Correlation and Simple Linear Regression

- Regression analysis is a statistical tool that utilizes the relation between two or more quantitative variables so that one variable can be predicted from the other, or others.
- Some Examples:
 - 1. Waistline and Weight.
 - 2. SAT score and First year college GPA.
 - 3. Number of customers and Revenue.
 - 4. Family income and Family expenditures.
- Functional Relation vs. Statistical Relation between two variables.
 - A functional relation between two variables is expressed by a mathematical formula. If X is the independent variable and Y the dependent variable, a functional relation is of the form:

$$Y = f(X).$$

That is, given a particular value of X, we get only one corresponding value Y.

- 1. For example, let x denote the number of printer cartridges that you order over the internet. Suppose each cartridge costs \$40 and there is a fixed shipping fee of \$10, determine the total cost y of ordering x cartridges.
- A statistical relation, unlike a functional relation, is not a perfect one. If X is the independent variable and Y the dependent variable, a statistical relation is of the form:

$$Y = f(x) + \epsilon.$$

In such cases, we call X an *explanatory variable* and Y a *response variable*.

1. For example, let x denote the distance that a person plans to jog and y the time that it will take this person to finish it. Consider his 22 jogging distances and times from last month shown in the table below. If this person plans to jog for 5.5 miles tomorrow, predict how long it will take him to finish the run.

	1	2	3	4	5	6	7	8	9	10	11
Distance (x)	2	2	3	3	2	2.5	2.5	3	3.5	3.5	4
Time (y)	25	22	35	36	23	30	31	35	41	40	49
	12	13	14	15	16	17	18	19	20	21	22
Distance (x)	12 4	13 4	14 4	$\frac{15}{4.5}$	$\frac{16}{4.5}$	$\frac{17}{5}$	$\frac{18}{5}$	19 5	$20 \\ 3.5$	21 3.5	22 4

- Scatterplots. A scatterplot (or scatter diagram) is a graph in which the paired (x, y) sample data are plotted with a horizontal x-axis and a vertical y-axis. Each individual (x, y) pair is plotted as a single point. Scatterplots are useful as they usually display the relationship between two quantitative variables.
 - Always plot the explanatory variable on the x-axis, while the response variable on the y-axis.
 - In examining a scatterplot, look for an overall pattern showing the **form**, **direction**, and **strength** of the relationship, and then for outliers or other deviations from this pattern.
 - * Form linear or not.
 - * Direction positive or negative association.
 - * Strength how close the points lie to the general pattern (usually a line).



- Correlation. A *correlation* exists between two variables when one of them is related to the other in some way.
- Linear Correlation. The linear correlation coefficient ρ measures the strength of the linear relationship between the paired x- and y-quantitative values in a sample. In Chapter 5, we defined the correlation coefficient $\rho(X, Y)$ by

$$\rho = \rho(X, Y) = \frac{Cov(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}.$$

The estimator for ρ is the sample correlation coefficient r:

=

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$
(1)

$$= \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$
(2)

$$= \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$
(3)

$$= \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}} \tag{4}$$

where,

$$SS_{xx} = \Sigma x^2 - \frac{1}{n} (\Sigma x)^2 = (n-1)s_x^2$$
(5)

$$SS_{yy} = \Sigma y^2 - \frac{1}{n} (\Sigma y)^2 = (n-1)s_y^2$$
(6)

$$SS_{xy} = \Sigma xy - \frac{1}{n} (\Sigma x) (\Sigma y)$$
(7)

• Sample plots with correlation values



• Properties of *r*.

- 1. The value of r does not depend on which of the two variables under study is labeled as x and which is labeled as y.
- **2.** The value of r is independent of the units in which x and y are measured.
- **3.** $-1 \le r \le 1$.
- Some guidelines in interpreting r.

Value of $ r $	Strength of linear relationship					
If $ r \ge .95$	Very Strong					
If $.85 \le r < .95$	Strong					
If $.65 \le r < .85$	Moderate to Strong					
If $.45 \le r < .65$	Moderate					
If $.25 \le r < .45$	Weak					
If $ r < .25$	Very weak/Close to none					

• Testing for the Absence of Correlation: Assuming that the data come from a bivariate normal distribution, when $H_o: \rho = 0$ is true, the test statistic

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

follows the *t*-distribution with (n-2) degrees of freedom.

• The table below displays data on age (in years) and price (in \$100) for a sample of 11 cars.

Age (x)	5	4	6	6	5	5	6	6	2	7	7
Price (y)	85	102	70	80	89	98	66	90	169	68	50

1. Determine the values of SS_{xx} , SS_{yy} , and SS_{xy} .

- **2.** Determine the correlation coefficient r.
- **3.** What can you say about the linear relationship of x and y? Is it a strong linear relationship.
- **4.** Test $H_o: \rho = 0$ against $H_1: \rho \neq 0$.