

# Glossary

## A

**Adjusted *R*-Squared:** A goodness-of-fit measure in multiple regression analysis that penalizes additional explanatory variables by using a degrees of freedom adjustment in estimating the error variance.

**Alternative Hypothesis:** The hypothesis against which the null hypothesis is tested.

**AR(1) Serial Correlation:** The errors in a time series regression model follow an AR(1) model.

**Asymptotic Bias:** *See* inconsistency.

**Asymptotic Confidence Interval:** A confidence interval that is approximately valid in large sample sizes.

**Asymptotic Normality:** The sampling distribution of a properly normalized estimator converges to the standard normal distribution.

**Asymptotic Properties:** Properties of estimators and test statistics that apply when the sample size grows without bound.

**Asymptotic Standard Error:** A standard error that is valid in large samples.

**Asymptotic *t* Statistic:** A *t* statistic that has an approximate standard normal distribution in large samples.

**Asymptotic Variance:** The square of the value by which we must divide an estimator in order to obtain an asymptotic standard normal distribution.

**Asymptotically Efficient:** For consistent estimators with asymptotically normal distributions, the estimator with the smallest asymptotic variance.

**Asymptotically Uncorrelated:** A time series process in which the correlation between random variables at two points in time tends to zero as the time interval between them increases. (*See also* weakly dependent.)

**Attenuation Bias:** Bias in an estimator that is always toward zero; thus, the expected value of an estimator with attenuation bias is less in magnitude than the absolute value of the parameter.

**Augmented Dickey-Fuller Test:** A test for a unit root that includes lagged changes of the variable as regressors.

**Autocorrelation:** *See* serial correlation.

**Autoregressive Conditional Heteroskedasticity (ARCH):** A model of dynamic heteroskedasticity where the variance of the error term, given past information, depends linearly on the past squared errors.

**Autoregressive Process of Order One [AR(1)]:** A time series model whose current value depends linearly on its most recent value plus an unpredictable disturbance.

**Auxiliary Regression:** A regression used to compute a test statistic—such as the test statistics for heteroskedasticity and serial correlation—or any other regression that does not estimate the model of primary interest.

**Average:** The sum of *n* numbers divided by *n*.

**Average Marginal Effect:** *See* average partial effect.

**Average Partial Effect:** For nonconstant partial effects, the partial effect averaged across the specified population.

**Average Treatment Effect:** A treatment, or policy, effect averaged across the population.

## B

**Balanced Panel:** A panel data set where all years (or periods) of data are available for all cross-sectional units.

**Base Group:** The group represented by the overall intercept in a multiple regression model that includes dummy explanatory variables.

**Base Period:** For index numbers, such as price or production indices, the period against which all other time periods are measured.

**Base Value:** The value assigned to the base period for constructing an index number; usually the base value is 1 or 100.

**Benchmark Group:** *See* base group.

**Bernoulli (or Binary) Random Variable:** A random variable that takes on the values zero or one.

**Best Linear Unbiased Estimator (BLUE):** Among all linear unbiased estimators, the one with the smallest variance. OLS is BLUE, conditional on the sample values of the explanatory variables, under the Gauss-Markov assumptions.

**Beta Coefficients:** *See* standardized coefficients.

- Bias:** The difference between the expected value of an estimator and the population value that the estimator is supposed to be estimating.
- Biased Estimator:** An estimator whose expectation, or sampling mean, is different from the population value it is supposed to be estimating.
- Biased Towards Zero:** A description of an estimator whose expectation in absolute value is less than the absolute value of the population parameter.
- Binary Response Model:** A model for a binary (dummy) dependent variable.
- Binary Variable:** *See* dummy variable.
- Binomial Distribution:** The probability distribution of the number of successes out of  $n$  independent Bernoulli trials, where each trial has the same probability of success.
- Bivariate Regression Model:** *See* simple linear regression model.
- BLUE:** *See* best linear unbiased estimator.
- Bootstrap:** A resampling method that draws random samples, with replacement, from the original data set.
- Bootstrap Standard Error:** A standard error obtained as the sample standard deviation of an estimate across all bootstrap samples.
- Breusch-Godfrey Test:** An asymptotically justified test for  $AR(p)$  serial correlation, with  $AR(1)$  being the most popular; the test allows for lagged dependent variables as well as other regressors that are not strictly exogenous.
- Breusch-Pagan Test:** A test for heteroskedasticity where the squared OLS residuals are regressed on the explanatory variables in the model.

## C

- Causal Effect:** A *ceteris paribus* change in one variable that has an effect on another variable.
- Censored Normal Regression Model:** The special case of the censored regression model where the underlying population model satisfies the classical linear model assumptions.
- Censored Regression Model:** A multiple regression model where the dependent variable has been censored above or below some known threshold.
- Central Limit Theorem (CLT):** A key result from probability theory which implies that the sum of independent random variables, or even weakly dependent random variables, when standardized by its standard deviation, has a distribution that tends to standard normal as the sample size grows.
- Ceteris Paribus:** All other relevant factors are held fixed.
- Chi-Square Distribution:** A probability distribution obtained by adding the squares of independent standard normal random variables. The number of terms in the sum equals the degrees of freedom in the distribution.
- Chi-Square Random Variable:** A random variable with a chi-square distribution.
- Chow Statistic:** An  $F$  statistic for testing the equality of regression parameters across different groups (say, men and women) or time periods (say, before and after a policy change).
- Classical Errors-in-Variables (CEV):** A measurement error model where the observed measure equals the actual variable plus an independent, or at least an uncorrelated, measurement error.
- Classical Linear Model:** The multiple linear regression model under the full set of classical linear model assumptions.
- Classical Linear Model (CLM) Assumptions:** The ideal set of assumptions for multiple regression analysis: for cross-sectional analysis, Assumptions MLR.1 through MLR.6, and for time series analysis, Assumptions TS.1 through TS.6. The assumptions include linearity in the parameters, no perfect collinearity, the zero conditional mean assumption, homoskedasticity, no serial correlation, and normality of the errors.
- Cluster Effect:** An unobserved effect that is common to all units, usually people, in the cluster.
- Cluster Sample:** A sample of natural clusters or groups that usually consist of people.
- Clustering:** The act of computing standard errors and test statistics that are robust to cluster correlation, either due to cluster sampling or to time series correlation in panel data.
- Cochrane-Orcutt (CO) Estimation:** A method of estimating a multiple linear regression model with  $AR(1)$  errors and strictly exogenous explanatory variables; unlike Prais-Winsten, Cochrane-Orcutt does not use the equation for the first time period.
- Coefficient of Determination:** *See*  $R$ -squared.
- Cointegration:** The notion that a linear combination of two series, each of which is integrated of order one, is integrated of order zero.
- Column Vector:** A vector of numbers arranged as a column.
- Composite Error Term:** In a panel data model, the sum of the time-constant unobserved effect and the idiosyncratic error.
- Conditional Distribution:** The probability distribution of one random variable, given the values of one or more other random variables.
- Conditional Expectation:** The expected or average value of one random variable, called the dependent or explained variable, that depends on the values of one or more other variables, called the independent or explanatory variables.
- Conditional Forecast:** A forecast that assumes the future values of some explanatory variables are known with certainty.
- Conditional Median:** The median of a response variable conditional on some explanatory variables.
- Conditional Variance:** The variance of one random variable, given one or more other random variables.

**Confidence Interval (CI):** A rule used to construct a random interval so that a certain percentage of all data sets, determined by the confidence level, yields an interval that contains the population value.

**Confidence Level:** The percentage of samples in which we want our confidence interval to contain the population value; 95% is the most common confidence level, but 90% and 99% are also used.

**Consistency:** An estimator converges in probability to the correct population value as the sample size grows.

**Consistent Estimator:** An estimator that converges in probability to the population parameter as the sample size grows without bound.

**Consistent Test:** A test where, under the alternative hypothesis, the probability of rejecting the null hypothesis converges to one as the sample size grows without bound.

**Constant Elasticity Model:** A model where the elasticity of the dependent variable, with respect to an explanatory variable, is constant; in multiple regression, both variables appear in logarithmic form.

**Contemporaneously Homoskedastic:** Describes a time series or panel data applications in which the variance of the error term, conditional on the regressors in the same time period, is constant.

**Contemporaneously Exogenous:** Describes a time series or panel data application in which a regressor is contemporaneously exogenous if it is uncorrelated with the error term in the same time period, although it may be correlated with the errors in other time periods.

**Continuous Random Variable:** A random variable that takes on any particular value with probability zero.

**Control Group:** In program evaluation, the group that does not participate in the program.

**Control Variable:** See explanatory variable.

**Corner Solution Response:** A nonnegative dependent variable that is roughly continuous over strictly positive values but takes on the value zero with some regularity.

**Correlated Random Effects:** An approach to panel data analysis where the correlation between the unobserved effect and the explanatory variables is modeled, usually as a linear relationship.

**Correlation Coefficient:** A measure of linear dependence between two random variables that does not depend on units of measurement and is bounded between  $-1$  and  $1$ .

**Count Variable:** A variable that takes on nonnegative integer values.

**Covariance:** A measure of linear dependence between two random variables.

**Covariance Stationary:** A time series process with constant mean and variance where the covariance between any two random variables in the sequence depends only on the distance between them.

**Covariate:** See explanatory variable.

**Critical Value:** In hypothesis testing, the value against which a test statistic is compared to determine whether or not the null hypothesis is rejected.

**Cross-Sectional Data Set:** A data set collected by sampling a population at a given point in time.

**Cumulative Distribution Function (cdf):** A function that gives the probability of a random variable being less than or equal to any specified real number.

**Cumulative Effect:** At any point in time, the change in a response variable after a permanent increase in an explanatory variable—usually in the context of distributed lag models.

## D

**Data Censoring:** A situation that arises when we do not always observe the outcome on the dependent variable because at an upper (or lower) threshold we only know that the outcome was above (or below) the threshold. (See also censored regression model.)

**Data Frequency:** The interval at which time series data are collected. Yearly, quarterly, and monthly are the most common data frequencies.

**Data Mining:** The practice of using the same data set to estimate numerous models in a search to find the “best” model.

**Davidson-MacKinnon Test:** A test that is used for testing a model against a nonnested alternative; it can be implemented as a  $t$  test on the fitted values from the competing model.

**Degrees of Freedom ( $df$ ):** In multiple regression analysis, the number of observations minus the number of estimated parameters.

**Denominator Degrees of Freedom:** In an  $F$  test, the degrees of freedom in the unrestricted model.

**Dependent Variable:** The variable to be explained in a multiple regression model (and a variety of other models).

**Derivative:** The slope of a smooth function, as defined using calculus.

**Descriptive Statistic:** A statistic used to summarize a set of numbers; the sample average, sample median, and sample standard deviation are the most common.

**Deseasonalizing:** The removing of the seasonal components from a monthly or quarterly time series.

**Detrending:** The practice of removing the trend from a time series.

**Diagonal Matrix:** A matrix with zeros for all off-diagonal entries.

**Dickey-Fuller Distribution:** The limiting distribution of the  $t$  statistic in testing the null hypothesis of a unit root.

**Dickey-Fuller (DF) Test:** A  $t$  test of the unit root null hypothesis in an AR(1) model. (See also augmented Dickey-Fuller test.)

**Difference in Slopes:** A description of a model where some slope parameters may differ by group or time period.

**Difference-in-Differences Estimator:** An estimator that arises in policy analysis with data for two time periods. One version of the estimator applies to independently pooled cross sections and another to panel data sets.

**Difference-Stationary Process:** A time series sequence that is  $I(0)$  in its first differences.

**Diminishing Marginal Effect:** The marginal effect of an explanatory variable becomes smaller as the value of the explanatory variable increases.

**Discrete Random Variable:** A random variable that takes on at most a finite or countably infinite number of values.

**Distributed Lag Model:** A time series model that relates the dependent variable to current and past values of an explanatory variable.

**Disturbance:** See error term.

**Downward Bias:** The expected value of an estimator is below the population value of the parameter.

**Dummy Dependent Variable:** See binary response model.

**Dummy Variable:** A variable that takes on the value zero or one.

**Dummy Variable Regression:** In a panel data setting, the regression that includes a dummy variable for each cross-sectional unit, along with the remaining explanatory variables. It produces the fixed effects estimator.

**Dummy Variable Trap:** The mistake of including too many dummy variables among the independent variables; it occurs when an overall intercept is in the model and a dummy variable is included for each group.

**Duration Analysis:** An application of the censored regression model where the dependent variable is time elapsed until a certain event occurs, such as the time before an unemployed person becomes reemployed.

**Durbin-Watson (DW) Statistic:** A statistic used to test for first order serial correlation in the errors of a time series regression model under the classical linear model assumptions.

**Dynamically Complete Model:** A time series model where no further lags of either the dependent variable or the explanatory variables help to explain the mean of the dependent variable.

## E

**Econometric Model:** An equation relating the dependent variable to a set of explanatory variables and unobserved disturbances, where unknown population parameters determine the ceteris paribus effect of each explanatory variable.

**Economic Model:** A relationship derived from economic theory or less formal economic reasoning.

**Economic Significance:** See practical significance.

**Elasticity:** The percentage change in one variable given a 1% ceteris paribus increase in another variable.

**Empirical Analysis:** A study that uses data in a formal econometric analysis to test a theory, estimate a relationship, or determine the effectiveness of a policy.

**Endogeneity:** A term used to describe the presence of an endogenous explanatory variable.

**Endogenous Explanatory Variable:** An explanatory variable in a multiple regression model that is correlated with the error term, either because of an omitted variable, measurement error, or simultaneity.

**Endogenous Sample Selection:** Nonrandom sample selection where the selection is related to the dependent variable, either directly or through the error term in the equation.

**Endogenous Variables:** In simultaneous equations models, variables that are determined by the equations in the system.

**Engle-Granger Test:** A test of the null hypothesis that two time series are not cointegrated; the statistic is obtained as the Dickey-Fuller statistic using OLS residuals.

**Engle-Granger Two-Step Procedure:** A two-step method for estimating error correction models whereby the cointegrating parameter is estimated in the first stage, and the error correction parameters are estimated in the second.

**Error Correction Model:** A time series model in first differences that also contains an error correction term, which works to bring two  $I(1)$  series back into long-run equilibrium.

**Error Term:** The variable in a simple or multiple regression equation that contains unobserved factors which affect the dependent variable. The error term may also include measurement errors in the observed dependent or independent variables.

**Error Variance:** The variance of the error term in a multiple regression model.

**Errors-in-Variables:** A situation where either the dependent variable or some independent variables are measured with error.

**Estimate:** The numerical value taken on by an estimator for a particular sample of data.

**Estimator:** A rule for combining data to produce a numerical value for a population parameter; the form of the rule does not depend on the particular sample obtained.

**Event Study:** An econometric analysis of the effects of an event, such as a change in government regulation or economic policy, on an outcome variable.

**Excluding a Relevant Variable:** In multiple regression analysis, leaving out a variable that has a nonzero partial effect on the dependent variable.

**Exclusion Restrictions:** Restrictions which state that certain variables are excluded from the model (or have zero population coefficients).

**Exogenous Explanatory Variable:** An explanatory variable that is uncorrelated with the error term.

**Exogenous Sample Selection:** A sample selection that either depends on exogenous explanatory variables or is independent of the error term in the equation of interest.

**Exogenous Variable:** Any variable that is uncorrelated with the error term in the model of interest.

**Expected Value:** A measure of central tendency in the distribution of a random variable, including an estimator.

**Experiment:** In probability, a general term used to denote an event whose outcome is uncertain. In econometric analysis, it denotes a situation where data are collected by randomly assigning individuals to control and treatment groups.

**Experimental Data:** Data that have been obtained by running a controlled experiment.

**Experimental Group:** *See* treatment group.

**Explained Sum of Squares (SSE):** The total sample variation of the fitted values in a multiple regression model.

**Explained Variable:** *See* dependent variable.

**Explanatory Variable:** In regression analysis, a variable that is used to explain variation in the dependent variable.

**Exponential Function:** A mathematical function defined for all values that has an increasing slope but a constant proportionate change.

**Exponential Smoothing:** A simple method of forecasting a variable that involves a weighting of all previous outcomes on that variable.

**Exponential Trend:** A trend with a constant growth rate.

## F

**F Distribution:** The probability distribution obtained by forming the ratio of two independent chi-square random variables, where each has been divided by its degrees of freedom.

**F Random Variable:** A random variable with an *F* distribution.

**F Statistic:** A statistic used to test multiple hypotheses about the parameters in a multiple regression model.

**Feasible GLS (FGLS) Estimator:** A GLS procedure where variance or correlation parameters are unknown and therefore must first be estimated. (*See also* generalized least squares estimator.)

**Finite Distributed Lag (FDL) Model:** A dynamic model where one or more explanatory variables are allowed to have lagged effects on the dependent variable.

**First Difference:** A transformation on a time series constructed by taking the difference of adjacent time periods, where the earlier time period is subtracted from the later time period.

**First-Differenced (FD) Equation:** In time series or panel data models, an equation where the dependent and independent variables have all been first differenced.

**First-Differenced (FD) Estimator:** In a panel data setting, the pooled OLS estimator applied to first differences of the data across time.

**First Order Autocorrelation:** For a time series process ordered chronologically, the correlation coefficient between pairs of adjacent observations.

**First Order Conditions:** The set of linear equations used to solve for the OLS estimates.

**Fitted Values:** The estimated values of the dependent variable when the values of the independent variables for each observation are plugged into the OLS regression line.

**Fixed Effect:** *See* unobserved effect.

**Fixed Effects Estimator:** For the unobserved effects panel data model, the estimator obtained by applying pooled OLS to a time-demeaned equation.

**Fixed Effects Model:** An unobserved effects panel data model where the unobserved effects are allowed to be arbitrarily correlated with the explanatory variables in each time period.

**Fixed Effects Transformation:** For panel data, the time-demeaned data.

**Forecast Error:** The difference between the actual outcome and the forecast of the outcome.

**Forecast Interval:** In forecasting, a confidence interval for a yet unrealized future value of a time series variable. (*See also* prediction interval.)

**Frisch-Waugh Theorem:** The general algebraic result that provides multiple regression analysis with its “partialling out” interpretation.

**Functional Form Misspecification:** A problem that occurs when a model has omitted functions of the explanatory variables (such as quadratics) or uses the wrong functions of either the dependent variable or some explanatory variables.

## G

**Gauss-Markov Assumptions:** The set of assumptions (Assumptions MLR.1 through MLR.5 or TS.1 through TS.5) under which OLS is BLUE.

**Gauss-Markov Theorem:** The theorem that states that, under the five Gauss-Markov assumptions (for cross-sectional or time series models), the OLS estimator is BLUE (conditional on the sample values of the explanatory variables).

**Generalized Least Squares (GLS) Estimator:** An estimator that accounts for a known structure of the error variance (heteroskedasticity), serial correlation pattern in the errors, or both, via a transformation of the original model.

**Geometric (or Koyck) Distributed Lag:** An infinite distributed lag model where the lag coefficients decline at a geometric rate.

**Goodness-of-Fit Measure:** A statistic that summarizes how well a set of explanatory variables explains a dependent or response variable.

**Granger Causality:** A limited notion of causality where past values of one series ( $x_t$ ) are useful for predicting future values of another series ( $y_t$ ), after past values of  $y_t$  have been controlled for.

**Growth Rate:** The proportionate change in a time series from the previous period. It may be approximated as the difference in logs or reported in percentage form.

## H

**Heckit Method:** An econometric procedure used to correct for sample selection bias due to incidental truncation or some other form of nonrandomly missing data.

**Heterogeneity Bias:** The bias in OLS due to omitted heterogeneity (or omitted variables).

**Heteroskedasticity:** The variance of the error term, given the explanatory variables, is not constant.

**Heteroskedasticity of Unknown Form:** Heteroskedasticity that may depend on the explanatory variables in an unknown, arbitrary fashion.

**Heteroskedasticity-Robust  $F$  Statistic:** An  $F$ -type statistic that is (asymptotically) robust to heteroskedasticity of unknown form.

**Heteroskedasticity-Robust  $LM$  Statistic:** An  $LM$  statistic that is robust to heteroskedasticity of unknown form.

**Heteroskedasticity-Robust Standard Error:** A standard error that is (asymptotically) robust to heteroskedasticity of unknown form.

**Heteroskedasticity-Robust  $t$  Statistic:** A  $t$  statistic that is (asymptotically) robust to heteroskedasticity of unknown form.

**Highly Persistent:** A time series process where outcomes in the distant future are highly correlated with current outcomes.

**Homoskedasticity:** The errors in a regression model have constant variance conditional on the explanatory variables.

**Hypothesis Test:** A statistical test of the null, or maintained, hypothesis against an alternative hypothesis.

## I

**Idempotent Matrix:** A (square) matrix where multiplication of the matrix by itself equals itself.

**Identification:** A population parameter, or set of parameters, can be consistently estimated.

**Identified Equation:** An equation whose parameters can be consistently estimated, especially in models with endogenous explanatory variables.

**Identity Matrix:** A square matrix where all diagonal elements are one and all off-diagonal elements are zero.

**Idiosyncratic Error:** In panel data models, the error that changes over time as well as across units (say, individuals, firms, or cities).

**Impact Elasticity:** In a distributed lag model, the immediate percentage change in the dependent variable given a 1% increase in the independent variable.

**Impact Multiplier:** *See* impact propensity.

**Impact Propensity:** In a distributed lag model, the immediate change in the dependent variable given a one-unit increase in the independent variable.

**Incidental Truncation:** A sample selection problem whereby one variable, usually the dependent variable, is only observed for certain outcomes of another variable.

**Inclusion of an Irrelevant Variable:** The including of an explanatory variable in a regression model that has a zero population parameter in estimating an equation by OLS.

**Inconsistency:** The difference between the probability limit of an estimator and the parameter value.

**Inconsistent:** Describes an estimator that does not converge (in probability) to the correct population parameter as the sample size grows.

**Independent Random Variables:** Random variables whose joint distribution is the product of the marginal distributions.

**Independent Variable:** *See* explanatory variable.

**Independently Pooled Cross Section:** A data set obtained by pooling independent random samples from different points in time.

**Index Number:** A statistic that aggregates information on economic activity, such as production or prices.

**Infinite Distributed Lag (IDL) Model:** A distributed lag model where a change in the explanatory variable can have an impact on the dependent variable into the indefinite future.

**Influential Observations:** *See* outliers.

**Information Set:** In forecasting, the set of variables that we can observe prior to forming our forecast.

**In-Sample Criteria:** Criteria for choosing forecasting models that are based on goodness-of-fit within the sample used to obtain the parameter estimates.

**Instrument:** *See* instrumental variable.

**Instrument Exogeneity:** In instrumental variables estimation, the requirement that an instrumental variable is uncorrelated with the error term.

**Instrument Relevance:** In instrumental variables estimation, the requirement that an instrumental variable helps to partially explain variation in the endogenous explanatory variable.

**Instrumental Variable (IV):** In an equation with an endogenous explanatory variable, an IV is a variable that does not appear in the equation, is uncorrelated with the error in the equation, and is (partially) correlated with the endogenous explanatory variable.

**Instrumental Variables (IV) Estimator:** An estimator in a linear model used when instrumental variables are available for one or more endogenous explanatory variables.

**Integrated of Order One [I(1)]:** A time series process that needs to be first-differenced in order to produce an I(0) process.

**Integrated of Order Zero [I(0)]:** A stationary, weakly dependent time series process that, when used in regression analysis, satisfies the law of large numbers and the central limit theorem.

**Interaction Effect:** In multiple regression, the partial effect of one explanatory variable depends on the value of a different explanatory variable.

**Interaction Term:** An independent variable in a regression model that is the product of two explanatory variables.

**Intercept:** In the equation of a line, the value of the  $y$  variable when the  $x$  variable is zero.

**Intercept Parameter:** The parameter in a multiple linear regression model that gives the expected value of the dependent variable when all the independent variables equal zero.

**Intercept Shift:** The intercept in a regression model differs by group or time period.

**Internet:** A global computer network that can be used to access information and download databases.

**Interval Estimator:** A rule that uses data to obtain lower and upper bounds for a population parameter. (*See also* confidence interval.)

**Inverse:** For an  $n \times n$  matrix, its inverse (if it exists) is the  $n \times n$  matrix for which pre- and post-multiplication by the original matrix yields the identity matrix.

**Inverse Mills Ratio:** A term that can be added to a multiple regression model to remove sample selection bias.

## J

**Joint Distribution:** The probability distribution determining the probabilities of outcomes involving two or more random variables.

**Joint Hypotheses Test:** A test involving more than one restriction on the parameters in a model.

**Jointly Insignificant:** Failure to reject, using an  $F$  test at a specified significance level, that all coefficients for a group of explanatory variables are zero.

**Jointly Statistically Significant:** The null hypothesis that two or more explanatory variables have zero population coefficients is rejected at the chosen significance level.

**Just Identified Equation:** For models with endogenous explanatory variables, an equation that is identified but would not be identified with one fewer instrumental variable.

## K

**Kurtosis:** A measure of the thickness of the tails of a distribution based on the fourth moment of the standardized random variable; the measure is usually compared to the value for the standard normal distribution, which is three.

## L

**Lag Distribution:** In a finite or infinite distributed lag model, the lag coefficients graphed as a function of the lag length.

**Lagged Dependent Variable:** An explanatory variable that is equal to the dependent variable from an earlier time period.

**Lagged Endogenous Variable:** In a simultaneous equations model, a lagged value of one of the endogenous variables.

**Lagrange Multiplier (LM) Statistic:** A test statistic with large-sample justification that can be used to test for omitted variables, heteroskedasticity, and serial correlation, among other model specification problems.

**Large Sample Properties:** *See* asymptotic properties.

**Latent Variable Model:** A model where the observed dependent variable is assumed to be a function of an underlying latent, or unobserved, variable.

**Law of Iterated Expectations:** A result from probability that relates unconditional and conditional expectations.

**Law of Large Numbers (LLN):** A theorem that says that the average from a random sample converges in probability to the population average; the LLN also holds for stationary and weakly dependent time series.

**Leads and Lags Estimator:** An estimator of a cointegrating parameter in a regression with I(1) variables, where the current, some past, and some future first differences in the explanatory variable are included as regressors.

**Least Absolute Deviations (LAD):** A method for estimating the parameters of a multiple regression model based on minimizing the sum of the absolute values of the residuals.

**Least Squares Estimator:** An estimator that minimizes a sum of squared residuals.

**Level-Level Model:** A regression model where the dependent variable and the independent variables are in level (or original) form.

**Level-Log Model:** A regression model where the dependent variable is in level form and (at least some of) the independent variables are in logarithmic form.

**Likelihood Ratio Statistic:** A statistic that can be used to test single or multiple hypotheses when the constrained and unconstrained models have been estimated by maximum likelihood. The statistic is twice the difference in the unconstrained and constrained log-likelihoods.

**Limited Dependent Variable (LDV):** A dependent or response variable whose range is restricted in some important way.

**Linear Function:** A function where the change in the dependent variable, given a one-unit change in an independent variable, is constant.

**Linear Probability Model (LPM):** A binary response model where the response probability is linear in its parameters.

**Linear Time Trend:** A trend that is a linear function of time.

**Linear Unbiased Estimator:** In multiple regression analysis, an unbiased estimator that is a linear function of the outcomes on the dependent variable.

**Linearly Independent Vectors:** A set of vectors such that no vector can be written as a linear combination of the others in the set.

**Log Function:** A mathematical function, defined only for strictly positive arguments, with a positive but decreasing slope.

**Logarithmic Function:** A mathematical function defined for positive arguments that has a positive, but diminishing, slope.

**Logit Model:** A model for binary response where the response probability is the logit function evaluated at a linear function of the explanatory variables.

**Log-Level Model:** A regression model where the dependent variable is in logarithmic form and the independent variables are in level (or original) form.

**Log-Likelihood Function:** The sum of the log-likelihoods, where the log-likelihood for each observation is the log of the density of the dependent variable given the explanatory variables; the log-likelihood function is viewed as a function of the parameters to be estimated.

**Log-Log Model:** A regression model where the dependent variable and (at least some of) the explanatory variables are in logarithmic form.

**Longitudinal Data:** *See* panel data.

**Long-Run Elasticity:** The long-run propensity in a distributed lag model with the dependent and independent variables in logarithmic form; thus, the long-run elasticity is the eventual percentage increase in the explained variable, given a permanent 1% increase in the explanatory variable.

**Long-Run Multiplier:** *See* long-run propensity.

**Long-Run Propensity (LRP):** In a distributed lag model, the eventual change in the dependent variable given a permanent, one-unit increase in the independent variable.

**Loss Function:** A function that measures the loss when a forecast differs from the actual outcome; the most common examples are absolute value loss and squared loss.

## M

**Marginal Effect:** The effect on the dependent variable that results from changing an independent variable by a small amount.

**Martingale:** A time series process whose expected value, given all past outcomes on the series, simply equals the most recent value.

**Martingale Difference Sequence:** The first difference of a martingale. It is unpredictable (or has a zero mean), given past values of the sequence.

**Matched Pair Sample:** A sample where each observation is matched with another, as in a sample consisting of a husband and wife or a set of two siblings.

**Matrix:** An array of numbers.

**Matrix Multiplication:** An algorithm for multiplying together two conformable matrices.

**Matrix Notation:** A convenient mathematical notation, grounded in matrix algebra, for expressing and manipulating the multiple regression model.

**Maximum Likelihood Estimation (MLE):** A broadly applicable estimation method where the parameter estimates are chosen to maximize the log-likelihood function.

**Maximum Likelihood Estimator:** An estimator that maximizes the (log of the) likelihood function.

**Mean:** *See* expected value.

**Mean Absolute Error (MAE):** A performance measure in forecasting, computed as the average of the absolute values of the forecast errors.

**Mean Independent:** The key requirement in multiple regression analysis, which says the unobserved error has a mean that does not change across subsets of the population defined by different values of the explanatory variables.

**Mean Squared Error (MSE):** The expected squared distance that an estimator is from the population value; it equals the variance plus the square of any bias.

**Measurement Error:** The difference between an observed variable and the variable that belongs in a multiple regression equation.

**Median:** In a probability distribution, it is the value where there is a 50% chance of being below the value and a 50% chance of being above it. In a sample of numbers, it is the middle value after the numbers have been ordered.

**Method of Moments Estimator:** An estimator obtained by using the sample analog of population moments; ordinary least squares and two stage least squares are both method of moments estimators.

**Micronumerosity:** A term introduced by Arthur Goldberger to describe properties of econometric estimators with small sample sizes.

**Minimum Variance Unbiased Estimator:** An estimator with the smallest variance in the class of all unbiased estimators.

**Missing at Random:** In multiple regression analysis, a missing data mechanism where the reason data are missing may be correlated with the explanatory variables but is independent of the error term.

**Missing Completely at Random (MCAR):** In multiple regression analysis, a missing data mechanism where the reason data are missing is statistically independent of the values of the explanatory variables as well as the unobserved error.

**Missing Data:** A data problem that occurs when we do not observe values on some variables for certain observations (individuals, cities, time periods, and so on) in the sample.

**Misspecification Analysis:** The process of determining likely biases that can arise from omitted variables, measurement error, simultaneity, and other kinds of model misspecification.

**Moving Average Process of Order One [MA(1)]:** A time series process generated as a linear function of the current value and one lagged value of a zero-mean, constant variance, uncorrelated stochastic process.

**Multicollinearity:** A term that refers to correlation among the independent variables in a multiple regression model; it is usually invoked when some correlations are “large,” but an actual magnitude is not well defined.

**Multiple Hypotheses Test:** A test of a null hypothesis involving more than one restriction on the parameters.

**Multiple Linear Regression (MLR) Model:** A model linear in its parameters, where the dependent variable is a function of independent variables plus an error term.

**Multiple Regression Analysis:** A type of analysis that is used to describe estimation of and inference in the multiple linear regression model.

**Multiple Restrictions:** More than one restriction on the parameters in an econometric model.

**Multiple-Step-Ahead Forecast:** A time series forecast of more than one period into the future.

**Multiplicative Measurement Error:** Measurement error where the observed variable is the product of the true unobserved variable and a positive measurement error.

**Multivariate Normal Distribution:** A distribution for multiple random variables where each linear combination of the random variables has a univariate (one-dimensional) normal distribution.

## N

***n*-R-Squared Statistic:** *See* Lagrange multiplier statistic.

**Natural Experiment:** A situation where the economic environment—sometimes summarized by an explanatory variable—exogenously changes, perhaps inadvertently, due to a policy or institutional change.

**Natural Logarithm:** *See* logarithmic function.

**Nominal Variable:** A variable measured in nominal or current dollars.

**Nonexperimental Data:** Data that have not been obtained through a controlled experiment.

**Nonlinear Function:** A function whose slope is not constant.

**Nonnested Models:** Two (or more) models where no model can be written as a special case of the other by imposing restrictions on the parameters.

**Nonrandom Sample:** A sample obtained other than by sampling randomly from the population of interest.

**Nonstationary Process:** A time series process whose joint distributions are not constant across different epochs.

**Normal Distribution:** A probability distribution commonly used in statistics and econometrics for modeling a population. Its probability distribution function has a bell shape.

**Normality Assumption:** The classical linear model assumption which states that the error (or dependent variable)

has a normal distribution, conditional on the explanatory variables.

**Null Hypothesis:** In classical hypothesis testing, we take this hypothesis as true and require the data to provide substantial evidence against it.

**Numerator Degrees of Freedom:** In an *F* test, the number of restrictions being tested.

## O

**Observational Data:** *See* nonexperimental data.

**OLS:** *See* ordinary least squares.

**OLS Intercept Estimate:** The intercept in an OLS regression line.

**OLS Regression Line:** The equation relating the predicted value of the dependent variable to the independent variables, where the parameter estimates have been obtained by OLS.

**OLS Slope Estimate:** A slope in an OLS regression line.

**Omitted Variable Bias:** The bias that arises in the OLS estimators when a relevant variable is omitted from the regression.

**Omitted Variables:** One or more variables, which we would like to control for, have been omitted in estimating a regression model.

**One-Sided Alternative:** An alternative hypothesis that states that the parameter is greater than (or less than) the value hypothesized under the null.

**One-Step-Ahead Forecast:** A time series forecast one period into the future.

**One-Tailed Test:** A hypothesis test against a one-sided alternative.

**Online Databases:** Databases that can be accessed via a computer network.

**Online Search Services:** Computer software that allows the Internet or databases on the Internet to be searched by topic, name, title, or keywords.

**Order Condition:** A necessary condition for identifying the parameters in a model with one or more endogenous explanatory variables: the total number of exogenous variables must be at least as great as the total number of explanatory variables.

**Ordinal Variable:** A variable where the ordering of the values conveys information but the magnitude of the values does not.

**Ordinary Least Squares (OLS):** A method for estimating the parameters of a multiple linear regression model. The ordinary least squares estimates are obtained by minimizing the sum of squared residuals.

**Outliers:** Observations in a data set that are substantially different from the bulk of the data, perhaps because of errors or because some data are generated by a different model than most of the other data.

**Out-of-Sample Criteria:** Criteria used for choosing forecasting models which are based on a part of the sample that was not used in obtaining parameter estimates.

**Over Controlling:** In a multiple regression model, including explanatory variables that should not be held fixed when studying the *ceteris paribus* effect of one or more other explanatory variables; this can occur when variables that are themselves outcomes of an intervention or a policy are included among the regressors.

**Overall Significance of a Regression:** A test of the joint significance of all explanatory variables appearing in a multiple regression equation.

**Overdispersion:** In modeling a count variable, the variance is larger than the mean.

**Overidentified Equation:** In models with endogenous explanatory variables, an equation where the number of instrumental variables is strictly greater than the number of endogenous explanatory variables.

**Overidentifying Restrictions:** The extra moment conditions that come from having more instrumental variables than endogenous explanatory variables in a linear model.

**Overspecifying a Model:** See inclusion of an irrelevant variable.

## P

***p*-Value:** The smallest significance level at which the null hypothesis can be rejected. Equivalently, the largest significance level at which the null hypothesis cannot be rejected.

**Pairwise Uncorrelated Random Variables:** A set of two or more random variables where each pair is uncorrelated.

**Panel Data:** A data set constructed from repeated cross sections over time. With a *balanced* panel, the same units appear in each time period. With an *unbalanced* panel, some units do not appear in each time period, often due to attrition.

**Parameter:** An unknown value that describes a population relationship.

**Parsimonious Model:** A model with as few parameters as possible for capturing any desired features.

**Partial Derivative:** For a smooth function of more than one variable, the slope of the function in one direction.

**Partial Effect:** The effect of an explanatory variable on the dependent variable, holding other factors in the regression model fixed.

**Partial Effect at the Average (PEA):** In models with non-constant partial effects, the partial effect evaluated at the average values of the explanatory variables.

**Percent Correctly Predicted:** In a binary response model, the percentage of times the prediction of zero or one coincides with the actual outcome.

**Percentage Change:** The proportionate change in a variable, multiplied by 100.

**Percentage Point Change:** The change in a variable that is measured as a percentage.

**Perfect Collinearity:** In multiple regression, one independent variable is an exact linear function of one or more other independent variables.

**Plug-In Solution to the Omitted Variables Problem:** A proxy variable is substituted for an unobserved omitted variable in an OLS regression.

**Point Forecast:** The forecasted value of a future outcome.

**Poisson Distribution:** A probability distribution for count variables.

**Poisson Regression Model:** A model for a count dependent variable where the dependent variable, conditional on the explanatory variables, is nominally assumed to have a Poisson distribution.

**Policy Analysis:** An empirical analysis that uses econometric methods to evaluate the effects of a certain policy.

**Pooled Cross Section:** A data configuration where independent cross sections, usually collected at different points in time, are combined to produce a single data set.

**Pooled OLS Estimation:** OLS estimation with independently pooled cross sections, panel data, or cluster samples, where the observations are pooled across time (or group) as well as across the cross-sectional units.

**Population:** A well-defined group (of people, firms, cities, and so on) that is the focus of a statistical or econometric analysis.

**Population Model:** A model, especially a multiple linear regression model, that describes a population.

**Population *R*-Squared:** In the population, the fraction of the variation in the dependent variable that is explained by the explanatory variables.

**Population Regression Function:** See conditional expectation.

**Positive Definite:** A symmetric matrix such that all quadratic forms, except the trivial one that must be zero, are strictly positive.

**Positive Semi-Definite:** A symmetric matrix such that all quadratic forms are nonnegative.

**Power of a Test:** The probability of rejecting the null hypothesis when it is false; the power depends on the values of the population parameters under the alternative.

**Practical Significance:** The practical or economic importance of an estimate, which is measured by its sign and magnitude, as opposed to its statistical significance.

**Prais-Winsten (PW) Estimation:** A method of estimating a multiple linear regression model with AR(1) errors and strictly exogenous explanatory variables; unlike Cochrane-Orcutt, Prais-Winsten uses the equation for the first time period in estimation.

**Predetermined Variable:** In a simultaneous equations model, either a lagged endogenous variable or a lagged exogenous variable.

**Predicted Variable:** See dependent variable.

**Prediction:** The estimate of an outcome obtained by plugging specific values of the explanatory variables into an estimated model, usually a multiple regression model.

**Prediction Error:** The difference between the actual outcome and a prediction of that outcome.

**Prediction Interval:** A confidence interval for an unknown outcome on a dependent variable in a multiple regression model.

**Predictor Variable:** *See* explanatory variable.

**Probability Density Function (pdf):** A function that, for discrete random variables, gives the probability that the random variable takes on each value; for continuous random variables, the area under the pdf gives the probability of various events.

**Probability Limit:** The value to which an estimator converges as the sample size grows without bound.

**Probit Model:** A model for binary responses where the response probability is the standard normal cdf evaluated at a linear function of the explanatory variables.

**Program Evaluation:** An analysis of a particular private or public program using econometric methods to obtain the causal effect of the program.

**Proportionate Change:** The change in a variable relative to its initial value; mathematically, the change divided by the initial value.

**Proxy Variable:** An observed variable that is related but not identical to an unobserved explanatory variable in multiple regression analysis.

**Pseudo *R*-Squared:** Any number of goodness-of-fit measures for limited dependent variable models.

## Q

**Quadratic Form:** A mathematical function where the vector argument both pre- and post-multiplies a square, symmetric matrix.

**Quadratic Functions:** Functions that contain squares of one or more explanatory variables; they capture diminishing or increasing effects on the dependent variable.

**Qualitative Variable:** A variable describing a nonquantitative feature of an individual, a firm, a city, and so on.

**Quasi-Demeaned Data:** In random effects estimation for panel data, it is the original data in each time period minus a fraction of the time average; these calculations are done for each cross-sectional observation.

**Quasi-Differenced Data:** In estimating a regression model with AR(1) serial correlation, it is the difference between the current time period and a multiple of the previous time period, where the multiple is the parameter in the AR(1) model.

**Quasi-Experiment:** *See* natural experiment.

**Quasi-Likelihood Ratio Statistic:** A modification of the likelihood ratio statistic that accounts for possible distributional misspecification, as in a Poisson regression model.

**Quasi-Maximum Likelihood Estimation (QMLE):** Maximum likelihood estimation where the log-likelihood function may not correspond to the actual conditional distribution of the dependent variable.

## R

***R*-Bar Squared:** *See* adjusted *R*-squared.

***R*-Squared:** In a multiple regression model, the proportion of the total sample variation in the dependent variable that is explained by the independent variable.

***R*-Squared Form of the *F* Statistic:** The *F* statistic for testing exclusion restrictions expressed in terms of the *R*-squareds from the restricted and unrestricted models.

**Random Coefficient (Slope) Model:** A multiple regression model where the slope parameters are allowed to depend on unobserved unit-specific variables.

**Random Effects Estimator:** A feasible GLS estimator in the unobserved effects model where the unobserved effect is assumed to be uncorrelated with the explanatory variables in each time period.

**Random Effects Model:** The unobserved effects panel data model where the unobserved effect is assumed to be uncorrelated with the explanatory variables in each time period.

**Random Sample:** A sample obtained by sampling randomly from the specified population.

**Random Sampling:** A sampling scheme whereby each observation is drawn at random from the population. In particular, no unit is more likely to be selected than any other unit, and each draw is independent of all other draws.

**Random Variable:** A variable whose outcome is uncertain.

**Random Vector:** A vector consisting of random variables.

**Random Walk:** A time series process where next period's value is obtained as this period's value, plus an independent (or at least an uncorrelated) error term.

**Random Walk with Drift:** A random walk that has a constant (or drift) added in each period.

**Rank Condition:** A sufficient condition for identification of a model with one or more endogenous explanatory variables.

**Rank of a Matrix:** The number of linearly independent columns in a matrix.

**Rational Distributed Lag (RDL) Model:** A type of infinite distributed lag model where the lag distribution depends on relatively few parameters.

**Real Variable:** A monetary value measured in terms of a base period.

**Reduced Form Equation:** A linear equation where an endogenous variable is a function of exogenous variables and unobserved errors.

**Reduced Form Error:** The error term appearing in a reduced form equation.

**Reduced Form Parameters:** The parameters appearing in a reduced form equation.

**Regressand:** *See* dependent variable.

**Regression Specification Error Test (RESET):** A general test for functional form in a multiple regression model; it is an *F* test of joint significance of the squares, cubes, and perhaps higher powers of the fitted values from the initial OLS estimation.

**Regression through the Origin:** Regression analysis where the intercept is set to zero; the slopes are obtained by minimizing the sum of squared residuals, as usual.

**Regressor:** *See* explanatory variable.

**Rejection Region:** The set of values of a test statistic that leads to rejecting the null hypothesis.

**Rejection Rule:** In hypothesis testing, the rule that determines when the null hypothesis is rejected in favor of the alternative hypothesis.

**Relative Change:** *See* proportionate change.

**Resampling Method:** A technique for approximating standard errors (and distributions of test statistics) whereby a series of samples are obtained from the original data set and estimates are computed for each subsample.

**Residual:** The difference between the actual value and the fitted (or predicted) value; there is a residual for each observation in the sample used to obtain an OLS regression line.

**Residual Analysis:** A type of analysis that studies the sign and size of residuals for particular observations after a multiple regression model has been estimated.

**Residual Sum of Squares:** *See* sum of squared residuals.

**Response Probability:** In a binary response model, the probability that the dependent variable takes on the value one, conditional on explanatory variables.

**Response Variable:** *See* dependent variable.

**Restricted Model:** In hypothesis testing, the model obtained after imposing all of the restrictions required under the null.

**Retrospective Data:** Data collected based on past, rather than current, information.

**Root Mean Squared Error (RMSE):** Another name for the standard error of the regression in multiple regression analysis.

**Row Vector:** A vector of numbers arranged as a row.

## S

**Sample Average:** The sum of  $n$  numbers divided by  $n$ ; a measure of central tendency.

**Sample Correlation:** For outcomes on two random variables, the sample covariance divided by the product of the sample standard deviations.

**Sample Correlation Coefficient:** An estimate of the (population) correlation coefficient from a sample of data.

**Sample Covariance:** An unbiased estimator of the population covariance between two random variables.

**Sample Regression Function (SRF):** *See* OLS regression line.

**Sample Selection Bias:** Bias in the OLS estimator which is induced by using data that arise from endogenous sample selection.

**Sample Standard Deviation:** A consistent estimator of the population standard deviation.

**Sample Variance:** An unbiased, consistent estimator of the population variance.

**Sampling Distribution:** The probability distribution of an estimator over all possible sample outcomes.

**Sampling Standard Deviation:** The standard deviation of an estimator, that is, the standard deviation of a sampling distribution.

**Sampling Variance:** The variance in the sampling distribution of an estimator; it measures the spread in the sampling distribution.

**Scalar Multiplication:** The algorithm for multiplying a scalar (number) by a vector or matrix.

**Scalar Variance-Covariance Matrix:** A variance-covariance matrix where all off-diagonal terms are zero and the diagonal terms are the same positive constant.

**Score Statistic:** *See* Lagrange multiplier statistic.

**Seasonal Dummy Variables:** A set of dummy variables used to denote the quarters or months of the year.

**Seasonality:** A feature of monthly or quarterly time series where the average value differs systematically by season of the year.

**Seasonally Adjusted:** Monthly or quarterly time series data where some statistical procedure—possibly regression on seasonal dummy variables—has been used to remove the seasonal component.

**Selected Sample:** A sample of data obtained not by random sampling but by selecting on the basis of some observed or unobserved characteristic.

**Self-Selection:** Deciding on an action based on the likely benefits, or costs, of taking that action.

**Semi-Elasticity:** The percentage change in the dependent variable given a one-unit increase in an independent variable.

**Sensitivity Analysis:** The process of checking whether the estimated effects and statistical significance of key explanatory variables are sensitive to inclusion of other explanatory variables, functional form, dropping of potentially outlying observations, or different methods of estimation.

**Sequentially Exogenous:** A feature of an explanatory variable in time series (or panel data) models where the error term in the current time period has a zero mean conditional on all current and past explanatory variables; a weaker version is stated in terms of zero correlations.

**Serial Correlation:** In a time series or panel data model, correlation between the errors in different time periods.

**Serial Correlation-Robust Standard Error:** A standard error for an estimator that is (asymptotically) valid whether or not the errors in the model are serially correlated.

**Serially Uncorrelated:** The errors in a time series or panel data model are pairwise uncorrelated across time.

**Short-Run Elasticity:** The impact propensity in a distributed lag model when the dependent and independent variables are in logarithmic form.

**Significance Level:** The probability of a Type I error in hypothesis testing.

- Simple Linear Regression Model:** A model where the dependent variable is a linear function of a single independent variable, plus an error term.
- Simultaneity:** A term that means at least one explanatory variable in a multiple linear regression model is determined jointly with the dependent variable.
- Simultaneity Bias:** The bias that arises from using OLS to estimate an equation in a simultaneous equations model.
- Simultaneous Equations Model (SEM):** A model that jointly determines two or more endogenous variables, where each endogenous variable can be a function of other endogenous variables as well as of exogenous variables and an error term.
- Skewness:** A measure of how far a distribution is from being symmetric, based on the third moment of the standardized random variable.
- Slope:** In the equation of a line, the change in the  $y$  variable when the  $x$  variable increases by one.
- Slope Parameter:** The coefficient on an independent variable in a multiple regression model.
- Smearing Estimate:** A retransformation method particularly useful for predicting the level of a response variable when a linear model has been estimated for the natural log of the response variable.
- Spreadsheet:** Computer software used for entering and manipulating data.
- Spurious Correlation:** A correlation between two variables that is not due to causality, but perhaps to the dependence of the two variables on another unobserved factor.
- Spurious Regression Problem:** A problem that arises when regression analysis indicates a relationship between two or more unrelated time series processes simply because each has a trend, is an integrated time series (such as a random walk), or both.
- Square Matrix:** A matrix with the same number of rows as columns.
- Stable AR(1) Process:** An AR(1) process where the parameter on the lag is less than one in absolute value. The correlation between two random variables in the sequence declines to zero at a geometric rate as the distance between the random variables increases, and so a stable AR(1) process is weakly dependent.
- Standard Deviation:** A common measure of spread in the distribution of a random variable.
- Standard Deviation of  $\hat{\beta}_j$ :** A common measure of spread in the sampling distribution of  $\hat{\beta}_j$ .
- Standard Error:** Generically, an estimate of the standard deviation of an estimator.
- Standard Error of  $\hat{\beta}_j$ :** An estimate of the standard deviation in the sampling distribution of  $\hat{\beta}_j$ .
- Standard Error of the Estimate:** *See* standard error of the regression.
- Standard Error of the Regression (SER):** In multiple regression analysis, the estimate of the standard deviation of the population error, obtained as the square root of the sum of squared residuals over the degrees of freedom.
- Standard Normal Distribution:** The normal distribution with mean zero and variance one.
- Standardized Coefficients:** Regression coefficients that measure the standard deviation change in the dependent variable given a one standard deviation increase in an independent variable.
- Standardized Random Variable:** A random variable transformed by subtracting off its expected value and dividing the result by its standard deviation; the new random variable has mean zero and standard deviation one.
- Static Model:** A time series model where only contemporaneous explanatory variables affect the dependent variable.
- Stationary Process:** A time series process where the marginal and all joint distributions are invariant across time.
- Statistical Inference:** The act of testing hypotheses about population parameters.
- Statistical Significance:** The importance of an estimate as measured by the size of a test statistic, usually a  $t$  statistic.
- Statistically Different from Zero:** *See* statistically significant.
- Statistically Insignificant:** Failure to reject the null hypothesis that a population parameter is equal to zero, at the chosen significance level.
- Statistically Significant:** Rejecting the null hypothesis that a parameter is equal to zero against the specified alternative, at the chosen significance level.
- Stochastic Process:** A sequence of random variables indexed by time.
- Stratified Sampling:** A nonrandom sampling scheme whereby the population is first divided into several non-overlapping, exhaustive strata, and then random samples are taken from within each stratum.
- Strict Exogeneity:** An assumption that holds in a time series or panel data model when the explanatory variables are strictly exogenous.
- Strictly Exogenous:** A feature of explanatory variables in a time series or panel data model where the error term at any time period has zero expectation, conditional on the explanatory variables in all time periods; a less restrictive version is stated in terms of zero correlations.
- Strongly Dependent:** *See* highly persistent.
- Structural Equation:** An equation derived from economic theory or from less formal economic reasoning.
- Structural Error:** The error term in a structural equation, which could be one equation in a simultaneous equations model.
- Structural Parameters:** The parameters appearing in a structural equation.
- Studentized Residuals:** The residuals computed by excluding each observation, in turn, from the estimation, divided by the estimated standard deviation of the error.

**Sum of Squared Residuals (SSR):** In multiple regression analysis, the sum of the squared OLS residuals across all observations.

**Summation Operator:** A notation, denoted by  $\Sigma$ , used to define the summing of a set of numbers.

**Symmetric Distribution:** A probability distribution characterized by a probability density function that is symmetric around its median value, which must also be the mean value (whenever the mean exists).

**Symmetric Matrix:** A (square) matrix that equals its transpose.

## T

***t* Distribution:** The distribution of the ratio of a standard normal random variable and the square root of an independent chi-square random variable, where the chi-square random variable is first divided by its *df*.

***t* Ratio:** See *t* statistic.

***t* Statistic:** The statistic used to test a single hypothesis about the parameters in an econometric model.

**Test Statistic:** A rule used for testing hypotheses where each sample outcome produces a numerical value.

**Text Editor:** Computer software that can be used to edit text files.

**Text (ASCII) File:** A universal file format that can be transported across numerous computer platforms.

**Time-Demeaned Data:** Panel data where, for each cross-sectional unit, the average over time is subtracted from the data in each time period.

**Time Series Data:** Data collected over time on one or more variables.

**Time Series Process:** See stochastic process.

**Time Trend:** A function of time that is the expected value of a trending time series process.

**Tobit Model:** A model for a dependent variable that takes on the value zero with positive probability but is roughly continuously distributed over strictly positive values. (See also corner solution response.)

**Top Coding:** A form of data censoring where the value of a variable is not reported when it is above a given threshold; we only know that it is at least as large as the threshold.

**Total Sum of Squares (SST):** The total sample variation in a dependent variable about its sample average.

**Trace of a Matrix:** For a square matrix, the sum of its diagonal elements.

**Transpose:** For any matrix, the new matrix obtained by interchanging its rows and columns.

**Treatment Group:** In program evaluation, the group that participates in the program.

**Trending Process:** A time series process whose expected value is an increasing or a decreasing function of time.

**Trend-Stationary Process:** A process that is stationary once a time trend has been removed; it is usually implicit that the detrended series is weakly dependent.

**True Model:** The actual population model relating the dependent variable to the relevant independent variables, plus a disturbance, where the zero conditional mean assumption holds.

**Truncated Normal Regression Model:** The special case of the truncated regression model where the underlying population model satisfies the classical linear model assumptions.

**Truncated Regression Model:** A linear regression model for cross-sectional data in which the sampling scheme entirely excludes, on the basis of outcomes on the dependent variable, part of the population.

**Two-Sided Alternative:** An alternative where the population parameter can be either less than or greater than the value stated under the null hypothesis.

**Two Stage Least Squares (2SLS) Estimator:** An instrumental variables estimator where the IV for an endogenous explanatory variable is obtained as the fitted value from regressing the endogenous explanatory variable on all exogenous variables.

**Two-Tailed Test:** A test against a two-sided alternative.

**Type I Error:** A rejection of the null hypothesis when it is true.

**Type II Error:** The failure to reject the null hypothesis when it is false.

## U

**Unbalanced Panel:** A panel data set where certain years (or periods) of data are missing for some cross-sectional units.

**Unbiased Estimator:** An estimator whose expected value (or mean of its sampling distribution) equals the population value (regardless of the population value).

**Uncentered *R*-squared:** The *R*-squared computed without subtracting the sample average of the dependent variable when obtaining the total sum of squares (SST).

**Unconditional Forecast:** A forecast that does not rely on knowing, or assuming values for, future explanatory variables.

**Uncorrelated Random Variables:** Random variables that are not linearly related.

**Underspecifying a Model:** See excluding a relevant variable.

**Unidentified Equation:** An equation with one or more endogenous explanatory variables where sufficient instrumental variables do not exist to identify the parameters.

**Unit Root Process:** A highly persistent time series process where the current value equals last period's value, plus a weakly dependent disturbance.

**Unobserved Effect:** In a panel data model, an unobserved variable in the error term that does not change over time. For cluster samples, an unobserved variable that is common to all units in the cluster.

**Unobserved Effects Model:** A model for panel data or cluster samples where the error term contains an unobserved effect.

**Unobserved Heterogeneity:** *See* unobserved effect.

**Unrestricted Model:** In hypothesis testing, the model that has no restrictions placed on its parameters.

**Upward Bias:** The expected value of an estimator is greater than the population parameter value.

## V

**Variance:** A measure of spread in the distribution of a random variable.

**Variance-Covariance Matrix:** For a random vector, the positive semi-definite matrix defined by putting the variances down the diagonal and the covariances in the appropriate off-diagonal entries.

**Variance-Covariance Matrix of the OLS Estimator:** The matrix of sampling variances and covariances for the vector of OLS coefficients.

**Variance Inflation Factor:** In multiple regression analysis under the Gauss-Markov assumptions, the term in the sampling variance affected by correlation among the explanatory variables.

**Variance of the Prediction Error:** The variance in the error that arises when predicting a future value of the dependent variable based on an estimated multiple regression equation.

**Vector Autoregressive (VAR) Model:** A model for two or more time series where each variable is modeled as a linear function of past values of all variables, plus disturbances that have zero means given all past values of the observed variables.

## W

**Wald Statistic:** A general test statistic for testing hypotheses in a variety of econometric settings; typically, the Wald statistic has an asymptotic chi-square distribution.

**Weak Instruments:** Instrumental variables that are only slightly correlated with the relevant endogenous explanatory variable or variables.

**Weakly Dependent:** A term that describes a time series process where some measure of dependence between random variables at two points in time—such as correlation—diminishes as the interval between the two points in time increases.

**Weighted Least Squares (WLS) Estimator:** An estimator used to adjust for a known form of heteroskedasticity, where each squared residual is weighted by the inverse of the (estimated) variance of the error.

**White Test:** A test for heteroskedasticity that involves regressing the squared OLS residuals on the OLS fitted values and on the squares of the fitted values; in its most general form, the squared OLS residuals are regressed on the explanatory variables, the squares of the explanatory variables, and all the nonredundant interactions of the explanatory variables.

**Within Estimator:** *See* fixed effects estimator.

**Within Transformation:** *See* fixed effects transformation.

## Y

**Year Dummy Variables:** For data sets with a time series component, dummy (binary) variables equal to one in the relevant year and zero in all other years.

## Z

**Zero Conditional Mean Assumption:** A key assumption used in multiple regression analysis that states that, given any values of the explanatory variables, the expected value of the error equals zero. (*See* Assumptions MLR.4, TS.3, and TS.3' in the text.)

**Zero Matrix:** A matrix where all entries are zero.

**Zero-One Variable:** *See* dummy variable.

## Numbers

2SLS. *See* two stage least squares

401(k) plans

asymptotic normality, 155–156

comparison of simple and multiple regression estimates, 70

statistical vs. practical significance, 121

WLS estimation, 259

## A

ability and wage

causality, 12

excluding ability from model, 78–83

IV for ability, 481

mean independence, 23

proxy variable for ability, 279–285

adaptive expectations, 353, 355

adjusted *R*-squareds, 181–184, 374

AFDC participation, 231

age

financial wealth and, 257–259, 263

smoking and, 261–262

aggregate consumption function, 511–514

air pollution and housing prices

beta coefficients, 175–176

logarithmic forms, 171–173

quadratic functions, 175–177

*t* test, 118

alcohol drinking, 230

alternative hypotheses

defined, 694

one-sided, 110–114, 695

two-sided, 114–115, 695

antidumping filings and chemical imports

AR(3) serial correlation, 381

dummy variables, 327–328

forecasting, 596, 597, 598

PW estimation, 384

seasonality, 336–338

apples, ecolabeled, 180–181

AR(1) models, consistency example, 350–351

testing for, after 2SLS estimation, 486

AR(1) serial correlation

correcting for, 381–387

testing for, 376–381

AR(2) models

EMH example, 352

forecasting example, 352, 397

ARCH model, 393–394

AR(*q*) serial correlation

correcting for, 386–387

testing for, 379–380

arrests

asymptotic normality, 155–156

average sentence length and, 249

goodness-of-fit, 72

heteroskedasticity-robust *LM* statistic, 249

linear probability model, 227–228

normality assumption and, 107

Poisson regression, 545–546

ASCII files, 609

assumptions

classical linear model (CLM), 106

establishing unbiasedness of OLS, 73–77, 317–320

homoskedasticity, 45–48, 82–83, 89, 363

matrix notation, 723–726

for multiple linear regressions, 73–77, 82, 89, 152

normality, 105–108, 322

for simple linear regressions, 40–45, 45–48

for time series regressions, 317–323, 348–354, 363

zero mean and zero correlation, 152

asymptotically uncorrelated sequences, 346–348

asymptotic bias, deriving, 153–154

asymptotic confidence interval, 157

asymptotic efficiency of OLS, 161–162

asymptotic normality of estimators, in general, 683–684

asymptotic normality of OLS

for multiple linear regressions, 156–158

for time series regressions, 351–354

asymptotic properties. *See* large sample properties

asymptotics, OLS. *See* OLS asymptotics  
 asymptotic sample properties of estimators, 681–684  
 asymptotic standard errors, 157  
 asymptotic  $t$  statistics, 157  
 asymptotic variance, 156  
 atime series data  
   applying 2SLS to, 485–486  
 attenuation bias, 291, 292  
 attrition, 441  
 augmented Dickey-Fuller test, 576  
 autocorrelation, 320–322. *See also* serial correlation  
 autoregressive conditional heteroskedacity (ARCH) model,  
   393–394  
 autoregressive model of order two [AR(2)]. *See* AR(2)  
   models  
 autoregressive process of order one [AR(1)], 347  
 auxiliary regression, 159  
 average, using summation operator, 629  
 average marginal effect (AME), 286, 532  
 average partial effect (APE), 286, 532, 540  
 average treatment effect, 410

## B

balanced panel, 420  
 baseball players' salaries  
   nonnested models, 183  
   testing exclusion restrictions, 127–132  
 base group, 208  
 base period  
   and value, 326  
 base value, 326  
 beer  
   price and demand, 185–186  
   taxes and traffic fatalities, 184  
 benchmark group, 208  
 Bernoulli random variables, 646–647  
 best linear unbiased estimator (BLUE), 89  
 beta coefficients, 169–170  
 between estimators, 435  
 bias  
   attenuation, 291, 292  
   heterogeneity, 413  
   omitted variable, 78–83  
   simultaneity, in OLS, 503–504  
 biased estimators, 677–678  
 biased toward zero, 80  
 binary random variable, 646  
 binary response models. *See* logit and probit models  
 binary variables. *See also* qualitative information  
   defined, 206  
   random, 646–647  
 binomial distribution, 651  
 birth weight  
   AFDC participation and, 231  
   asymptotic standard error, 158  
   data scaling, 166–168  
    $F$  statistic, 133–134  
   IV estimation, 470–471  
 bivariate linear regression model. *See* simple regression  
   model  
 BLUE (best linear unbiased estimator), 89  
 bootstrap standard error, 204  
 Breusch-Godfrey test, 381  
 Breusch-Pagan test  
   for heteroskedasticity, 251

## C

calculus, differential, 640–642  
 campus crimes,  $t$  test, 116–117  
 causality, 10–14  
 cdf (cumulative distribution functions), 648–649  
 censored regression models, 547–552  
 Center for Research in Security Prices (CRSP), 608  
 central limit theorem, 684  
 CEO salaries  
   in multiple regressions  
     motivation for multiple regression, 63–64  
     nonnested models, 183–184  
     predicting, 192, 193–194  
     writing in population form, 74  
   returns on equity and  
     fitted values and residuals, 32  
     goodness-of-fit, 35  
     OLS Estimates, 29–30  
   sales and, constant elasticity model, 39  
 ceteris paribus, 10–14, 66, 67–68  
 chemical firms, nonnested models, 183  
 chemical imports. *See* antidumping filings and chemical  
   imports  
 chi-square distribution  
   critical values table, 749  
   discussions, 669, 717  
 Chow tests  
   differences across groups, 223–224  
   heteroskedasticity and, 247–248  
   for panel data, 423–424  
   for structural change across time, 406  
 cigarettes. *See* smoking  
 city crimes. *See also* crimes  
   law enforcement and, 13  
   panel data, 9–10  
 classical errors-in-variables (CEV), 290  
 classical linear model (CLM) assumptions, 106  
 clear-up rate, distributed lag estimation, 416–417

- clusters, 449–450
  - effect, 449
  - sample, 449
- Cochrane-Orcutt (CO) estimation, 383, 391, 395
- coefficient of determination. *See* *R*-squareds
- cointegration, 580–584
- college admission, omitting unobservables, 285
- college GPA
  - beta coefficients, 169–170
  - fitted values and intercept, 68
  - gender and, 221–224
  - goodness-of-fit, 71
  - heteroskedasticity-robust *F* statistic, 247–248
  - interaction effect, 178–179
  - interpreting equations, 66
  - with measurement error, 292
  - partial effect, 67
  - population regression function, 23
  - predicted, 187–188, 189
  - with single dummy variable, 209–210
  - t* test, 115
- college proximity, as IV for education, 473–474
- colleges, junior vs. four-year, 124–127
- collinearity, perfect, 74–76
- column vectors, 709
- commute time and freeway width, 702–703
- compact discs, demand for, 732
- complete cases estimator, 293
- composite error, 413
  - term, 441
- Compustat, 608
- computer ownership
  - college GPA and, 209–210
  - determinants of, 267
- computers, grants to buy
  - reducing error variance, 185–186
  - R*-squared size, 180–181
- computer usage and wages
  - with interacting terms, 218
  - proxy variable in, 282–283
- conceptual framework, 615
- conditional distributions
  - features, 652–658
  - overview, 649, 651–653
- conditional expectations, 661–665
- conditional forecasts, 587
- conditional median, 300–302
- conditional variances, 665
- confidence intervals
  - 95%, rule of thumb for, 691
  - asymptotic, 157
  - asymptotic, for nonnormal populations, 692–693
  - hypothesis testing and, 701–702
  - interval estimation and, 687–693
  - main discussions, 122–123, 687–688
  - for mean from normally distributed population, 689–691
  - for predictions, 186–189
- consistency of estimators, in general, 681–683
- consistency of OLS
  - in multiple regressions, 150–154
  - sampling selection and, 553–554
  - in time series regressions, 348–351, 372
- consistent tests, 703
- constant dollars, 326
- constant elasticity model, 39, 75, 638
- constant terms, 21
- consumer price index (CPI), 323
- consumption. *See under* family income
- contemporaneously exogenous variables, 318
- continuous random variables, 648–649
- control group, 210
- control variables, 21. *See also* independent variables
- corner solution response, 525
- corrected *R*-squareds, 181–184
- correlated random effects, 445–447
- correlation, 22–23
  - coefficients, 659–660
- count variables, 543–547
- county crimes, multi-year panel data, 422–423
- covariances, 658–659
  - stationary processes, 345–346
- covariates, 21
- CPI (consumer price index), 323
- crimes. *See also* arrests
  - on campuses, *t* test, 116–117
  - in cities, law enforcement and, 13
  - in cities, panel data, 9–10
  - clear-up rate, 416–417
  - in counties, multi-year panel data, 422–423
  - earlier data, use of, 283–284
  - econometric model of, 4–5
  - economic model of, 3, 160, 275–277
  - functional form misspecification, 275–277
  - housing prices and, beta coefficients, 175–176
  - LM statistic, 160
  - prison population and, SEM, 515–516
  - unemployment and, two-period panel data, 412–417
- criminologists, 607
- critical values
  - discussions, 110, 695
  - tables of, 743–749
- crop yields and fertilizers
  - causality, 11, 12
  - simple equation, 21–22
- cross-sectional analysis, 612

cross-sectional data. *See also* panel data; pooled cross sections; regression analysis  
 Gauss-Markov assumptions and, 82, 354  
 main discussion, 5–7  
 time series data vs., 312–313  
 CRSP (Center for Research in Security Prices), 608  
 cumulative areas under standard normal distribution, 743–744  
 cumulative distribution functions (cdf), 648–649  
 cumulative effect, 316  
 current dollars, 326  
 cyclical unemployment, 353

## D

data  
 collection, 608–611  
 economic, types of, 5–12  
 experimental vs. nonexperimental, 2  
 frequency, 7  
 data issues. *See also* misspecification  
 measurement error, 287–292  
 missing data, 293–294  
 multicollinearity, 83–86, 293–294  
 nonrandom samples, 294–295  
 outliers and influential observations, 296–300  
 random slopes, 285–287  
 unobserved explanatory variables, 279–285  
 data mining, 613  
 data scaling, effects on OLS statistics, 166–170  
 Davidson-MacKinnon test, 278  
 deficits. *See* interest rates  
 degrees of freedom (*df*)  
 chi-square distributions with *n*, 669  
 for fixed effects estimator, 436  
 for OLS estimators, 88  
 dependent variables. *See also* regression analysis; *specific event studies*  
 defined, 21  
 measurement error in, 289–292  
 derivatives, 635  
 descriptive statistics, 629  
 deseasonalizing data, 337  
 detrending, 334–335  
 diagonal matrices, 710  
 Dickey-Fuller distribution, 575  
 Dickey-Fuller (DF) test, 575–578  
 augmented, 576  
 difference-in-differences estimator, 408, 410  
 difference in slopes, 218–224  
 difference-stationary processes, 358  
 differencing  
 panel data  
 with more than two periods, 420–425  
 two-period, 412–417  
 serial correlation and, 387–388

differential calculus, 640–642  
 diminishing marginal effects, 635  
 discrete random variables, 646–647  
 disturbance terms, 4, 21, 63  
 disturbance variances, 45  
 downward bias, 80  
 drug usage, 230  
 drunk driving laws and fatalities, 419  
 dummy variables. *See also* qualitative information; year  
 dummy variables  
 defined, 206  
 regression, 438–439  
 trap, 208  
 duration analysis, 549–551  
 Durbin-Watson test, 378–379, 381  
 dynamically complete models, 360–363

## E

Eagle-Granger test, 581–582  
 earnings of veterans, IV estimation, 469  
*EconLit*, 606, 607  
 econometric analysis in projects, 611–614  
 econometric models, 4–5. *See also* econometric models  
 econometrics, 1–2. *See also specific topics*  
 economic growth and government policies, 7  
 economic models, 2–5  
 economic significance. *See* practical significance  
 economic vs. statistical significance, 120–124,  
 702–703  
 economists, types of, 606–607  
 education  
 birth weight and, 133–134  
 fertility and  
 2SLS, 487  
 with discrete dependent variables, 231–232  
 independent cross sections, 404–405  
 gender wage gap and, 405–406  
 IV for, 463, 473–474  
 logarithmic equation, 639  
 return to  
 2SLS, 477  
 differencing, 448  
 fixed effects estimation, 438  
 independent cross sections, 405–406  
 IQ and, 281–282  
 IV estimation, 467–469  
 testing for endogeneity, 482  
 testing overidentifying restrictions, 482  
 wages and. *See under* wages  
 return to education, over time, 405–406  
 smoking and, 261–262  
 women and, 225–227. *See also under* women in  
 labor force

- efficiency
- asymptotic, 161–162
  - of estimators in general, 679–680
  - of OLS with serially correlated errors, 373–374
- efficient markets hypothesis (EMH)
- asymptotic analysis example, 352–353
  - heteroskedasticity and, 393
- elasticity, 39, 637–638
- elections. *See* voting outcomes
- EMH. *See* efficient markets hypothesis (EMH)
- empirical analysis
- data collection, 608–611
  - econometric analysis, 611–614
  - literature review, 607–608
  - posing question, 605–607
  - sample projects, 621–625
  - steps in, 2–5
  - writing paper, 614–621
- employment and unemployment. *See also* wages
- arrests and, 227–228
  - crimes and, 412–417
  - enterprise zones and, 422
  - estimating average rate, 675
  - forecasting, 589, 591–592, 594
  - inflation and. *See under* inflation
  - in Puerto Rico
    - logarithmic form, 323–324
    - time series data, 7–8
  - women and. *See* women in labor force
- endogenous explanatory variables. *See also* instrumental variables; simultaneous equations models; two stage least squares
- defined, 76, 274
  - in logit and probit models, 536
  - sample selection and, 557
  - testing for, 481–482
- endogenous sample selection, 294
- Engle-Granger two-step procedure, 586
- enrollment, *t* test, 116–117
- enterprise zones
- business investments and, 696–697
  - unemployment and, 422
- error correction models, 584–586
- errors-in-variables problem, 479–481, 512
- error terms, 4, 21, 63
- error variances
- adding regressors to reduce, 185–186
  - defined, 45, 83
  - estimating, 48–50
- estimated GLS. *See* feasible GLS
- estimation and estimators. *See also* first differencing; fixed effects; instrumental variables; logit and probit models; OLS (ordinary least squares); random effects; Tobit model
- advantages of multiple over simple regression, 60–64
  - asymptotic sample properties of, 681–684
  - changing independent variables simultaneously, 68
  - defined, 675
  - difference-in-differences, 408–410
  - finite sample properties of, 675–680
  - LAD, 300–302
  - language of, 90–91
  - method of moments approach, 25–26
  - misspecifying models, 78–83
  - sampling distributions of OLS estimators, 105–108
- event studies, 325, 327–328
- Excel, 610
- excluding relevant variables, 78–83
- exclusion restrictions, 127
- for 2SLS, 475
  - general linear, 136–137
  - Lagrange multiplier (LM) statistic, 158–160
  - overall significance of regressions, 135
  - for SEM, 510–511
  - testing, 127–132
- exogenous explanatory variables, 76
- exogenous sample selection, 294, 553
- expectations augmented Phillips curve, 353–354, 377, 378
- expectations hypothesis, 14
- expected values, 652–654, 716
- experience
- wage and
    - causality, 12
    - interpreting equations, 67
    - motivation for multiple regression, 61
    - omitted variable bias, 81
    - partial effect, 642
    - quadratic functions, 173–175, 636
    - women and, 225–227
- experimental data, 2
- experimental group, 210
- experiments, defined, 645
- explained sum of squares (SSE), 34, 70
- explained variables. *See also* dependent variables
- defined, 21
- explanatory variables, 21. *See also* independent variables
- exponential function, 639
- exponential smoothing, 587
- exponential trends, 330–331
- F**
- family income. *See also* savings
- birth weight and
    - asymptotic standard error, 158
    - data scaling, 166–168

- family income (*continued*)
- college GPA and, 292
  - consumption and
    - motivation for multiple regression, 62, 63
    - perfect collinearity and, 75
- farmers and pesticide usage, 185
- F* distribution
- critical values table, 746–748
  - discussions, 670, 671, 717
- FDL (finite distributed lag) models, 314–316, 350, 416–417
- feasible GLS
- with heteroskedasticity and AR(1) serial correlations, 395
  - main discussion, 258–263
  - OLS vs., 385–386
- Federal Bureau of Investigation, 608
- fertility rate
- education and, 487
  - forecasting, 597
  - over time, 404–405
  - tax exemption and
    - with binary variables, 324–325
    - cointegration, 582–583
    - FDL model, 314–316
    - first differences, 363–364
    - serial correlation, 362
    - trends, 333
- fertility studies, with discrete dependent variable, 231–232
- fertilizers
- land quality and, 23
  - soybean yields and
    - causality, 11, 12
    - simple equation, 21–22
- final exam scores
- interaction effect, 178–179
  - skipping classes and, 464–465
- financial wealth
- nonrandom sampling, 294–295
  - WLS estimation, 257–259, 263
- finite distributed lag (FDL) models, 314–316, 350
- finite sample properties
- of estimators, 675–680
  - of OLS in matrix form, 723–726
- firm sales. *See* sales
- first-differenced equations, 414
- first-differenced estimator, 414
- first differencing
- defined, 414
  - fixed effects vs, 439–440
  - I(1) time series and, 358
  - panel data, pitfalls in, 423–424
- first order autocorrelation, 359
- first order conditions, 27, 65, 642, 721
- fitted values. *See also* OLS (ordinary least squares)
- in multiple regressions, 68–69
  - in simple regressions, 27, 32
- fixed effects
- defined, 413
  - dummy variable regression, 438–439
  - estimation, 435–441
  - first differencing vs., 439–440
  - random effects vs, 444–445
  - transformation, 435
  - with unbalanced panels, 440–441
- forecast error, 586
- forecasting
- multiple-step-ahead, 592–594
  - one-step-ahead, 588
  - overview and definitions, 586–587
  - trending, seasonal, and integrated processes, 594–598
  - types of models used for, 587–588
- forecast intervals, 588
- free throw shooting, 651–652
- freeway width and commute time, 702–703
- frequency, data, 7
- frequency distributions, 401(k) plans, 155
- F* statistics. *See also* *F* tests
- defined, 129
  - heteroskedasticity-robust, 247–248
- F* tests. *See also* Chow tests; *F* statistics
- F* and *t* statistics, 132–133
  - functional form misspecification and, 275–279
  - general linear restrictions, 136–137
  - LM tests and, 160
  - overall significance of regressions, 135
  - p*-values for, 134–135
  - reporting regression results, 137–138
  - R*-squared form, 133–134
  - testing exclusion restrictions, 127–132
- functional forms
- in multiple regressions
    - with interaction terms, 177–179
    - logarithmic, 171–173
    - misspecification, 275–279
    - quadratic, 173–177
  - in simple regressions, 36–40
  - in time series regressions, 323–324
- G**
- Gaussian distribution, 665
- Gauss-Markov assumptions
- for multiple linear regressions, 73–77, 82
  - for simple linear regressions, 40–44, 45–48
  - for time series regressions, 319–322

- Gauss-Markov Theorem  
 for multiple linear regressions, 89–90  
 for OLS in matrix form, 725–726
- GDL (geometric distributed lag), 571–572
- GDP. *See* gross domestic product (GDP)
- gender  
 oversampling, 295  
 wage gap, 405–406
- gender gap  
 independent cross sections, 405–406  
 panel data, 405–406
- generalized least squares (GLS) estimators  
 for AR(1) models, 383–387  
 with heteroskedasticity and AR(1) serial correlations, 395  
 when heteroskedasticity function must be estimated, 258–263  
 when heteroskedasticity is known up to a multiplicative constant, 255–256
- geometric distributed lag (GDL), 571–572
- GLS estimators. *See* generalized least squares (GLS) estimators
- Goldberger, Arthur, 85
- goodness-of-fit. *See also* predictions; *R*-squareds  
 change in unit of measurement and, 37  
 in multiple regressions, 70–71  
 overemphasizing, 184–185  
 percent correctly predicted, 227, 530  
 in simple regressions, 35–36  
 in time series regressions, 374
- Google Scholar*, 606
- government policies  
 economic growth and, 6, 8–9
- GPA. *See* college GPA
- Granger, Clive W. J., 150
- Granger causality, 590
- gross domestic product (GDP)  
 data frequency for, 7  
 government policies and, 6  
 high persistence, 355–357  
 in real terms, 326  
 seasonal adjustment of, 336  
 unit root test, 578
- growth rate, 331
- gun control laws, 230
- H**
- HAC standard errors, 389
- Hartford School District, 190
- Hausman test*, 262, 444
- Head Start participation, 230
- Heckit method, 556
- heterogeneity bias, 413
- heteroskedasticity. *See also* weighted least squares estimation  
 2SLS with, 484–485  
 consequences of, for OLS, 243–244  
 defined, 45  
 HAC standard errors, 389  
 heteroskedasticity-robust procedures, 244–249  
 linear probability model and, 265–267  
 robust *F* statistic, 247  
 robust *LM* Statistic, 248  
 robust *t* statistic, 246  
 for simple linear regressions, 45–48  
 testing for, 249–254  
 for time series regressions, 363  
 in time series regressions, 391–395  
 of unknown form, 244  
 in wage equation, 46
- highly persistent time series  
 deciding whether I(0) or I(1), 359–360  
 description of, 354–363  
 transformations on, 358–360
- histogram, 401(k) plan participation, 155
- homoskedasticity  
 for IV estimation, 466–467  
 for multiple linear regressions, 82–83, 89  
 for OLS in matrix form, 724  
 for time series regressions, 319–322, 351–352
- hourly wages. *See* wages
- housing prices and expenditures  
 general linear restrictions, 136–137  
 heteroskedasticity  
 BP test, 251–252  
 White test, 252–254
- incinerators and  
 inconsistency in OLS, 153  
 pooled cross sections, 407–411
- income and, 631
- inflation, 572–574
- investment and  
 computing *R*-squared, 334–335  
 spurious relationship, 332–333  
 over controlling, 185  
 with qualitative information, 211  
 RESET, 278–279  
 savings and, 502
- hypotheses. *See also* hypothesis testing  
 about single linear combination of parameters, 124–127  
 after 2SLS estimation, 479  
 expectations, 14  
 language of classical testing, 120  
 in logit and probit models, 529–530  
 multiple linear restrictions. *See F* tests  
 residual analysis, 190  
 stating, in empirical analysis, 4

- hypothesis testing  
 about mean in normal population, 695–696  
 asymptotic tests for nonnormal populations, 698  
 computing and using  $p$ -values, 698–700  
 confidence intervals and, 701–702  
 in matrix form, Wald statistics for, 730–731  
 overview and fundamentals, 693–695  
 practical vs. statistical significance, 702–703
- I(0) and I(1) processes, 359–360
- idempotent matrices, 715
- identification  
 defined, 465  
 in systems with three or more equations, 510–511  
 in systems with two equations, 504–510
- identified equation, 505
- identity matrices, 710
- idiosyncratic error, 413
- IDL (infinite distributed lag models), 569–574
- IIP (index of industrial production), 326–327
- impact propensity/multiplier, 315
- incidental truncation, 553, 554–558
- incinerators and housing prices  
 inconsistency in OLS, 153  
 pooled cross sections, 407–411
- including irrelevant variables, 77–78
- income. *See also* wages  
 family. *See* family income  
 housing expenditure and, 631  
 PIH, 513–514  
 savings and. *See under* savings
- inconsistency in OLS, deriving, 153–154
- inconsistent estimators, 681
- independence, joint distributions and, 649–651
- independently pooled cross sections. *See also* pooled cross sections  
 across time, 403–407  
 defined, 402
- independent variables. *See also* regression analysis;  
*specific event studies*  
 changing simultaneously, 68  
 defined, 21  
 measurement error in, 289–291  
 in misspecified models, 78–83  
 random, 650  
 simple vs. multiple regression, 61–64
- index numbers, 324–327
- industrial production, index of (IIP), 326–327
- infant mortality rates, outliers, 299–300
- inference  
 in multiple regressions  
 confidence intervals, 122–124  
 statistical, with IV estimator, 466–469  
 in time series regressions, 322–323, 373–374
- infinite distributed lag models, 569–574
- inflation  
 from 1948 to 2003, 313  
 openness and, 508, 509–510  
 random walk model for, 355  
 unemployment and  
 expectations augmented Phillips curve, 353–354  
 forecasting, 589  
 static Phillips curve, 314, 322–323  
 unit root test, 577
- influential observations, 296–300
- information set, 587
- in-sample criteria, 591
- instrumental variables  
 computing  $R$ -squared after estimation, 471  
 in multiple regressions, 471–475  
 overview and definitions, 462, 463, 465  
 properties, with poor instrumental variable, 469–471  
 in simple regressions, 462–471  
 solutions to errors-in-variables problems, 479–481  
 statistical inference, 466–469
- integrated of order zero/one processes, 358–360
- integrated processes, forecasting, 594–598
- interaction effect, 177–179
- interaction terms, 217–218
- intercept parameter, 21
- intercepts. *See also* OLS estimators; regression analysis  
 change in unit of measurement and, 36–37  
 defined, 21, 630  
 in regressions on a constant, 51  
 in regressions through origin, 50–51
- intercept shifts, 207
- interest rates  
 differencing, 387–388  
 inference under CLM assumptions, 323  
 T-bill. *See* T-bill rates
- internet services, 606
- interval estimation, 674, 687–688
- inverse Mills ratio, 538
- inverse of matrix, 713
- IQ  
 ability and, 279–283, 284–285  
 nonrandom sampling, 294–295
- irrelevant variables, including, 77–78
- IV. *See* instrumental variables

## J

JEL. *See Journal of Economic Literature (JEL)*

job training

sample model

as self-selection problem, 3

worker productivity and

program evaluation, 229

as self-selection problem, 230

joint distributions

features of, 652–658

independence and, 649–651

joint hypotheses tests, 127

jointly statistically significant/insignificant, 130

joint probability, 649

*Journal of Economic Literature (JEL)*, 606

junior colleges vs. universities, 124–127

just identified equations, 511

## K

Koyck distributed lag, 571–572

kurtosis, 658

## L

labor economists, 605, 607

labor force. *See* employment and unemployment;

women in labor force

labor supply and demand, 500–501

labor supply function, 639

LAD (least absolute deviations) estimation, 300–302

lag distribution, 315

lagged dependent variables

as proxy variables, 283–284

serial correlation and, 374–375

lagged endogenous variables, 591–592

lagged explanatory variables, 316

Lagrange multiplier (LM) statistics

heteroskedasticity-robust, 248–249. *See also*

heteroskedasticity

main discussion, 158–160

land quality and fertilizers, 23

large sample properties, 681–683

latent variable models, 526

law enforcement

city crime levels and (causality), 13

murder rates and (SEM), 501–502

law of iterated expectations, 664

law of large numbers, 682

law school rankings

as dummy variables, 216–217

residual analysis, 190

leads and lags estimators, 584

least absolute deviations (LAD) estimation, 300–302

least squares estimator, 686

likelihood ratio statistic, 529

limited dependent variables

censored and truncated regression models, 547–552

corner solution response. *See* Tobit model

count response, Poisson regression for, 543–547

overview, 524–525

sample selection corrections, 554–558

linear functions, 630–631

linear independence, 714

linear in parameters assumption

for OLS in matrix form, 723–724

for simple linear regressions, 40, 44

for time series regressions, 317–318

linearity and weak dependence assumption,

348–349

linear probability model (LPM). *See also* limited dependent variables

heteroskedasticity and, 265–266

main discussion, 224–229

linear regression model, 40, 64

linear relationship among independent variables, 83–86

linear time trends, 330

literature review, 607–608

loan approval rates

$F$  and  $t$  statistics, 150

multicollinearity, 85

program evaluation, 230

logarithms

in multiple regressions, 171–173

natural, overview, 736–739

predicting  $y$  when  $\log(y)$  is dependent, 191–193

qualitative information and, 211–212

real dollars and, 327

in simple regressions, 37–39

in time series regressions, 323–324

log function, 636

logit and probit models

interpreting estimates, 530–536

maximum likelihood estimation of, 528–529

specifying, 525–528

testing multiple hypotheses, 529–530

log-likelihood functions, 529

longitudinal data. *See* panel data

long-run elasticity, 324

long-run multiplier. *See* long-run propensity (LRP)

long-run propensity (LRP), 316

loss functions, 586

LRP (long-run propensity), 316

lunch program and math performance, 44–45

## M

- macroeconomists, 606
- MAE (mean absolute error), 591
- marginal effect, 630
- marital status. *See* qualitative information
- martingale difference sequence, 574
- martingale functions, 587
- matched pair samples, 449
- mathematical statistics. *See* statistics
- math performance and lunch program, 44–45
- matrices. *See also* OLS in matrix form
  - addition, 710
  - basic definitions, 709–710
  - differentiation of linear and quadratic forms, 715
  - idempotent, 715
  - linear independence and rank of, 714
  - moments and distributions of random vectors, 716–717
  - multiplication, 711–712
  - operations, 710–713
  - quadratic forms and positive definite, 714–715
- matrix notation, 721
- maximum likelihood estimation, 528–529, 685–686
- MCAR (missing completely at random), 293
- mean, using summation operator, 629–630
- mean absolute error (MAE), 591
- mean independence, 23
- mean squared error (MSE), 680
- measurement error
  - IV solutions to, 479–481
  - men, return to education, 468
  - properties of OLS under, 287–292
- measures of association, 658
- measures of central tendency, 655–657
- measures of variability, 656
- median, 630, 655
- method of moments approach, 25–26, 685
- micronumerosity, 85
- military personnel survey, oversampling in, 295
- minimum variance unbiased estimators, 106, 686, 727
- minimum wages
  - causality, 13
  - employment/unemployment and
    - AR(1) serial correlation, testing for, 377–378
    - detrending, 334–335
    - logarithmic form, 323–324
    - SC-robust standard error, 391
  - in Puerto Rico, effects of, 7–8
- minorities and loans. *See* loan approval rates
- missing at random, 294
- missing completely at random (MCAR), 293
- missing data, 293–294
- misspecification
  - in empirical projects, 613
  - functional form, 275–279
  - unbiasedness and, 78–83
  - variances, 86–87
- motherhood, teenage, 448–449
- moving average process of order one [MA(1)], 346
- MSE (mean squared error), 680
- multicollinearity
  - 2SLS and, 477
  - among explanatory variables, 293
  - main discussion, 83–86
- multiple hypotheses tests, 127
- multiple linear regression (MLR) model, 63
- multiple regression analysis. *See also* data issues; estimation
  - and estimators; heteroskedasticity; hypotheses; OLS (ordinary least squares); predictions; *R*-squareds
  - adding regressors to reduce error variance, 185–186
  - advantages over simple regression, 60–64
  - confidence intervals, 122–124
  - interpreting equations, 67
  - null hypothesis, 108
  - omitted variable bias, 78–83
  - over controlling, 184–185
- multiple regressions. *See also* qualitative information
  - beta coefficients, 169
  - hypotheses with more than one parameter, 124–127
  - misspecified functional forms, 275
  - motivation for multiple regression, 61, 62
  - nonrandom sampling, 294–295
  - normality assumption and, 107
  - productivity and, 360
  - quadratic functions, 173–177
  - with qualitative information
    - of baseball players, race and, 220–221
    - computer usage and, 218
    - with different slopes, 218–221
    - education and, 218–220
    - gender and, 207–211, 212–214, 218–221
    - with interacting terms, 218
    - law school rankings and, 216–217
    - with  $\log(y)$  dependent variable, 213–214
    - marital status and, 219–220
    - with multiple dummy variables, 212–213
    - with ordinal variables, 215–217
    - physical attractiveness and, 216–217
  - random effects model, 443–444
  - random slope model, 285
  - reporting results, 137–138
  - t* test, 110
  - with unobservables, general approach, 284–285
  - with unobservables, using proxy, 279–285
  - working individuals in 1976, 6

multiple restrictions, 127  
 multiple-step-ahead forecast, 587, 592–594  
 multiplicative measurement error, 289  
 multivariate normal distribution, 716–717  
 municipal bond interest rates, 214–215  
 murder rates  
   SEM, 501–502  
   static Phillips curve, 314

**N**

natural experiments, 410, 469  
 natural logarithms, 736–739. *See also* logarithms  
 netted out, 69  
 nominal dollars, 326  
 nominal *vs.* real, 326  
 nonexperimental data, 2  
 nonlinear functions, 634–640  
 nonlinearities, incorporating in simple regressions, 37–39  
 nonnested models  
   choosing between, 182–184  
   functional form misspecification and, 278–279  
 nonrandom samples, 294–295, 553  
 nonstationary time series processes, 345–346  
 no perfect collinearity assumption  
   form, 723  
   for multiple linear regressions, 74–76, 77  
   for time series regressions, 318, 349  
 normal distribution, 665–669  
 normality assumption  
   for multiple linear regressions, 105–108  
   for time series regressions, 322  
 normality of errors assumption, 726  
 normality of estimators in general, asymptotic, 683–684  
 normality of OLS, asymptotic  
   in multiple regressions, 154–160  
   in time series regressions, 351–354  
 normal sampling distributions  
   for multiple linear regressions, 107–108  
   for time series regressions, 322–323  
 no serial correlation assumption. *See also* serial correlation  
   for OLS in matrix form, 724–725  
   for time series regressions, 320–322, 351–352  
*n-R*-squared statistic, 159  
 null hypothesis, 108–110, 694. *See also* hypotheses  
 numerator degrees of freedom, 129

**O**

observational data, 2  
 OLS (ordinary least squares)  
   cointegration and, 583–584  
   comparison of simple and multiple regression estimates, 69–70  
   consistency. *See* consistency of OLS

logit and probit *vs.*, 533–535  
 in multiple regressions  
   algebraic properties, 64–72  
   computational properties, 64–66, 64–72  
   effects of data scaling, 166–170  
   fitted values and residuals, 68  
   goodness-of-fit, 70–71  
   interpreting equations, 65–66  
   Lagrange multiplier (LM) statistic, 158–160  
   measurement error and, 287–292  
   normality, 154–160  
   partialling out, 69  
   regression through origin, 73  
   statistical properties, 73–81  
 Poisson *vs.*, 545, 546–547  
 in simple regressions  
   algebraic properties, 32–34  
   defined, 27  
   deriving estimates, 24–32  
   statistical properties, 45–50  
   units of measurement, changing, 36–37  
 simultaneity bias in, 503–504  
 in time series regressions  
   correcting for serial correlation, 383–386  
   FGLS *vs.*, 385–386  
   finite sample properties, 317–323  
   normality, 351–354  
   SC-robust standard errors, 388–391  
   with serially correlated errors, properties of, 373–375  
 Tobit *vs.*, 540–542  
 OLS and Tobit estimates, 540–542  
 OLS asymptotics  
   in matrix form, 728–731  
   in multiple regressions  
     consistency, 150–154  
     efficiency, 161–162  
     overview, 149–150  
   in time series regressions  
     consistency, 348–354  
 OLS estimators. *See also* heteroskedasticity  
   defined, 40  
   in multiple regressions  
     efficiency of, 89–90  
     variances of, 81–89  
   sampling distributions of, 105–108  
 in simple regressions  
   expected value of, 73–81  
   unbiasedness of, 40–45, 77  
   variances of, 45–48  
 in time series regressions  
   sampling distributions of, 322–323  
   unbiasedness of, 317–323  
   variances of, 320–322

- OLS in matrix form
    - asymptotic analysis, 728–731
    - finite sample properties, 723–726
    - overview, 720–722
    - statistical inference, 726–728
    - Wald statistics for testing multiple hypotheses, 730–731
  - OLS intercept estimates, defined, 65–66
  - OLS regression line. *See also* OLS (ordinary least squares)
    - defined, 28
    - in multiple regressions, 65
  - OLS slope estimates, defined, 65
  - omitted variable bias. *See also* instrumental variables
    - general discussions, 78–83
    - using proxy variables, 279–285
  - one-sided alternatives, 695
  - one-step-ahead forecasts, 586, 588
  - one-tailed tests, 110, 696. *See also* *t* tests
  - online databases, 609
  - online search services, 607–608
  - order condition, 479, 507
  - ordinal variables, 214–217
  - outliers
    - guarding against, 300–302
    - main discussion, 296–300
  - out-of-sample criteria, 591
  - overall significance of regressions, 135
  - over controlling, 184–185
  - overdispersion, 545
  - overidentified equations, 511
  - overidentifying restrictions, testing, 482–485
  - overspecifying the model, 78
- P**
- pairwise uncorrelated random variables, 660–661
  - panel data
    - applying 2SLS to, 487–488
    - applying methods to other structures, 448–450
    - correlated random effects, 445–447
    - differencing with more than two periods, 420–425
    - fixed effects, 435–441
    - independently pooled cross sections vs, 403
    - organizing, 417
    - overview, 9–10
    - pitfalls in first differencing, 424
    - random effects, 441–445
    - simultaneous equations models with, 514–516
    - two-period, analysis, 417–419
    - two-period, policy analysis with, 417–419
    - unbalanced, 440–441
  - Panel Study of Income Dynamics, 608
  - parameters
    - defined, 4, 674
    - estimation, general approach to, 684–686
    - partial derivatives, 641
    - partial effect, 66, 67–68
    - partial effect at average (PEA), 531–532
    - partially out, 69
    - partitioned matrix multiplication, 712–713
    - pdf (probability density functions), 647
    - percentage point change, 634
    - percentages, 633–634
      - change, 633
    - percent correctly predicted, 227, 530
    - perfect collinearity, 74–76
    - permanent income hypothesis, 513–514
    - pesticide usage, over controlling, 185
    - physical attractiveness and wages, 215–216
    - pizzas, expected revenue, 654
    - plug-in solution
      - to the omitted variables problem, 280
    - point estimates, 674
    - point forecasts, 588
    - poisson distribution, 544, 545
    - poisson regression model, 543–547
    - policy analysis
      - with pooled cross sections, 407–412
      - with qualitative information, 210, 229–231
      - with two-period panel data, 417–419
    - pooled cross sections. *See also* independently pooled cross sections
      - applying 2SLS to, 487–488
      - overview, 8
      - policy analysis with, 407–412
    - population, defined, 674
    - population model, defined, 73
    - population regression function (PRF), 23
    - population *R*-squareds, 181
    - positive definite and semi-definite matrices,
      - defined, 715
    - poverty rate
      - in absence of suitable proxies, 285
      - excluding from model, 80
    - power of test, 694
    - practical significance, 120
    - practical vs. statistical significance, 120–124, 702–703
    - Prais-Winsten (PW) estimation, 383–384, 386, 390
    - predetermined variables, 592
    - predicted variables, 21. *See also* dependent variables
    - prediction error, 188
    - predictions
      - confidence intervals for, 186–189
      - with heteroskedasticity, 264–266
      - residual analysis, 190
      - for *y* when  $\log(y)$  is dependent, 191–193
    - predictor variables, 23. *See also* dependent variables
    - price index, 326–327

prisons  
 population and crime rates, 515–516  
 recidivism, 549–551

probability. *See also* conditional distributions; joint distributions  
 features of distributions, 652–658  
 independence, 649–651  
 joint, 649  
 normal and related distributions, 665–669  
 overview, 645  
 random variables and their distributions, 645–649

probability density function (pdf), 647

probability limits, 681–683

probit model. *See* logit and probit models

productivity. *See* worker productivity

program evaluation, 210, 229–231

projects. *See* empirical analysis

property taxes and housing prices, 8

proportions, 733–734

proxy variables, 279–285

pseudo *R*-squareds, 531

public finance study researchers, 606

Puerto Rico, employment in  
 detrending, 334–335  
 logarithmic form, 323–324  
 time series data, 7–8

*p*-values  
 computing and using, 698–700  
 for *F* tests, 134–135  
 for *t* tests, 118–120

## Q

QMLE (quasi-maximum likelihood estimation), 728

quadratic form for matrices, 714–715, 716

quadratic function, 634–636

quadratic time trends, 331

qualitative information. *See also* linear probability model (LPM)  
 in multiple regressions  
 allowing for different slopes, 218–221  
 binary dependent variable, 224–229  
 describing, 205–206  
 discrete dependent variables, 231–232  
 interactions among dummy variables, 217  
 with  $\log(y)$  dependent variable, 211–212  
 multiple dummy independent variables, 212–217  
 ordinal variables, 214–217  
 overview, 205  
 policy analysis and program evaluation, 229–231  
 proxy variables, 282–283  
 single dummy independent variable, 206–212

testing for differences in regression functions across groups, 221–224  
 in time series regressions  
 main discussion, 324–329  
 seasonal, 336–338

quantile regression, 302

quasi-demeaned data, 442

quasi-differenced data, 382, 390

quasi-experiment, 410

quasi- (natural) experiments, 410, 469

quasi-likelihood ratio statistic, 546

quasi-maximum likelihood estimation (QMLE), 545, 728

## R

$R^2_j$ , 83–86

race  
 arrests and, 229  
 baseball player salaries and, 220–221  
 discrimination in hiring  
 asymptotic confidence interval, 692–693  
 hypothesis testing, 698  
*p*-value, 701

random coefficient model, 285–287

random effects  
 correlated, 445–447  
 estimator, 442  
 fixed effects vs., 444–445  
 main discussion, 441–445

random sampling  
 assumption  
 for multiple linear regressions, 74  
 for simple linear regressions, 40–41, 42, 44  
 cross-sectional data and, 5–7  
 defined, 675

random slope model, 285–287

random variables, 645–649

random vectors, 716

random walks, 354

rank condition, 479, 497, 506–507

rank of matrix, 714

rational distributed lag models, 572–574

R&D and sales  
 confidence intervals, 123–124  
 nonnested models, 182–184  
 outliers, 296–298

RDL (rational distributed lag models), 572–574

real dollars, 326

recidivism, duration analysis, 549–551

reduced form equations, 473, 504

reduced form error, 504

reduced form parameters, 504

regressands, 21. *See also* dependent variables

regression analysis, 50–51. *See also* multiple regression analysis; simple regression model; time series data

regression specification error test (RESET), 277–278

regression through origin, 50–52

regressors, 21, 185–186. *See also* independent variables

rejection region, 695

rejection rule, 110. *See also* *t* tests

relative change, 633

relative efficiency, 679–680

relevant variables, excluding, 78–83

reporting multiple regression results, 137–138

resampling method, 203

rescaling, 166–168

RESET (regression specification error test), 277

residual analysis, 190

residuals. *See also* OLS (ordinary least squares)

- in multiple regressions, 68, 297–298
- in simple regressions, 27, 32, 48
- studentized, 297–298

residual sum of squares (SSR). *See* sum of squared residuals

response probability, 225, 525

response variables, 21. *See also* dependent variables

REST (regression specification error test), 277–278

restricted model, 128–129. *See also* *F* tests

retrospective data, 2

returns on equity and CEO salaries

- fitted values and residuals, 32
- goodness-of-fit, 35
- OLS Estimates, 29–30

RMSE (root mean squared error), 50, 88, 591

robust regression, 302

rooms and housing prices

- beta coefficients, 175–176
- interaction effect, 177–179
- quadratic functions, 175–177
- residual analysis, 190

root mean squared error (RMSE), 50, 88, 591

row vectors, 709

*R*-squareds. *See also* predictions

- adjusted, 181–184, 374
- after IV estimation, 471
- change in unit of measurement and, 37
- in fixed effects estimation, 437, 438–439
- for *F* statistic, 133–134
- in multiple regressions, main discussion, 70–73
- for probit and logit models, 531
- for PW estimation, 383–384
- in regressions through origin, 50–51, 73
- in simple regressions, 35–36
- size of, 180–181
- in time series regressions, 374
- trending dependent variables and, 334–335
- uncentered, 214

## S

salaries. *See* CEO salaries; income; wages

sales

- CEO salaries and
  - constant elasticity model, 39
  - nonnested models, 183–184
  - motivation for multiple regression, 63–64
  - R&D and. *See* R&D and sales

sales tax increase, 634

sample average, 675

sample correlation coefficient, 685

sample covariance, 685

sample regression function (SRF), 28, 65

sample selection corrections, 553–558

sample standard deviation, 683

sample variation in the explanatory variable

- assumption, 42, 44

sampling, nonrandom, 293–300

sampling distributions

- defined, 676
- of OLS estimators, 105–108

sampling standard deviation, 693

sampling variances

- of estimators in general, 678–679
- of OLS estimators
  - for multiple linear regressions, 82, 83
  - for simple linear regressions, 47–48

savings

- housing expenditures and, 502
- income and
  - heteroskedasticity, 254–256
  - scatterplot, 25
- measurement error in, 289
- with nonrandom sample, 294–295

scalar multiplication, 710

scalar variance-covariance matrices, 724

scatterplots

- R&D and sales, 297–298
- savings and income, 25
- wage and education, 27

school lunch program and math performance, 44–45

school size and student performance, 113–114

score statistic, 158–160

scrap rates and job training

- 2SLS, 487
- confidence interval, 700–701
- confidence interval and hypothesis testing, 702
- fixed effects estimation, 436–437
- measurement error in, 289
- program evaluation, 229
- p*-value, 700–701
- statistical vs. practical significance, 121–122

- two-period panel data, 418
- unbalanced panel data, 441
- seasonal dummy variables, 337
- seasonality
  - forecasting, 594–598
  - serial correlation and, 381
  - of time series, 336–338
- seasonally adjusted patterns, 336
- selected samples, 553
- self-selection problems, 230
- SEM. *See* simultaneous equations models
- semi-elasticity, 39, 639
- sensitivity analysis, 613
- sequential exogeneity, 363
- serial correlation
  - correcting for, 381–387
  - differencing and, 387–389
  - heteroskedasticity and, 395
  - lagged dependent variables and, 374–375
  - no serial correlation assumption, 320–322, 351–354
  - properties of OLS with, 373–375
  - testing for, 376–381
- serial correlation-robust standard errors, 388–391
- serially uncorrelation, 360
- short-run elasticity, 324
- significance level, 110
- simple linear regression model, 20
- simple regression model, 20–24. *See also* OLS (ordinary least squares)
  - incorporating nonlinearities in, 37–39
  - IV estimation, 462–471
  - multiple regression *vs.*, 60–63
  - regression on a constant, 51
  - regression through origin, 50–51
- simultaneity bias, 504
- simultaneous equations models
  - bias in OLS, 503–504
  - identifying and estimating structural equations, 504–510
  - overview and nature of, 449–503
  - with panel data, 514–516
  - systems with more than two equations, 510–511
  - with time series, 511–514
- skewness, 658
- sleeping *vs.* working tradeoff, 415–416
- slopes. *See also* OLS estimators; regression analysis
  - change in unit of measurement and, 36–37, 39
  - defined, 21, 630
  - parameter, 21
  - qualitative information and, 218–221
  - random, 285–287
  - in regressions on a constant, 51
  - in regressions through origin, 50–51
- smearing estimates, 191
- smoking
  - birth weight and
    - asymptotic standard error, 158
    - data scaling, 166–170
  - cigarette taxes and consumption, 411–412
  - demand for cigarettes, 261–262
  - IV estimation, 470
  - measurement error, 292
- Social Sciences Citation Index*, 606
- soybean yields and fertilizers
  - causality, 11, 12
  - simple equation, 21–22
- specification search, 613
- spreadsheets, 610
- spurious regression, 332–333, 578–580
- square matrices, 709–710
- SRF (sample regression function), 28, 65
- SSE (explained sum of squares), 34, 70–71
- SSR (residual sum of squares). *See* sum of squared residuals
- SST (total sum of squares), 34, 70–71
- SSTj (total sample variation in  $x_j$ ), 83
- stable AR(1) processes, 347
- standard deviation
  - of  $\hat{\beta}_j$ , 89–90
  - defined, 45, 657
  - estimating, 49
  - properties of, 657
- standard error of the regression (SER), 50, 88
- standard errors
  - asymptotic, 157
  - of  $\hat{\beta}_j$ , 88
  - heteroskedasticity-robust, 246–247
  - of OLS estimators, 87–89
  - of  $\hat{\beta}_1$ , 50
  - serial correlation-robust, 388–391
- standardized coefficients, 169–170
- standardized random variables, 657–658
- standardized test scores
  - beta coefficients, 169
  - collinearity, 74–75
  - interaction effect, 178–179
  - motivation for multiple regression, 61, 62
  - omitted variable bias, 80, 81
  - omitting unobservables, 285
  - residual analysis, 190
- standard normal distribution, 666–668, 743–744
- static models, 314, 350
- static Phillips curve, 314, 322–323, 377, 378, 386
- stationary time series processes, 345–346
- statistical inference
  - with IV estimator, 466–469
  - for OLS in matrix form, 726–728

statistical significance  
   defined, 115  
   economic/practical significance vs., 120–124  
   economic/practical significance vs., 702  
   joint, 130  
 statistical tables, 743–749  
 statistics. *See also* hypothesis testing  
   asymptotic properties of estimators, 681–684  
   finite sample properties of estimators, 675–680  
   interval estimation and confidence  
     intervals, 687–693  
   notation, 703  
   overview and definitions, 674–675  
   parameter estimation, general approaches  
     to, 684–686  
 stepwise regression, 614  
 stochastic process, 313, 345  
 stock prices and trucking regulations, 325  
 stock returns, 393, 394. *See also* efficient markets  
   hypothesis (EMH)  
 stratified sampling, 295  
 strict exogeneity assumption, 414–420, 570  
 strictly exogenous variables  
   correcting for, 381–387  
   serial correlation  
     testing for, 376–381  
 strict stationarity, 345  
 strongly dependent time series. *See* highly persistent time series  
 structural equations  
   definitions, 471, 500, 501, 504  
   identifying and estimating, 504–510  
 structural error, 501  
 structural parameters, 504  
 student enrollment, *t* test, 116–117  
 studentized residuals, 298  
 student performance. *See also* college GPA; final exam scores; standardized test scores  
   in math, lunch program and, 44–45  
   school expenditures and, 85  
   school size and, 113–114  
 style hints for empirical papers, 619–621  
 summation operator, 628–630  
 sum of squared residuals. *See also* OLS (ordinary least squares)  
   in multiple regressions, 70–71  
   in simple regressions, 34  
 supply shock, 353  
 Survey of Consumer Finances, 608  
 symmetric matrices, 712  
 systematic part, defined, 24  
 system estimation methods, 511

## T

tables, statistical, 743–749  
 tax exemption. *See under* fertility rate  
 T-bill rates  
   cointegration, 580–584  
   error correction model, 585  
   inflation, deficits. *See under* interest rates  
   random walk characterization of, 355, 356  
   unit root test, 576  
*t* distribution  
   critical values table, 745  
   discussions, 108–110, 660–670, 717  
   for standardized estimators, 108–110  
 teachers, salary-pension tradeoff, 137–138  
 teenage motherhood, 448–449  
 tenure. *See also* wages  
   interpreting equations, 67  
   motivation for multiple regression, 63–64  
 testing overidentifying restrictions, 482–485  
 test scores, as indicators of ability, 481  
 test statistic, 695  
 text editor, 609  
 text files and editors, 608–609  
 theorems  
   asymptotic efficiency of OLS, 162  
     for time series regressions, 351–354  
   consistency of OLS  
     for multiple linear regressions, 150–154  
     for time series regressions, 348–351  
   Gauss-Markov  
     for multiple linear regressions, 89–90  
     for time series regressions, 320–322  
   normal sampling distributions, 107–108  
   for OLS in matrix form  
     Gauss-Markov, 725–726  
     statistical inference, 726–728  
     unbiasedness, 726  
     variance-covariance matrix of OLS  
       estimator, 724–725  
   sampling variances of OLS estimators  
     for simple linear regressions, 47–48  
     for time series regressions, 320–322  
   unbiased estimation of  $s^2$   
     for multiple linear regressions, 88–89  
     for time series regressions, 321  
   unbiasedness of OLS  
     for multiple linear regressions, 77  
     for time series regressions, 317–320  
   theoretical framework, 615  
   three stage least squares, 511  
   time-demeaned data, 435

- time series data
    - absence of serial correlation, 360–363
    - applying 2SLS to, 485–486
    - cointegration, 580–584
    - dynamically complete models, 360–363
    - error correction models, 584–586
    - examples of models, 313–316
    - functional forms, 323–324
    - heteroskedasticity in, 391–395
    - highly persistent. *See* highly persistent time series
    - homoskedasticity assumption for, 363–364
    - infinite distributed lag models, 569–574
    - nature of, 312–313
    - OLS. *See under* OLS (ordinary least squares); OLS estimators
    - overview, 7–8
    - in panel data, 9–10
    - in pooled cross sections, 8–9
    - with qualitative information. *See under* qualitative information
    - seasonality, 336–338
    - simultaneous equations models with, 511–514
    - spurious regression, 578–580
    - stationary and nonstationary, 345–346
    - unit roots, testing for, 574–579
    - weakly dependent, 346–348
  - time trends. *See* trends
  - time-varying error, 413
  - tobit model
    - interpreting estimates, 537–542
    - overview, 536–537
    - specification issues in, 543
  - top coding, 548
  - total sample variation in  $x_j$ , 83
  - total sum of squares (SST), 34, 70–71
  - trace of matrix, 713
  - traffic fatalities
    - beer taxes and, 184
  - training grants. *See also* job training program evaluation, 229
    - single dummy variable, 210–211
  - transpose of matrix, 712
  - treatment group, 210
  - trends
    - characterizing trending time series, 329–332
    - detrending, 334–335
    - forecasting, 594–598
    - high persistence *vs.*, 352
    - $R$ -squared and trending dependent variable, 334–335
    - seasonality and, 337–338
    - time, 329
    - using trending variables, 332–333
  - trend-stationary processes, 348
  - trucking regulations and stock prices, 325
  - true model, defined, 74
  - truncated normal regression model, 551
  - truncated regression models, 548, 551–552
  - $t$  statistics. *See also*  $t$  tests
    - asymptotic, 157
    - defined, 109, 696
    - $F$  statistic and, 132–133
    - heteroskedasticity-robust, 246–247
  - $t$  tests. *See also*  $t$  statistics
    - for AR(1) serial correlation, 376–378
    - null hypothesis, 108–110
    - one-sided alternatives, 110–114
    - other hypotheses about  $b_j$ , 116–118
    - overview, 108–110
    - $p$ -values for, 118–120
    - two-sided alternatives, 114–115
  - two-period panel data
    - analysis, 417–419
    - policy analysis with, 417–419
  - two-sided alternatives, 695–696
  - two stage least squares
    - applied to pooled cross sections and panel data, 487–488
    - applied to time series data, 485–486
    - with heteroskedasticity, 485–486
    - multiple endogenous explanatory variables, 478–479
    - for SEM, 508–510, 511
    - single endogenous explanatory variable, 475–477
    - testing multiple hypotheses after estimation, 479
    - testing for endogeneity, 481–482
  - two-tailed tests, 115, 697. *See also*  $t$  tests
  - Type I/II error, 694
- ## U
- $u$  (“unobserved” term)
    - CEV assumption and, 292
    - foregoing specifying models with, 284–285
    - general discussions, 4–5, 21–23
    - in time series regressions, 319
    - using proxy variables for, 279–285
  - unanticipated inflation, 353
  - unbalanced panels, 440–441
  - unbiased estimation of  $s^2$ 
    - for multiple linear regressions, 88–89
    - for simple linear regressions, 49
    - for time series regressions, 321
  - unbiasedness
    - in general, 677–678
    - of OLS
      - in matrix form, 724

unbiasedness (*continued*)  
 in multiple regressions, 77  
 for simple linear regressions, 43–44  
 in simple regressions, 40–44  
 in time series regressions, 317–323,  
 373–375  
 of  $\hat{\sigma}^2$ , 726

uncentered *R*-squareds, 214

unconditional forecasts, 587

uncorrelated random variables, 660

underspecifying the model, 78–83

unemployment. *See* employment and unemployment

unidentified equations, 511

unit roots  
 forecasting processes with, 597–598  
 testing for, 574–579  
 gross domestic product (GDP), 578  
 inflation, 577  
 process, 355, 358

units of measurement, effects of changing, 36–37,  
 166–168

universities vs. junior colleges, 124–127

unobserved effects/heterogeneity, 413, 435. *See also* fixed effects

“unobserved” terms. *See* *u* (“unobserved” term)

unrestricted model, 128–129. *See also* *F* tests

unsystematic part, defined, 24

upward bias, 80, 81

utility maximization, 2

## V

variables. *See also* dependent variables; independent variables; *specific types*  
 dummy, 206. *See also* qualitative information  
 in multiple regressions, 61–64  
 seasonal dummy, 337  
 in simple regressions, 20–21

variance-covariance matrices, 716, 724–725

variance inflation factor (VIF), 86

variance of prediction error, 188

variances  
 conditional, 665  
 of OLS estimators  
 in multiple regressions, 81–89  
 in simple regressions, 45–50  
 in time series regressions, 320–322  
 overview and properties of, 656–657, 660–661  
 of prediction error, 189

VAR model, 589, 597–598

vector autoregressive model, 589, 597–598

vectors, defined, 709

veterans, earnings of, 469

voting outcomes  
 campaign expenditures and deriving OLS estimate, 31  
 economic performance and, 328–329  
 perfect collinearity, 75–76

## W

wages  
 causality, 13–14  
 education and  
 2SLS, 488  
 conditional expectation, 661–665  
 heteroskedasticity, 46–47  
 independent cross sections, 405–406  
 nonlinear relationship, 37–39  
 OLS estimates, 30–31  
 partial effect, 641  
 rounded averages, 33  
 scatterplot, 27  
 simple equation, 22  
 experience and. *See under* experience  
 with heteroskedasticity-robust standard errors, 246–247  
 labor supply and demand, 500–501  
 labor supply function, 639  
 multiple regressions. *See also* qualitative information  
 homoskedasticity, 82–83

Wald test/statistics, 529–530, 537, 730–731

weak instruments, 471

weakly dependent time series, 346–348

wealth. *See* financial wealth

weighted least squares estimation  
 linear probability model, 265–267  
 overview, 254  
 prediction and prediction intervals, 264–265  
 for time series regressions, 390, 393–394  
 when assumed heteroskedasticity function is wrong,  
 262–264  
 when heteroskedasticity function must be estimated,  
 258–263  
 when heteroskedasticity is known up to a multiplicative  
 constant, 254–259

White test for heteroskedasticity, 252–254

within estimators, 435. *See also* fixed effects

within transformation, 435

women in labor force  
 heteroskedasticity, 265–267  
 LPM, logit, and probit estimates, 533–535  
 return to education  
 2SLS, 477  
 IV estimation, 467  
 testing for endogeneity, 482  
 testing overidentifying restrictions, 482  
 sample selection correction, 556–557

women's fertility. *See* fertility rate  
worker compensation laws and weeks out of work, 411  
worker productivity  
  job training and  
    program evaluation, 229  
  sample model, 4  
  in U.S., trend in, 331  
  wages and, 360  
working vs. sleeping tradeoff, 415–416  
working women. *See* women in labor force  
writing empirical papers, 614–621  
  conceptual (or theoretical) framework, 615  
  conclusions, 618–619  
  data description, 617–618  
  econometric models and estimation methods, 615–617  
  introduction, 614–615  
  results section, 618  
  style hints, 619–621

## Y

year dummy variables  
  in fixed effects model, 436–438  
  pooling independent cross sections across time,  
    403–407  
  in random effects model, 443–444

## Z

zero conditional mean assumption  
  homoskedasticity vs., 45  
  for multiple linear regressions, 62–63, 76–77  
  for OLS in matrix form, 724  
  for simple linear regressions, 23–24, 42, 44  
  for time series regressions, 318–319, 349  
zero mean and zero correlation assumption, 152  
zero-one variables, 206. *See also* qualitative information