# $\rangle\rangle\rangle$

The nature of the relationship between variables can take many forms, ranging from simple mathematical functions to extremely complicated ones. The simplest relationship consists of a straight-line or linear relationship (linear function).

## This is an example plot of linear function:





Suppose that a variable Y is a linear function of another variable X, with unknown parameters  $\beta_0$  and  $\beta_1$  that we wish to estimate.

Suppose that we have a sample of 4 observations with *X* values as shown.



If the relationship were an exact one, the observations would lie on a straight line and we would have no trouble obtaining accurate estimates of  $\beta_0$  and  $\beta_1$ . When all empirical pairs of X-Y points lie on a straight line – it is called a functional or deterministic relationship.



In practice, most economic relationships are not exact and the actual values of Y are different from those corresponding to the straight line.



To allow for such divergences, we will write the model as  $Y = \beta_0 + \beta_1 X + e$ , where *e* is a disturbance term.



Each value of Y thus has a nonrandom component,  $\beta_0 + \beta_1 X$ , and a random component, e. The first observation has been decomposed into these two components.



In practice we can see only the *P* points.



Obviously, we can use the *P* points to draw a line which is an approximation to the line  $Y = \beta_0 + \beta_1 X$ . If we write this line  $\hat{Y} = b_0 + b_1 X$ ,  $b_0$  is an estimate of  $\beta_0$  and  $b_1$  is an estimate of  $\beta_1$ .

## **Population Linear Regression**

Population regression line is a straight line that describes the dependence of the average value (conditional mean) of one variable on the other



However, we have obtained data from only a random sample of the population. For a sample,  $b_0$  and  $b_1$  can be used as estimates (estimators) of the respective population parameters  $\beta_0$  and  $\beta_1$ 

$$\hat{y}_i = b_0 + b_1 x_i + e_i$$

The intercept  $b_0$  and the slope  $b_1$  are the coefficients of the regression line. The slope  $b_1$  is the change in Y (increase, if >0, and decrease, if <0) associated with a unit change in X. The intercept is the value of Y when X=0; it's the point at which the population regression line intersects the Y axis. In some cases the intercept has no real-world meaning (for example when X is the class size, Y is the test score – the intercept is the predicted value of test scores when there are no students in the class!).

Random error contains all the other factors besides X that determine the value of the dependent variable Y, for a specific observation.



The line is called the fitted model and the values of *Y* predicted by it are called the fitted values of *Y*. They are given by the heights of the *R* points.



The discrepancies between the actual and fitted values of Y are known as the residuals.

Least squares criterion:

Minimize SSE (residual sum of squares), where

$$SSE = \sum_{i=1}^{n} e_i^2 = e_1^2 + \dots + e_n^2$$

To begin with, we will draw the fitted line so as to minimize the sum of the squares of the residuals, *SSE*. This is described as the least squares criterion.

Least squares criterion:

Minimize SSE (residual sum of squares), where

$$SSE = \sum_{i=1}^{n} e_i^2 = e_1^2 + \dots + e_n^2$$

Why not minimize

$$\sum_{i=1}^{n} e_{i} = e_{1} + \dots + e_{n}$$

Why the squares of the residuals? Why not just minimize the sum of the residuals?



The answer is that you would get an apparently perfect fit by drawing a horizontal line through the mean value of *Y*. The sum of the residuals would be zero.



You must prevent negative residuals from cancelling positive ones, and one way to do this is to use the squares of the residuals.

Since  $\hat{y}_i = b_0 + b_1 x_i$  we are minimizing  $SSE = \sum e_i^2 = \sum (y_i - (b_0 + b_1 x_i))^2$ 

which has two unknowns,  $b_0$  and  $b_1$ . A mathematical technique which determines the values of  $b_0$  and  $b_1$  that best fit the observed data is known as the Ordinary Least Squares method (OLS).

Ordinary Least Squares is a procedure that selects the best fit line given a set of data points, by minimizing the sum of the squared deviations of the points from a line. That is, if  $\hat{y}_i = b_0 + b_1 x_i$  is the equation of the best line to fit through the data then in order to get this best line, using the least squares criteria, for each value data point  $(\mathbf{x}_i, \mathbf{y}_i)$  if  $e_i = y_i - \hat{y}_i$  where  $\hat{y} = b_0 + b_1 x_i$ , then  $e_i$  is the amount of deviation of the data point from the line. The least squares criteria minimizes, finds the slope  $\mathbf{b}_1$  and the y-intercept  $\mathbf{b}_0$  from the data, that minimizes the sum of the square deviations,  $\sum_{i=1}^n e_i^2$ 

For the mathematically curious, I provide a condensed derivation of the coefficients.

To minimize  $SSE = \sum e_i^2 = \sum (y_i - (b_0 + b_1 x_i))^2$  determine the

partial derivatives with respect to  $b_0$  and with respect to  $b_1$ . These are:

$$f_{b0} = 2\sum (y - b_0 - b_1 x)(-1)$$
  
$$f_{b1} = 2\sum (y - b_0 - b_1 x)(-x)$$

Setting  $f_{b0} = f_{b1} = 0$  and solving for  $b_0$  and  $b_1$  results in equations given below.

$$\sum_{i=1}^{n} y_{i} = nb_{0} + b_{1} \sum_{i=1}^{n} x_{i}$$
$$\sum_{i=1}^{n} x_{i} y_{i} = b_{0} \sum_{i=1}^{n} x_{i} + b_{1} \sum_{i=1}^{n} x_{i}^{2}$$

Since there are two equations with two unknown, we can solve these equations simultaneously for  $b_0$  and  $b_1$  as follows:

$$b_{1} = \frac{n \sum_{i=1}^{n} x_{i} y_{i} - \sum_{i=1}^{n} x_{i} \sum_{i=1}^{n} y_{i}}{n \sum_{i=1}^{n} x_{i}^{2} - (\sum_{i=1}^{n} x_{i})^{2}} \qquad b_{0} = \overline{Y} - b_{1} \overline{X}$$

ONLY FOR REGRESSION MODELS WITH ONE INDEPENDENT VARIABLE!

We also note that the regression line always goes through the mean (  $_{\overline{X},\overline{Y}}$  ).

In matrix notation OLS may be written as: Y = Xb + e The normal equations in matrix form are now

$$\mathbf{X}^{\mathsf{T}} \mathbf{Y} = \mathbf{X}^{\mathsf{T}} \mathbf{X} \mathbf{b}$$

And when we solve it for b we get:  $b = (X^TX)^{-1}X^TY$ 

where Y is a column vector of the Y values and X is a matrix containing a column of ones (to pick up the intercept) followed by a column of the X variable containing the observations on it and b is a vector containing the estimators of regression parameters.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \qquad \qquad X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \qquad \qquad b = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$$

We can state as follows:

$$X^{T}X = \begin{bmatrix} n & \sum_{i=1}^{n} x_{i} \\ \sum_{i=1}^{n} x_{i} & \sum_{i=1}^{n} x_{i}^{2} \end{bmatrix} \qquad \qquad X^{T}Y = \begin{bmatrix} \sum_{i=1}^{n} y_{i} \\ \sum_{i=1}^{n} y_{i}x_{i} \end{bmatrix}$$

How to inverse X<sup>T</sup>X? 1. matrix determinant  $deX^T X = n \cdot \sum x^2 - (\sum x)^2$ 2. minor matrix  $\min X^T X = \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix} = \begin{bmatrix} \sum x^2 & \sum x \\ \sum x & n \end{bmatrix}$ 

3. cofactor matrix 
$$(X^T X)D = \begin{bmatrix} \sum x^2 \cdot (-1)^{1+1} & \sum x \cdot (-1)^{1+2} \\ \sum x \cdot (-1)^{2+1} & n \cdot (-1)^{2+2} \end{bmatrix}$$

4. inverse matrix 
$$(X^T X)^{-1} = \frac{1}{\det X^T X} \begin{bmatrix} \sum x^2 & -\sum x \\ -\sum x & n \end{bmatrix}$$

#### EXAMPLE

In this problem we were looking at the way home size is effected by the family income. We will use this model to try to predict the value of the dependent variable based on the independent variable. Also, the slope will help us to understand how the Y variable changes for each unit change in the X variable.

Assume a real-estate developer is interested in determining the relationship between family income (X, in thousand of dollars) of the local resident and the square footage of their homes (Y, in hundreds of square feet). A random sample of ten families is obtained with the following results:

X	22	26	45	37	28	50	56	34	60	40
Y	16	17	26	24	22	21	32	18	30	20

### For our example (X-family income; Y-home size)

Family	x	y y								
1	22	16			1	22			16	
2	26	17			1	26			17	2 2
3	45	26			1	45			26	
4	37	24			1	37			24	
5	28	22		X	1	28		Y	22	
6	50	21			1	50			21	С.
7	56	32			1	56			32	3 2
8	34	18			1	34			18	
9	60	30			1	60			30	3 
10	40	20			1	40			20	
ΧT	1	1	1	1	1	1	1	1	1	1
	22	26	45	37	28	50	56	34	60	40
X™X	10	398		X <sup>⊤</sup> Y	226					
	398	17330			9522					

det X™X so it's nonsingular		14896	X —family income (\$ 000) Y — home size (hundreds of square feet)							
Matrix of minors							f the family has no			
Min X <sup>T</sup> X	17330	398	E	ors	ir	] ] income, the home size is				
398		10	<b>b</b> o	8,51		equal to 851 square				
Cofactor matrix			<b>b</b> 1	0,35			qual to obti square reet			
(X <sup>⊤</sup> X)D	17330	-398			If the family income increase 1000 \$					
	-398	10								
Inverse matrix					the ho	me s	ize will increase 35			
(X <sup>T</sup> X) <sup>-1</sup>	1,163	-0,03			square feet on average					
	-0,03	0,0007			square reet, on average.					

The regression equation can be stated as follows:

$$\hat{y}_i = 8,51 + 0,35x_i + e_i$$



## Let's try another example:

## X – commercial time (minutes) Y – sales (\$ hundred thousand)

2332	.02562	1	`				N:)	1		
Х	У									
1	9			1	1			9		
5	20			1	5			20		
6	22			1	6			22		
5	15			1	5			15		
5	17		X =	1	5		Y =	17		
9	30			1	9			30		
3	18			1	3			18		
7	25			1	7			25		
3	10			1	3			10		
6	20			1	6			20		
Хт =	1	1	1	1	1	1	1	1	1	1
	1	5	6	5	5	9	3	7	3	6
							E	stimator	s	
XTX =	10	50		Х <sup>т</sup> Υ =	186		bo	5,56		
	50	296			1050		b <sub>1</sub>	2,61		

## $\hat{y} = 5.56 + 2.61x$

How are the slope and intercept interpreted for this example?

X - commercial time (minutes)

Y - sales (\$ hundred thousand)

**Intercept:**  $b_0 = 5.56$  When we have no commercials, we would expect \$556,000 in sales. **Slope:**  $b_1 = 2.61$  For each minute increase in commercial time we have a \$261,000 increase in sales.

Notice that when  $x = \bar{x} = 5$   $\hat{y} = 5.56 + 2.61(5) = 18.61 = \bar{y}$ . As we said, the least squares regression line always goes through the points  $(\bar{x}, \bar{y})$ 



## **REGRESSION MODEL WITH TWO EXPLANATORY VARIABLES**

 $\mathbf{Y} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{X}_1 + \boldsymbol{\beta}_2 \boldsymbol{X}_2 + \mathbf{e}_i$ 

This sequence provides a geometrical interpretation of a multiple regression model with two explanatory variables.

- Y-weekly salary (\$)
- X1 length of employment (in months)

```
X2 – age (in years)
```



Specifically, we will look at weekly salary function model where weekly salary, Y, depend on length of employment X1, and age, X2.

 $\mathbf{Y} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{X}_1 + \boldsymbol{\beta}_2 \boldsymbol{X}_2 + \mathbf{e}_i$ 

Y-weekly salary (\$)

X1 – length of employment (in months)

```
X2 – age (in years)
```



The model has three dimensions, one each for *Y*, *X1*, and X2. The starting point for investigating the determination of *Y* is the intercept,  $\beta_0$ .

 $\mathbf{Y} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{X}_1 + \boldsymbol{\beta}_2 \boldsymbol{X}_2 + \mathbf{e}_i$ 

Y-weekly salary (\$)

X1 – length of employment (in months)

```
X2 – age (in years)
```



Literally the intercept gives weekly salary for those respondents who have no age (??) and no length of employment (??). Hence a literal interpretation of  $\beta_0$  would be unwise.

 $\mathbf{Y} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{X}_1 + \boldsymbol{\beta}_2 \boldsymbol{X}_2 + \mathbf{e}_i$ 

Y – weekly salary (\$)

- X1 length of employment (in months)
- X2 age (in years)



The next term on the right side of the equation gives the effect of X1. A one month of employment increase in X1 causes weekly salary to increase by  $\beta_1$  dollars, holding X2 constant.



Similarly, the third term gives the effect of variations in X2. A one year of age increase in X2 causes weekly salary to increase by  $\beta_2$  dollars, holding X1 constant.



Different combinations of X1 and X2 give rise to values of weekly salary which lie on the plane shown in the diagram, defined by the equation  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ . This is the nonrandom component of the model.



The final element of the model is the error term, e. This causes the actual values of Y to deviate from the plane. In this observation, e happens to have a positive value.


A sample consists of a number of observations generated in this way. Note that the interpretation of the model does not depend on whether *X1* and *X2* are correlated or not.

# $\mathbf{Y} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{X}_1 + \boldsymbol{\beta}_2 \boldsymbol{X}_2 + \mathbf{e}_i$ $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$ $\beta_0 + \beta_1 X_1 + \beta_2 X_2$ е $\beta_0 + \beta_2 X_2$ combined effect of X1 pure X2 effect and X2 pure X1 effect $\beta_0 + \beta_1 X_1$ $\beta_0$ Y **X**2 X1

However we do assume that the effects of X1 and X2 on salary are additive. The impact of a difference in X1 on salary is not affected by the value of X2, or vice versa.

$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i}$$

Slope coefficients are interpreted as partial slope/partial regression coefficients:

- D b<sub>1</sub> = average change in Y associated with a unit change in X<sub>1</sub>, with the other independent variables held constant (all else equal);
- $\Box$  b<sub>2</sub> = average change in Y associated with a unit change in X<sub>2</sub>, with the other independent variables held constant (all else equal).

$$Y_{i} = \beta_{0} + \beta_{1}X_{1i} + \beta_{2}X_{2i} + e_{i}$$

$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i}$$

The regression coefficients are derived using the same least squares principle used in simple regression analysis. The fitted value of Y in observation *i* depends on our choice of  $b_0$ ,  $b_1$ , and  $b_2$ .

$$Y_{i} = \beta_{0} + \beta_{1}X_{1i} + \beta_{2}X_{2i} + e_{i}$$

$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i}$$

$$e_i = Y_i - \hat{Y}_i = Y_i - b_0 - b_1 X_{1i} - b_2 X_{2i}$$

The residual  $e_i$  in observation *i* is the difference between the actual and fitted values of *Y*.

 $SSE = \sum e_i^2 = \sum (Y_i - b_0 - b_1 X_{1i} - b_2 X_{2i})^2$ 

We define SSE, the sum of the squares of the residuals, and choose  $b_0$ ,  $b_1$ , and  $b_2$  so as to minimize it.

$$SSE = \sum e_i^2 = \sum (Y_i - b_0 - b_1 X_{1i} - b_2 X_{2i})^2$$
  

$$= \sum (Y_i^2 + b_0^2 + b_1^2 X_{1i}^2 + b_2^2 X_{2i}^2 - 2b_0 Y_i - 2b_1 X_{1i} Y_i$$
  

$$- 2b_2 X_{2i} Y_i + 2b_0 b_1 X_{1i} + 2b_0 b_2 X_{1i} + 2b_1 b_2 X_{1i} X_{2i})$$
  

$$= \sum Y_i^2 + nb_0^2 + b_1^2 \sum X_{1i}^2 + b_2^2 \sum X_{2i}^2 - 2b_0 \sum Y_i$$
  

$$- 2b_1 \sum X_{1i} Y_i - 2b_2 \sum X_{2i} Y_i + 2b_0 b_1 \sum X_{1i}$$
  

$$+ 2b_0 b_2 \sum X_{2i} + 2b_1 b_2 \sum X_{1i} X_{2i}$$
  

$$\frac{\partial SSE}{\partial b_0} = 0 \qquad \frac{\partial SSE}{\partial b_1} = 0 \qquad \frac{\partial SSE}{\partial b_2} = 0$$

First we expand SSE as shown, and then we use the first order conditions for minimizing it.

$$b_0 = \overline{Y} - b_1 \overline{X}_1 - b_2 \overline{X}_2$$

$$b_{1} = \frac{\text{Cov}(X_{1}, Y)\text{Var}(X_{2}) - \text{Cov}(X_{2}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

$$b_{2} = \frac{\text{Cov}(X_{2}, Y)\text{Var}(X_{1}) - \text{Cov}(X_{1}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

We thus obtain three equations in three unknowns. Solving for  $b_0$ ,  $b_1$ , and  $b_2$ , we obtain the expressions shown above.

$$b_0 = \overline{Y} - b_1 \overline{X}_1 - b_2 \overline{X}_2$$

$$b_{1} = \frac{\text{Cov}(X_{1}, Y)\text{Var}(X_{2}) - \text{Cov}(X_{2}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

$$b_{2} = \frac{\text{Cov}(X_{2}, Y)\text{Var}(X_{1}) - \text{Cov}(X_{1}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

The expression for  $b_0$  is a straightforward extension of the expression for it in simple regression analysis.

$$b_0 = \overline{Y} - b_1 \overline{X}_1 - b_2 \overline{X}_2$$

$$b_{1} = \frac{\text{Cov}(X_{1}, Y)\text{Var}(X_{2}) - \text{Cov}(X_{2}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

$$b_{2} = \frac{\text{Cov}(X_{2}, Y)\text{Var}(X_{1}) - \text{Cov}(X_{1}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

However, the expressions for the slope coefficients are considerably more complex than that for the slope coefficient in simple regression analysis.

$$b_0 = \overline{Y} - b_1 \overline{X}_1 - b_2 \overline{X}_2$$

$$b_{1} = \frac{\text{Cov}(X_{1}, Y)\text{Var}(X_{2}) - \text{Cov}(X_{2}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

$$b_{2} = \frac{\text{Cov}(X_{2}, Y)\text{Var}(X_{1}) - \text{Cov}(X_{1}, Y)\text{Cov}(X_{1}, X_{2})}{\text{Var}(X_{1})\text{Var}(X_{2}) - [\text{Cov}(X_{1}, X_{2})]^{2}}$$

For the general case when there are many explanatory variables, ordinary algebra is inadequate. It is necessary to switch to matrix algebra.

In matrix notation OLS may be written as:

Y = Xb + e

The normal equations in matrix form are now

 $\mathbf{X}^{\mathsf{T}} \mathbf{Y} = \mathbf{X}^{\mathsf{T}} \mathbf{X} \mathbf{b}$ 

And when we solve it for b we get:

 $\mathbf{b} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{Y}$ 

where Y is a column vector of the Y values and X is a matrix containing a column of ones (to pick up the intercept) followed by a column of the X variables containing the observations on them and *b* is a vector containing the estimators of regression parameters.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \qquad X = \begin{bmatrix} 1 & X_{11} & X_{21} \\ 1 & X_{12} & X_{22} \\ \dots & \dots & \dots \\ 1 & X_{1n} & X_{2n} \end{bmatrix} \qquad b = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

#### MATRIX ALGEBRA: SUMMARY

A **vector** is a collection of n numbers or elements, collected either in a column (a **column vec**tor) or in a row (a **row vector**).

A **matrix** is a collection, or array, of numbers of elements in which the elements are laid out in columns and rows. The dimension of matrix is  $n \times m$  where *n* is the number of rows and *m* is the number of columns.

#### Types of matrices

A matrix is said to be **square** if the number of rows equals the number of columns. A square matrix is said to be **symmetric** if its (i, j) element equals its (j, i) element. A **diagonal** matrix is a square matrix in which all the off-diagonal elements equal zero, that is, if the square matrix A is diagonal, then  $a_{ij} = 0$  for  $i \neq j$ .

The **transpose** of a matrix switches the rows and the columns. That is, the transpose of a matrix turns the  $n \times m$  matrix A into the  $m \times n$  matrix denoted by  $A^{T}$ , where the (i, j) element of A becomes the (j, i) element of  $A^{T}$ ; said differently, the transpose of a matrix A turns the rows of A into the columns of  $A^{T}$ . The **inverse** of the matrix A is defined as the matrix for which  $A^{-1}A=1$ . If in fact the inverse matrix  $A^{-1}$  exists, then A is said to be **invertible** or **nonsingular**.

#### Vector and matrix multiplication

The matrices A and B can be multiplied together if they are conformable, that is, if the number of columns of A equals the number of rows of B. In general, matrix multiplication does not commute, that is, in general  $AB \neq BA$ .

# Data for weekly salary based upon the length of employment and age of employees of a large industrial corporation are shown in the table.

	Weekly	Length of	Age
Employee	salary	employment	(X2,
	(\$)	(X1, months)	years)
1	639	330	46
2	746	569	65
3	670	375	57
4	518	113	47
5	602	215	41
6	612	343	59
7	548	252	45
8	591	348	57
9	552	352	55
10	529	256	61
11	456	87	28
12	674	337	51
13	406	42	28
14	529	129	37
15	528	216	46
16	592	327	56

Calculate the OLS estimates for regression coefficients for the available sample. Comment on your results.

Y-weekly salary (\$) X1 –length of employment (months) X2-age (years)

i	Y	<b>X</b> 1	Х <sub>2</sub>
1	639	330	46
2	746	569	65
3	670	375	57
4	518	113	47
5	602	215	41
6	612	343	59
7	548	252	45
8	591	348	57
9	552	352	55
10	529	256	61
11	456	87	28
12	674	337	51
13	406	42	28
14	529	129	37
15	528	216	46
16	592	327	56

11111	<u> </u>	330	46
()))))	1	569	65
11111	1	375	57
11111	1	113	47
11111	1	215	41
Х	1	343	59
11111	1	252	45
11111	1	348	57
11111	1	352	55
11111	1	256	61
11111	1	87	28
11111	1	337	51
11111	1	42	28
11111	1	129	37
(1111)	1	216	46
11111	1	327	56

	639
	746
	670
	518
	602
Y	612
11111	548
	591
	552
	529
	456
	674
	406
	529
	528
	592

	1	Y	1	1	1	1	1	T	1		1	Ţ	T	Ţ	1	1
XT	330	569	375	113	215	343	252	348	352	256	87	337	42	129	216	327
	46	65	57	47	41	59	45	57	55	61	28	51	28	37	46	56

		<u></u>	//////////////////////////////////////	<u>_</u>						
	16	4291	779					9192		
X <sup>T</sup> X	4291	1417105	227875				X <sup>T</sup> Y	2617701		
111111111	779	227875	39771					457709		
( <u>)))))))))</u> ))										
det	2105037674									
						HUUUHH	+///////			
min11	1417105	227875		min12	4291	227875	4///////	min13	4291	1417105
	227875	39771			779	39771	<u>1</u> ///////		779	227875
det	4432667330			det	-6857264			det	-126113170	
min21	4291	779	//////////	min22	16	779		min23	16	4291
	227875	39771			779	39771	AIIIIIIA		779	227875
det	-6857264			det	29495			det	303311	T///////
				111111	11111111		<u>_</u>			
min31	4291	779		min32	16	779		min33	16	4291
	1417105	227875			4291	227875			4291	1417105
det	-126113170			det	303311			det	4260999	Ι
	matrix of minors					cofactor ma	atrix		]	
<u>(1111111111</u> 1)			11111111111	<u>1</u> //////		11111111111	<u>((((((((()))</u>		_//////////////////////////////////////	

	4432667330	-6857264	-126113170	())
minors XIX	-6857264	29495	303311	111
	-126113170	303311	4260999	

N		4432667330	6857264	-126113170
	(X <sup>⊺</sup> X)D	6857264	29495	-303311
	11111111	-126113170	-303311	4260999



Y-weekly salary (\$) X1 –length of employment (months) X2-age (years)

Our regression equation with two predictors (X1, X2):

$$\hat{y}_i = 461,85 + 0,671 \cdot X_1 - 1,383 \cdot X_2$$



Y

These are our data points in 3dimensional space (graph drawn using *Statistica 6.0*)

¥1



X2

Y

X1

•

Data points with the regression surface (*Statistica 6.0*)

X2

Y

X1

Data points with the regression surface (Statistica 6.0) after rotation.

There are times when a variable of interest in a regression cannot possibly be considered quantitative. An example is the variable gender.

Although this variable may be considered important in predicting a quantitative dependent variable, it cannot be regarded as quantitative.

The best course of action in such case is to take separate samples of males and females and conduct two separate regression analyses.

The results for the males can be compared with the results for the females to see if the same predictor variables and the same regression coefficients results. If a large sample size is not possible, a dummy variable can be employed to introduce qualitative variable into the analysis.

A DUMMY VARIABLE IN A REGRESSION ANALYSIS IS A QUALITATIVE OR CATEGORICAL VARIABLE THAT IS USED AS A PREDICTOR VARIABLE.

# For example, a male could be designated with the code 0 and the female could be coded as 1. Each person sampled could then be measured as

either a 0 or a 1 for the variable gender, and this variable, along with the quantitative variables for the persons, could be entered into a multiple regression program and analyzed. **Example 1** Returning to real-estate developer, we noticed that all the houses in the population were from three neighborhoods, A, B, and C.

X1 - family income (\$ 000)

X2 - family size

X3 - neighborhood

Family	Y	X1	X2	<b>X</b> 3
1	16	22	2	В
2	17	26	2	С
3	26	45	3	А
4	24	37	4	С
5	22	28	4	В
6	21	50	3	С
7	32	56	6	В
8	18	34	3	В
9	30	60	5	А
10	20	40	3	А

Using these data, we can construct the necessary dummy variables and determine whether they contribute significantly to the prediction of home size (Y).

One way to code neighborhoods would be to define:

 $X_{3} \begin{cases} 0 \text{ if neighborhood } A \\ 1 \text{ if neighborhood } B \\ 2 \text{ if neighborhood } C \end{cases}$ 

However, this type of coding has many problems. First, because 0 < 1 < 2, the codes imply that neighborhood A is smaller then neighborhood B, which is smaller then neighborhood C. A better procedure is to use the necessary number of dummy variables to represent the neighborhood. To represent the three neighborhoods, we use two dummy variables, by letting

$$X_{3} = \begin{cases} 1 & if house is in A \\ 0 & otherwise \end{cases}$$

 $X_{4} = \begin{cases} 1 & if house is in B \\ 0 & otherwise \end{cases}$ 

What happened to neighborhood C? It is not necessary to develop a third dummy variable. IT IS VERY IMPORTANT THAT YOU NOT INCLUDE IT!! If you attempted to use three such dummy variables in your model, you would receive a message in your computer output informing you that no solution exists for this model.

# Why?

One predictor variable is a linear combination (including a constant term) of one or more other predictors, then mathematically no solution exists for the least squares coefficients. To arrive at a usable equation, any such predictor variable must not be included. We don't lose any information - this excluded category is the reference system. The coefficients are the measure of the categories included in comparison to this one excluded.

# The final array of data is

Family	Y	X1	<b>X</b> 2	X3 (A)	X4 (B)
1	16	22	2	0	1
2	17	26	2	0	0
3	26	45	3	1	0
4	24	37	4	0	0
5	22	28	4	0	1
6	21	50	3	0	0
7	32	56	6	0	1
8	18	34	3	0	1
9	30	60	5	1	0
10	20	40	3	1	0

 $\hat{Y} = 7,772 + 0,082x_1 + 3,27x_2 + 1,613x_3 - 0,9x_4$ (2,5573) (0,1059) (0,987) (1,8009) (1,8414)

 If family income increases 1000\$ the average home size will increase about 0,082 hundred of square feet (holding family size constant)

 If family size increases 1 person the average home size will increase about 3,27 hundred of square feet (holding family income constant)  $\hat{Y} = 7,772 + 0,082x_1 + 3,27x_2 + 1,613x_3 - 0,9x_4$ (2,5573) (0,1059) (0,987) (1,8009) (1,8414)

The houses located in neighborhood A are 1,613
 hundred of square feet bigger then houses from
 neighborhood C.

• The houses located in neighborhood B are 0,9 hundred of square feet smaller then houses from neighborhood C.

## Example 2

Joanne Herr, an analyst for the Best Foods grocery chain, wanted to know whether three stores have the same average dollar amount per purchase or not. Stores can be thought of a single qualitative variable set at 3 levels – A, B, and C. A model can be set up to predict the dollar amount per purchase:

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + e_i$$

where

Y^- expected dollar amount per purchase

 $x_{1} = \begin{cases} 1 \text{ if the purchase is made in store } A \\ 0 \text{ if the purchase is not made in store } A \end{cases}$ 

 $x_{2} = \begin{cases} 1 \text{ if the purchase is made in store } B \\ 0 \text{ if the purchase is not made in store } B \end{cases}$ 

# The data

purchase	///////////////////////////////////////	
(dollars)	Store A	Store B
Y	X1	X2
12,05	1	0
23,94	1	0
14,63	1	0
25,78	1	0
17,52	1	0
18,45	1	0
15,17	0	1
18, 52	0	1
19,57	0	1
21,4	0	1
13, 59	0	1
20, 57	0	1
9,48	0	0
6,92	0	0
10,47	0	0
7,63	0	0
11,9	0	0
5,92	0	0
<u>273,51</u>		

The variables X1 and X2 are dummy variables representing purchases in store A or B, respectively. Note that the three levels of the qualitative variable have been described with only two variables.

$$b = (X^{T}X)^{-1}X^{T}Y$$
$$b = \begin{bmatrix} 8,72\\ 10,01\\ 9,42 \end{bmatrix}$$

The regression equation

$$\hat{Y} = 8,72 + 10,01x_1 + 9,42x_2$$
$$\hat{Y} = 8,72 + 10,01x_1 + 9,42x_2$$

the average dollar amount per purchase is for store A
is 10,01\$ higher comparing with store C

the average dollar amount per purchase is for store B
is 9,42\$ higher comparing with store C

## always compare to the excluded category!!

## Store A $\hat{Y} = 8,72 + 10,01 \cdot 1 + 9,42 \cdot 0 = 18,73$ \$ Store B $\hat{Y} = 8,72 + 10,01 \cdot 0 + 9,42 \cdot 1 = 18,14$ \$ Store C $\hat{Y} = 8,72 + 10,01 \cdot 0 + 9,42 \cdot 0 = 8,72$ \$