiated when called for) they would not be very useful. The real advantage of these variables is the ability to define *interaction effects* for both integer-valued and continuous variables. For instance, consider the indicator `foreign` in the `auto` dataset. We may use a new operator, `#`, to define an interaction:

```
regress mpg i.rep78 i.foreign i.rep78#i.foreign
```

All combinations of the two categorical variables will be defined, and included in the regression as appropriate (omitting base levels and cells with no observations).

In fact, we can specify this model more simply: rather than

```
regress mpg i.rep78 i.foreign i.rep78#i.foreign
```

we can use the *factorial interaction* operator, ##:

```
regress mpg i.rep78##i.foreign
```

which will provide exactly the same regression, producing all first-level and second-level interactions. Interactions are not limited to pairs of variables; up to eight factor variables may be included.



Furthermore, factor variables may be interacted with continuous variables to produce analysis of covariance models. The continuous variables are signalled by the new `c .` operator:

```
regress mpg i.foreign i.foreign#c.displacement
```

which essentially estimates two regression lines: one for domestic cars, one for foreign cars. Again, the factorial operator could be used to estimate the same model:

```
regress mpg i.foreign##c.displacement
```

As we will see in discussing marginal effects, it is very advantageous to use this syntax to describe interactions, both among categorical variables and between categorical variables and continuous variables. Indeed, it is likewise useful to use the same syntax to describe squared (and cubed...) terms:

```
regress mpg i.foreign c.displacement c.displacement#c.displacement
```

In this model, we allow for an intercept shift for `foreign`, but constrain the slopes to be equal across foreign and domestic cars. However, by using this syntax, we may ask Stata to calculate the marginal effect  $\partial \text{mpg} / \partial \text{displacement}$ , taking account of the squared term as well, as Stata understands the mathematics of the specification in this explicit form.

# Computing marginal effects

With the introduction of factor variables in Stata 11, a powerful new command has been added: `margins`, which supersedes earlier versions' `mf` and `adjust` commands. Those commands remain available, but the new command has many advantages. Like those commands, `margins` is used after an estimation command.

In the simplest case, `margins` applied after a simple one-way ANOVA estimated with `regress i.rep78`, with `margins i.rep78`, merely displays the conditional means for each category of `rep78`.

```
. regress mpg i.rep78
```

Source	SS	df	MS	Number of obs = 69		
Model	549.415777	4	137.353944	F( 4, 64) = 4.91		
Residual	1790.78712	64	27.9810488	Prob > F = 0.0016		
Total	2340.2029	68	34.4147485	R-squared = 0.2348		
				Adj R-squared = 0.1869		
				Root MSE = 5.2897		
mpg	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
rep78						
2	-1.875	4.181884	-0.45	0.655	-10.22927	6.479274
3	-1.566667	3.863059	-0.41	0.686	-9.284014	6.150681
4	.6666667	3.942718	0.17	0.866	-7.209818	8.543152
5	6.363636	4.066234	1.56	0.123	-1.759599	14.48687
_cons	21	3.740391	5.61	0.000	13.52771	28.47229

```
. margins i.rep78
```

Adjusted predictions

Number of obs = 69

Model VCE : OLS

Expression : Linear prediction, predict()

	Delta-method					
	Margin	Std. Err.	z	P> z	[95% Conf. Interval]	
rep78						
1	21	3.740391	5.61	0.000	13.66897	28.33103
2	19.125	1.870195	10.23	0.000	15.45948	22.79052
3	19.43333	.9657648	20.12	0.000	17.54047	21.3262
4	21.66667	1.246797	17.38	0.000	19.22299	24.11034
5	27.36364	1.594908	17.16	0.000	24.23767	30.4896

We now estimate a model including both displacement and its square:

```
. regress mpg i.foreign c.displacement c.displacement#c.displacement
```

Source	SS	df	MS	Number of obs = 74		
Model	1416.01205	3	472.004018	F( 3, 70) = 32.16		
Residual	1027.44741	70	14.6778201	Prob > F = 0.0000		
Total	2443.45946	73	33.4720474	R-squared = 0.5795		
				Adj R-squared = 0.5615		
				Root MSE = 3.8312		

mpg	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
1.foreign	-2.88953	1.361911	-2.12	0.037	-5.605776	-.1732833
displacement	-.1482539	.0286111	-5.18	0.000	-.2053169	-.0911908
c. displacement# c. displacement	.0002116	.0000583	3.63	0.001	.0000953	.0003279
_cons	41.40935	3.307231	12.52	0.000	34.81328	48.00541

margins can then properly evaluate the regression function for domestic and foreign cars at selected levels of displacement:

```
. margins i.foreign, at(displacement=(100 300))
```

Adjusted predictions                      Number of obs       =            74

Model VCE : OLS

Expression : Linear prediction, `predict()`

```
1._at      : displacement      =      100
```

2. at : displacement = 300

		Delta-method				
		Margin	Std. Err.	z	P> z	[95% Conf. Interval]
_at#foreign						
	1 0	28.69991	1.216418	23.59	0.000	26.31578 31.08405
	1 1	25.81038	.8317634	31.03	0.000	24.18016 27.44061
	2 0	15.97674	.7014015	22.78	0.000	14.60201 17.35146
	2 1	13.08721	1.624284	8.06	0.000	9.903668 16.27074

In earlier versions of Stata, calculation of marginal effects in this model required some programming due to the nonlinear term displacement. Using `margins, dydx`, that is now simple. Furthermore, and most importantly, the default behavior of `margins` is to calculate average marginal effects (AMEs) rather than marginal effects at the average (MAE) or at some other point in the space of the regressors. In Stata 10, the user-written command `margeff` (Tamas Bartus, on the SSC Archive) was required to compute AMEs.

Current econometric practice favors the use of AMEs: the computation of each observation's marginal effect with respect to an explanatory factor, averaged over the estimation sample, to the computation of MAEs (which reflect an average individual: e.g. a family with 2.3 children).



We illustrate by computing average marginal effects (AMEs) for the prior regression:

```
. margins, dydx(foreign displacement)
Average marginal effects          Number of obs   =          74
Model VCE      : OLS
Expression    : Linear prediction, predict()
dy/dx w.r.t.  : 1.foreign displacement
```

	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
1.foreign displacement	-2.88953	1.361911	-2.12	0.034	-5.558827	-.2202327
	-.0647596	.007902	-8.20	0.000	-.0802473	-.049272

Note: dy/dx for factor levels is the discrete change from the base level.

## Alternatively, we may compute elasticities or semi-elasticities:

```
. margins, eyex(displacement) at(displacement=(100(100)400))
Average marginal effects                                Number of obs   =           74
Model VCE      : OLS
Expression     : Linear prediction, predict()
ey/ex w.r.t.   : displacement
1._at         : displacement      =           100
2._at         : displacement      =           200
3._at         : displacement      =           300
4._at         : displacement      =           400
```

	Delta-method					
	ey/ex	Std. Err.	z	P> z	[95% Conf. Interval]	
displacement						
_at						
1	-.3813974	.0537804	-7.09	0.000	-.486805	-.2759898
2	-.6603459	.0952119	-6.94	0.000	-.8469578	-.473734
3	-.4261477	.193751	-2.20	0.028	-.8058926	-.0464028
4	.5613844	.4817784	1.17	0.244	-.3828839	1.505653

Consider a model where we specify a factorial interaction between categorical and continuous covariates:

```
regress mpg i.foreign i.rep78##c.displacement
```

In this specification, each level of `rep78` has its own intercept and slope, whereas `foreign` only shifts the intercept term.

We may compute elasticities or semi-elasticities with the `over` option of `margins` for all combinations of `foreign` and `rep78`:

```
. margins, eyex(displacement) over(foreign rep78)
```

Average marginal effects	Number of obs	=	69
Model VCE	: OLS		
Expression	: Linear prediction, predict()		
ey/ex w.r.t.	: displacement		
over	: foreign rep78		

Model VCE : OLS

ey/ex w.r.t. : displacement

```
over          : foreign rep78
```

		Delta-method				
		ey/ex	Std. Err.	z	P> z	[95% Conf. Interval]
displacement						
foreign#						
rep78						
0	1	-.7171875	.5342	-1.34	0.179	-1.7642 .3298253
0	2	-.5953046	.219885	-2.71	0.007	-1.026271 -.1643379
0	3	-.4620597	.0999242	-4.62	0.000	-.6579077 -.2662118
0	4	-.6327362	.1647866	-3.84	0.000	-.955712 -.3097604
0	5	-.8726071	.0983042	-8.88	0.000	-1.06528 -.6799345
1	3	-.128192	.0228214	-5.62	0.000	-.1729213 -.0834628
1	4	-.1851193	.0380458	-4.87	0.000	-.2596876 -.110551
1	5	-1.689962	.3125979	-5.41	0.000	-2.302642 -1.077281

The `margins` command has many other capabilities which we will not discuss here. The lengthy reference manual article on `margins` is a useful reference.

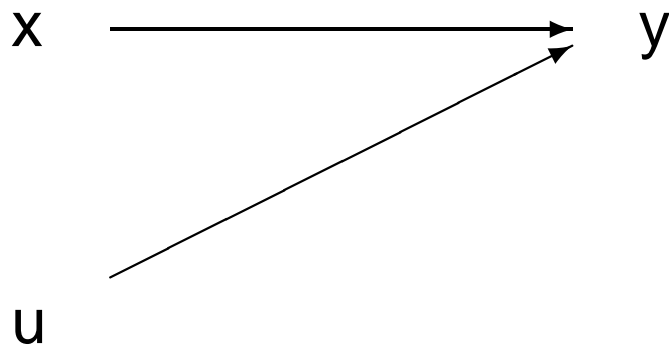
# Regression with instrumental variables

What are instrumental variables (IV) methods? Most widely known as a solution to *endogenous regressors*: explanatory variables correlated with the regression error term, IV methods provide a way to nonetheless obtain consistent parameter estimates.

First let us consider a path diagram illustrating the problem addressed by IV methods. We can use ordinary least squares (OLS) regression to consistently estimate a model of the following sort.

**Standard regression:**  $y = xb + u$

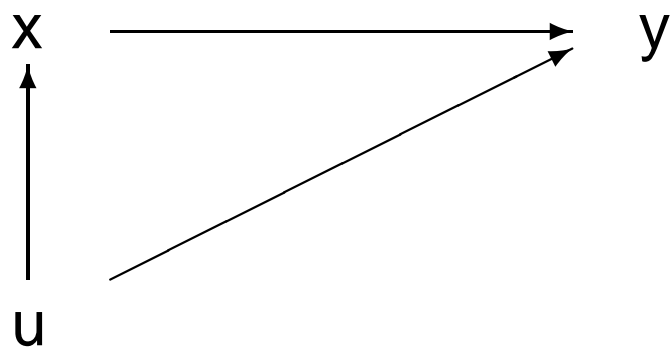
*no association between  $x$  and  $u$ ; OLS consistent*



However, OLS regression breaks down in the following circumstance:

**Endogeneity:**  $y = xb + u$

*correlation between  $x$  and  $u$ ; OLS inconsistent*



The correlation between  $x$  and  $u$  (or the failure of the zero conditional mean assumption  $E[u|x] = 0$ ) can be caused by any of several factors.



We have stated the problem as that of *endogeneity*: the notion that two or more variables are jointly determined in the behavioral model. This arises naturally in the context of a *simultaneous equations model* such as a supply-demand system in economics, in which price and quantity are jointly determined in the market for that good or service.

A shock or disturbance to either supply or demand will affect both the equilibrium price and quantity in the market, so that by construction both variables are correlated with any shock to the system. OLS methods will yield inconsistent estimates of any regression including both price and quantity, however specified.

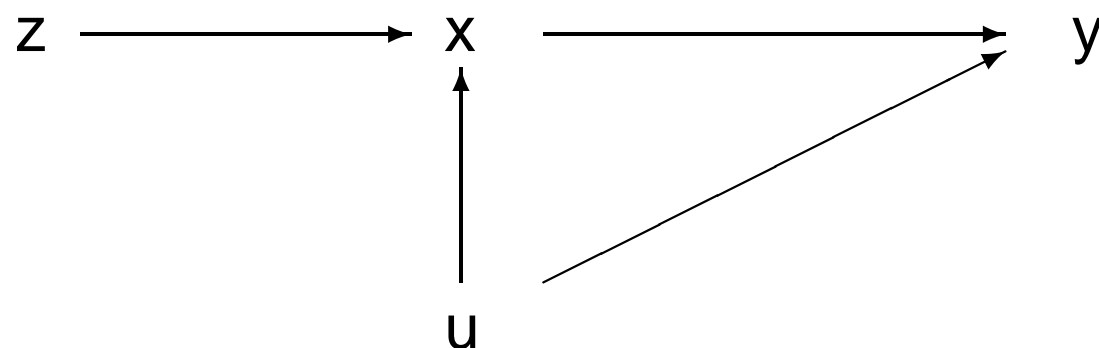
In a macroeconomic context, many of the behavioral equations that we might specify for consumption, investment, money demand, and so on are likely to contain endogenous regressors. In a consumption function, a shock to consumption or saving will also affect the level of GDP, and thus disposable income.

In this context, the zero conditional mean assumption cannot hold, even in terms of weak exogeneity of the regressors. OLS is no longer an appropriate estimation method, and we must rely upon other estimators to produce consistent estimates.

The solution provided by IV methods may be viewed as:

**Instrumental variables regression:**  $y = xb + u$

*z uncorrelated with  $u$ , correlated with  $x$*



The additional variable  $z$  is termed an *instrument* for  $x$ . In general, we may have many variables in  $x$ , and more than one  $x$  correlated with  $u$ . In that case, we shall need at least that many variables in  $z$ .

To deal with the problem of *endogeneity* in a supply-demand system, a candidate  $z$  will affect (e.g.) the quantity supplied of the good, but not directly impact the demand for the good. An example for an agricultural commodity might be temperature or rainfall: clearly exogenous to the market, but likely to be important in the production process.

For the model of macro consumption, we might use autonomous government expenditure or the level of exports as an instrument. Those components of GDP are clearly related to the level of GDP and disposable income, but they are not directly affected by consumption shocks.

# But why should we not always use IV?

First, It may be difficult to find variables that can serve as valid instruments. Many variables that have an effect on included endogenous variables also have a direct effect on the dependent variable. Chris Sims' critique of macro modelers employing 'incredible identifying restrictions' should be taken seriously, as identification requires that certain variables not appear in the equation to be estimated.

Second, IV estimators are innately *biased*, and their finite-sample properties are often problematic. Thus, most of the justification for the use of IV is asymptotic. Performance in small samples may be poor.

Third, the precision of IV estimates is lower than that of OLS estimates (least squares is just that). In the presence of *weak instruments* (excluded instruments only weakly correlated with included endogenous regressors) the loss of precision will be severe, and IV estimates may be no improvement over OLS. This suggests we need a test to determine whether a particular regressor must be treated as endogenous in order to produce consistent estimates.

# The IV–GMM estimator

To discuss the implementation of IV estimators and test statistics, we consider a more general framework: an instrumental variables estimator implemented using the Generalized Method of Moments (GMM). As we will see, conventional IV estimators such as two-stage least squares (2SLS) are special cases of this IV-GMM estimator.

The model:

$$y = X\beta + u, \quad u \sim (0, \Omega)$$

with  $X$  ( $N \times k$ ) and define a matrix  $Z$  ( $N \times \ell$ ) where  $\ell \geq k$ . This is the Generalized Method of Moments IV (IV-GMM) estimator.

The  $\ell$  instruments give rise to a set of  $\ell$  moments:

$$g_i(\beta) = Z_i' u_i = Z_i'(y_i - x_i\beta), \quad i = 1, N$$

where each  $g_i$  is an  $\ell$ -vector. The method of moments approach considers each of the  $\ell$  moment equations as a sample moment, which we may estimate by averaging over  $N$ :

$$\bar{g}(\beta) = \frac{1}{N} \sum_{i=1}^N z_i(y_i - x_i\beta) = \frac{1}{N} Z' u$$

The GMM approach chooses an estimate that solves  $\bar{g}(\hat{\beta}_{GMM}) = 0$ .



If  $\ell = k$ , the equation to be estimated is said to be *exactly identified* by the *order condition* for identification: that is, there are as many excluded instruments as included right-hand endogenous variables. The method of moments problem is then  $k$  equations in  $k$  unknowns, and a unique solution exists, equivalent to the standard IV estimator:

$$\hat{\beta}_{IV} = (Z'X)^{-1}Z'y$$

In the case of *overidentification* ( $\ell > k$ ) we may define a set of  $k$  instruments:

$$\hat{X} = Z'(Z'Z)^{-1}Z'X = P_ZX$$

which gives rise to the *two-stage least squares* (2SLS) estimator

$$\hat{\beta}_{2SLS} = (\hat{X}'X)^{-1}\hat{X}'y = (X'P_ZX)^{-1}X'P_Zy$$

which despite its name is computed by this single matrix equation.

In the 2SLS method with overidentification, the  $\ell$  available instruments are “boiled down” to the  $k$  needed by defining the  $P_Z$  matrix. In the IV-GMM approach, that reduction is not necessary. All  $\ell$  instruments are used in the estimator. Furthermore, a *weighting matrix* is employed so that we may choose  $\hat{\beta}_{GMM}$  so that the elements of  $\bar{g}(\hat{\beta}_{GMM})$  are as close to zero as possible. With  $\ell > k$ , not all  $\ell$  moment conditions can be exactly satisfied, so a criterion function that weights them appropriately is used to improve the efficiency of the estimator.

The GMM estimator minimizes the criterion

$$J(\hat{\beta}_{GMM}) = N \bar{g}(\hat{\beta}_{GMM})' W \bar{g}(\hat{\beta}_{GMM})$$

where  $W$  is a  $\ell \times \ell$  symmetric weighting matrix.

Solving the set of FOCs, we derive the IV-GMM estimator of an overidentified equation:

$$\hat{\beta}_{GMM} = (X'ZWZ'X)^{-1}X'ZWZ'y$$

which will be identical for all  $W$  matrices which differ by a factor of proportionality. The *optimal* weighting matrix, as shown by Hansen (1982), chooses  $W = S^{-1}$  where  $S$  is the covariance matrix of the moment conditions to produce the most *efficient* estimator:

$$S = E[Z'uu'Z] = \lim_{N \rightarrow \infty} N^{-1}[Z'\Omega Z]$$

With a consistent estimator of  $S$  derived from 2SLS residuals, we define the feasible IV-GMM estimator as

$$\hat{\beta}_{FEGMM} = (X'Z \hat{S}^{-1}Z'X)^{-1}X'Z \hat{S}^{-1}Z'y$$

where *FEGMM* refers to the *feasible efficient* GMM estimator.

# IV-GMM and the distribution of $u$

The derivation makes no mention of the form of  $\Omega$ , the variance-covariance matrix (vce) of the error process  $u$ . If the errors satisfy all classical assumptions are *i.i.d.*,  $S = \sigma_u^2 I_N$  and the optimal weighting matrix is proportional to the identity matrix. The IV-GMM estimator is merely the standard IV (or 2SLS) estimator.

# IV-GMM robust estimates

If there is heteroskedasticity of unknown form, we usually compute *robust* standard errors in any Stata estimation command to derive a consistent estimate of the *vce*. In this context,

$$\hat{S} = \frac{1}{N} \sum_{i=1}^N \hat{u}_i^2 \mathbf{z}_i' \mathbf{z}_i$$

where  $\hat{u}$  is the vector of residuals from any consistent estimator of  $\beta$  (e.g., the 2SLS residuals). For an overidentified equation, the IV-GMM estimates computed from this estimate of  $S$  will be more efficient than 2SLS estimates.

# IV-GMM cluster-robust estimates

If errors are considered to exhibit arbitrary intra-cluster correlation in a dataset with  $M$  clusters, we may derive a *cluster-robust* IV-GMM estimator using

$$\hat{S} = \sum_{j=1}^M \hat{u}_j' \hat{u}_j$$

where

$$\hat{u}_j = (y_j - x_j \hat{\beta}) X' Z (Z' Z)^{-1} z_j$$

The IV-GMM estimates employing this estimate of  $S$  will be both robust to arbitrary heteroskedasticity and intra-cluster correlation, equivalent to estimates generated by Stata's `cluster(varname)` option. For an overidentified equation, IV-GMM cluster-robust estimates will be more efficient than 2SLS estimates.

# IV-GMM HAC estimates

The IV-GMM approach may also be used to generate *HAC standard errors*: those robust to arbitrary heteroskedasticity and autocorrelation. Although the best-known *HAC* approach in econometrics is that of Newey and West, using the Bartlett kernel (per Stata's `newey`), that is only one choice of a *HAC* estimator that may be applied to an IV-GMM problem.

Baum–Schaffer–Stillman's `ivreg2` (from the SSC Archive) and Stata 10's `ivregress` provide several choices for kernels. For some kernels, the kernel *bandwidth* (roughly, number of lags employed) may be chosen automatically in either command.

# Example of IV and IV-GMM estimation

We illustrate various forms of the IV estimator with a Phillips curve equation fit to quarterly US data from the `usmacro1` dataset. The model should not be taken seriously, as its specification is for pedagogical purposes. We first fit the relationship with the standard 2SLS estimator, using Baum–Schaffer–Stillman’s `ivreg2` command. You could fit the same equation with `ivregress 2sls`.

We model the year-over-year rate of inflation in a wage measure (average hourly earnings in manufacturing) as a function of the current unemployment rate. To deal with potential endogeneity of the unemployment rate, we use lags 2–4 of the unemployment rate as instruments. We first fit the equation through 1973q4, prior to the first oil shock. Some of the standard `ivreg2` output, relating to weak instruments, has been edited on the following slides.



```
. ivreg2 wageinfl (unemp = L(2/4).unemp) if tin(,1973q4)
```

IV (2SLS) estimation

Estimates efficient for homoskedasticity only  
 Statistics consistent for homoskedasticity only

Total (centered) SS	=	158.1339335	Number of obs	=	56
Total (uncentered) SS	=	1362.450328	F( 1, 54)	=	4.95
Residual SS	=	142.674146	Prob > F	=	0.0303
			Centered R2	=	0.0978
			Uncentered R2	=	0.8953
			Root MSE	=	1.596

wageinfl	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
unemp	-.6012813	.265382	-2.27	0.023	-1.121421	-.0811421
_cons	7.610898	1.329598	5.72	0.000	5.004934	10.21686

Underidentification test (Anderson canon. corr. LM statistic):	32.622
Chi-sq(3) P-val =	0.0000

Sargan statistic (overidentification test of all instruments):	0.046
Chi-sq(2) P-val =	0.9771

Instrumented: unemp

Excluded instruments: L2.unemp L3.unemp L4.unemp

We may fit this equation with different assumptions about the error process. The estimates above assume *i.i.d.* errors. We may also compute robust standard errors in the 2SLS context.

We then apply IV-GMM with robust standard errors. As the equation is overidentified, the IV-GMM estimates will differ, and will be more efficient than the robust 2SLS estimates.

Last, we may estimate the equation with IV-GMM and HAC standard errors, using the default Bartlett kernel (as employed by Newey–West) and a bandwidth of 5 quarters. This corresponds to four lags in the `newey` command.

```
. estimates table IID Robust IVGMM IVGMM_HAC, b(%9.4f) se(%5.3f) t(%5.2f) ///
> title(Alternative IV estimates of pre-1974 Phillips curve) stat(rmse)
Alternative IV estimates of pre-1974 Phillips curve
```

Variable	IID	Robust	IVGMM	IVGMM_HAC
unemp	-0.6013	-0.6013	-0.6071	-0.6266
	0.265	0.219	0.217	0.295
	-2.27	-2.75	-2.80	-2.13
_cons	7.6109	7.6109	7.6320	7.7145
	1.330	1.018	1.007	1.363
	5.72	7.48	7.58	5.66
rmse	1.5962	1.5962	1.5966	1.5982

legend: b/se/t

Note that the coefficients' point estimates change when IV-GMM is employed, and that their  $t$ -statistics are larger than those of robust IV. The point estimates are also altered when IV-GMM with HAC VCE is computed. As expected, 2SLS yields the smallest RMS error.

# Tests of overidentifying restrictions

If and only if an equation is *overidentified*, with more excluded instruments than included endogenous variables, we may test whether the excluded instruments are appropriately independent of the error process. That test should always be performed when it is possible to do so, as it allows us to evaluate the validity of the instruments.

A test of *overidentifying restrictions* regresses the residuals from an IV or 2SLS regression on all instruments in  $Z$ . Under the null hypothesis that all instruments are uncorrelated with  $u$ , the test has a large-sample  $\chi^2(r)$  distribution where  $r$  is the number of overidentifying restrictions.

Under the assumption of *i.i.d.* errors, this is known as a *Sargan test*, and is routinely produced by `ivreg2` for IV and 2SLS estimates. After `ivregress`, the command `estat overid` provides the test.

If we have used IV-GMM estimation in `ivreg2`, the test of overidentifying restrictions becomes the Hansen  $J$  statistic: the GMM criterion function. Although  $J$  will be identically zero for any exactly-identified equation, it will be positive for an overidentified equation. If it is “too large”, doubt is cast on the satisfaction of the moment conditions underlying GMM.

The test in this context is known as the *Hansen test* or *J test*, and is calculated by `ivreg2` when the `gmm2s` option is employed.

The Sargan–Hansen test of overidentifying restrictions should be performed routinely in any overidentified model estimated with instrumental variables techniques. Instrumental variables techniques are powerful, but if a strong rejection of the null hypothesis of the Sargan–Hansen test is encountered, you should strongly doubt the validity of the estimates.

For instance, consider a variation of the IV-GMM model estimated above (with robust standard errors) and focus on the test of overidentifying restrictions provided by the Hansen  $J$  statistic.

In this form of the model, estimated through 1979q4, we also include the growth rate of the monetary base, `basegro`, as an excluded instrument. The model is overidentified by three degrees of freedom, as there is one endogenous regressor and four excluded instruments. We see that the  $J$  statistic clearly rejects its null, casting doubt on our choice of instruments.

```
. ivreg2 wageinfl (unemp = L(2/4).unemp basegro) if tin(,1979q4), robust gmm2s
2-Step GMM estimation
```

Estimates efficient for arbitrary heteroskedasticity  
 Statistics robust to heteroskedasticity

Total (centered) SS	=	414.4647455	Number of obs	=	80
Total (uncentered) SS	=	3075.230877	F( 1, 78)	=	22.46
Residual SS	=	377.8689419	Prob > F	=	0.0000
			Centered R2	=	0.0883
			Uncentered R2	=	0.8771
			Root MSE	=	2.173

wageinfl	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
unemp	.7228864	.1506083	4.80	0.000	.4276996	1.018073
_cons	2.229875	.8310032	2.68	0.007	.6011386	3.858611

Underidentification test (Kleibergen-Paap rk LM statistic): 27.693  
 Chi-sq(4) P-val = 0.0000

Hansen J statistic (overidentification test of all instruments): 30.913  
 Chi-sq(3) P-val = 0.0000

Instrumented: unemp  
 Excluded instruments: L2.unemp L3.unemp L4.unemp basegro

We reestimate the model, retaining money base growth as an exogenous variable, but including it in the estimated equation rather than applying an exclusion restriction. The resulting  $J$  statistic now fails to reject its null.



```
. ivreg2 wageinfl (unemp = L(2/4).unemp) basegro if tin(,1979q4), robust gmm2s
2-Step GMM estimation
```

Estimates efficient for arbitrary heteroskedasticity  
 Statistics robust to heteroskedasticity

		Number of obs =	80
		F( 2, 77) =	122.14
		Prob > F =	0.0000
Total (centered) SS	=	414.4647455	
Total (uncentered) SS	=	3075.230877	
Residual SS	=	100.724328	
		Centered R2 =	0.7570
		Uncentered R2 =	0.9672
		Root MSE =	1.122

wageinfl	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
unemp	.3350836	.0796765	4.21	0.000	.1789206	.4912466
basegro	.7582774	.0592661	12.79	0.000	.6421181	.8744368
_cons	-.346625	.5022148	-0.69	0.490	-1.330948	.6376979

Underidentification test (Kleibergen-Paap rk LM statistic): 29.279  
 Chi-sq(3) P-val = 0.0000

Hansen J statistic (overidentification test of all instruments): 1.147  
 Chi-sq(2) P-val = 0.5635

Instrumented: unemp  
 Included instruments: basegro  
 Excluded instruments: L2.unemp L3.unemp L4.unemp

It is important to understand that the Sargan–Hansen test of overidentifying restrictions is a joint test of the hypotheses that the instruments, excluded and included, are independently distributed of the error process *and* that they are properly excluded from the model.

Note as well that all exogenous variables in the equation—excluded and included—appear in the set of instruments  $Z$ . In the context of single-equation IV estimation, they must. You cannot pick and choose which instruments appear in which ‘first-stage’ regressions.

# Testing a subset of overidentifying restrictions

We may be quite confident of some instruments' independence from  $u$  but concerned about others. In that case a *GMM distance* or *C* test may be used. The `orthog( )` option of `ivreg2` tests whether a *subset* of the model's overidentifying restrictions appear to be satisfied.

This is carried out by calculating two Sargan–Hansen statistics: one for the full model and a second for the model in which the listed variables are (a) considered endogenous, if included regressors, or (b) dropped, if excluded regressors. In case (a), the model must still satisfy the order condition for identification. The difference of the two Sargan–Hansen statistics, often termed the *GMM distance* or Hayashi *C statistic*, will be distributed  $\chi^2$  under the null hypothesis that the specified orthogonality conditions are satisfied, with d.f. equal to the number of those conditions.

We perform the  $C$  test on the estimated equation by challenging the exogeneity of `basegro`. Is it properly considered exogenous? The `orthog( )` option reestimates the equation, treating it as endogenous, and evaluates the difference in the  $J$  statistics from the two models. Considering `basegro` as exogenous is essentially imposing one more orthogonality condition on the GMM estimation problem.

```
. ivreg2 wageinfl (unemp = L(2/4).unemp) basegro if tin(,1979q4), ///
  robust gmm2s orthog(basegro)
```

### 2-Step GMM estimation

Estimates efficient for arbitrary heteroskedasticity

Statistics robust to heteroskedasticity

Total (centered) SS	=	414.4647455	Number of obs	=	80
Total (uncentered) SS	=	3075.230877	F( 2, 77)	=	122.14
Residual SS	=	100.724328	Prob > F	=	0.0000
			Centered R2	=	0.7570
			Uncentered R2	=	0.9672
			Root MSE	=	1.122

wageinfl	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
unemp	.3350836	.0796765	4.21	0.000	.1789206	.4912466
basegro	.7582774	.0592661	12.79	0.000	.6421181	.8744368
_cons	-.346625	.5022148	-0.69	0.490	-1.330948	.6376979

Hansen J statistic (overidentification test of all instruments): 1.147  
Chi-sq(2) P-val = 0.5635

-orthog- option:

Hansen J statistic (eqn. excluding suspect orthog. conditions): 0.620  
Chi-sq(1) P-val = 0.4312

C statistic (exogeneity/orthogonality of suspect instruments): 0.528  
Chi-sq(1) P-val = 0.4676

Instruments tested: basegro

It appears that `basegro` may be considered exogenous in this specification.

A variant on this strategy is implemented by the `endog( )` option of `ivreg2`, in which one or more variables considered endogenous can be tested for exogeneity. The *C* test in this case will consider whether the null hypothesis of their exogeneity is supported by the data.

If all endogenous regressors are included in the `endog( )` option, the test is essentially a test of whether IV methods are required to estimate the equation. If OLS estimates of the equation are consistent, they should be preferred. In this context, the test is equivalent to a (*Durbin–Wu–*)*Hausman test* comparing IV and OLS estimates, as implemented by Stata's `hausman` command with the `sigmaless` option. Using `ivreg2`, you need not estimate and store both models to generate the test's verdict.

For instance, with the model above, we might question whether IV techniques are needed. We can conduct the C test via:

```
ivreg2 wageinfl (unemp = L(2/4).unemp) basegro if tin(,1979q4), ///  
robust gmm2s endog(unemp)
```

where the `endog(unemp)` option tests the null hypothesis that the variable can be treated as exogenous in this model, rather than as an endogenous variable.

```
. ivreg2 wageinfl (unemp = L(2/4).unemp) basegro if tin(,1979q4), robust gmm2s
> endog(unemp)
```

### 2-Step GMM estimation

Estimates efficient for arbitrary heteroskedasticity

Statistics robust to heteroskedasticity

Total (centered) SS	=	414.4647455	Number of obs	=	80
Total (uncentered) SS	=	3075.230877	F( 2, 77)	=	122.14
Residual SS	=	100.724328	Prob > F	=	0.0000
			Centered R2	=	0.7570
			Uncentered R2	=	0.9672
			Root MSE	=	1.122

wageinfl	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
unemp	.3350836	.0796765	4.21	0.000	.1789206	.4912466
basegro	.7582774	.0592661	12.79	0.000	.6421181	.8744368
_cons	-.346625	.5022148	-0.69	0.490	-1.330948	.6376979

Hansen J statistic (overidentification test of all instruments): 1.147  
Chi-sq(2) P-val = 0.5635

-endog- option:

Endogeneity test of endogenous regressors: 1.505  
Chi-sq(1) P-val = 0.2200

Regressors tested: unemp

Instrumented: unemp



In this context, it appears that we could safely estimate this equation with OLS techniques, as the P-value for the C test of endogenous regressors of 0.2200 does not reject the null hypothesis.

There are a number of other diagnostic tools that may be employed in instrumental variables estimation. Although time constraints prevents their thorough discussion, full details can be found in the Baum–Schaffer–Stillman *Stata Journal* articles.

# The weak instruments problem

Instrumental variables methods rely on two assumptions: the excluded instruments are distributed independently of the error process, and they are sufficiently correlated with the included endogenous regressors. Tests of overidentifying restrictions address the *first* assumption, although we should note that a rejection of their null may be indicative that the exclusion restrictions for these instruments may be inappropriate. That is, some of the instruments have been improperly excluded from the regression model's specification.

The specification of an instrumental variables model asserts that the excluded instruments affect the dependent variable only *indirectly*, through their correlations with the included endogenous variables. If an excluded instrument exerts both direct and indirect influences on the dependent variable, the exclusion restriction should be rejected. This can be readily tested by including the variable as a regressor, as we did above with `basegro`.

To test the *second* assumption—that the excluded instruments are sufficiently correlated with the included endogenous regressors—we should consider the goodness-of-fit of the “first stage” regressions relating each endogenous regressor to the entire set of instruments.

It is important to understand that the theory of single-equation (“limited information”) IV estimation requires that all columns of  $X$  are conceptually regressed on all columns of  $Z$  in the calculation of the estimates. We cannot meaningfully speak of “this variable is an instrument for that regressor” or somehow restrict which instruments enter which first-stage regressions. Stata’s `ivregress` or `ivreg2` will not let you do that because such restrictions only make sense in the context of estimating an entire system of equations by full-information methods (for instance, with `reg3`).

The `first` and `ffirst` options of `ivreg2` (or the `first` option of `ivregress`) present several useful diagnostics that assess the first-stage regressions. If there is a single endogenous regressor, these issues are simplified, as the instruments either explain a reasonable fraction of that regressor's variability or not. With multiple endogenous regressors, diagnostics are more complicated, as each instrument is being called upon to play a role in each first-stage regression.

With sufficiently weak instruments, the asymptotic identification status of the equation is called into question. An equation identified by the order and rank conditions in a finite sample may still be *effectively unidentified* or it numerically unidentified.

As Staiger and Stock (*Econometrica*, 1997) show, the weak instruments problem can arise even when the first-stage  $t$ - and  $F$ -tests are significant at conventional levels in a large sample. In the worst case, the bias of the IV estimator is the same as that of OLS, IV becomes inconsistent, and instrumenting only aggravates the problem.

Beyond the informal “rule-of-thumb” diagnostics such as  $F > 10$ , `ivreg2` computes several statistics that can be used to critically evaluate the strength of instruments. We can write the first-stage regressions as

$$X = Z\Pi + v$$

With  $X_1$  as the endogenous regressors,  $Z_1$  the excluded instruments and  $Z_2$  as the included instruments, this can be partitioned as

$$X_1 = [Z_1 Z_2] [\Pi'_{11} \Pi'_{12}]' + v_1$$

The rank condition for identification states that the  $L \times K_1$  matrix  $\Pi_{11}$  must be of full column rank.

We do not observe the true  $\Pi_{11}$ , so we must replace it with an estimate. Anderson's (John Wiley, 1984) approach to testing the rank of this matrix (or that of the full  $\Pi$  matrix) considers the *canonical correlations* of the  $X$  and  $Z$  matrices. If the equation is to be identified, all  $K$  of the canonical correlations will be significantly different from zero.

The squared canonical correlations can be expressed as eigenvalues of a matrix. Anderson's CC test considers the null hypothesis that the minimum canonical correlation is zero. Under the null, the test statistic is distributed  $\chi^2$  with  $(L - K + 1)$  d.f., so it may be calculated even for an exactly-identified equation. Failure to reject the null suggests the equation is unidentified. `ivreg2` routinely reports this Lagrange Multiplier (LM) statistic. In the first example of 2SLS shown above, you see the Anderson canonical correlation statistic as a test for underidentification.



The C–D statistic is a closely related test of the rank of a matrix. While the Anderson CC test is a LR test, the C–D test is a Wald statistic, with the same asymptotic distribution. The C–D statistic plays an important role in Stock and Yogo's work (see below). Both the Anderson and C–D tests are reported by `ivreg2` with the `first` option.

Research by Kleibergen and Paap (KP) (*J. Econometrics*, 2006) has developed a robust version of a test for the rank of a matrix: e.g. testing for *underidentification*. The statistic has been implemented by Kleibergen and Schaffer as command `ranktest`, which is part of the `ivreg2` package. If non-*i.i.d.* errors are assumed, the `ivreg2` output contains the K–P `rk` statistic in place of the Anderson canonical correlation statistic as a test of underidentification, as you can see in the first IV-GMM example above.

The canonical correlations may also be used to test a set of instruments for *redundancy* by considering their statistical significance in the first stage regressions. This can be calculated, in robust form, as a K–P LM test. The `redundant( )` option of `ivreg2` allows a set of excluded instruments to be tested for relevance, with the null hypothesis that they do not contribute to the asymptotic efficiency of the equation.

Stock and Yogo (Camb. U. Press festschrift, 2005) propose testing for weak instruments by using the  $F$ -statistic form of the C–D statistic. Their null hypothesis is that the estimator is weakly identified in the sense that it is subject to bias that the investigator finds unacceptably large.

Their test comes in two flavors: maximal relative bias (relative to the bias of OLS) and maximal size. The former test has the null that instruments are weak, where weak instruments are those that can lead to an asymptotic relative bias greater than some level  $b$ . This test uses the finite sample distribution of the IV estimator, and can only be calculated where the appropriate moments exist (when the equation is suitably overidentified: the  $m^{\text{th}}$  moment of an IV estimator exists iff  $m < (L - K + 1)$ ). The test is routinely reported in `ivreg2` and `ivregress` output when it can be calculated, with the relevant critical values calculated by Stock and Yogo.

The second test proposed by Stock and Yogo is based on the performance of the Wald test statistic for the endogenous regressors. Under weak identification, the test rejects too often. The test statistic is based on the rejection rate  $r$  tolerable to the researcher if the true rejection rate is 5%. Their tabulated values consider various values for  $r$ . To be able to reject the null that the size of the test is unacceptably large (versus 5%), the Cragg–Donald  $F$  statistic must exceed the tabulated critical value.

The Stock–Yogo test statistics, like others discussed above, assume *i.i.d.* errors. The Cragg–Donald  $F$  can be robustified in the absence of *i.i.d.* errors by using the Kleibergen–Paap  $\text{rk}$  statistic, which `ivreg2` reports in that circumstance.

# LIML and GMM-CUE

OLS and IV estimators are special cases of *k-class estimators*: OLS with  $k = 0$  and IV with  $k = 1$ . Limited-information maximum likelihood (LIML) is another member of this class, with  $k$  chosen optimally in the estimation process. Like any ML estimator, LIML is invariant to normalization. In an equation with two endogenous variables, it does not matter whether you specify  $y_1$  or  $y_2$  as the left-hand variable.

One of the other virtues of the LIML estimator is that it has been found to be more resistant to weak instruments problems than the IV estimator. On the down side, it makes the distributional assumption of normally distributed (and *i.i.d.*) errors. `ivreg2` produces LIML estimates with the `liml` option, and `liml` is a subcommand for official Stata's `ivregress`.

If the *i.i.d.* assumption of LIML is not reasonable, you may use the GMM equivalent: the *continuously updated* GMM estimator, or CUE estimator. In `ivreg2`, the `cue` option combined with `robust`, `cluster` and/or `bw( )` options specifies that non-*i.i.d.* errors are to be modeled. GMM-CUE requires numerical optimization, and may require many iterations to converge.

`ivregress` provides an iterated GMM estimator, which is not the same estimator as GMM-CUE.

# Testing for *i.i.d.* errors in IV

In the context of an equation estimated with instrumental variables, the standard diagnostic tests for heteroskedasticity and autocorrelation are generally not valid.

In the case of heteroskedasticity, Pagan and Hall (*Econometric Reviews*, 1983) showed that the Breusch–Pagan or Cook–Weisberg tests (`estat hettest`) are generally not usable in an IV setting. They propose a test that will be appropriate in IV estimation where heteroskedasticity may be present in more than one structural equation. Mark Schaffer's `ivhettest`, part of the `ivreg2` suite, performs the Pagan–Hall test under a variety of assumptions on the indicator variables. It will also reproduce the Breusch–Pagan test if applied in an OLS context.

In the same token, the Breusch–Godfrey statistic used in the OLS context (`estat bgodfrey`) will generally not be appropriate in the presence of endogenous regressors, overlapping data or conditional heteroskedasticity of the error process. Cumby and Huizinga (*Econometrica*, 1992) proposed a generalization of the BG statistic which handles each of these cases.

Their test is actually more general in another way. Its null hypothesis of the test is that the regression error is a moving average of known order  $q \geq 0$  against the general alternative that autocorrelations of the regression error are nonzero at lags greater than  $q$ . In that context, it can be used to test that autocorrelations beyond any  $q$  are zero. Like the BG test, it can test multiple lag orders. The C–H test is available as Baum and Schaffer's `ivactest` routine, part of the `ivreg2` suite.



For more details on IV and IV-GMM, please see

- Enhanced routines for instrumental variables/GMM estimation and testing. Baum, C.F., Schaffer, M.E., Stillman, S., *Stata Journal* 7:4, 2007.
- *An Introduction to Modern Econometrics Using Stata*, Baum, C.F., Stata Press, 2006 (particularly Chapter 8).
- Instrumental variables and GMM: Estimation and testing. Baum, C.F., Schaffer, M.E., Stillman, S., *Stata Journal* 3:1–31, 2003.

Both of the *Stata Journal* papers are freely downloadable from <http://stata-journal.com>.

# Nonlinear least squares estimation

Besides the capabilities for maximum likelihood estimation of one or several equations via the `ml` suite of commands, Stata provides facilities for single-equation nonlinear least squares estimation with `nl` and the estimation of nonlinear systems of equations with `nlSUR`.

The `nl` and `nlSUR` commands may be invoked in three ways: interactively, using a “programmed substitutable expression”, and using a “function evaluator program”. We discuss the first and third methods here. The function evaluator program is quite similar to likelihood function evaluators for `ml` (maximum likelihood estimation).

In the interactive mode, you specify the nonlinear least squares expression, including starting values if necessary, on the command line. For example, consider the two-factor CES production function:

$$\ln Q_i = \beta_0 - \frac{1}{\rho} \ln \left( \delta K_i^{-\rho} + (1 - \delta) L_i^{-\rho} \right) + u_i$$

with the parameters  $\beta_0, \rho, \delta$ .

This could be estimated with:

```
nl (lnQ={b0}-1/{rho=1}*ln({delta=0.5}*K^(-1*{rho}) +
    (1-{delta})*L^(-1*{rho})))
```

Note that the parameters are enclosed in `{ }`, with initial values given if needed. The entire equation must be enclosed by `( )`. You may use options such as `robust` and `cluster(varname)` with `nl`.

The standard apparatus for any estimation command is available after invoking `nl`. For instance, we might want to calculate the elasticity of substitution for the CES function, defined as  $\sigma = 1/(1 + \rho)$ . The `nlcom` command can provide point and interval estimates of this expression via the delta method:

```
nlcom (sigma: 1 / ( 1 + [rho]_b[_cons] ))
```

where we refer to the “constant” in the `rho` equation, and label the resulting expression `sigma` in the output.

After `nl`, all of the standard results from any estimation command are available for further use; `ereturn list` for details.

If you want to use `nl` extensively for a particular problem, it makes sense to develop a *function evaluator program*. That program is quite similar to any Stata `ado`-file or `ml` program. It must be named `nlfunc.ado`, where *func* is a name of your choice: e.g., `nlces.ado` for a CES function evaluator.

The stylized function evaluator program contains:

```
program nlfunc
    version 11
    syntax varlist(min=n max=n) if, at(name)
    // extract vars from varlist
    // extract params as scalars from at matrix
    // fill in dependent variable with replace
end
```

To use the program `nlces`, call it with the `nl` command, but only include the unique part of its name, followed by `@`:

```
nl ces @ lnQ cap lab, parameters(b0 rho delta) ///  
    initial(b0 0 rho 1 delta 0.5)
```

You could restrict analysis to a subsample with the *if exp* qualifier:

```
nl ces @ lnQ cap lab if industry==33, ...
```

Note that the `nlsur` command estimates systems of *seemingly unrelated* nonlinear equations, just as `sureg` estimates systems of seemingly unrelated linear equations. In that context, `nlsur` cannot be used to estimate a system of simultaneous nonlinear equations. The `gmm` command, as we now discuss, could be used for that purpose, as could Stata's maximum likelihood commands (`ml`).

# Programs for GMM estimation

There are various Stata commands, official and user-written, that make use of Generalized Method of Moments (GMM) estimation. In Stata version 11, there is a general-purpose GMM command, `gmm`, that can be used to solve GMM estimation problems of any type. Like the `nl` (nonlinear-least squares) command, `gmm` can be used interactively, but it is likely to be used in its *function evaluator program* form. In that form, just as with `ml` or the programmed version of `nl`, you write a program specifying the estimation problem.



The function evaluator program, or moment-evaluator program, is passed a *varlist* containing the moments to be evaluated for each observation. Your program replaces the elements of the *varlist* with the ‘error part’ of the moment conditions. For instance, if we were to solve an OLS regression problem with GMM we might write a moment-evaluator program as:

```

. program gmm_reg
1.      version 11
2.      syntax varlist if, at(name)
3.      qui {
4.          tempvar xb
5.          gen double `xb' = x1*`at'[1,1] + x2*`at'[1,2] + ///
>          x3*`at'[1,3] + `at'[1,4] `if'
6.          replace `varlist' = y - `xb' `if'
7.      }
8. end

```

where we have specified that the regression has three explanatory variables and a constant term, with variable  $y$  as the dependent variable. The row vector `at()` contains the current values of the estimated parameters. The contents of *varlist* are replaced with the discrepancies,  $y - X\beta$ , defined by those parameters. A *varlist* is used as `gmm` can handle multiple-equation problems.

To perform the estimation (using the standard `auto` dataset), we specify the parameters and instruments to be used in the `gmm` command:

```
. gen y = price
. gen x1 = weight
. gen x2 = length
. gen x3 = turn
. gmm gmm_reg, nequations(1) parameters(b1 b2 b3 b0) ///
> instruments(weight length turn) onestep nolog
```

Final GMM criterion  $Q(b) = 2.43e-16$

GMM estimation

Number of parameters = 4

Number of moments = 4

Initial weight matrix: Unadjusted

Number of obs = 74

	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
/b1	5.382135	1.719276	3.13	0.002	2.012415	8.751854
/b2	-66.17856	57.56738	-1.15	0.250	-179.0086	46.65143
/b3	-318.2055	171.6618	-1.85	0.064	-654.6564	18.24543
/b0	14967.64	6012.23	2.49	0.013	3183.881	26751.39

Instruments for equation 1: weight length turn \_cons

This may seem unusual syntax, but we are just stating that we want to use the regressors as instruments for themselves in solving the GMM problem, as under the hypothesis of  $E[u|X] = 0$ , the appropriate moment conditions can be written as  $EX'u = 0$ .

Inspection of the parameters and their standard errors shows that these estimates match those from `regress, robust` for the same model. It is quite unnecessary to use GMM in this context, of course, but it illustrates the way in which you may set up a GMM problem.

To perform linear instrumental variables, we can use the same moment-evaluator program and merely alter the instrument list:

```
. webuse hsng2, clear
(1980 Census housing data)

. gen y = rent
. gen x1 = hsngval
. gen x2 = pcturban
. gen x3 = popden
. gmm gmm_reg, nequations(1) parameters(b1 b2 b3 b0) ///
> instruments(pcturban popden faminc reg2-reg4) onestep nolog
```

Final GMM criterion  $Q(b) = 150.8821$

GMM estimation

Number of parameters = 4

Number of moments = 7

Initial weight matrix: Unadjusted

Number of obs = 50

	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
/b1	.0022538	.0006785	3.32	0.001	.000924	.0035836
/b2	.0281637	.5017214	0.06	0.955	-.9551922	1.01152
/b3	.0006083	.0012742	0.48	0.633	-.0018891	.0031057
/b0	122.6632	17.26189	7.11	0.000	88.83052	156.4959

Instruments for equation 1: pcturban popden faminc reg2 reg3 reg4 \_cons

These estimates match those produced by `ivregress 2sls, robust`.

Let us consider solving a nonlinear estimation problem: a binomial probit model, using GMM rather than the usual ML estimator. The moment-evaluator program:

```
. program gmm_probit
  1.         version 11
  2.         syntax varlist if, at(name)
  3.         qui {
  4.             tempvar xb
  5.             gen double `xb' = x1*`at'[1,1] + x2*`at'[1,2] + ///
>                x3*`at'[1,3] + `at'[1,4] `if'
  6.             replace `varlist' = y - normal(`xb')
  7.         }
  8. end
```

To perform the estimation, we specify the parameters and instruments to be used in the `gmm` command:

```
. webuse hsng2, clear
(1980 Census housing data)
. gen y = (region >= 3)
. gen x1 = hsngval
. gen x2 = pcturban
. gen x3 = popden
. gmm gmm_probit, nequations(1) parameters(b1 b2 b3 b0) ///
> instruments(pcturban hsngval popden) onestep nolog
```

Final GMM criterion  $Q(b) = 3.18e-21$

GMM estimation

Number of parameters = 4

Number of moments = 4

Initial weight matrix: Unadjusted

Number of obs = 50

	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
/b1	.0000198	.0000146	1.35	0.177	-8.92e-06	.0000484
/b2	.0139055	.0177526	0.78	0.433	-.020889	.0487001
/b3	-.0003561	.0001142	-3.12	0.002	-.0005799	-.0001323
/b0	-1.136154	.9463889	-1.20	0.230	-2.991042	.7187345

Instruments for equation 1: pcturban hsngval popden \_cons

Inspection of the parameters shows that these estimates are quite similar to those from `probit` for the same model. However, whereas `probit` requires the assumption of *i.i.d.* errors, GMM does not; the standard errors produced by `gmm` are robust to arbitrary heteroskedasticity.

As in the case of our linear regression estimation example, we can use the same moment-evaluator program to estimate an instrumental-variables probit model, similar to that estimated by `ivprobit`. Unlike that ML command, though, we need not make any distributional assumptions about the error process in order to use GMM.



```
. gmm gmm_probit, nequations(1) parameters(b1 b2 b3 b0) ///
> instruments(pcturban popden rent hsnggrew) onestep nolog
```

Final GMM criterion Q(b) = .0470836

GMM estimation

Number of parameters = 4

Number of moments = 5

Initial weight matrix: Unadjusted

Number of obs = 50

	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
/b1	-6.25e-06	.0000203	-0.31	0.758	-.000046	.0000335
/b2	.0370542	.0333466	1.11	0.266	-.028304	.1024124
/b3	-.0014897	.0013724	-1.09	0.278	-.0041795	.0012001
/b0	-.538059	1.278787	-0.42	0.674	-3.044435	1.968317

Instruments for equation 1: pcturban popden rent hsnggrew \_cons

Although the examples of `gmm` moment-evaluator programs we have shown here largely duplicate the functionality of existing Stata commands, they should illustrate that the general-purpose `gmm` command may be used to solve estimation problems not amenable to any existing commands, or indeed to a maximum-likelihood approach. In that sense, familiarity with `gmm` capabilities is likely to be quite helpful if you face challenging estimation problems in your research.