

Chapter 2

The Simple Linear Regression Model: Specification and Estimation

Walter R. Paczkowski
Rutgers University

Chapter Contents

- 2.1 An Economic Model
- 2.2 An Econometric Model
- 2.3 Estimating the Regression Parameters
- 2.4 Assessing the Least Squares Estimators
- 2.5 The Gauss-Markov Theorem
- 2.6 The Probability Distributions of the Least Squares Estimators
- 2.7 Estimating the Variance of the Error Term
- 2.8 Estimating Nonlinear Relationships
- 2.9 Regression with Indicator Variables

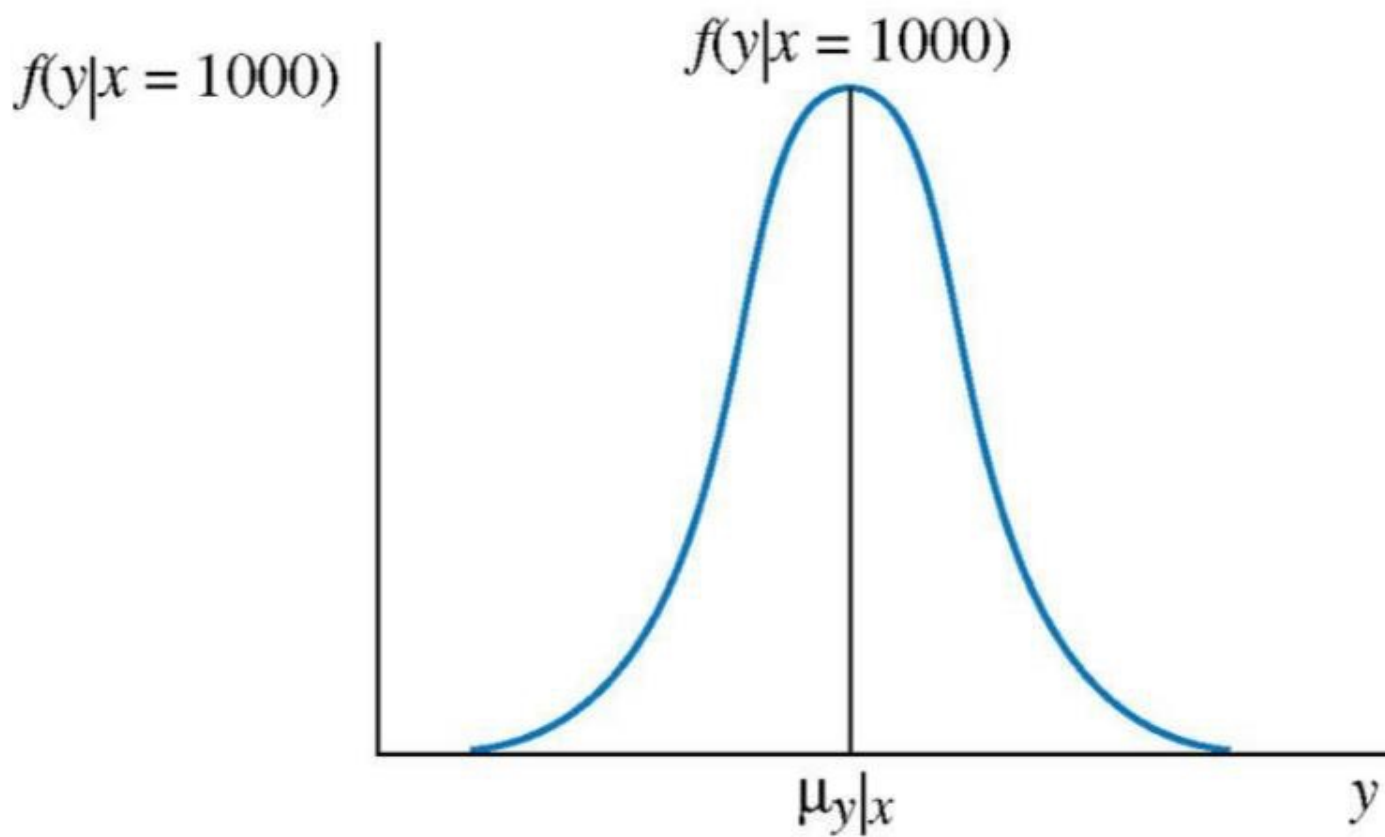
2.1

An Economic Model

- As economists we are usually more interested in studying relationships between variables
 - Economic theory tells us that expenditure on economic goods depends on income
 - Consequently we call y the “dependent variable” and x the independent” or “explanatory” variable
 - In econometrics, we recognize that real-world expenditures are **random variables**, and we want to use data to learn about the relationship

- The *pdf* is a conditional probability density function since it is “conditional” upon an x
 - The conditional mean, or expected value, of y is $E(y|x)$
 - The expected value of a random variable is called its “mean” value, which is really a contraction of population mean, the center of the probability distribution of the random variable
 - This is not the same as the sample mean, which is the arithmetic average of numerical values

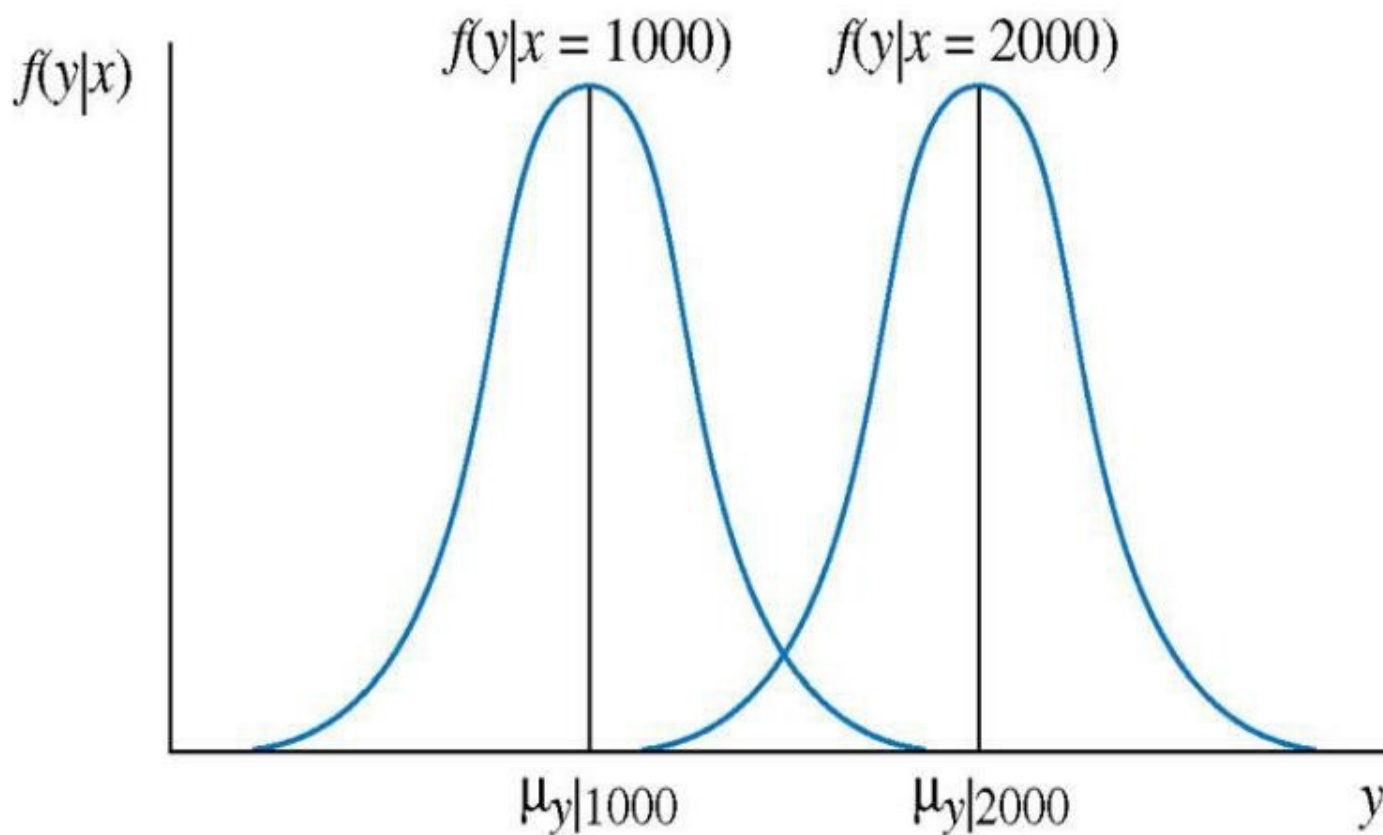
Figure 2.1a Probability distribution of food expenditure y given
income $x = \$1000$



(a)

- The conditional variance of y is σ^2 which measures the dispersion of y about its mean $\mu_{y|x}$
 - The parameters $\mu_{y|x}$ and σ^2 , if they were known, would give us some valuable information about the population we are considering

Figure 2.1b Probability distributions of food expenditures y
given incomes $x = \$1000$ and $x = \$2000$



(b)

- In order to investigate the relationship between expenditure and income we must build an economic model and then a corresponding econometric model that forms the basis for a quantitative or empirical economic analysis
 - This econometric model is also called a **regression model**

- The simple regression function is written as

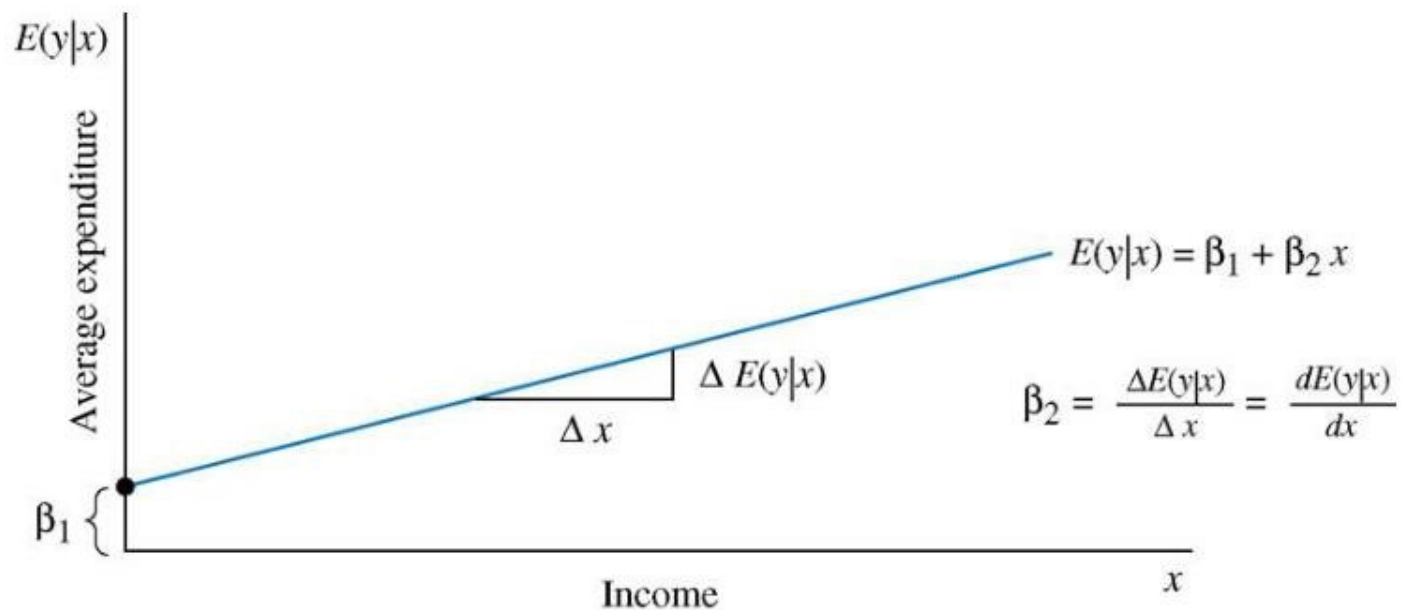
$$E(y | x) = \mu_y = \beta_1 + \beta_2 x$$

Eq. 2.1

where β_1 is the intercept and β_2 is the slope

- It is called simple regression not because it is easy, but because there is only one explanatory variable on the right-hand side of the equation

Figure 2.2 The economic model: a linear relationship between average per person food expenditure and income



- The slope of the regression line can be written as:

Eq. 2.2

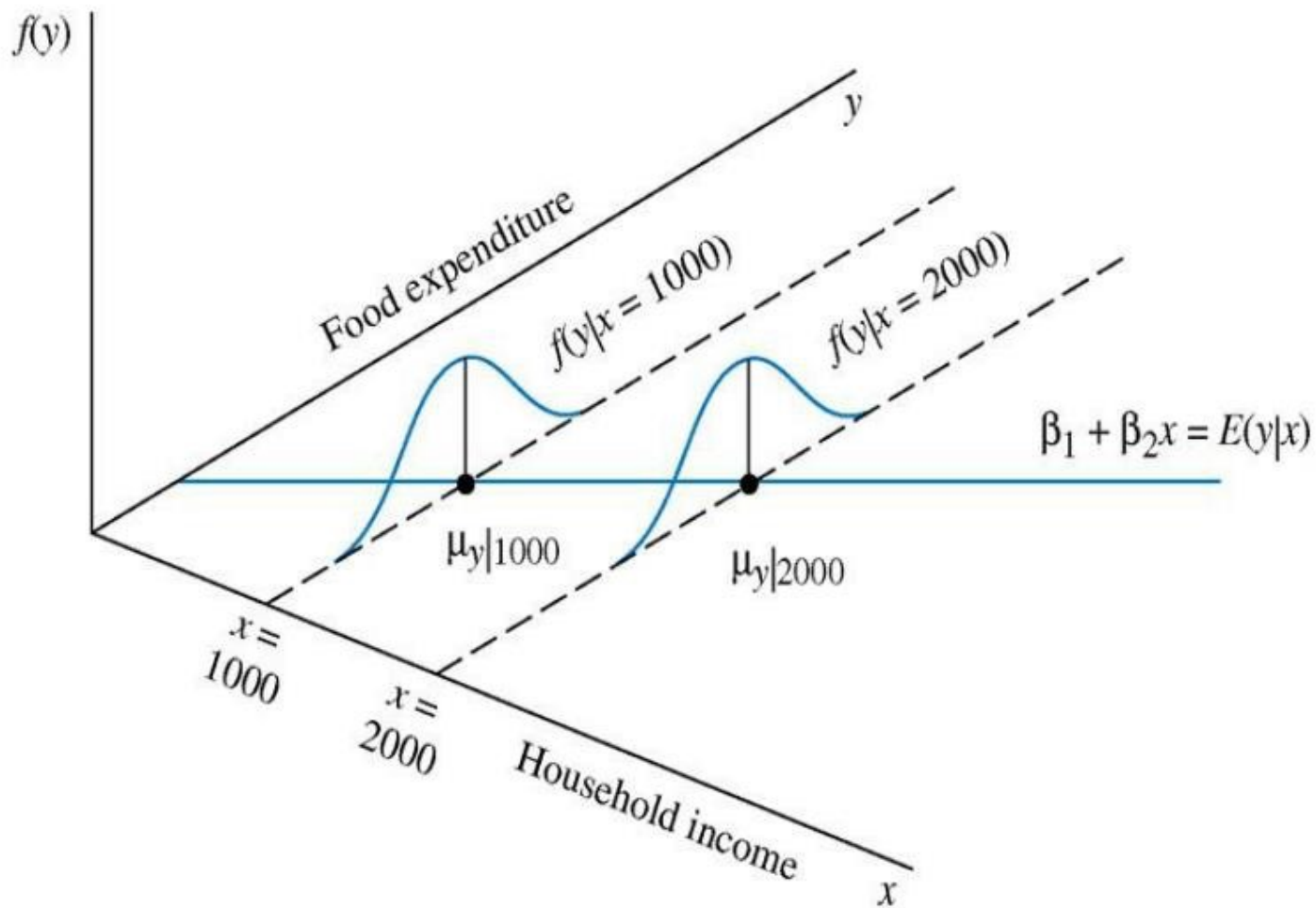
$$\beta_2 = \frac{\Delta E(y | x)}{\Delta x} = \frac{dE(y | x)}{dx}$$

where “ Δ ” denotes “change in” and “ $dE(y|x)/dx$ ” denotes the derivative of the expected value of y given an x value

2.2

An Econometric Model

Figure 2.3 The probability density function for y at two levels of income



- There are several key assumptions underlying the simple linear regression
 - More will be added later

Assumption 1:

The mean value of y , for each value of x , is given by the *linear regression*

$$E(y | x) = \beta_1 + \beta_2 x$$

Assumption 2:

For each value of x , the values of y are distributed about their mean value, following probability distributions that all have the same variance

$$\text{var}(y | x) = \sigma^2$$

Assumption 3:

The sample values of y are all *uncorrelated*, and have zero *covariance*, implying that there is no linear association among them

$$\text{cov}(y_i, y_j) = 0$$

This assumption can be made stronger by assuming that the values of y are all statistically independent

Assumption 4:

The variable x is not random, and must take at least two different values

Assumption 5:

(optional) The values of y are *normally distributed* about their mean for each value of x

$$y \sim N(\beta_1 + \beta_2 x, \sigma^2)$$

■ The random error term is defined as

Eq. 2.3

$$e = y - E(y | x) = y - \beta_1 - \beta_2 x$$

– Rearranging gives

Eq. 2.4

$$y = \beta_1 + \beta_2 x + e$$

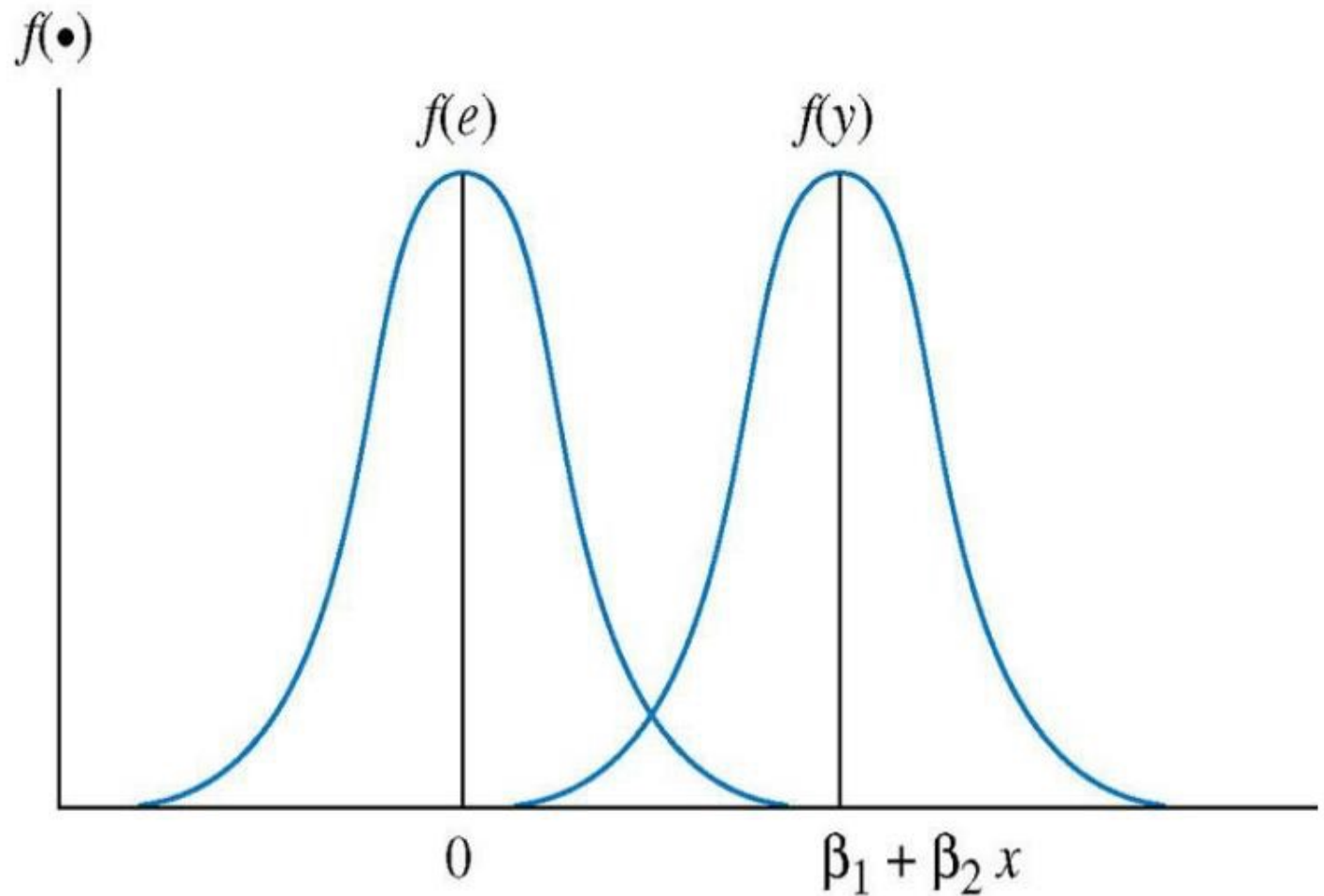
where y is the dependent variable and x is the independent variable

- The expected value of the error term, given x , is

$$E(e | x) = E(y | x) - \beta_1 - \beta_2 x = 0$$

The mean value of the error term, given x , is zero

Figure 2.4 Probability density functions for e and y



Assumption SR1:

The value of y , for each value of x , is:

$$y = \beta_1 + \beta_2 x + e$$

Assumption SR2:

The expected value of the random error e is:

$$E(e) = 0$$

This is equivalent to assuming that

$$E(y) = \beta_1 + \beta_2 x$$

Assumption SR3:

The variance of the random error e is:

$$\text{var}(e) = \sigma^2 = \text{var}(y)$$

The random variables y and e have the same variance because they differ only by a constant.

Assumption SR4:

The covariance between any pair of random errors, e_i and e_j is:

$$\text{cov}(e_i, e_j) = \text{cov}(y_i, y_j) = 0$$

The stronger version of this assumption is that the random errors e are statistically independent, in which case the values of the dependent variable y are also statistically independent

Assumption SR5:

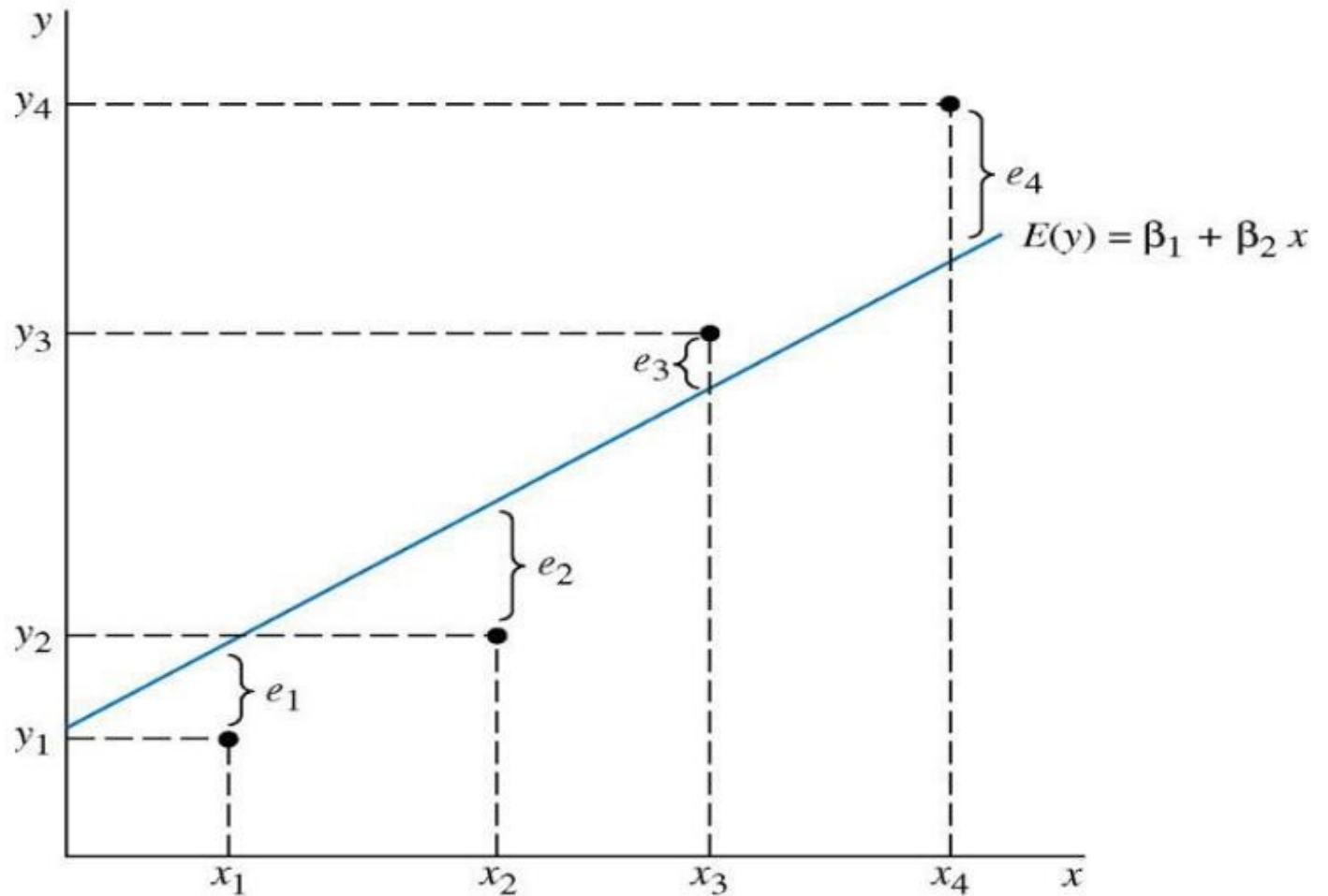
The variable x is not random, and must take at least two different values

Assumption SR6:

(optional) The values of e are *normally distributed* about their mean if the values of y are normally distributed, and *vice versa*

$$e \sim N(0, \sigma^2)$$

Figure 2.5 The relationship among y , e and the true regression line



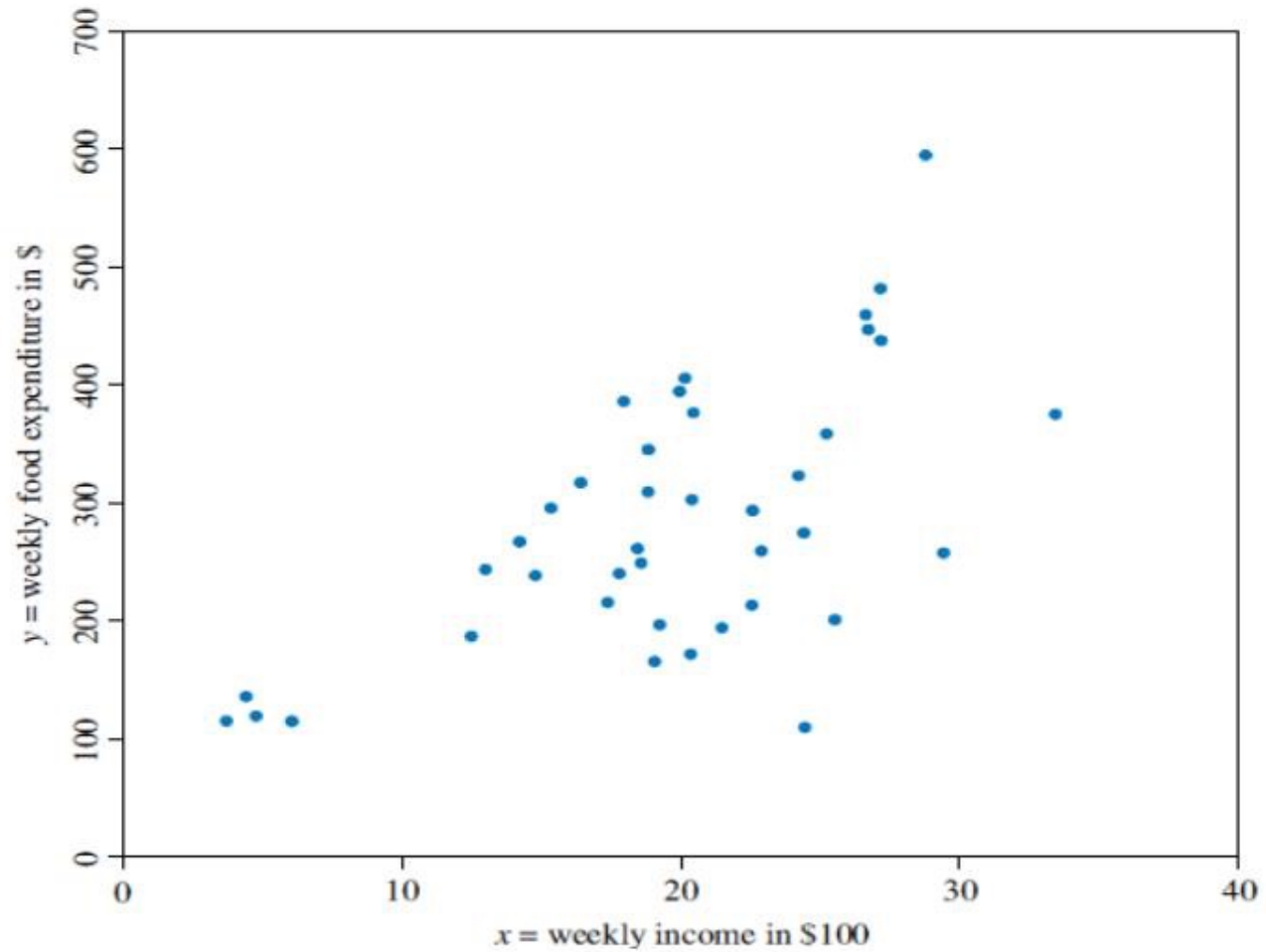
2.3

Estimating the Regression Parameters

Table 2.1 Food Expenditure and Income Data

Observation (household)	Food expenditure (\$)	Weekly income (\$100)
i	y_i	x_i
1	115.22	3.69
2	135.98	4.39
	\vdots	
39	257.95	29.40
40	375.73	33.40
Summary statistics		
Sample mean	283.5735	19.6048
Median	264.4800	20.0300
Maximum	587.6600	33.4000
Minimum	109.7100	3.6900
Std. Dev.	112.6752	6.8478

Figure 2.6 Data for food expenditure example



- The fitted regression line is:

Eq. 2.5

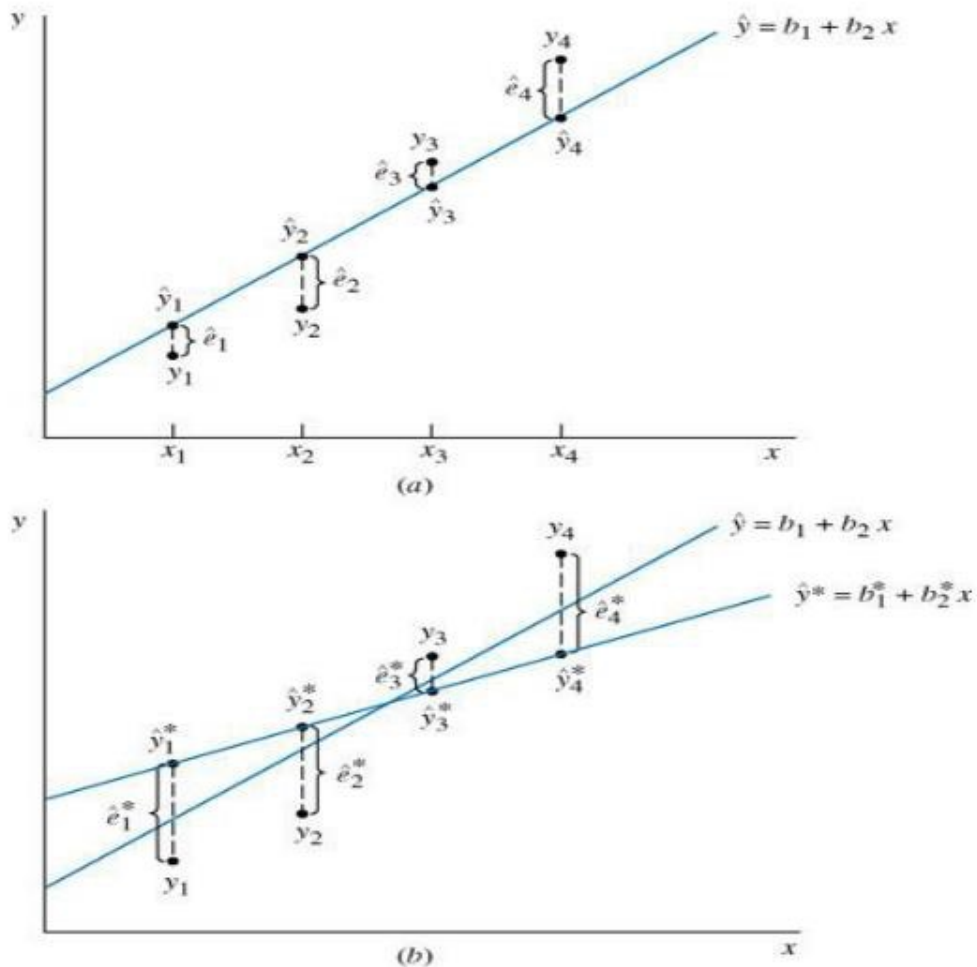
$$\hat{y}_i = b_1 + b_2 x_i$$

The least squares residual is:

Eq. 2.6

$$\hat{e}_i = y_i - \hat{y}_i = y_i - b_1 - b_2 x_i$$

Figure 2.7 The relationship among y , \hat{e} and the fitted regression line



- Suppose we have another fitted line:

$$\hat{y}_i^* = b_1^* + b_2^* x_i$$

The least squares line has the smaller sum of squared residuals:

$$\text{if } SSE = \sum_{i=1}^N \hat{e}_i^2 \text{ and } SSE^* = \sum_{i=1}^N \hat{e}_i^{*2} \text{ then } SSE < SSE^*$$

- Least squares estimates for the unknown parameters β_1 and β_2 are obtained by minimizing the sum of squares function:

$$S(\beta_1, \beta_2) = \sum_{i=1}^N (y_i - \beta_1 - \beta_2 x_i)^2$$

THE LEAST SQUARES ESTIMATORS

Eq. 2.7

$$b_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Eq. 2.8

$$b_1 = \bar{y} - b_2 \bar{x}$$

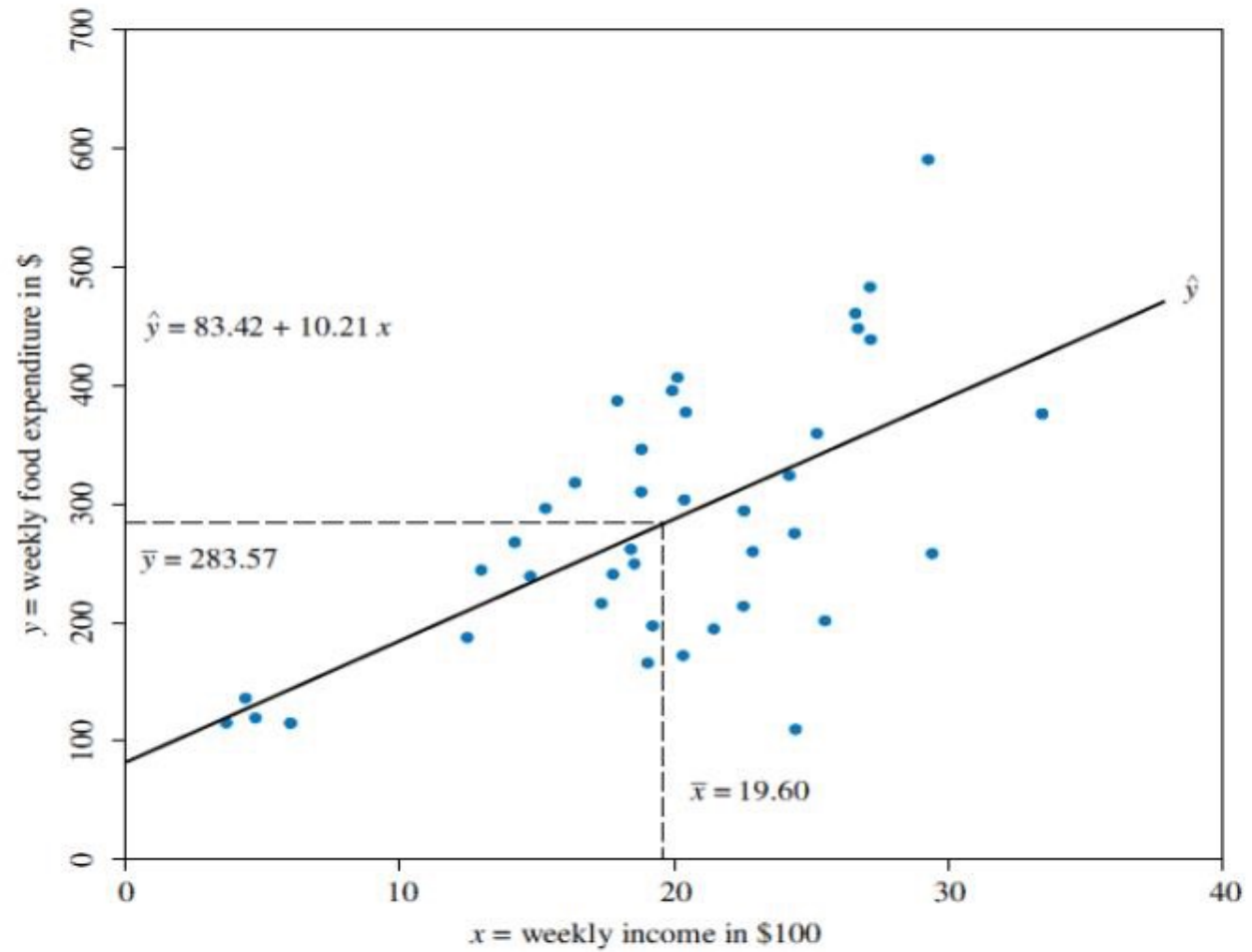
$$b_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{18671.2684}{1828.7876} = 10.2096$$

$$b_1 = \bar{y} - b_2\bar{x} = 283.5735 - (10.2096)(19.6048) = 83.4160$$

A convenient way to report the values for b_1 and b_2 is to write out the *estimated* or *fitted* regression line:

$$\hat{y}_i = 83.42 + 10.21x_i$$

Figure 2.8 The fitted regression line



- The value $b_2 = 10.21$ is an estimate of β_2 , the amount by which weekly expenditure on food per household increases when household weekly income increases by \$100. Thus, we estimate that if income goes up by \$100, expected weekly expenditure on food will increase by approximately \$10.21
 - Strictly speaking, the intercept estimate $b_1 = 83.42$ is an estimate of the weekly food expenditure on food for a household with zero income

- Income elasticity is a useful way to characterize the responsiveness of consumer expenditure to changes in income. The elasticity of a variable y with respect to another variable x is:

$$\varepsilon = \frac{\text{percentage change in } y}{\text{percentage change in } x} = \frac{\Delta y}{\Delta x} \frac{x}{y}$$

In the linear economic model given by Eq. 2.1 we have shown that

$$\beta_2 = \frac{\Delta E(y)}{\Delta x}$$

Eq. 2.9

- The elasticity of mean expenditure with respect to income is:

$$\varepsilon = \frac{\Delta E(y)/E(y)}{\Delta x/x} = \frac{\Delta E(y)}{\Delta x} \frac{x}{E(y)} = \beta_2 \frac{x}{E(y)}$$

A frequently used alternative is to calculate the elasticity at the “point of the means” because it is a representative point on the regression line.

$$\hat{\varepsilon} = b_2 \frac{\bar{x}}{\bar{y}} = 10.21 \times \frac{19.60}{283.57} = 0.71$$

- Suppose that we wanted to predict weekly food expenditure for a household with a weekly income of \$2000. This prediction is carried out by substituting $x = 20$ into our estimated equation to obtain:

$$\hat{y} = 83.42 + 10.21x_i = 83.42 + 10.21(20) = 287.61$$

We predict that a household with a weekly income of \$2000 will spend \$287.61 per week on food

Figure 2.9 EViews Regression Output

Dependent Variable: *FOOD_EXP*

Method: Least Squares

Sample: 1 40

Included observations: 40

	Coefficient	Std. Error	t-Statistic	Prob.
<i>C</i>	83.41600	43.41016	1.921578	0.0622
<i>INCOME</i>	10.20964	2.093264	4.877381	0.0000
R-squared	0.385002	Mean dependent var		283.5735
Adjusted R-squared	0.368818	S.D. dependent var		112.6752
S.E. of regression	89.51700	Akaike info criterion		11.87544
Sum squared resid	304505.2	Schwarz criterion		11.95988
Log likelihood	-235.5088	Hannan-Quinn criter		11.90597
F-statistic	23.78884	Durbin-Watson stat		1.893880
Prob(F-statistic)	0.000019			

- The simple regression model can be applied to estimate the parameters of many relationships in economics, business, and the social sciences
 - The applications of regression analysis are fascinating and useful

2.4

Assessing the Least Squares Fit

- We call b_1 and b_2 the *least squares estimators*.
 - We can investigate the properties of the estimators b_1 and b_2 , which are called their sampling properties, and deal with the following important questions:
 1. If the least squares estimators are random variables, then what are their expected values, variances, covariances, and probability distributions?
 2. How do the least squares estimators compare with other procedures that might be used, and how can we compare alternative estimators?

■ The estimator b_2 can be rewritten as:

Eq. 2.10

$$b_2 = \sum_{i=1}^N w_i y_i$$

where

Eq. 2.11

$$w_i = \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2}$$

It could also be written as:

Eq. 2.12

$$b_2 = \beta_2 + \sum w_i e_i$$

- We will show that if our model assumptions hold, then $E(b_2) = \beta_2$, which means that the estimator is **unbiased**. We can find the expected value of b_2 using the fact that the expected value of a sum is the sum of the expected values:

$$\begin{aligned} E(b_2) &= E(b_2 + \sum w_i e_i) = E(\beta_2 + w_1 e_1 + w_2 e_2 + \dots + w_N e_N) \\ &= E(\beta_2) + E(w_1 e_1) + E(w_2 e_2) + \dots + E(w_N e_N) \\ &= E(\beta_2) + \sum E(w_i e_i) \\ &= \beta_2 + \sum w_i E(e_i) \\ &= \beta_2 \end{aligned}$$

using $E(e_i) = 0$ and $E(w_i e_i) = w_i E(e_i)$

Eq. 2.13

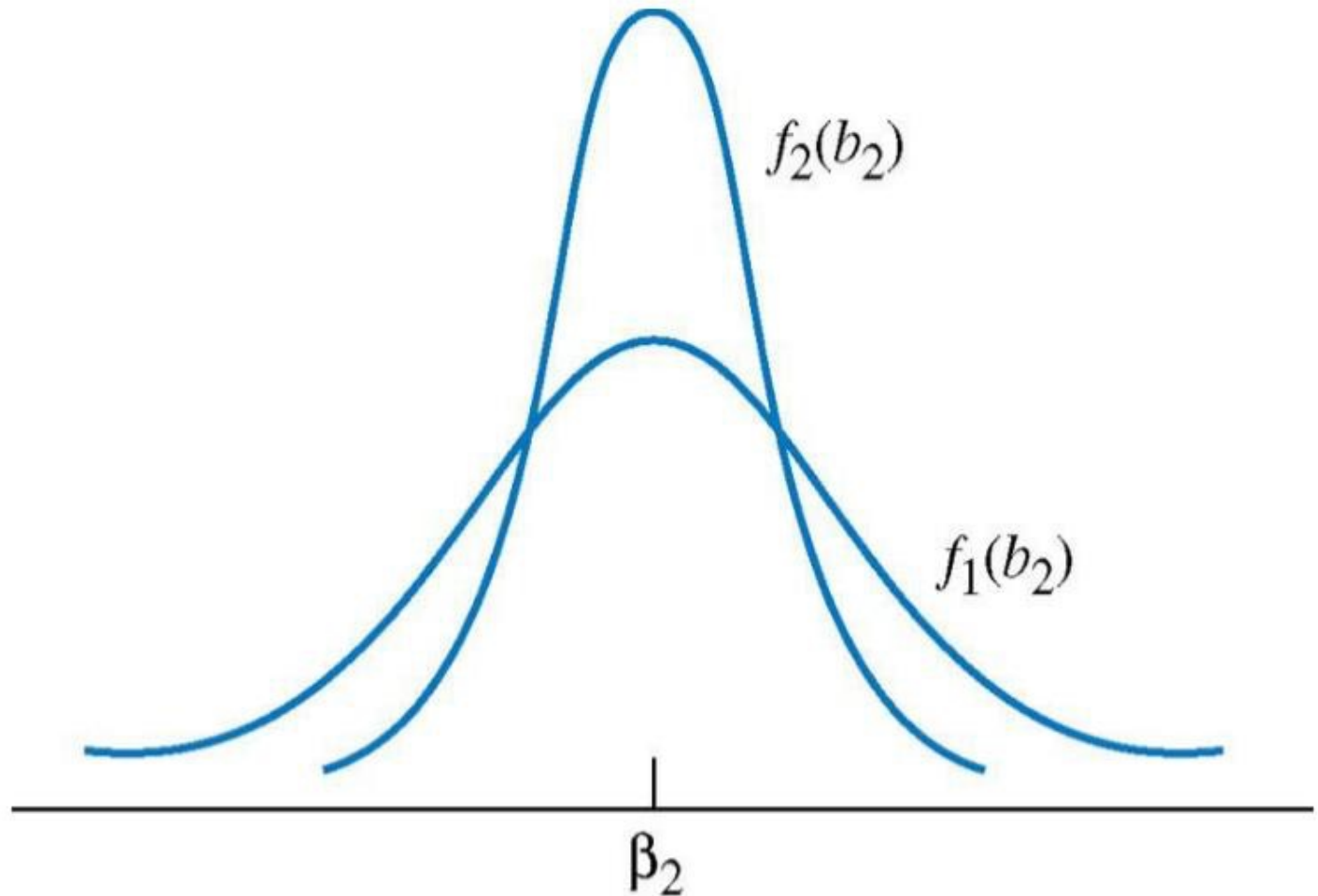
- The property of unbiasedness is about the average values of b_1 and b_2 if many samples of the same size are drawn from the same population
 - If we took the averages of estimates from many samples, these averages would approach the true parameter values b_1 and b_2
 - Unbiasedness does not say that an estimate from any one sample is close to the true parameter value, and thus we cannot say that an estimate is unbiased
 - We can say that the least squares estimation procedure (or the least squares estimator) is unbiased

Table 2.2 Estimates from 10 Samples

Sample	b_1	b_2
1	131.69	6.48
2	57.25	10.88
3	103.91	8.14
4	46.50	11.90
5	84.23	9.29
6	26.63	13.55
7	64.21	10.93
8	79.66	9.76
9	97.30	8.05
10	95.96	7.77

Figure 2.10 Two possible probability density functions for b_2

The variance of b_2 is defined as $\text{var}(b_2) = E[b_2 - E(b_2)]^2$



- If the regression model assumptions SR1-SR5 are correct (assumption SR6 is not required), then the variances and covariance of b_1 and b_2 are:

Eq. 2.14

$$\text{var}(b_1) = \sigma^2 \left[\frac{\sum x_i^2}{N \sum (x_i - \bar{x})^2} \right]$$

Eq. 2.15

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

Eq. 2.16

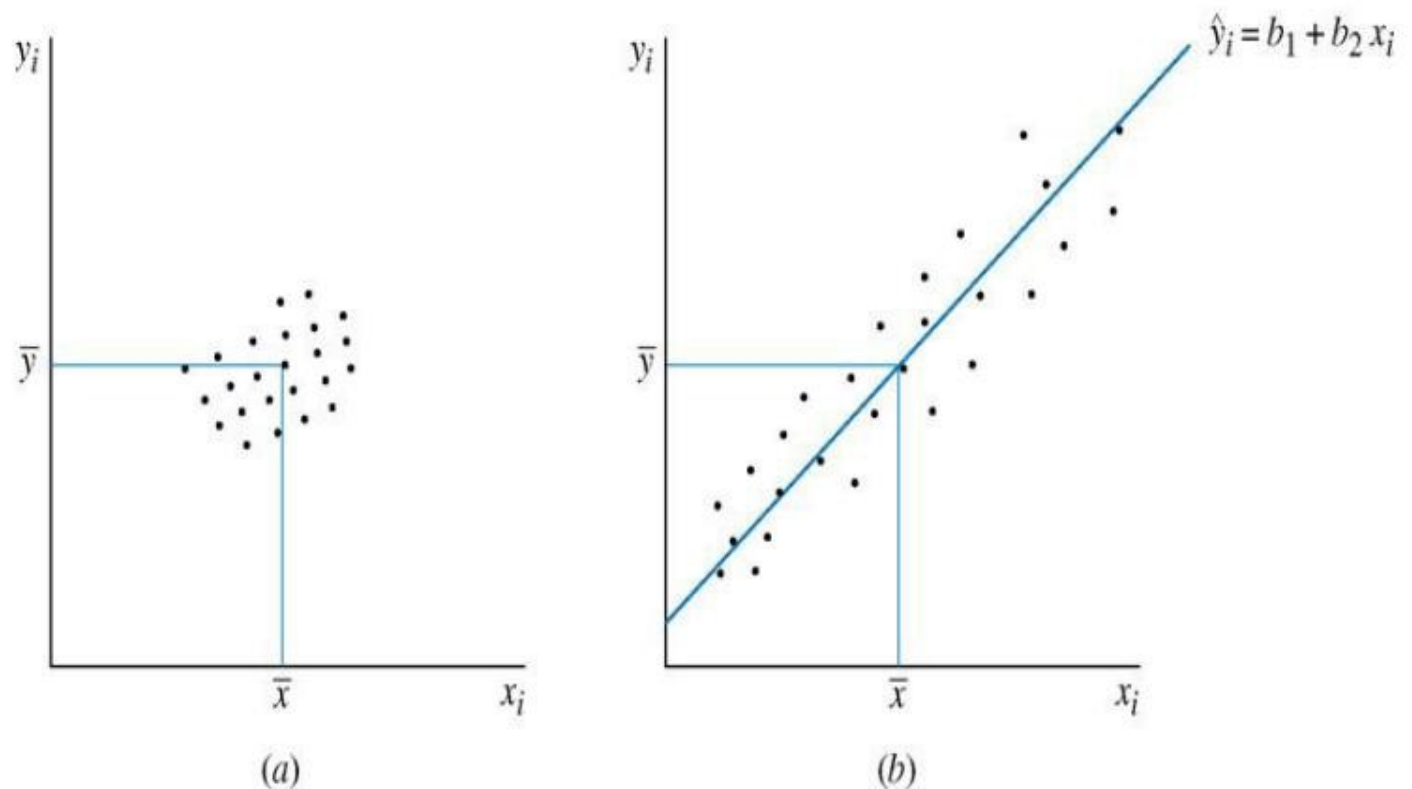
$$\text{cov}(b_1, b_2) = \sigma^2 \left[\frac{-\bar{x}}{\sum (x_i - \bar{x})^2} \right]$$

MAJOR POINTS ABOUT THE VARIANCES AND COVARIANCES OF b_1 AND b_2

1. The *larger* the variance term σ^2 , the *greater* the uncertainty there is in the statistical model, and the *larger* the variances and covariance of the least squares estimators.
2. The *larger* the sum of squares, $\sum(x_i - \bar{x})^2$, the *smaller* the variances of the least squares estimators and the more *precisely* we can estimate the unknown parameters.
3. The larger the sample size N , the *smaller* the variances and covariance of the least squares estimators.
4. The larger the term $\sum x_i^2$, the larger the variance of the least squares estimator b_1 .
5. The absolute magnitude of the covariance *increases* the larger in magnitude is the sample mean \bar{x} , and the covariance has a *sign* opposite to that of \bar{x} .

Figure 2.11 The influence of variation in the explanatory variable x on precision of estimation (a) Low x variation, low precision (b) High x variation, high precision

The variance of b_2 is defined as $\text{var}(b_2) = E[b_2 - E(b_2)]^2$



2.5

The Gauss-Markov Theorem

GAUSS-MARKOV THEOREM

Under the assumptions SR1-SR5 of the linear regression model, the estimators b_1 and b_2 have the smallest variance of all linear and unbiased estimators of b_1 and b_2 . They are the **Best Linear Unbiased Estimators (BLUE) of b_1 and b_2**

MAJOR POINTS ABOUT THE GAUSS-MARKOV THEOREM

1. The estimators b_1 and b_2 are “best” when compared to similar estimators, those which are linear and unbiased. The Theorem does *not* say that b_1 and b_2 are the best of all *possible* estimators.
2. The estimators b_1 and b_2 are best within their class because they have the minimum variance. When comparing two linear and unbiased estimators, we *always* want to use the one with the smaller variance, since that estimation rule gives us the higher probability of obtaining an estimate that is close to the true parameter value.
3. In order for the Gauss-Markov Theorem to hold, assumptions SR1-SR5 must be true. If any of these assumptions are *not* true, then b_1 and b_2 are *not* the best linear unbiased estimators of β_1 and β_2 .

MAJOR POINTS ABOUT THE GAUSS-MARKOV THEOREM

4. The Gauss-Markov Theorem does *not* depend on the assumption of normality (assumption SR6).
5. In the simple linear regression model, if we want to use a linear and unbiased estimator, then we have to do no more searching. The estimators b_1 and b_2 are the ones to use. This explains why we are studying these estimators and why they are so widely used in research, not only in economics but in all social and physical sciences as well.
6. The Gauss-Markov theorem applies to the least squares estimators. It *does not* apply to the least squares *estimates* from a single sample.

2.6

The Probability Distributions of the Least Squares Estimators

- *If we make the normality assumption (assumption SR6 about the error term) then the least squares estimators are normally distributed:*

Eq. 2.17

$$b_1 \sim N\left(\beta_1, \frac{\sigma^2 \sum x_i^2}{N \sum (x_i - \bar{x})^2}\right)$$

Eq. 2.18

$$b_2 \sim N\left(\beta_2, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$$

A CENTRAL LIMIT THEOREM

If assumptions SR1-SR5 hold, and if the sample size N is *sufficiently large*, then the least squares estimators have a distribution that approximates the normal distributions shown in Eq. 2.17 and Eq. 2.18

2.7

Estimating the Variance of the Error Term

- The variance of the random error e_i is:

$$\text{var}(e_i) = \sigma^2 = E[e_i - E(e_i)]^2 = E(e_i)^2$$

if the assumption $E(e_i) = 0$ is correct.

Since the “expectation” is an average value we might consider estimating σ^2 as the average of the squared errors:

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{N}$$

where the error terms are $e_i = y_i - \beta_1 - \beta_2 x_i$

- The least squares residuals are obtained by replacing the unknown parameters by their least squares estimates:

$$\hat{e}_i = y_i - \hat{y}_i = y_i - b_1 - b_2 x_i$$

$$\sigma^2 = \frac{\sum \hat{e}_i^2}{N}$$

There is a simple modification that produces an unbiased estimator, and that is:

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{N - 2}$$

so that:

$$E(\hat{\sigma}^2) = \sigma^2$$

Eq. 2.19

- Replace the unknown error variance σ^2 in Eq. 2.14 – Eq. 2.16 by $\hat{\sigma}^2$ to obtain:

Eq. 2.20

$$\text{var}(b_1) = \hat{\sigma}^2 \left[\frac{\sum x_i^2}{N \sum (x_i - \bar{x})^2} \right]$$

Eq. 2.21

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

Eq. 2.22

$$\text{cov}(b_1, b_2) = \hat{\sigma}^2 \left[\frac{-\bar{x}}{\sum (x_i - \bar{x})^2} \right]$$

- The square roots of the estimated variances are the “standard errors” of b_1 and b_2 :

Eq. 2.23

$$se(b_1) = \sqrt{\text{var}(b_1)}$$

Eq. 2.24

$$se(b_2) = \sqrt{\text{var}(b_2)}$$

Table 2.3 Least Squares Residuals

x	y	\hat{y}	$\hat{e} = y - \hat{y}$
3.69	115.22	121.09	-5.87
4.39	135.98	128.24	7.74
4.75	119.34	131.91	-12.57
6.03	114.96	144.98	-30.02
12.47	187.05	210.73	-23.68

$$\hat{\sigma}^2 = \frac{\sum \hat{e}_i^2}{N-2} = \frac{304505.2}{38} = 8013.29$$

- The estimated variances and covariances for a regression are arrayed in a rectangular array, or *matrix*, with variances on the diagonal and covariances in the “off-diagonal” positions.

$$\begin{bmatrix} \text{var}(b_1) & \text{cov}(b_1, b_2) \\ \text{cov}(b_1, b_2) & \text{var}(b_2) \end{bmatrix}$$

- For the food expenditure data the estimated covariance matrix is:

	C	Income
C	1884.442	-85.90316
Income	-85.90316	4.381752

- The standard errors of b_1 and b_2 are measures of the sampling variability of the least squares estimates b_1 and b_2 in repeated samples.
 - The estimators are random variables. As such, they have probability distributions, means, and variances.
 - In particular, if assumption SR6 holds, and the random error terms e_i are normally distributed, then:

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

- The estimator variance, $\text{var}(b_2)$, or its square root, $\sigma_{b_2} = \sqrt{\text{var}(b_2)}$ which we might call the true standard deviation of b_2 , measures the sampling variation of the estimates b_2
 - The bigger σ_{b_2} is the more variation in the least squares estimates b_2 we see from sample to sample. If σ_{b_2} is large then the estimates might change a great deal from sample to sample
 - If σ_{b_2} is small relative to the parameter b_2 , we know that the least squares estimate will fall near b_2 with high probability

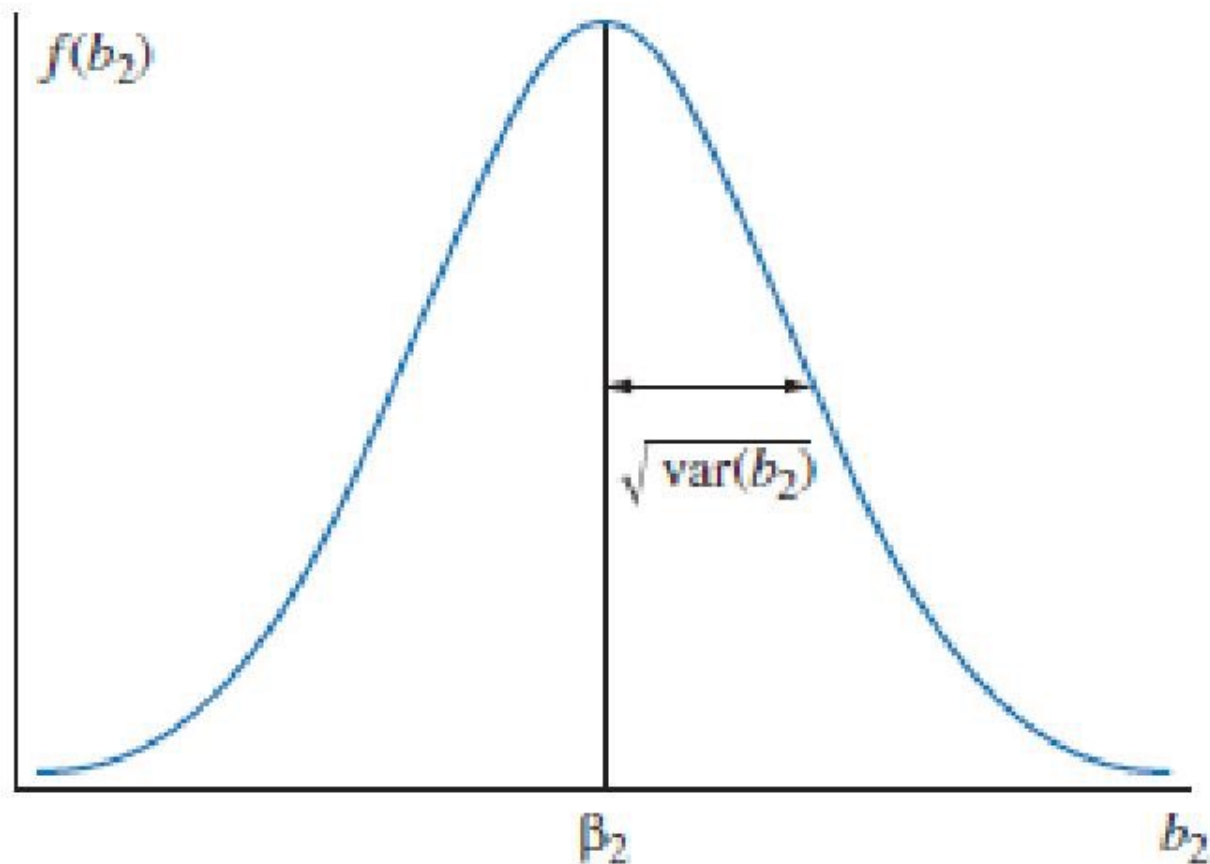
- The question we address with the standard error is “*How much variation about their means do the estimates exhibit from sample to sample?*”

- We estimate σ^2 , and then estimate σ_{b_2} using:

$$\begin{aligned} \text{se}(b_2) &= \sqrt{\text{var}(b_2)} \\ &= \sqrt{\frac{\hat{\sigma}^2}{\sum (x_i - \bar{x})^2}} \end{aligned}$$

- The standard error of b_2 is thus an estimate of what the standard deviation of many estimates b_2 would be in a very large number of samples, and is an indicator of the width of the *pdf* of b_2 shown in Figure 2.12

Figure 2.12 The probability density function of the least squares estimator b_2 .



2.8

Estimating Nonlinear Relationships

THE WORLD IS NOT LINEAR

Economic variables are not always related by straight-line relationships; in fact, many economic relationships are represented by curved lines, and are said to display *curvilinear forms*.

Fortunately, the simple linear regression model $y = \beta_1 + \beta_2 + e$ is much more flexible than it looks at first glance, because the variables y and x can be transformations, involving logarithms, squares, cubes or reciprocals, of the basic economic variables, or they can be indicator variables that take only the values zero and one.

Including these possibilities means the simple linear regression model can be used to account for nonlinear relationships between variables

- Consider the linear model of house prices:

Eq. 2.25

$$PRICE = \beta_1 + \beta_2 SQFT + e$$

where $SQFT$ is the square footage.

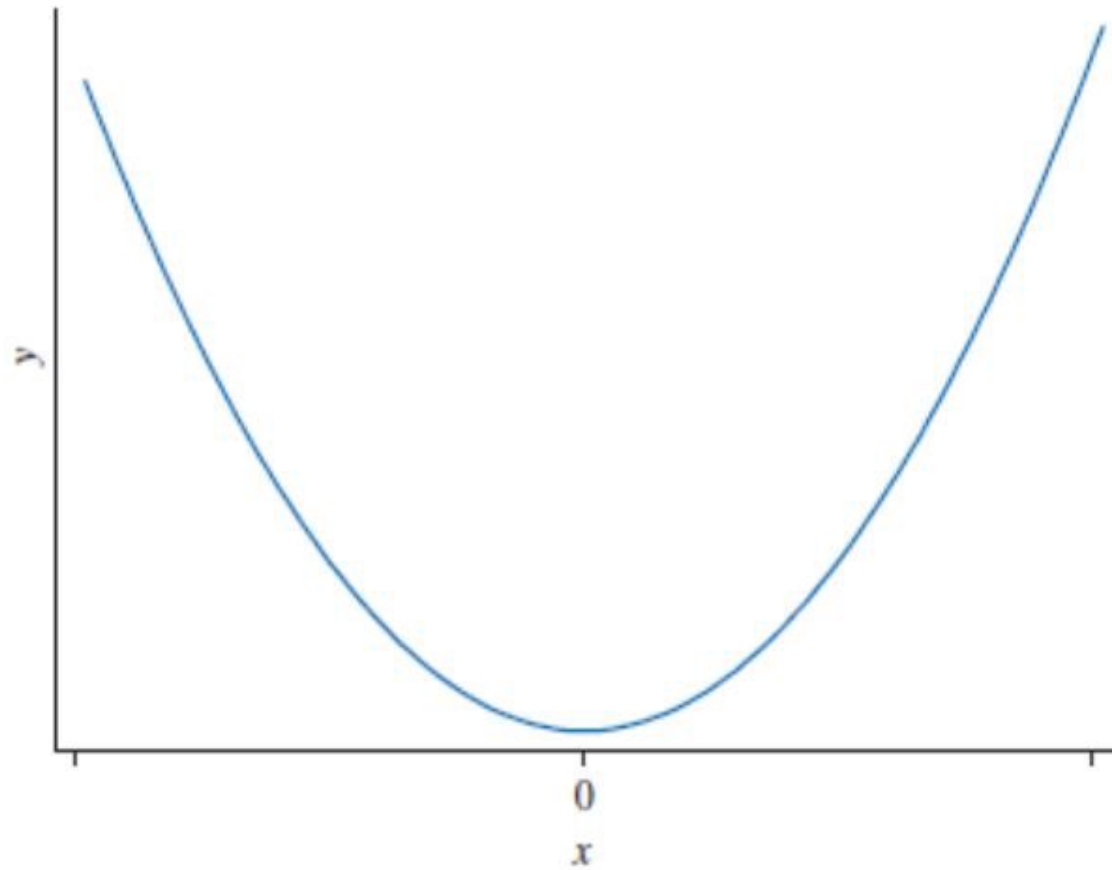
- It may be reasonable to assume that larger and more expensive homes have a higher value for an additional square foot of living area than smaller, less expensive, homes

- We can build this into our model in two ways:
 1. a quadratic equation in which the explanatory variable is $SQFT^2$
 2. a loglinear equation in which the dependent variable is $\ln(PRICE)$
- In each case we will find that the slope of the relationship between $PRICE$ and $SQFT$ is not constant, but changes from point to point.

- The quadratic function $y = \beta_1 + \beta_2 x^2$ is a parabola
 - The elasticity, or the percentage change in y given a 1% change in x , is:

$$\begin{aligned}\varepsilon &= \text{slope} \times x / y \\ &= 2bx^2 / y\end{aligned}$$

Figure 2.13 A quadratic function



- A quadratic model for house prices includes the squared value of $SQFT$, giving:

Eq. 2.26

$$PRICE = \alpha_1 + \alpha_2 SQFT^2 + e$$

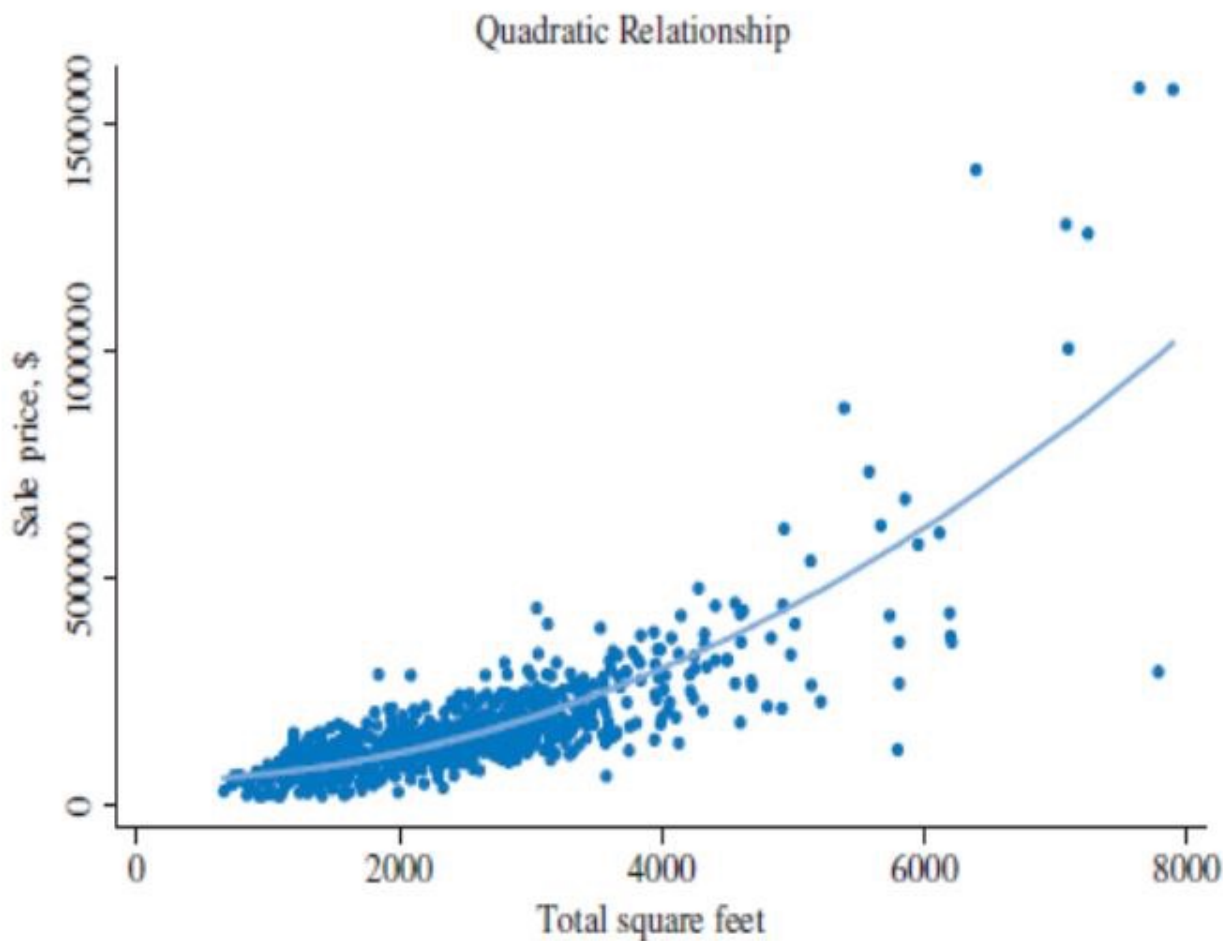
- The slope is:

Eq. 2.27

$$\frac{d(PRICE)}{dSQFT} = 2\hat{\alpha}_2 SQFT$$

- If $\hat{\alpha}_2 > 0$, then larger houses will have larger slope, and a larger estimated price per additional square foot

Figure 2.14 A fitted quadratic relationship



- For 1080 houses sold in Baton Rouge, LA during mid-2005, the estimated quadratic equation is:

$$PRICE = 55776.56 + 0.0154SQFT^2$$

- The estimated slope is:

$$slope = 2(0.0154)SQFT$$

- The elasticity is:

$$\begin{aligned}\hat{\varepsilon} &= slope \times \frac{SQFT}{PRICE} \\ &= (2\hat{\alpha}_2 SQFT) \times \frac{SQFT}{PRICE}\end{aligned}$$

- To compute an estimate we must select values for *SQFT* and *PRICE*
 - A common approach is to choose a point on the fitted relationship
 - That is, we choose a value for *SQFT* and choose for price the corresponding fitted value

- For houses of 2000, 4000 and 6000 square feet, the estimated elasticities are:

1.05 using $\widehat{\text{PRICE}} = \$117,461.77$

1.63 using $\widehat{\text{PRICE}} = \$302,517.39$

1.82 using $\widehat{\text{PRICE}} = \$610,943.42$

respectively

For a 2000-square-foot house, we estimate that a 1% increase in house size will increase price by 1.05%

■ The log-linear equation $\ln(y) = a + bx$ has a logarithmic term on the left-hand side of the equation and an untransformed (linear) variable on the right-hand side

– Both its slope and elasticity change at each point and are the same sign as b

• The slope is:

$$dy/dx = by$$

– The elasticity, the percentage change in y given a 1% increase in x , at a point on this curve is:

$$\varepsilon = slope \times x/y = bx$$

- Using the slope expression, we can solve for a **semi-elasticity**, which tells us the percentage change in y given a 1-unit increase in x :

$$\eta = \frac{100(dy/dx)}{dx} = 100b$$

Eq. 2.28

- Consider again the model for the price of a house as a function of the square footage, but now written in semi-log form:

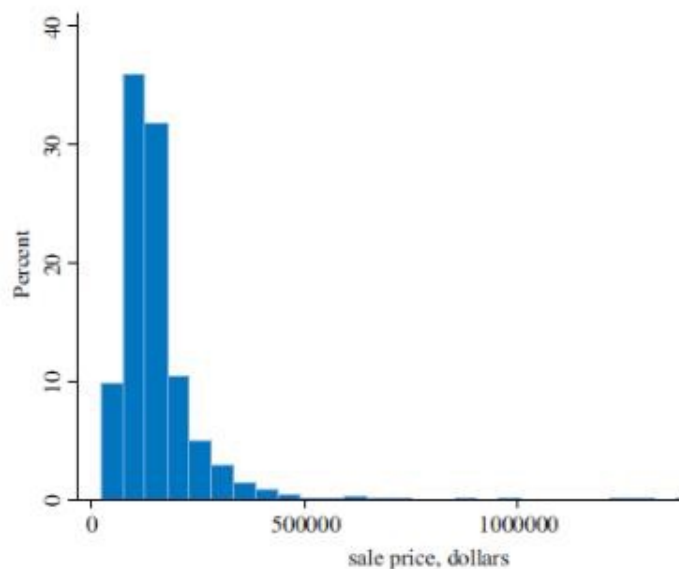
Eq. 2.29

$$\ln(PRICE) = \gamma_1 + \gamma_2 SQFT + e$$

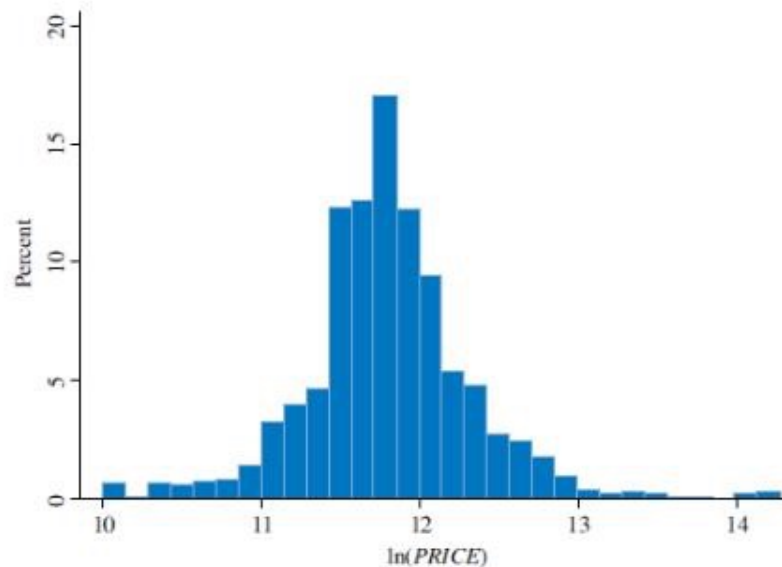
- This logarithmic transformation can regularize data that is skewed with a long tail to the right

Figure 2.16 (a) Histogram of PRICE (b) Histogram of $\ln(\text{PRICE})$

2.8.4
Using a Log-Linear
Model



(a)



(b)

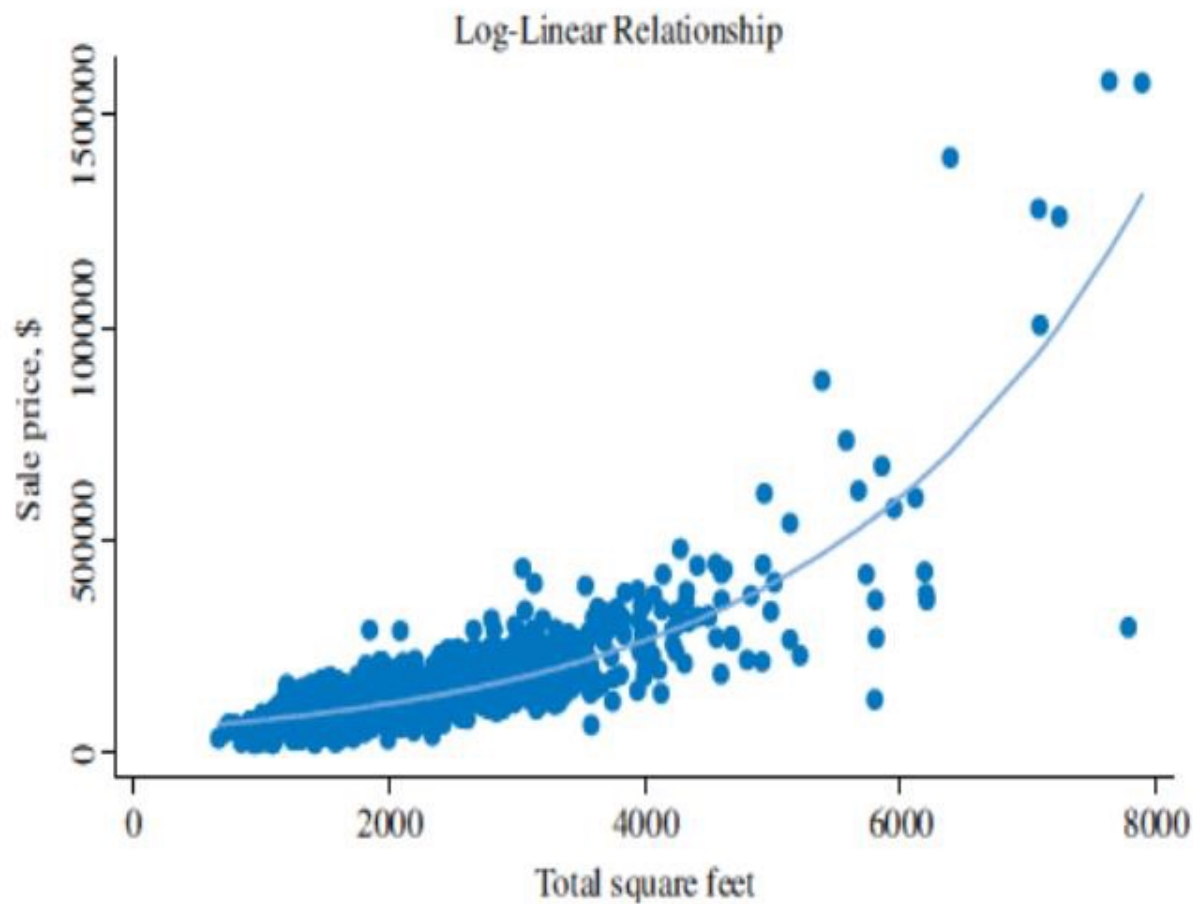
- Using the Baton Rouge data, the fitted log-linear model is:

$$\ln(PRICE) = 10.8386 + 0.0004113SQFT$$

- To obtain predicted price take the anti-logarithm, which is the exponential function:

$$PRICE = \exp\left[\ln(PRICE)\right] = \exp(10.8386 + 0.0004113SQFT)$$

Figure 2.17 The fitted log-linear model



- The slope of the log-linear model is:

$$\frac{d(PRICE)}{dSQFT} = \hat{\gamma}_2 PRICE = 0.0004113 PRICE$$

For a house with a predicted *PRICE* of \$100,000, the estimated increase in *PRICE* for an additional square foot of house area is \$41.13, and for a house with a predicted *PRICE* of \$500,000, the estimated increase in *PRICE* for an additional square foot of house area is \$205.63

■ The estimated elasticity is:

$$\hat{\varepsilon} = \hat{\gamma}_2 SQFT = 0.0004113 SQFT$$

- For a house with 2000-square-feet, the estimated elasticity is 0.823:
 - A 1% increase in house size is estimated to increase selling price by 0.823%
- For a house with 4000 square feet, the estimated elasticity is 1.645:
 - A 1% increase in house size is estimated to increase selling price by 1.645%

- Using the “semi-elasticity” defined in Eq. 2.28 we can say that, for a one-square-foot increase in size, we estimate a price increase of 0.04%
 - Or, perhaps more usefully, we estimate that a 100-square-foot increase will increase price by approximately 4%.

- We should do our best to choose a functional form that is:
 - consistent with economic theory
 - that fits the data well
 - that is such that the assumptions of the regression model are satisfied

- In real-world problems it is sometimes difficult to achieve all these goals
 - Furthermore, we will never truly know the correct functional relationship, no matter how many years we study econometrics
 - The truth is out there, but we will never know it
 - In applications of econometrics we must simply do the best we can to choose a satisfactory functional form

2.9

Regression with Indicator Variables

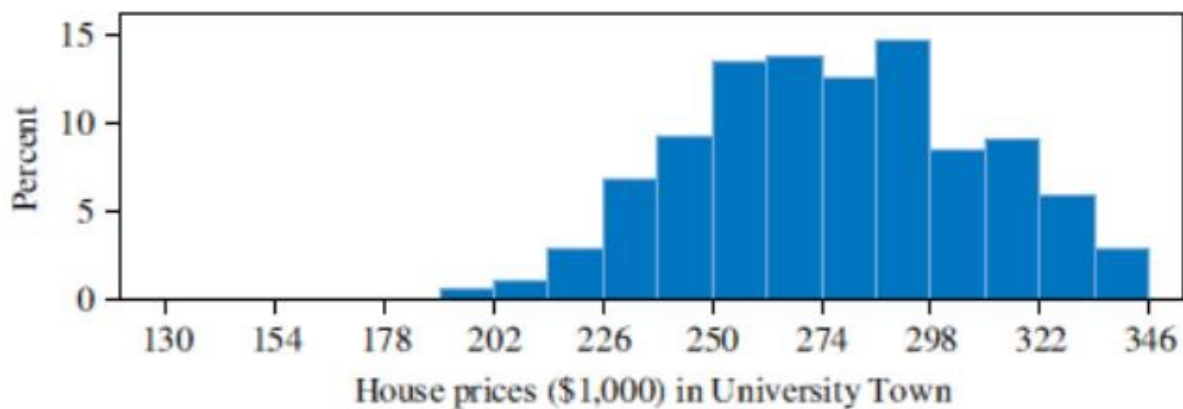
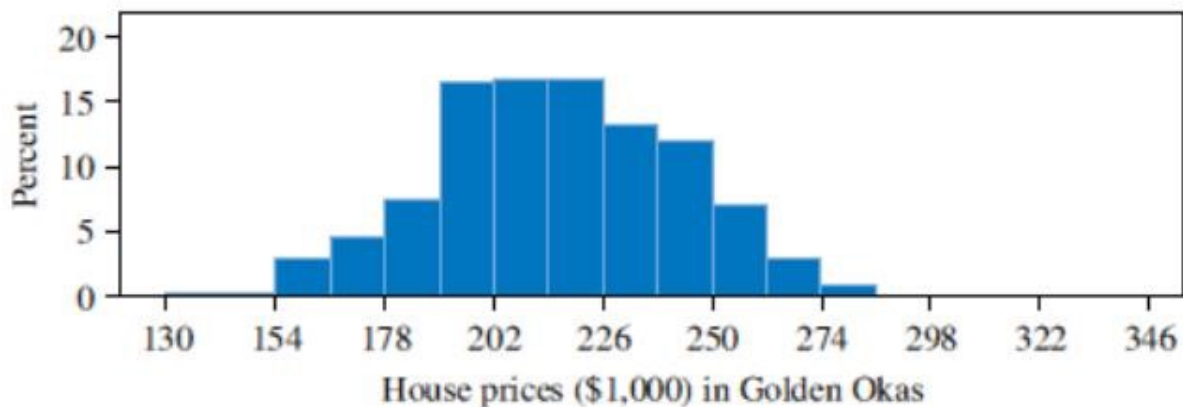
- An indicator variable is a binary variable that takes the values zero or one; it is used to represent a nonquantitative characteristic, such as gender, race, or location

$$UTOWN = \begin{cases} 1 & \text{house is in University Town} \\ 0 & \text{house is in Golden Oaks} \end{cases}$$

$$PRICE = \beta_1 + \beta_2 UTOWN + e$$

– How do we model this?

Figure 2.18 Distributions of house prices



- When an indicator variable is used in a regression, it is important to write out the regression function for the different values of the indicator variable

$$E(PRICE) = \begin{cases} \beta_1 + \beta_2 & \text{if } UTOWN = 1 \\ \beta_1 & \text{if } UTOWN = 0 \end{cases}$$

- The estimated regression is:

$$\begin{aligned} PRICE &= b_1 + b_2 UTOWN \\ &= 215.7325 + 61.5091 UTOWN \\ &= \begin{cases} 277.2416 & \text{if } UTOWN = 1 \\ 215.7325 & \text{if } UTOWN = 0 \end{cases} \end{aligned}$$

- The least squares estimators b_1 and b_2 in this indicator variable regression can be shown to be:

$$b_1 = \overline{PRICE}_{\text{Golden Oaks}}$$

$$b_2 = \overline{PRICE}_{\text{University Town}} - \overline{PRICE}_{\text{Golden Oaks}}$$

- In the simple regression model, an indicator variable on the right-hand side gives us a way to estimate the differences between population means

Key Words

- assumptions
- asymptotic
- B.L.U.E.
- biased estimator
- degrees of freedom
- dependent variable
- deviation from the mean form
- econometric model
- economic model
- elasticity
- Gauss-Markov Theorem
- heteroskedastic
- homoskedastic
- independent variable
- least squares estimates
- least squares estimators
- least squares principle
- least squares residuals
- linear estimator
- prediction
- random error term
- regression model
- regression parameters
- repeated sampling
- sampling precision
- sampling properties
- scatter diagram
- simple linear regression function
- specification error
- unbiased estimator

Appendices

- 2A Derivation of the Least Squares Estimates
- 2B Deviation from the Mean Form of b_2
- 2C b_2 is a Linear Estimator
- 2D Derivation of Theoretical Expression for b_2
- 2E Deriving the Variance of b_2
- 2F Proof of the Gauss-Markov Theorem

Eq. 2A.1

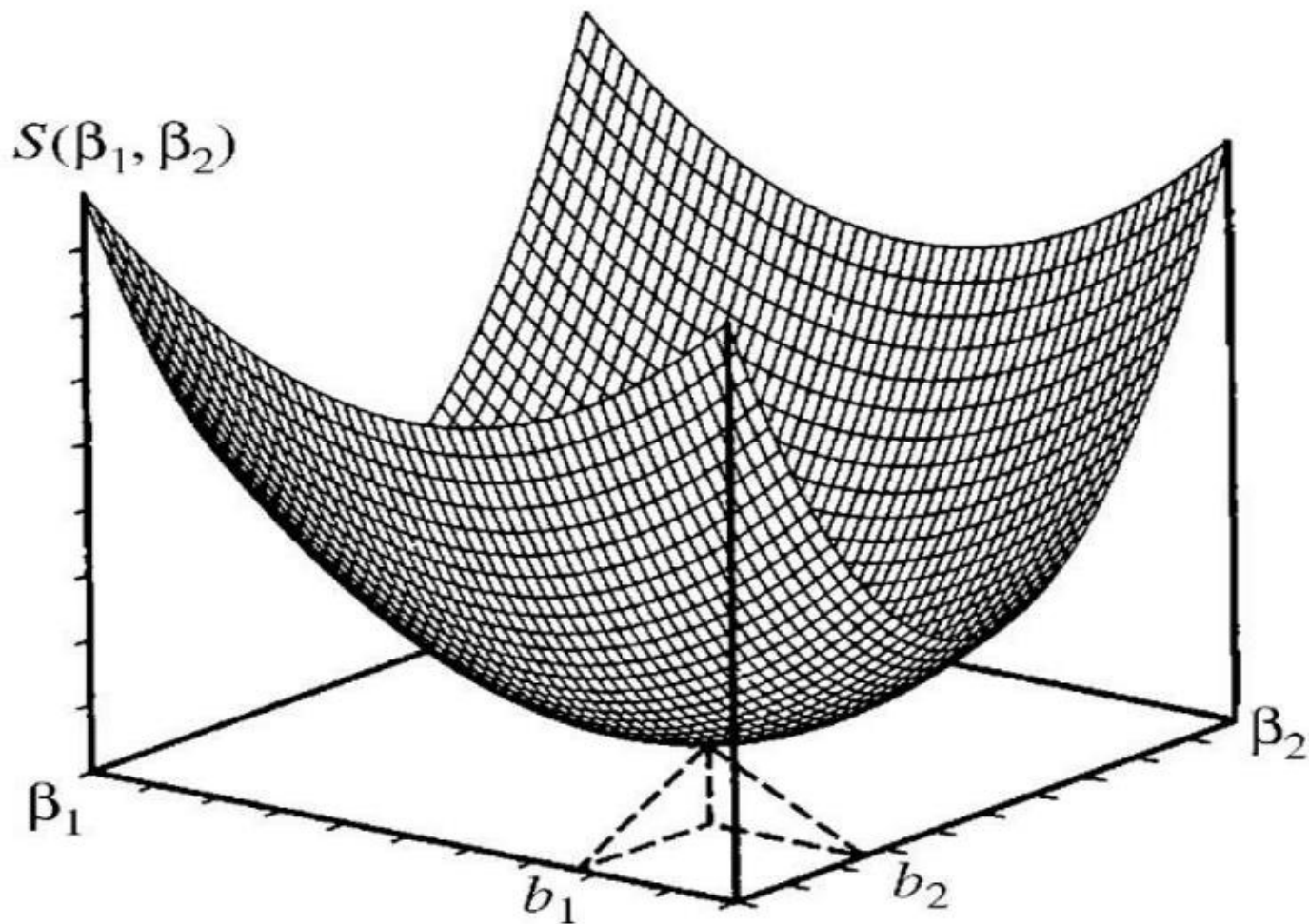
$$S(\beta_1, \beta_2) = \sum_{i=1}^N (y_i - \beta_1 - \beta_2 x_i)^2$$

Eq. 2A.2

$$\frac{\partial S}{\partial \beta_1} = 2N\beta_1 - 2\sum y_i + 2\left(\sum x_i\right)\beta_2$$

$$\frac{\partial S}{\partial \beta_2} = 2\left(\sum x_i^2\right)\beta_2 - 2\sum x_i y_i + 2\left(\sum x_i\right)\beta_1$$

Figure 2A.1 The sum of squares function and the minimizing values b_1 and b_2



Set the derivatives equal to zero to get two equations:

$$2\left[\sum y_i - Nb_1 - \left(\sum x_i\right)b_2\right] = 0$$

$$2\left[\sum x_i y_i - \left(\sum x_i\right)b_1 - \left(\sum x_i^2\right)b_2\right] = 0$$

Simplify these to:

Eq. 2A.3

$$Nb_1 + \left(\sum x_i\right)b_2 = \sum y_i$$

Eq. 2A.4

$$\left(\sum x_i\right)b_1 + \left(\sum x_i^2\right)b_2 = \sum x_i y_i$$

Solving the two equations simultaneously, we get for b_2 :

Eq. 2A.5

$$b_2 = \frac{N \sum x_i y_i - \sum x_i \sum y_i}{N \sum x_i^2 - (\sum x_i)^2}$$

We can rewrite the equation 2A.5 by noting the following:

Eq. 2B.1

$$\begin{aligned}
 \sum (x_i - \bar{x})^2 &= \sum x_i^2 - 2\bar{x} \sum x_i + N\bar{x}^2 \\
 &= \sum x_i^2 - 2\bar{x} \left(N \frac{1}{N} \sum x_i \right) + N\bar{x}^2 \\
 &= \sum x_i^2 - 2N\bar{x}^2 + N\bar{x}^2 \\
 &= \sum x_i^2 - N\bar{x}^2
 \end{aligned}$$

Also note that:

Eq. 2B.2

$$\begin{aligned}
 \sum (x_i - \bar{x})^2 &= \sum x_i^2 - N\bar{x}^2 \\
 &= \sum x_i^2 - \bar{x} \sum x_i \\
 &= \sum x_i^2 - \frac{(\sum x_i)^2}{N}
 \end{aligned}$$

Finally, we have:

$$\begin{aligned}\sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum x_i y_i - N\bar{x}\bar{y} \\ &= \sum x_i y_i - \frac{\sum x_i \sum y_i}{N}\end{aligned}$$

Eq. 2B.3

We can rewrite b_2 in deviation from the mean form as:

$$b_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

First note that we will always have:

$$\sum (x_i - \bar{x}) = 0$$

Now rewrite our formula for b_2 and use this fact:

$$\begin{aligned} b_2 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x})y_i - \bar{y}(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} \\ &= \sum \left[\frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right] y_i \\ &= \sum w_i y_i \end{aligned}$$

To obtain Eq. 2.12, replace y_i in Eq. 2.11 by $y_i = \beta_1 + \beta_2 x_i + e_i$ and simplify:

$$\begin{aligned} b_2 &= \sum w_i y_i \\ &= \sum w_i (\beta_1 + \beta_2 x_i + e_i) \\ &= \sum w_i \beta_1 + \beta_2 \sum w_i x_i + \sum w_i e_i \\ &= \beta_2 + \sum w_i e_i \end{aligned}$$

For this, we used the facts that:

$$\begin{aligned} \sum w_i &= \sum \left[\frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right] \\ &= \frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) \\ &= 0 \end{aligned}$$

$$\sum w_i x_i = 1$$

$$\beta_2 \sum w_i x_i = \beta_2$$

$$\sum (x_i - \bar{x}) = 0$$

We can show that $\sum w_i x_i = 1$ by using:

$$\begin{aligned} \sum (x_i - \bar{x})^2 &= \sum (x_i - \bar{x})(x_i - \bar{x}) \\ &= \sum (x_i - \bar{x})x_i - \bar{x} \sum (x_i - \bar{x}) \\ &= \sum (x_i - \bar{x})x_i \end{aligned}$$

so that:

$$\begin{aligned} \sum w_i x_i &= \frac{\sum (x_i - \bar{x})x_i}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x})x_i}{\sum (x_i - \bar{x})x_i} \\ &= 1 \end{aligned}$$

First note that:

$$b_2 = \beta_2 + \sum w_i e_i$$

and that:

$$\text{var}(b_2) = E[b_2 - E(b_2)]^2$$

$$\begin{aligned}
 \text{var}(b_2) &= E\left[\beta_2 + \sum w_i e_i - \beta_2\right]^2 \\
 &= E\left[\sum w_i e_i\right]^2 \\
 &= E\left[\sum w_i^2 e_i^2 + 2\sum_{i \neq j} \sum w_i w_j e_i e_j\right] \\
 &= \sum w_i^2 E(e_i^2) + 2\sum_{i \neq j} \sum w_i w_j E(e_i e_j) \\
 &= \sigma^2 \sum w_i^2 \\
 &= \frac{\sigma^2}{\sum (x_i - \bar{x})^2}
 \end{aligned}$$

$$\sigma^2 = \text{var}(e_i) = E[e_i - E(e_i)]^2 = E[e_i - 0]^2 = E(e_i^2)$$

$$\text{cov}(e_i, e_j) = E[(e_i - E(e_i))(e_j - E(e_j))] = E(e_i e_j) = 0$$

$$\sum w_i^2 = \sum \left[\frac{(x_i - \bar{x})^2}{\left\{ \sum (x_i - \bar{x})^2 \right\}^2} \right] = \frac{\sum (x_i - \bar{x})^2}{\left\{ \sum (x_i - \bar{x})^2 \right\}^2} = \frac{1}{\sum (x_i - \bar{x})^2}$$

$$\text{var}(aX + bY) = a^2 \text{var}(X) + b^2 \text{var}(Y) + 2 \text{cov}(X, Y)$$

$$\text{var}(b_2) = \text{var}\left(\beta_2 + \sum w_i e_i\right) \quad [\text{since } \beta_2 \text{ is a constant}]$$

$$= \sum w_i^2 \text{var}(e_i) + \sum_{i \neq j} \sum w_i w_j \text{cov}(e_i, e_j) \quad [\text{generalizing the variance rule}]$$

$$= \sum w_i^2 \text{var}(e_i) \quad [\text{using } \text{cov}(e_i, e_j) = 0]$$

$$= \sigma^2 \sum w_i^2 \quad [\text{using } \text{var}(e_i) = \sigma^2]$$

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

To begin, let $b_2^* = \sum k_i y_i$ be any other linear estimator of β_1 . Also, suppose $k_i = w_i + c_i$. Then:

$$\begin{aligned} b_2^* &= \sum k_i y_i = \sum (w_i + c_i) y_i = \sum (w_i + c_i) (\beta_1 + \beta_2 x_i + e_i) \\ &= \sum (w_i + c_i) \beta_1 + \sum (w_i + c_i) \beta_2 x_i + \sum (w_i + c_i) e_i \\ &= \beta_1 \sum w_i + \beta_1 \sum c_i + \beta_2 \sum w_i x_i + \beta_2 \sum c_i x_i + \sum (w_i + c_i) e_i \\ &= \beta_1 \sum c_i + \beta_2 + \beta_2 \sum c_i x_i + \sum (w_i + c_i) e_i \end{aligned}$$

Eq. 2F.1

Now:

Eq. 2F.2

$$\begin{aligned} E(b_2^*) &= \beta_1 \sum c_i + \beta_2 + \beta_2 \sum c_i x_i + \sum (w_i + c_i) E(e_i) \\ &= \beta_1 \sum c_i + \beta_2 + \beta_2 \sum c_i x_i \end{aligned}$$

Eq. 2F.3

$$\sum c_i = 0 \text{ and } \sum c_i x_i = 0$$

Eq. 2F.4

$$b_2^* = \sum k_i y_i = \beta_2 + \sum (w_i + c_i) e_i$$

For now, observe that :

$$\sum c_i w_i = \sum \left[\frac{c_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \right] = \frac{1}{\sum (x_i - \bar{x})^2} \sum c_i x_i - \frac{\bar{x}}{\sum (x_i - \bar{x})^2} \sum c_i = 0$$

Now we can write:

$$\text{var}(b_2^*) = \text{var} \left[\beta_2 + \sum (w_i + c_i) e_i \right]$$

$$= \sum (w_i + c_i)^2 \text{var}(e_i)$$

$$= \sigma^2 \sum (w_i + c_i)^2$$

$$= \sigma^2 \sum w_i^2 + \sigma^2 \sum c_i^2$$

$$= \text{var}(b_2) + \sigma^2 \sum c_i^2$$

$$\geq \text{var}(b_2)$$

- Monte Carlo simulation experiments use random number generators to replicate the random way that data are obtained
 - In Monte Carlo simulations we specify a data generation process and create samples of artificial data
 - Then we “try out” estimation methods on the data we have created
 - We create many samples of size N and examine the **repeated sampling properties** of the estimators
 - In this way, we can study how statistical procedures behave under ideal, as well as not so ideal, conditions

- The data generation process for the simple linear regression model is given by:

$$\begin{aligned}y_i &= E(y_i | x_i) + e_i \\ &= b_1 + b_2 x_i + e_i \\ i &= 1, \dots, N\end{aligned}$$

- Each value of the dependent variable y_i is obtained, or generated, by adding a random error e_i to the regression function $E(y_i | x_i)$
- To simulate values of y_i we create values for the systematic portion of the regression relationship $E(y_i | x_i)$ and add to it the random error e_i

- To create the variables for the regression function, we must:
 1. Select sample size N
 2. Choose x_i values
 3. Choose β_1 and β_2

- To be consistent with assumptions SR2–SR4 the random errors should have mean zero, a constant variance and be uncorrelated with one another
 - We can generate random numbers
 - Of course the computer-generated numbers cannot be truly random, because they are generated by a computer code
 - The random numbers created by computer software are “pseudorandom,” in that they behave like random numbers

Figure 2G.1 The true probability density functions of the data

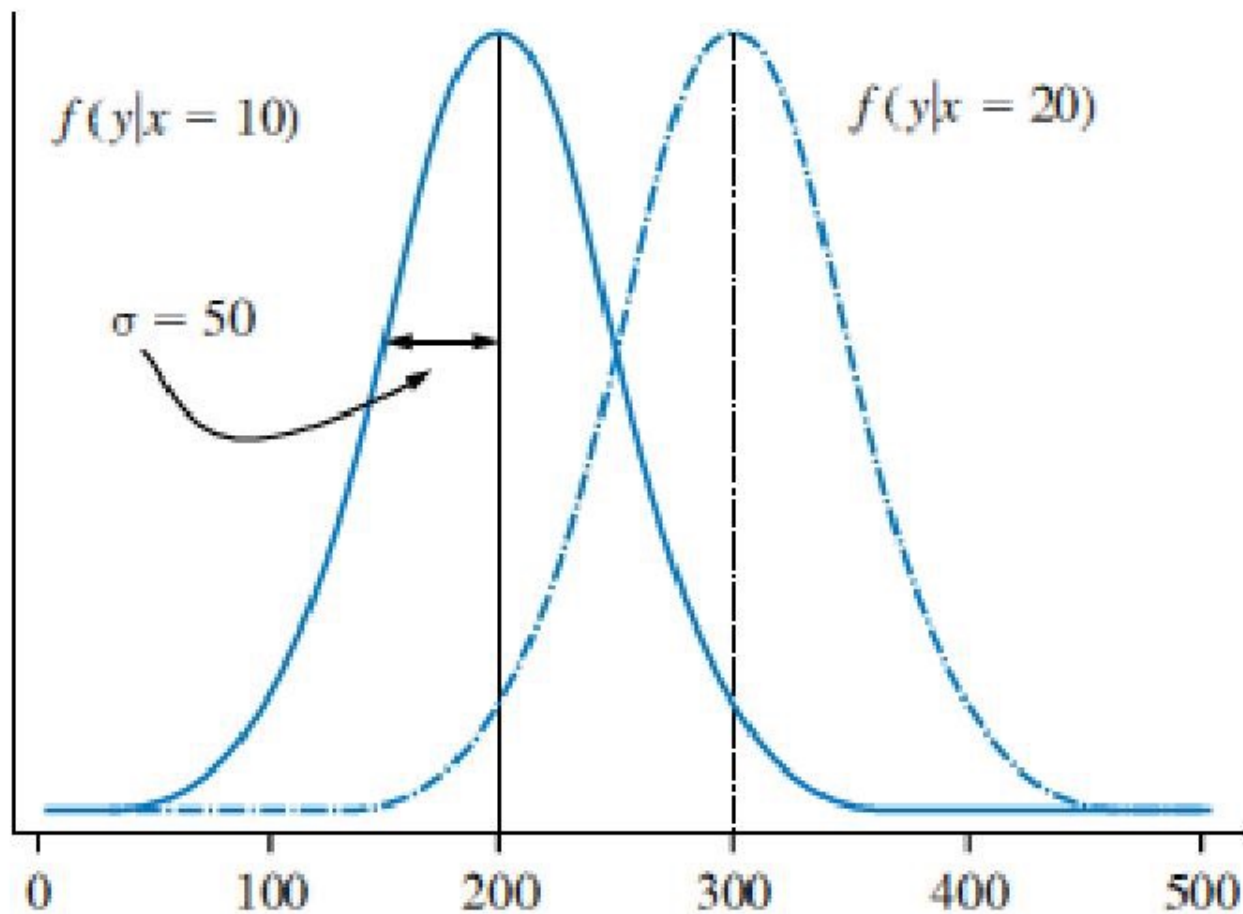
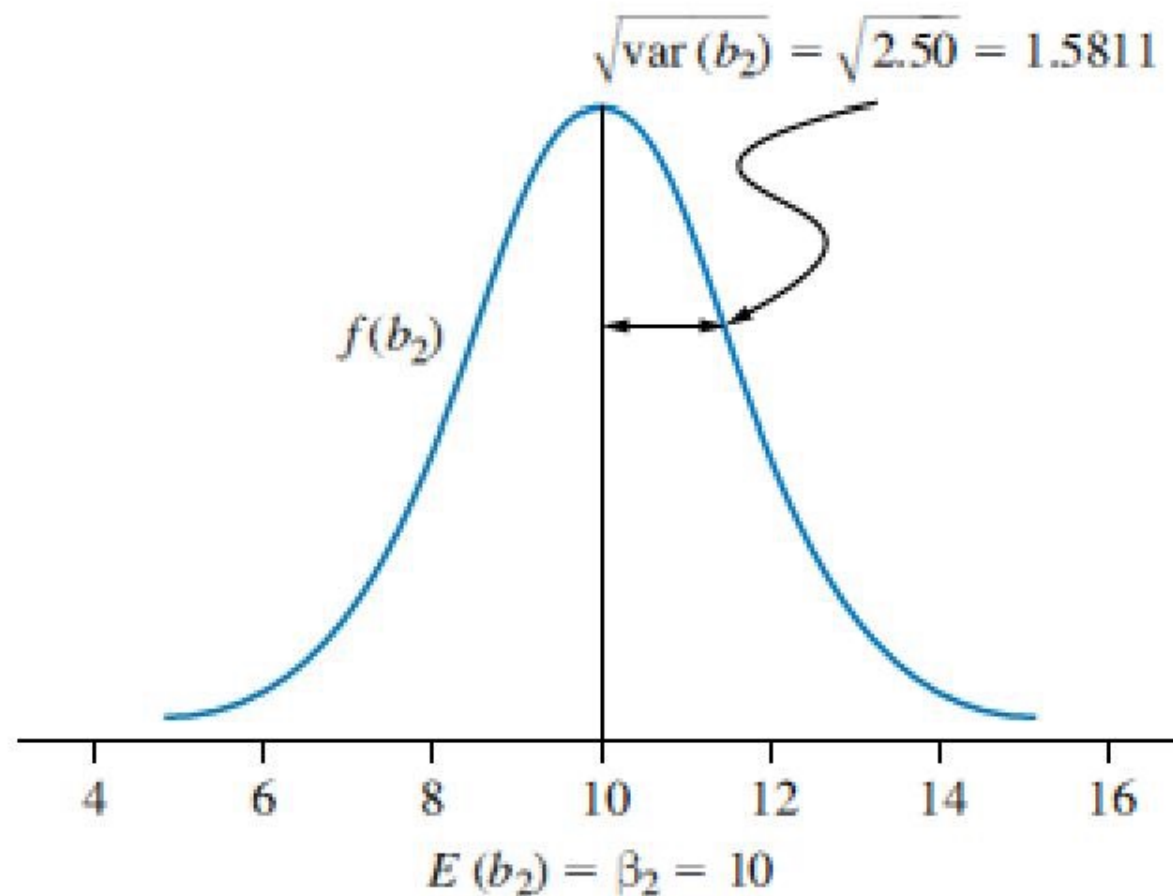


Figure 2G.2 The true probability density functions of the estimator b_2



$$\hat{y}_i = 75.7679 + 11.9683x_i$$
$$(se)(25.7928) \quad (1.6313)$$

$$\hat{\sigma} = 51.5857$$

Variance-Covariance Matrix

	b₁	b₂
b₁	665.2699	-39.9162
b₂	-39.9162	2.6611

- What do we hope to achieve with a Monte Carlo experiment?
 - We would like to verify that under SR1–SR5 the least squares estimators are unbiased
 - We would like to verify that under SR1–SR5 the least squares estimators have sampling variances given by Eq. 2.14 and Eq. 2.16
 - We would like to verify that the estimator of the error variance Eq. 2.19 is unbiased
 - Because we have assumed the random errors are normal, SR6, we expect the least squares estimates to have a normal distribution.

Table 2G.1 Summary of 1,000 Monte Carlo Samples

	Mean	Variance	Std. Dev.	Minimum	Maximum	1st Pct.	99th Pct.
b_1 (100)	99.7581	575.3842	23.9872	25.8811	174.6061	42.1583	156.0710
b_2 (10)	10.0143	2.3174	1.5223	5.1401	14.9928	6.3811	13.5620
$\hat{\sigma}^2$ (2,500)	2489.935	329909.9	574.3778	1024.191	5200.785	1360.764	4031.641

Figure 2G.3 The sampling distribution of b_2 in 1000 Monte Carlo samples

