

Prediction, Goodness-of-Fit, and Modeling Issues

Chapter 4

Prepared by Vera Tabakova, East Carolina University

Chapter 4:
Prediction, Goodness-of-Fit, and Modeling Issues

- 4.1 Least Squares Prediction
- 4.2 Measuring Goodness-of-Fit
- 4.3 Modeling Issues
- 4.4 Log-Linear Models

Principles of Econometrics, 3rd Edition Slide 4-2

4.1 Least Squares Prediction

$$y_0 = \beta_1 + \beta_2 x_0 + e_0$$
(4.1)

where e_0 is a random error. We assume that $E(y_0) = \beta_1 + \beta_2 x_0$ and $E(e_0) = 0$. We also assume that $\text{var}(e_0) = \sigma^2$ and $\text{cov}(e_0, e_i) = 0 \quad i = 1, 2, \dots, N$

$$\hat{y}_0 = b_1 + b_2 x_0$$
(4.2)

Principles of Econometrics, 3rd Edition Slide 4-3

4.1 Least Squares Prediction

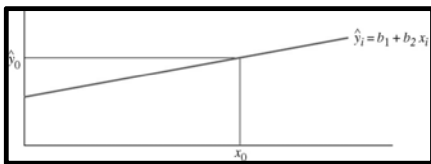


Figure 4.1 A point prediction

Principles of Econometrics, 3rd Edition

Slide 4-4

4.1 Least Squares Prediction

$$f = y_0 - \hat{y}_0 = (\beta_1 + \beta_2 x_0 + e_0) - (b_1 + b_2 x_0) \quad (4.3)$$

$$\begin{aligned} E(f) &= \beta_1 + \beta_2 x_0 + E(e_0) - [E(b_1) + E(b_2) x_0] \\ &= \beta_1 + \beta_2 x_0 + 0 - [\beta_1 + \beta_2 x_0] = 0 \end{aligned}$$

$$\text{var}(f) = \sigma^2 \left[1 + \frac{1}{N} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \quad (4.4)$$

Principles of Econometrics, 3rd Edition

Slide 4-5

4.1 Least Squares Prediction

The variance of the forecast error is smaller when

- i. the overall uncertainty in the model is smaller, as measured by the variance of the random errors ;
- ii. the sample size N is larger;
- iii. the variation in the explanatory variable is larger; and
- iv. the value of x_0 is small.

Principles of Econometrics, 3rd Edition

Slide 4-6

4.1 Least Squares Prediction

$$\widehat{\text{var}}(f) = \hat{\sigma}^2 \left[1 + \frac{1}{N} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

$$\text{se}(f) = \sqrt{\widehat{\text{var}}(f)} \tag{4.5}$$

$$\hat{y}_0 \pm t_c \text{se}(f) \tag{4.6}$$

Principles of Econometrics, 3rd Edition Slide 4-7

4.1 Least Squares Prediction

Figure 4.2 Point and interval prediction

Principles of Econometrics, 3rd Edition Slide 4-8

4.1.1 Prediction in the Food Expenditure Model

$$\hat{y}_0 = b_1 + b_2 x_0 = 83.4160 + 10.2096(20) = 287.6089$$

$$\widehat{\text{var}}(f) = \hat{\sigma}^2 \left[1 + \frac{1}{N} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

$$= \hat{\sigma}^2 + \frac{\hat{\sigma}^2}{N} + (x_0 - \bar{x})^2 \frac{\hat{\sigma}^2}{\sum (x_i - \bar{x})^2}$$

$$= \hat{\sigma}^2 + \frac{\hat{\sigma}^2}{N} + (x_0 - \bar{x})^2 \widehat{\text{var}}(b_2)$$

$$\hat{y}_0 \pm t_c \text{se}(f) = 287.6069 \pm 2.0244(90.6328) = [104.1323, 471.0854]$$

Principles of Econometrics, 3rd Edition Slide 4-9

4.2 Measuring Goodness-of-Fit

$$y_i = \beta_1 + \beta_2 x_i + e_i \quad (4.7)$$

$$y_i = E(y_i) + e_i \quad (4.8)$$

$$y_i = \hat{y}_i + \hat{e}_i \quad (4.9)$$

$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + \hat{e}_i \quad (4.10)$$

Principles of Econometrics, 3rd Edition

Slide 4-10

4.2 Measuring Goodness-of-Fit

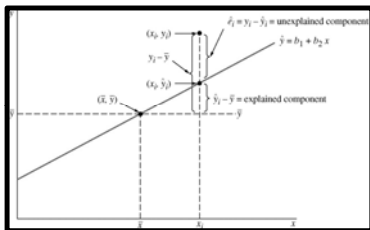


Figure 4.3 Explained and unexplained components of y_i

Principles of Econometrics, 3rd Edition

Slide 4-11

4.2 Measuring Goodness-of-Fit

$$\hat{\sigma}_y^2 = \frac{\sum (y_i - \bar{y})^2}{N-1}$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum \hat{e}_i^2 \quad (4.11)$$

Principles of Econometrics, 3rd Edition

Slide 4-12

4.2 Measuring Goodness-of-Fit

- $\sum (y_i - \bar{y})^2$ = total sum of squares = *SST*: a measure of *total variation* in *y* about the sample mean.
- $\sum (\hat{y}_i - \bar{y})^2$ = sum of squares due to the regression = *SSR*: that part of total variation in *y*, about the sample mean, that is explained by, or due to, the regression. Also known as the “explained sum of squares.”
- $\sum e_i^2$ = sum of squares due to error = *SSE*: that part of total variation in *y* about its mean that is not explained by the regression. Also known as the unexplained sum of squares, the residual sum of squares, or the sum of squared errors.
- $SST = SSR + SSE$

Principles of Econometrics, 3rd Edition

Slide 4-13

4.2 Measuring Goodness-of-Fit

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad (4.12)$$

- The closer R^2 is to one, the closer the sample values y_i are to the fitted regression equation $\hat{y}_i = b_1 + b_2x_i$. If $R^2 = 1$, then all the sample data fall exactly on the fitted least squares line, so $SSE = 0$, and the model fits the data “perfectly.” If the sample data for y and x are uncorrelated and show no linear association, then the least squares fitted line is “horizontal,” so that $SSR = 0$ and $R^2 = 0$.

Principles of Econometrics, 3rd Edition

Slide 4-14

4.2.1 Correlation Analysis

$$\rho_{xy} = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)}\sqrt{\text{var}(y)}} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x\hat{\sigma}_y} \quad (4.13)$$

$$r_{xy} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x\hat{\sigma}_y} \quad (4.14)$$

$$\begin{aligned} \hat{\sigma}_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) / (N - 1) \\ \hat{\sigma}_x &= \sqrt{\sum (x_i - \bar{x})^2 / (N - 1)} \\ \hat{\sigma}_y &= \sqrt{\sum (y_i - \bar{y})^2 / (N - 1)} \end{aligned} \quad (4.15)$$

Principles of Econometrics, 3rd Edition

Slide 4-15

4.2.2 Correlation Analysis and R^2

$$r_{xy}^2 = R^2$$

$$R^2 = r_{xy}^2$$

R^2 measures the linear association, or goodness-of-fit, between the sample data and their predicted values. Consequently R^2 is sometimes called a measure of "goodness-of-fit."

Principles of Econometrics, 3rd Edition

Slide 4-16

4.2.3 The Food Expenditure Example

$$SST = \sum (y_i - \bar{y})^2 = 495132.160$$

$$SSE = \sum (y_i - \hat{y}_i)^2 = \sum e_i^2 = 304505.176$$

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{304505.176}{495132.160} = .385$$

$$r_{xy} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y} = \frac{478.75}{(6.848)(112.675)} = .62$$

Principles of Econometrics, 3rd Edition

Slide 4-17

4.2.4 Reporting the Results

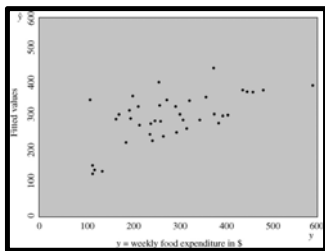


Figure 4.4 Plot of predicted \hat{y} against y

Principles of Econometrics, 3rd Edition

Slide 4-18

4.2.4 Reporting the Results

- *FOOD_EXP* = weekly food expenditure by a household of size 3, in dollars
- *INCOME* = weekly household income, in \$100 units

$$FOOD_EXP = 83.42 + 10.21 INCOME \quad R^2 = .385$$

(se) (43.41)* (2.09)***

* indicates significant at the 10% level
 ** indicates significant at the 5% level
 *** indicates significant at the 1% level

Principles of Econometrics, 3rd Edition Slide 4-19

4.3 Modeling Issues

- 4.3.1 The Effects of Scaling the Data
- Changing the scale of *x*:
 $y = \beta_1 + \beta_2 x + e = \beta_1 + (c\beta_2)(x/c) + e = \beta_1 + \beta_2^* x^* + e$
 where $\beta_2^* = c\beta_2$ and $x^* = x/c$
- Changing the scale of *y*:
 $y/c = (\beta_1/c) + (\beta_2/c)x + (e/c)$ or $y^* = \beta_1^* + \beta_2^* x + e^*$

Principles of Econometrics, 3rd Edition Slide 4-20

4.3.2 Choosing a Functional Form

Variable transformations:

- Power: if *x* is a variable then x^p means raising the variable to the power *p*; examples are quadratic (x^2) and cubic (x^3) transformations.
- The natural logarithm: if *x* is a variable then its natural logarithm is $\ln(x)$.
- The reciprocal: if *x* is a variable then its reciprocal is $1/x$.

Principles of Econometrics, 3rd Edition Slide 4-21

4.3.2 Choosing a Functional Form

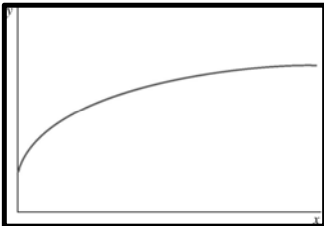


Figure 4.5 A nonlinear relationship between food expenditure and income

Principles of Econometrics, 3rd Edition *Slide 4-22*

4.3.2 Choosing a Functional Form

- The **log-log model**
 $\ln(y) = \beta_1 + \beta_2 \ln(x)$
 The parameter β_2 is the elasticity of y with respect to x .
- The **log-linear model**
 $\ln(y_i) = \beta_1 + \beta_2 x_i$
 A one-unit increase in x leads to (approximately) a $100 \beta_2$ percent change in y .
- The **linear-log model**
 $y = \beta_1 + \beta_2 \ln(x)$ or $\frac{\Delta y}{100(\Delta x/x)} = \frac{\beta_2}{100}$
 A 1% increase in x leads to a $\beta_2/100$ unit change in y .

Principles of Econometrics, 3rd Edition *Slide 4-23*

4.3.3 The Food Expenditure Model

- The reciprocal model is
 $FOOD_EXP = \beta_1 + \beta_2 \frac{1}{INCOME} + e$
- The linear-log model is
 $FOOD_EXP = \beta_1 + \beta_2 \ln(INCOME) + e$

Principles of Econometrics, 3rd Edition *Slide 4-24*

4.3.3 The Food Expenditure Model

Remark: Given this array of models, that involve different transformations of the dependent and independent variables, and some of which have similar shapes, what are some guidelines for choosing a functional form?

1. Choose a shape that is consistent with what economic theory tells us about the relationship.
2. Choose a shape that is sufficiently flexible to "fit" the data
3. Choose a shape so that assumptions SR1-SR6 are satisfied, ensuring that the least squares estimators have the desirable properties described in Chapters 2 and 3.

Principles of Econometrics, 3rd Edition

Slide 4-25

4.3.4 Are the Regression Errors Normally Distributed?

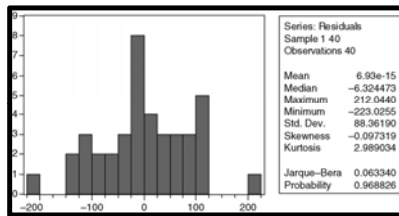


Figure 4.6 EViews output: residuals histogram and summary statistics for food expenditure example

Principles of Econometrics, 3rd Edition

Slide 4-26

4.3.4 Are the Regression Errors Normally Distributed?

- The Jarque-Bera statistic is given by

$$JB = \frac{N}{6} \left(S^2 + \frac{(K-3)^2}{4} \right)$$

where N is the sample size, S is skewness, and K is kurtosis.

- In the food expenditure example

$$JB = \frac{40}{6} \left(-0.097^2 + \frac{(2.99-3)^2}{4} \right) = .063$$

Principles of Econometrics, 3rd Edition

Slide 4-27

4.3.5 Another Empirical Example

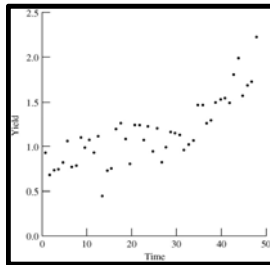


Figure 4.7 Scatter plot of wheat yield over time

Principles of Econometrics, 3rd Edition

Slide 4-28

4.3.5 Another Empirical Example

$$YIELD_t = \beta_1 + \beta_2 TIME_t + e_t$$

$$\widehat{YIELD}_t = 638 + .0210 TIME_t \quad R^2 = .649$$

(se) (.064) (.0022)

Principles of Econometrics, 3rd Edition

Slide 4-29

4.3.5 Another Empirical Example

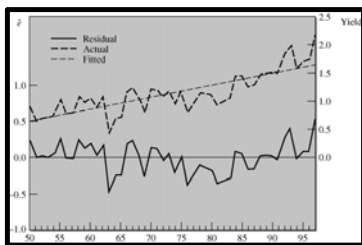


Figure 4.8 Predicted, actual and residual values from straight line

Principles of Econometrics, 3rd Edition

Slide 4-30

4.3.5 Another Empirical Example

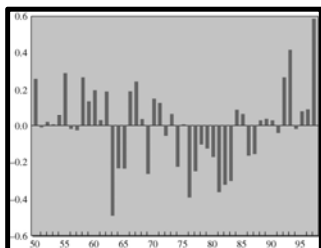


Figure 4.9 Bar chart of residuals from straight line

Principles of Econometrics, 3rd Edition

Slide 4-31

4.3.5 Another Empirical Example

$$YIELD_t = \beta_1 + \beta_2 TIME_t^3 + e_t$$

$$TIMECUBE = TIME^3 / 1000000$$

$$\widehat{YIELD}_t = 0.874 + 9.68 TIMECUBE_t \quad R^2 = 0.751$$

(se) (.036) (.082)

Principles of Econometrics, 3rd Edition

Slide 4-32

4.3.5 Another Empirical Example

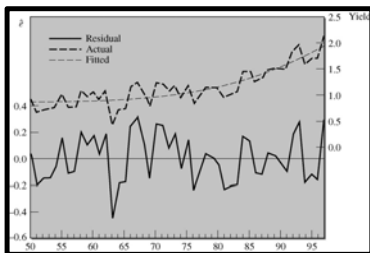


Figure 4.10 Fitted, actual and residual values from equation with cubic term

Principles of Econometrics, 3rd Edition

Slide 4-33

4.4 Log-Linear Models

- 4.4.1 The Growth Model

$$\ln(YIELD_t) = \ln(YIELD_0) + \ln(1+g)t$$

$$= \beta_1 + \beta_2 t$$

$$\widehat{\ln(YIELD_t)} = -.3434 + .0178t$$

(se) (.0584) (.0021)

Principles of Econometrics, 3rd Edition Slide 4-34

4.4 Log-Linear Models

- 4.4.2 A Wage Equation

$$\ln(WAGE) = \ln(WAGE_0) + \ln(1+r)EDUC$$

$$= \beta_1 + \beta_2 EDUC$$

$$\widehat{\ln(WAGE)} = .7884 + .1038 \times EDUC$$

(se) (.0849) (.0063)

Principles of Econometrics, 3rd Edition Slide 4-35

4.4 Log-Linear Models

- 4.4.3 Prediction in the Log-Linear Model

$$\hat{y}_n = \exp(\widehat{\ln(y)}) = \exp(b_1 + b_2 x)$$

$$\hat{y}_c = E(y) = \exp(b_1 + b_2 x + \hat{\sigma}^2/2) = \hat{y}_n e^{\hat{\sigma}^2/2}$$

$$\widehat{\ln(WAGE)} = .7884 + .1038 \times EDUC = .7884 + .1038 \times 12 = 2.0335$$

$$\hat{y}_c = E(y) = \hat{y}_n e^{\hat{\sigma}^2/2} = 7.6408 \times 1.1276 = 8.6161$$

Principles of Econometrics, 3rd Edition Slide 4-36

4.4 Log-Linear Models

- 4.4.4 A Generalized R^2 Measure

$$R_g^2 = [\text{corr}(y, \hat{y})]^2 = r_{y,\hat{y}}^2$$

$$R_g^2 = [\text{corr}(y, \hat{y}_e)]^2 = .4739^2 = .2246$$

R^2 values tend to be small with microeconomic, cross-sectional data, because the variations in individual behavior are difficult to fully explain.

Principles of Econometrics, 3rd Edition Slide 4-37

4.4 Log-Linear Models

- 4.4.5 Prediction Intervals in the Log-Linear Model

$$\left[\exp(\ln(\hat{y}) - t_{\alpha} \text{se}(f)), \exp(\ln(\hat{y}) + t_{\alpha} \text{se}(f)) \right]$$

$$\left[\exp(2.0335 - 1.96 \times .4905), \exp(2.0335 + 1.96 \times .4905) \right] = [2.9184, 20.0046]$$

Principles of Econometrics, 3rd Edition Slide 4-38

Keywords

- coefficient of determination
- correlation
- data scale
- forecast error
- forecast standard error
- functional form
- goodness-of-fit
- growth model
- Jarque-Bera test
- kurtosis
- least squares predictor
- linear model
- linear relationship
- linear-log model
- log-linear model
- log-log model
- log-normal distribution
- prediction
- prediction interval
- R^2
- residual
- skewness

Principles of Econometrics, 3rd Edition Slide 4-39

Chapter 4 Appendices

- **Appendix 4A** Development of a Prediction Interval
- **Appendix 4B** The Sum of Squares Decomposition
- **Appendix 4C** The Log-Normal Distribution

Principles of Econometrics, 3rd Edition

Slide 4-40

Appendix 4A Development of a Prediction Interval

$$f = y_0 - \hat{y}_0 = (\beta_1 + \beta_2 x_0 + e_0) - (b_1 + b_2 x_0)$$

$$\begin{aligned} \text{var}(\hat{y}_0) &= \text{var}(b_1 + b_2 x_0) = \text{var}(b_1) + x_0^2 \text{var}(b_2) + 2x_0 \text{cov}(b_1, b_2) \\ &= \frac{\sigma^2 \sum x_i^2}{N \sum (x_i - \bar{x})^2} + x_0^2 \frac{\sigma^2}{\sum (x_i - \bar{x})^2} + 2x_0 \sigma^2 \frac{-\bar{x}}{\sum (x_i - \bar{x})^2} \end{aligned}$$

Principles of Econometrics, 3rd Edition

Slide 4-41

Appendix 4A Development of a Prediction Interval

$$\begin{aligned} \text{var}(\hat{y}_0) &= \left[\frac{\sigma^2 \sum x_i^2}{N \sum (x_i - \bar{x})^2} - \left\{ \frac{\sigma^2 N \bar{x}^2}{N \sum (x_i - \bar{x})^2} \right\} \right] + \left[\frac{\sigma^2 x_0^2}{\sum (x_i - \bar{x})^2} + \frac{\sigma^2 (-2x_0 \bar{x})}{\sum (x_i - \bar{x})^2} + \left\{ \frac{\sigma^2 N \bar{x}^2}{N \sum (x_i - \bar{x})^2} \right\} \right] \\ &= \sigma^2 \left[\frac{\sum x_i^2 - N \bar{x}^2}{N \sum (x_i - \bar{x})^2} + \frac{x_0^2 - 2x_0 \bar{x} + \bar{x}^2}{\sum (x_i - \bar{x})^2} \right] \\ &= \sigma^2 \left[\frac{\sum (x_i - \bar{x})^2}{N \sum (x_i - \bar{x})^2} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \\ &= \sigma^2 \left[\frac{1}{N} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \end{aligned}$$

Principles of Econometrics, 3rd Edition

Slide 4-42

Appendix 4A
Development of a Prediction Interval

$$\frac{f}{\sqrt{\text{var}(f)}} \sim N(0,1)$$

$$\widehat{\text{var}}(f) = \hat{\sigma}^2 \left[1 + \frac{1}{N} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

$$\frac{f}{\sqrt{\text{var}(f)}} = \frac{y_0 - \hat{y}_0}{\text{se}(f)} \sim t_{(N-2)} \quad (4A.1)$$

$$P(-t_c \leq t \leq t_c) = 1 - \alpha \quad (4A.2)$$

Principles of Econometrics, 3rd Edition Slide 4-43

Appendix 4A
Development of a Prediction Interval

$$P[-t_c \leq \frac{y_0 - \hat{y}_0}{\text{se}(f)} \leq t_c] = 1 - \alpha$$

$$P[\hat{y}_0 - t_c \text{se}(f) \leq y_0 \leq \hat{y}_0 + t_c \text{se}(f)] = 1 - \alpha$$

Principles of Econometrics, 3rd Edition Slide 4-44

Appendix 4B
The Sum of Squares Decomposition

$$(y_i - \bar{y})^2 = [(\hat{y}_i - \bar{y}) + \hat{e}_i]^2 = (\hat{y}_i - \bar{y})^2 + \hat{e}_i^2 + 2(\hat{y}_i - \bar{y})\hat{e}_i$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum \hat{e}_i^2 + 2\sum (\hat{y}_i - \bar{y})\hat{e}_i$$

$$\begin{aligned} \sum (\hat{y}_i - \bar{y})\hat{e}_i &= \sum \hat{y}_i \hat{e}_i - \bar{y} \sum \hat{e}_i = \sum (b_1 + b_2 x_i) \hat{e}_i - \bar{y} \sum \hat{e}_i \\ &= b_1 \sum \hat{e}_i + b_2 \sum x_i \hat{e}_i - \bar{y} \sum \hat{e}_i \end{aligned}$$

Principles of Econometrics, 3rd Edition Slide 4-45

Appendix 4B
The Sum of Squares Decomposition

$$\sum \hat{e}_i = \sum (y_i - b_1 - b_2 x_i) = \sum y_i - Nb_1 - b_2 \sum x_i = 0$$

$$\sum x_i \hat{e}_i = \sum x_i (y_i - b_1 - b_2 x_i) = \sum x_i y_i - b_1 \sum x_i - b_2 \sum x_i^2 = 0$$

$$\sum (\hat{y}_i - \bar{y}) \hat{e}_i = 0$$

If the model contains an intercept it is guaranteed that $SST = SSR + SSE$.
If, however, the model does not contain an intercept, then $\sum \hat{e}_i \neq 0$ and $SST \neq SSR + SSE$.

Principles of Econometrics, 3rd Edition Slide 4-46

Appendix 4C
The Log-Normal Distribution

Suppose that the variable y has a normal distribution, with mean μ and variance σ^2 .
If we consider $w = e^y$ then $y = \ln(w) \sim N(\mu, \sigma^2)$ is said to have a **log-normal** distribution.

$$E(w) = e^{\mu + \sigma^2/2}$$

$$\text{var}(w) = e^{2\mu + \sigma^2} (e^{\sigma^2} - 1)$$

Principles of Econometrics, 3rd Edition Slide 4-47

Appendix 4C
The Log-Normal Distribution

Given the log-linear model $\ln(y) = \beta_1 + \beta_2 x + e$
If we assume that $e \sim N(0, \sigma^2)$

$$E(y_i) = E(e^{\beta_1 + \beta_2 x_i + e_i}) = E(e^{\beta_1 + \beta_2 x_i} e^{e_i}) = e^{\beta_1 + \beta_2 x_i} E(e^{e_i}) = e^{\beta_1 + \beta_2 x_i} e^{\sigma^2/2} = e^{\beta_1 + \beta_2 x_i + \sigma^2/2}$$

$$\widehat{E}(y_i) = e^{\beta_1 + \beta_2 x_i + \hat{\sigma}^2/2}$$

Principles of Econometrics, 3rd Edition Slide 4-48

Appendix 4C The Log-Normal Distribution

The growth and wage equations:

$$\beta_2 = \ln(1+r) \quad \text{and} \quad r = e^{\beta_2} - 1$$

$$b_2 \sim N(\beta_2, \text{var}(b_2) = \hat{\sigma}^2 / \sum (x_i - \bar{x})^2)$$

$$E[e^{\beta_2}] = e^{\beta_2 + \text{var}(b_2)/2} \quad \hat{r} = e^{\beta_2 + \text{var}(b_2)/2} - 1$$

$$\widehat{\text{var}}(b_2) = \hat{\sigma}^2 / \sum (x_i - \bar{x})^2$$

Principles of Econometrics, 3rd Edition

Slide 4-49
