

## 6.7: Sample Correlation and Regression

### Descriptive Theory

Recall the basic model of statistics: we have a population of objects of interest, and we have various measurements (variables) that we make on these objects. We select objects from the population and record the variables for the objects in the sample; these become our data. Our first discussion is from a purely descriptive point of view. That is, we do not assume that the data are generated by an underlying probability distribution. But as always, remember that the data themselves define a probability distribution, namely the *empirical distribution* that assigns equal probability to each data point.

Suppose that  $x$  and  $y$  are real-valued variables for a population, and that  $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$  is an observed sample of size  $n$  from  $(x, y)$ . We will let  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  denote the sample from  $x$  and  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  the sample from  $y$ . In this section, we are interested in statistics that are *measures of association* between the  $\mathbf{x}$  and  $\mathbf{y}$ , and in finding the line (or other curve) that best fits the data.

Recall that the sample means are

$$m(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n x_i, \quad m(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n y_i \quad (6.7.1)$$

and the sample variances are

$$s^2(\mathbf{x}) = \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})]^2, \quad s^2(\mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n [y_i - m(\mathbf{y})]^2 \quad (6.7.2)$$

### Scatterplots

Often, the first step in *exploratory data analysis* is to draw a graph of the points; this is called a *scatterplot* and can give a visual sense of the statistical relationship between the variables.

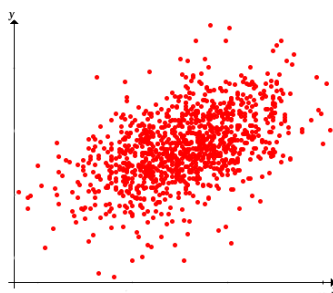


Figure 6.7.1: A scatterplot

In particular, we are interested in whether the cloud of points seems to show a linear trend or whether some nonlinear curve might fit the cloud of points. We are interested in the extent to which one variable  $x$  can be used to predict the other variable  $y$ .

### Definitions

Our next goal is to define statistics that measure the association between the  $\mathbf{x}$  and  $\mathbf{y}$  data.

The *sample covariance* is defined to be

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] \quad (6.7.3)$$

Assuming that the data vectors are not constant, so that the standard deviations are positive, the *sample correlation* is defined to be

$$r(\mathbf{x}, \mathbf{y}) = \frac{s(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})s(\mathbf{y})} \quad (6.7.4)$$

Note that the sample covariance is an average of the product of the deviations of the  $x$  and  $y$  data from their means. Thus, the physical unit of the sample covariance is the product of the units of  $x$  and  $y$ . Correlation is a standardized version of covariance. In particular, correlation is dimensionless (has no physical units), since the covariance in the numerator and the product of the standard deviations in the denominator have the same units (the product of the units of  $x$  and  $y$ ). Note also that covariance and correlation have the same *sign*: positive, negative, or zero. In the first case, the data  $\mathbf{x}$  and  $\mathbf{y}$  are said to be *positively correlated*; in the second case  $\mathbf{x}$  and  $\mathbf{y}$  are said to be *negatively correlated*; and in the third case  $\mathbf{x}$  and  $\mathbf{y}$  are said to be *uncorrelated*.

To see that the sample covariance is a measure of association, recall first that the point  $(m(\mathbf{x}), m(\mathbf{y}))$  is a measure of the center of the bivariate data. Indeed, if each point is the location of a unit mass, then  $(m(\mathbf{x}), m(\mathbf{y}))$  is the *center of mass* as defined in physics. Horizontal and vertical lines through this center point divide the plane into four quadrants. The product deviation  $[x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$  is positive in the first and third quadrants and negative in the second and fourth quadrants. After we study linear regression [below](#), we will have a much deeper sense of what covariance measures.

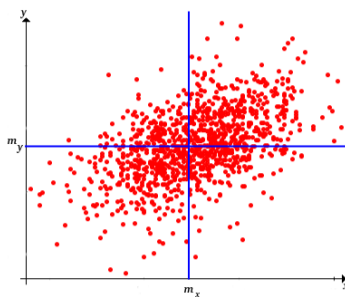


Figure 6.7.2: Scatterplot with means

You may be perplexed that we average the product deviations by dividing by  $n - 1$  rather than  $n$ . The best explanation is that in the probability model discussed [below](#), the sample covariance is an unbiased estimator of the distribution covariance. However, the mode of averaging can also be understood in terms of *degrees of freedom*, as was done for sample variance. Initially, we have  $2n$  degrees of freedom in the bivariate data. We lose two by computing the sample means  $m(\mathbf{x})$  and  $m(\mathbf{y})$ . Of the remaining  $2n - 2$  degrees of freedom, we lose  $n - 1$  by computing the product deviations. Thus, we are left with  $n - 1$  degrees of freedom total. As is typical in statistics, we average not by dividing by the number of terms in the sum but rather by the number of degrees of freedom in those terms. However, from a purely descriptive point of view, it would also be reasonable to divide by  $n$ .

Recall that there is a natural probability distribution associated with the data, namely the empirical distribution that gives probability  $\frac{1}{n}$  to each data point  $(x_i, y_i)$ . (Thus, if these points are distinct this is the discrete uniform distribution on the data.) The sample means are simply the expected values of this bivariate distribution, and except for a constant multiple (dividing by  $n - 1$  rather than  $n$ ), the sample variances are simply the variances of this bivariate distribution. Similarly, except for a constant multiple (again dividing by  $n - 1$  rather than  $n$ ), the sample covariance is the covariance of the bivariate distribution and the sample correlation is the correlation of the bivariate distribution. All of the following results in our discussion of descriptive statistics are actually special cases of more general results for probability distributions.

### Properties of Covariance

The next few exercises establish some essential properties of sample covariance. As usual, bold symbols denote samples of a fixed size  $n$  from the corresponding population variables (that is, vectors of length  $n$ ), while symbols in regular type denote real numbers. Our first result is a formula for sample covariance that is sometimes better than the definition for computational purposes. To state the result succinctly, let  $\mathbf{xy} = (x_1 y_1, x_2 y_2, \dots, x_n y_n)$  denote the sample from the product variable  $xy$ .

The sample covariance can be computed as follows:

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n x_i y_i - \frac{n}{n-1} m(\mathbf{x}) m(\mathbf{y}) = \frac{n}{n-1} [m(\mathbf{xy}) - m(\mathbf{x}) m(\mathbf{y})] \quad (6.7.5)$$

Proof

Note that

$$\sum_{i=1}^n [(x_i - m(\mathbf{x}))][y_i - m(\mathbf{y})] = \sum_{i=1}^n [x_i y_i - x_i m(\mathbf{y}) - y_i m(\mathbf{x}) + m(\mathbf{x}) m(\mathbf{y})] \quad (6.7.6)$$

$$= \sum_{i=1}^n x_i y_i - m(\mathbf{y}) \sum_{i=1}^n x_i - m(\mathbf{x}) \sum_{i=1}^n y_i + n m(\mathbf{x}) m(\mathbf{y}) \quad (6.7.7)$$

$$= \sum_{i=1}^n x_i y_i - n m(\mathbf{y}) m(\mathbf{x}) - n m(\mathbf{x}) m(\mathbf{y}) + n m(\mathbf{x}) m(\mathbf{y}) \quad (6.7.8)$$

$$= \sum_{i=1}^n x_i y_i - n m(\mathbf{x}) m(\mathbf{y}) \quad (6.7.9)$$

The following theorem gives another formula for the sample covariance, one that does not require the computation of intermediate statistics.

The sample covariance can be computed as follows:

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)(y_i - y_j) \quad (6.7.10)$$

Proof

Note that

$$\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)(y_i - y_j) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n [x_i - m(\mathbf{x}) + m(\mathbf{x}) - x_j][y_i - m(\mathbf{y}) + m(\mathbf{y}) - y_j] \quad (6.7.11)$$

$$= \sum_{i=1}^n \sum_{j=1}^n [(x_i - m(\mathbf{x}))[y_i - m(\mathbf{y})] + [x_i - m(\mathbf{x})][m(\mathbf{y}) - y_j] + [m(\mathbf{x}) - x_j][y_i - m(\mathbf{y})] + [m(\mathbf{x}) - x_j][m(\mathbf{y}) - y_j]] \quad (6.7.12)$$

We compute the sums term by term. The first is

$$n \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] \quad (6.7.13)$$

The second two sums are 0. The last sum is

$$n \sum_{j=1}^n [m(\mathbf{x}) - x_j][m(\mathbf{y}) - y_j] = n \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] \quad (6.7.14)$$

Dividing the entire sum by  $2n(n-1)$  results in  $\text{cov}(\mathbf{x}, \mathbf{y})$ .

As the name suggests, sample covariance generalizes sample variance.

$$s(\mathbf{x}, \mathbf{x}) = s^2(\mathbf{x}).$$

In light of the previous theorem, we can now see that the [first computational formula](#) and the [second computational formula](#) above generalize the computational formulas for sample variance. Clearly, sample covariance is *symmetric*.

$$s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}).$$

Sample covariance is linear in the first argument with the second argument fixed.

If  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$  are data vectors from population variables  $x$ ,  $y$ , and  $z$ , respectively, and if  $c$  is a constant, then

1.  $s(\mathbf{x} + \mathbf{y}, \mathbf{z}) = s(\mathbf{x}, \mathbf{z}) + s(\mathbf{y}, \mathbf{z})$
2.  $s(c\mathbf{x}, \mathbf{y}) = cs(\mathbf{x}, \mathbf{y})$

Proof

1. Recall that  $m(\mathbf{x} + \mathbf{y}) = m(\mathbf{x}) + m(\mathbf{y})$ . Hence

$$s(\mathbf{x} + \mathbf{y}, \mathbf{z}) = \frac{1}{n-1} \sum_{i=1}^n [x_i + y_i - m(\mathbf{x} + \mathbf{y})][z_i - m(\mathbf{z})] \quad (6.7.15)$$

$$= \frac{1}{n-1} \sum_{i=1}^n ([x_i - m(\mathbf{x})] + [y_i - m(\mathbf{y})])[z_i - m(\mathbf{z})] \quad (6.7.16)$$

$$= \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})][z_i - m(\mathbf{z})] + \frac{1}{n-1} \sum_{i=1}^n [y_i - m(\mathbf{y})][z_i - m(\mathbf{z})] \quad (6.7.17)$$

$$= s(\mathbf{x}, \mathbf{z}) + s(\mathbf{y}, \mathbf{z}) \quad (6.7.18)$$

2. Recall that  $m(c\mathbf{x}) = cm(\mathbf{x})$ . Hence

$$s(c\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n [cx_i - m(c\mathbf{x})][y_i - m(\mathbf{y})] \quad (6.7.19)$$

$$= \frac{1}{n-1} \sum_{i=1}^n [cx_i - cm(\mathbf{x})][y_i - m(\mathbf{y})] = cs(\mathbf{x}, \mathbf{y}) \quad (6.7.20)$$

By symmetry, sample covariance is also linear in the second argument with the first argument fixed, and hence is *bilinear*. The general version of the bilinear property is given in the following theorem:

Suppose that  $\mathbf{x}_i$  is a data vector from a population variable  $x_i$  for  $i \in \{1, 2, \dots, k\}$  and that  $\mathbf{y}_j$  is a data vector from a population variable  $y_j$  for  $j \in \{1, 2, \dots, l\}$ . Suppose also that  $a_1, a_2, \dots, a_k$  and  $b_1, b_2, \dots, b_l$  are constants. Then

$$s\left(\sum_{i=1}^k a_i \mathbf{x}_i, \sum_{j=1}^l b_j \mathbf{y}_j\right) = \sum_{i=1}^k \sum_{j=1}^l a_i b_j s(\mathbf{x}_i, \mathbf{y}_j) \quad (6.7.21)$$

A special case of the bilinear property provides a nice way to compute the sample variance of a sum.

$$s^2(\mathbf{x} + \mathbf{y}) = s^2(\mathbf{x}) + 2s(\mathbf{x}, \mathbf{y}) + s^2(\mathbf{y}).$$

Proof

From the preceding results,

$$s^2(\mathbf{x} + \mathbf{y}) = s(\mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y}) = s(\mathbf{x}, \mathbf{x}) + s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{x}) + s(\mathbf{y}, \mathbf{y}) \quad (6.7.22)$$

$$= s^2(\mathbf{x}) + 2s(\mathbf{x}, \mathbf{y}) + s^2(\mathbf{y}) \quad (6.7.23)$$

The generalization of this result to sums of three or more vectors is completely straightforward: namely, the sample variance of a sum is the sum of all of the pairwise sample covariances. Note that the sample variance of a sum can be greater than, less than, or equal to the sum of the sample variances, depending on the sign and magnitude of the pure covariance term. In particular, if the vectors are pairwise uncorrelated, then the variance of the sum is the sum of the variances.

If  $\mathbf{c}$  is a constant data set then  $s(\mathbf{x}, \mathbf{c}) = 0$ .

Proof

This follows directly from the definition. If  $c_i = c$  for each  $i$ , then  $m(\mathbf{c}) = c$  and hence  $c_i - m(\mathbf{c}) = 0$  for each  $i$ .

Combining the result in the last exercise with the [bilinear property](#), we see that covariance is unchanged if constants are added to the data sets. That is, if  $\mathbf{c}$  and  $\mathbf{d}$  are constant vectors then  $s(\mathbf{x} + \mathbf{c}, \mathbf{y} + \mathbf{d}) = s(\mathbf{x}, \mathbf{y})$ .

### Properties of Correlation

A few simple properties of correlation are given next. Most of these follow easily from the corresponding properties of covariance. First, recall that the standard scores of  $x_i$  and  $y_i$  are, respectively,

$$u_i = \frac{x_i - m(\mathbf{x})}{s(\mathbf{x})}, \quad v_i = \frac{y_i - m(\mathbf{y})}{s(\mathbf{y})} \quad (6.7.24)$$

The standard scores from a data set are dimensionless quantities that have mean 0 and variance 1.

The correlation between  $\mathbf{x}$  and  $\mathbf{y}$  is the covariance of their standard scores  $\mathbf{u}$  and  $\mathbf{v}$ . That is,  $r(\mathbf{x}, \mathbf{y}) = s(\mathbf{u}, \mathbf{v})$ .

Proof

In vector notation, note that

$$\mathbf{u} = \frac{1}{s(\mathbf{x})}[\mathbf{x} - m(\mathbf{x})], \quad \mathbf{v} = \frac{1}{s(\mathbf{y})}[\mathbf{y} - m(\mathbf{y})] \quad (6.7.25)$$

Hence the result follows immediately from properties of covariance:

$$s(\mathbf{u}, \mathbf{v}) = \frac{1}{s(\mathbf{x})s(\mathbf{y})}s(\mathbf{x}, \mathbf{y}) = r(\mathbf{x}, \mathbf{y}) \quad (6.7.26)$$

Correlation is symmetric.

$$r(\mathbf{x}, \mathbf{y}) = r(\mathbf{y}, \mathbf{x}).$$

Unlike covariance, correlation is unaffected by multiplying one of the data sets by a positive constant (recall that this can always be thought of as a change of scale in the underlying variable). On the other hand, multiplying a data set by a negative constant changes the sign of the correlation.

If  $c \neq 0$  is a constant then

1.  $r(c\mathbf{x}, \mathbf{y}) = r(\mathbf{x}, \mathbf{y})$  if  $c > 0$
2.  $r(c\mathbf{x}, \mathbf{y}) = -r(\mathbf{x}, \mathbf{y})$  if  $c < 0$

Proof

By definition and from the [scaling property of covariance](#),

$$r(c\mathbf{x}, \mathbf{y}) = \frac{s(c\mathbf{x}, \mathbf{y})}{s(c\mathbf{x})s(\mathbf{y})} = \frac{cs(\mathbf{x}, \mathbf{y})}{|c|s(\mathbf{x})s(\mathbf{y})} = \frac{c}{|c|}r(\mathbf{x}, \mathbf{y}) \quad (6.7.27)$$

and of course,  $c/|c| = 1$  if  $c > 0$  and  $c/|c| = -1$  if  $c < 0$ .

Like covariance, correlation is unaffected by adding constants to the data sets. Adding a constant to a data set often corresponds to a *change of location*.

If  $\mathbf{c}$  and  $\mathbf{d}$  are constant vectors then  $r(\mathbf{x} + \mathbf{c}, \mathbf{y} + \mathbf{d}) = r(\mathbf{x}, \mathbf{y})$ .

Proof

This result follows directly from the corresponding properties of covariance and standard deviation:

$$r(\mathbf{x} + \mathbf{c}, \mathbf{y} + \mathbf{d}) = \frac{s(\mathbf{x} + \mathbf{c}, \mathbf{y} + \mathbf{d})}{s(\mathbf{x} + \mathbf{c})s(\mathbf{y} + \mathbf{d})} = \frac{s(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})s(\mathbf{y})} = r(\mathbf{x}, \mathbf{y}) \quad (6.7.28)$$

The last couple of properties reinforce the fact that correlation is a standardized measure of association that is not affected by changing the units of measurement. In the first Challenger data set, for example, the variables of interest are temperature at time of launch (in degrees Fahrenheit) and O-ring erosion (in millimeters). The correlation between these variables is of critical importance. If we were to measure temperature in degrees Celsius and O-ring erosion in inches, the correlation between the two variables would be unchanged.

The most important properties of correlation arise from studying the line that best fits the data, our next topic.

### Linear Regression

We are interested in finding the line  $y = a + bx$  that best fits the sample points  $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$ . This is a basic and important problem in many areas of mathematics, not just statistics. We think of  $x$  as the *predictor variable* and  $y$  as the *response variable*. Thus, the term *best* means that we want to find the line (that is, find the coefficients  $a$  and  $b$ ) that minimizes the average of the squared errors between the actual  $y$  values in our data and the predicted  $y$  values:

$$\text{mse}(a, b) = \frac{1}{n-1} \sum_{i=1}^n [y_i - (a + bx_i)]^2 \quad (6.7.29)$$

Note that the minimizing value of  $(a, b)$  would be the same if the function were simply the sum of the squared errors, or if we averaged by dividing by  $n$  rather than  $n-1$ , or if we used the square root of any of these functions. Of course that actual *minimum value* of the function would be different if we changed the function, but again, not the point  $(a, b)$  where the minimum occurs. Our particular choice of mse as the error function is best for statistical purposes. Finding  $(a, b)$  that minimize mse is a standard problem in calculus.

The graph of mse is a paraboloid opening upward. The function mse is minimized when

$$b(\mathbf{x}, \mathbf{y}) = \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} \quad (6.7.30)$$

$$a(\mathbf{x}, \mathbf{y}) = m(\mathbf{y}) - b(\mathbf{x}, \mathbf{y})m(\mathbf{x}) = m(\mathbf{y}) - \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})}m(\mathbf{x}) \quad (6.7.31)$$

Proof

We can tell from the algebraic form of mse that the graph is a paraboloid opening upward. To find the unique point that minimizes mse, note that

$$\frac{\partial}{\partial a} \text{mse}(a, b) = \frac{1}{n-1} \sum 2[y_i - (a + bx_i)](-1) = \frac{2}{n-1} \left[ -\sum_{i=1}^n y_i + na + b \sum_{i=1}^n x_i \right] \quad (6.7.32)$$

$$\frac{\partial}{\partial b} \text{mse}(a, b) = \frac{1}{n-1} \sum 2[y_i - (a + bx_i)](-x_i) = \frac{2}{n-1} \left[ -\sum_{i=1}^n x_i y_i + a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \right] \quad (6.7.33)$$

Solving  $\frac{\partial}{\partial a} \text{mse}(a, b) = 0$ , gives  $a = m(\mathbf{y}) - bm(\mathbf{x})$ . Substituting this into  $\frac{\partial}{\partial b} \text{mse}(a, b) = 0$  and solving for  $b$  gives

$$b = \frac{n[m(\mathbf{x}\mathbf{y}) - m(\mathbf{x})m(\mathbf{y})]}{n[m(\mathbf{x}^2) - m^2(\mathbf{x})]} \quad (6.7.34)$$

Dividing the numerator and denominator in the last expression by  $n-1$  and using the computational formula above, we see that  $b = s(\mathbf{x}, \mathbf{y})/s^2(\mathbf{x})$ .

Of course, the optimal values of  $a$  and  $b$  are *statistics*, that is, functions of the data. Thus the *sample regression line* is

$$\mathbf{y} = m(\mathbf{y}) + \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} [\mathbf{x} - m(\mathbf{x})] \quad (6.7.35)$$

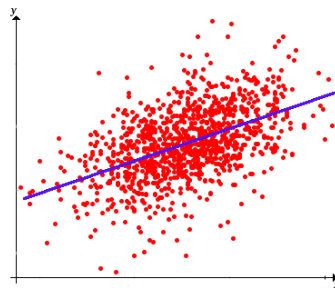


Figure 6.7.3: Scatterplot with regression line

Note that the regression line passes through the point  $(m(\mathbf{x}), m(\mathbf{y}))$ , the center of the sample of points.

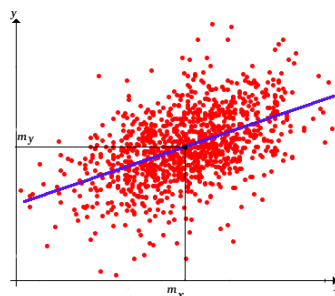


Figure 6.7.4: The regression line passes through the center

The minimum mean square error is

$$\text{mse} [a(\mathbf{x}, \mathbf{y}), b(\mathbf{x}, \mathbf{y})] = s(\mathbf{y})^2 [1 - r^2(\mathbf{x}, \mathbf{y})] \quad (6.7.36)$$

Proof

This follows from substituting  $a(\mathbf{x}, \mathbf{y})$   $b(\mathbf{x}, \mathbf{y})$  into mse and simplifying.

Sample correlation and covariance satisfy the following properties.

1.  $-1 \leq r(\mathbf{x}, \mathbf{y}) \leq 1$
2.  $-s(\mathbf{x})s(\mathbf{y}) \leq s(\mathbf{x}, \mathbf{y}) \leq s(\mathbf{x})s(\mathbf{y})$
3.  $r(\mathbf{x}, \mathbf{y}) = -1$  if and only if the sample points lie on a line with negative slope.
4.  $r(\mathbf{x}, \mathbf{y}) = 1$  if and only if the sample points lie on a line with positive slope.

Proof

Note that  $\text{mse} \geq 0$  and hence from the previous theorem, we must have  $r^2(\mathbf{x}, \mathbf{y}) \leq 1$ . This is equivalent to part (a), which in turn, from the definition of sample correlation, is equivalent to part (b). For parts (c) and (d), note that  $\text{mse}(a, b) = 0$  if and only if  $y_i = a + bx_i$  for each  $i$ , and moreover,  $b(\mathbf{x}, \mathbf{y})$  has the same sign as  $r(\mathbf{x}, \mathbf{y})$ .

Thus, we now see in a deeper way that the sample covariance and correlation measure the degree of linearity of the sample points. Recall from our discussion of measures of center and spread that the constant  $a$  that minimizes

$$\text{mse}(a) = \frac{1}{n-1} \sum_{i=1}^n (y_i - a)^2 \quad (6.7.37)$$

is the sample mean  $m(\mathbf{y})$ , and the minimum value of the mean square error is the sample variance  $s^2(\mathbf{y})$ . Thus, the difference between this value of the mean square error and the one above, namely  $s^2(\mathbf{y})r^2(\mathbf{x}, \mathbf{y})$  is the reduction in the variability of the  $y$  data when the linear term in  $x$  is added to the predictor. The fractional reduction is  $r^2(\mathbf{x}, \mathbf{y})$ ,

and hence this statistics is called the (sample) *coefficient of determination*. Note that if the data vectors  $\mathbf{x}$  and  $\mathbf{y}$  are uncorrelated, then  $x$  has no value as a predictor of  $y$ ; the regression line in this case is the horizontal line  $y = m(\mathbf{y})$  and the mean square error is  $s^2(\mathbf{y})$ .

The choice of predictor and response variables is important.

The sample regression line with predictor variable  $x$  and response variable  $y$  is not the same as the sample regression line with predictor variable  $y$  and response variable  $x$ , except in the extreme case  $r(\mathbf{x}, \mathbf{y}) = \pm 1$  where the sample points all lie on a line.

## Residuals

The difference between the actual  $y$  value of a data point and the value predicted by the regression line is called the *residual* of that data point. Thus, the residual corresponding to  $(x_i, y_i)$  is  $d_i = y_i - \hat{y}_i$  where  $\hat{y}_i$  is the regression line at  $x_i$ :

$$\hat{y}_i = m(\mathbf{y}) + \frac{s(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})^2} [x_i - m(\mathbf{x})] \quad (6.7.38)$$

Note that the predicted value  $\hat{y}_i$  and the residual  $d_i$  are *statistics*, that is, functions of the data  $(\mathbf{x}, \mathbf{y})$ , but we are suppressing this in the notation for simplicity.

The residuals sum to 0:  $\sum_{i=1}^n d_i = 0$ .

Proof

This follows from the definition, and is a restatement of the fact that the regression line passes through the center of the data set  $(m(\mathbf{x}), m(\mathbf{y}))$ .

Various plots of the residuals can help one understand the relationship between the  $x$  and  $y$  data. Some of the more common are given in the following definition:

### Residual plots

1. A plot of  $(i, d_i)$  for  $i \in \{1, 2, \dots, n\}$ , that is, a plot of indices versus residuals.
2. A plot of  $(x_i, d_i)$  for  $i \in \{1, 2, \dots, n\}$ , that is, a plot of  $x$  values versus residuals.
3. A plot of  $(d_i, y_i)$  for  $i \in \{1, 2, \dots, n\}$ , that is, a plot of residuals versus actual  $y$  values.
4. A plot of  $(d_i, \hat{y}_i)$  for  $i \in \{1, 2, \dots, n\}$ , that is a plot of residuals versus predicted  $y$  values.
5. A histogram of the residuals  $(d_1, d_2, \dots, d_n)$ .

## Sums of Squares

For our next discussion, we will re-interpret the minimum mean square error formula above. Here are the new definitions:

### Sums of squares

1.  $\text{sst}(\mathbf{y}) = \sum_{i=1}^n [y_i - m(\mathbf{y})]^2$  is the *total sum of squares*.
2.  $\text{ssr}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n [\hat{y}_i - m(\mathbf{y})]^2$  is the *regression sum of squares*
3.  $\text{sse}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  is the *error sum of squares*.

Note that  $\text{sst}(\mathbf{y})$  is simply  $n - 1$  times the variance  $s^2(\mathbf{y})$  and is the total of the sums of the squares of the deviations of the  $y$  values from the mean of the  $y$  values. Similarly,  $\text{sse}(\mathbf{x}, \mathbf{y})$  is simply  $n - 1$  times the minimum mean square error given above. Of course,  $\text{sst}(\mathbf{y})$  has  $n - 1$  degrees of freedom, while  $\text{sse}(\mathbf{x}, \mathbf{y})$  has  $n - 2$  degrees of freedom and  $\text{ssr}(\mathbf{x}, \mathbf{y})$  a single degree of freedom. The total sum of squares is the sum of the regression sum of squares and the error sum of squares:

The sums of squares are related as follows:

1.  $\text{ssr}(\mathbf{x}, \mathbf{y}) = r^2(\mathbf{x}, \mathbf{y}) \text{sst}(\mathbf{y})$
2.  $\text{sst}(\mathbf{y}) = \text{ssr}(\mathbf{x}, \mathbf{y}) + \text{sse}(\mathbf{x}, \mathbf{y})$

Proof

By definition of  $\text{sst}$  and  $r$ , we see that  $r^2(\mathbf{x}, \mathbf{y}) \text{sst}(\mathbf{y}) = s^2(\mathbf{x}, \mathbf{y}) / s^2(\mathbf{x})$ . But from the regression equation,

$$[\hat{y}_i - m(\mathbf{y})]^2 = \frac{s^2(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} [x_i - m(\mathbf{x})]^2 \quad (6.7.39)$$

Summing over  $i$  gives

$$\text{ssr}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n [\hat{y}_i - m(\mathbf{y})]^2 = \frac{s^2(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} \quad (6.7.40)$$

Hence  $\text{ssr}(\mathbf{x}, \mathbf{y}) = r^2(\mathbf{x}, \mathbf{y}) \text{sst}(\mathbf{y})$ . Finally, multiplying the result above by  $n - 1$  gives  $\text{sse}(\mathbf{x}, \mathbf{y}) = \text{sst}(\mathbf{y}) - r^2(\mathbf{x}, \mathbf{y}) \text{sst}(\mathbf{y}) = \text{sst}(\mathbf{y}) - \text{ssr}(\mathbf{x}, \mathbf{y})$ .

Note that  $r^2(\mathbf{x}, \mathbf{y}) = \text{ssr}(\mathbf{x}, \mathbf{y}) / \text{sst}(\mathbf{y})$ , so once again,  $r^2(\mathbf{x}, \mathbf{y})$  is the coefficient of determination—the proportion of the variability in the  $y$  data explained by the  $x$  data. We can average  $\text{sse}$  by dividing by its degrees of freedom and then take the square root to obtain a standard error:

The *standard error of estimate* is

$$\text{se}(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{\text{sse}(\mathbf{x}, \mathbf{y})}{n - 2}} \quad (6.7.41)$$

This really is a *standard* error in the same sense as a *standard* deviation. It's an average of the errors of sorts, but in the root mean square sense.

Finally, it's important to note that linear regression is a much more powerful idea than might first appear, and in fact the term *linear* can be a bit misleading. By applying various transformations to  $y$  or  $x$  or both, we can fit a variety of two-parameter curves to the given data  $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$ . Some of the most common transformations are explored in the exercises below.

## Probability Theory

We continue our discussion of sample covariance, correlation, and regression but now from the more interesting point of view that the variables are random. Specifically, suppose that we have a basic random experiment, and that  $X$  and  $Y$  are real-valued random variables for the experiment. Equivalently,  $(X, Y)$  is a random vector taking values in  $\mathbb{R}^2$ . Let  $\mu = \mathbb{E}(X)$  and  $\nu = \mathbb{E}(Y)$  denote the distribution means,  $\sigma^2 = \text{var}(X)$  and  $\tau^2 = \text{var}(Y)$  the distribution variances, and let  $\delta = \text{cov}(X, Y)$  denote the distribution covariance, so that the distribution correlation is

$$\rho = \text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\text{sd}(X)\text{sd}(Y)} = \frac{\delta}{\sigma\tau} \quad (6.7.42)$$

We will also need some higher order moments. Let  $\sigma_4 = \mathbb{E}[(X - \mu)^4]$ ,  $\tau_4 = \mathbb{E}[(Y - \nu)^4]$ , and  $\delta_2 = \mathbb{E}[(X - \mu)^2(Y - \nu)^2]$ . Naturally, we assume that all of these moments are finite.

Now suppose that we run the basic experiment  $n$  times. This creates a compound experiment with a sequence of independent random vectors  $((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$  each with the same distribution as  $(X, Y)$ . In statistical terms, this is a random sample of size  $n$  from the distribution of  $(X, Y)$ . The statistics discussed in previous section are well defined but now they are all random variables. We use the notation established previously, except that we use our usual convention of denoting random variables with capital letters. Of course, the deterministic properties and relations established [above](#) still hold. Note that  $\mathbf{X} = (X_1, X_2, \dots, X_n)$  is a random sample of size  $n$  from the distribution of  $X$  and  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$  is a random sample of size  $n$  from the distribution of  $Y$ . The main purpose of this subsection is to study the relationship between various statistics from  $\mathbf{X}$  and  $\mathbf{Y}$ , and to study statistics that are natural estimators of the distribution covariance and correlation.

### The Sample Means

Recall that the sample means are

$$M(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n X_i, \quad M(\mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n Y_i \quad (6.7.43)$$

From the sections on the law of large numbers and the central limit theorem, we know a great deal about the distributions of  $M(\mathbf{X})$  and  $M(\mathbf{Y})$  *individually*. But we need to know more about the *joint distribution*.

The covariance and correlation between  $M(\mathbf{X})$  and  $M(\mathbf{Y})$  are

1.  $\text{cov}[M(\mathbf{X}), M(\mathbf{Y})] = \delta/n$
2.  $\text{cor}[M(\mathbf{X}), M(\mathbf{Y})] = \rho$

Proof

Part (a) follows from the bilinearity of the covariance operator:

$$\text{cov}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{j=1}^n Y_j\right) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i, Y_j) \quad (6.7.44)$$

By independence, the terms in the last sum are 0 if  $i \neq j$ . For  $i = j$  the terms are  $\text{cov}(X, Y) = \delta$ . There are  $n$  such terms so  $\text{cov}[M(\mathbf{X}), M(\mathbf{Y})] = \delta/n$ . For part (b), recall that  $\text{var}[M(\mathbf{X})] = \sigma^2/n$  and  $\text{var}[M(\mathbf{Y})] = \tau^2/n$ . Hence

$$\text{cor}[M(\mathbf{X}), M(\mathbf{Y})] = \frac{\delta/n}{(\sigma/\sqrt{n})(\tau/\sqrt{n})} = \frac{\delta}{\sigma\tau} = \rho \quad (6.7.45)$$

Note that the correlation between the sample means is the same as the correlation of the underlying sampling distribution. In particular, the correlation does not depend on the sample size  $n$ .

### The Sample Variances

Recall that special versions of the sample variances, in the unlikely event that the distribution means are known, are

$$W^2(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2, \quad W^2(\mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \nu)^2 \quad (6.7.46)$$

Once again, we have studied these statistics individually, so our emphasis now is on the joint distribution.

The covariance and correlation between  $W^2(\mathbf{X})$  and  $W^2(\mathbf{Y})$  are

1.  $\text{cov}[W^2(\mathbf{X}), W^2(\mathbf{Y})] = (\delta_2 - \sigma^2\tau^2)/n$
2.  $\text{cor}[W^2(\mathbf{X}), W^2(\mathbf{Y})] = (\delta_2 - \sigma^2\tau^2)/\sqrt{(\sigma_4 - \sigma^4)(\tau_4 - \tau^4)}$

Proof

For part (a), we use the bilinearity of the covariance operator to obtain

$$\text{cov}[W^2(\mathbf{X}), W^2(\mathbf{Y})] = \text{cov}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2, \frac{1}{n} \sum_{j=1}^n (Y_j - \nu)^2\right) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \text{cov}[(X_i - \mu)^2, (Y_j - \nu)^2] \quad (6.7.47)$$

By independence, the terms in the last sum are 0 when  $i \neq j$ . When  $i = j$  the terms are

$$\text{cov}[(X - \mu)^2(Y - \nu)^2] = \mathbb{E}[(X - \mu)^2(Y - \nu)^2] - \mathbb{E}[(X - \mu)^2]\mathbb{E}[(Y - \nu)^2] = \delta_2 - \sigma^2\tau^2 \quad (6.7.48)$$

There are  $n$  such terms, so  $\text{cov}[W^2(\mathbf{X}), W^2(\mathbf{Y})] = (\delta_2 - \sigma^2\tau^2)/n$ . Part (b) follows from part (a) and the variances of  $W^2(\mathbf{X})$  and  $W^2(\mathbf{Y})$  from the section on Sample Variance.

Note that the correlation does not depend on the sample size  $n$ . Next, recall that the standard versions of the sample variances are

$$S^2(\mathbf{X}) = \frac{1}{n-1} \sum_{i=1}^n [X_i - M(\mathbf{X})]^2, \quad S^2(\mathbf{Y}) = \frac{1}{n-1} \sum_{i=1}^n [Y_i - M(\mathbf{Y})]^2 \quad (6.7.49)$$

The covariance and correlation of the sample variances are

1.  $\text{cov}[S^2(\mathbf{X}), S^2(\mathbf{Y})] = (\delta_2 - \sigma^2 \tau^2) / n + 2\delta^2 / [n(n-1)]$
2.  $\text{cor}[S^2(\mathbf{X}), S^2(\mathbf{Y})] = [(n-1)(\delta_2 - \sigma^2 \tau^2) + 2\delta^2] / \sqrt{[(n-1)\sigma_4 - (n-3)\sigma^4][(n-1)\tau_4 - (n-3)\tau^4]}$

Proof

Recall that

$$S^2(\mathbf{X}) = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)^2, \quad S^2(\mathbf{Y}) = \frac{1}{2n(n-1)} \sum_{k=1}^n \sum_{l=1}^n (Y_k - Y_l)^2 \quad (6.7.50)$$

Hence using the bilinearity of the covariance operator we have

$$\text{cov}[S^2(\mathbf{X}), S^2(\mathbf{Y})] = \frac{1}{4n^2(n-1)^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \text{cov}[(X_i - X_j)^2, (Y_k - Y_l)^2] \quad (6.7.51)$$

We compute the covariances in this sum by considering disjoint cases:

- $\text{cov}[(X_i - X_j)^2, (Y_k - Y_l)^2] = 0$  if  $i = j$  or if  $k = l$ , and there are  $2n^3 - n^2$  such terms.
- $\text{cov}[(X_i - X_j)^2, (Y_k - Y_l)^2] = 0$  by independence if  $i, j, k, l$  are distinct, and there are  $n(n-1)(n-2)(n-3)$  such terms.
- $\text{cov}[(X_i - X_j)^2, (Y_k - Y_l)^2] = 2\delta_2 - 2\sigma^2 \tau^2 + 4\delta^2$  if  $i \neq j$  and  $\{k, l\} = \{i, j\}$ , and there are  $2n(n-1)$  such terms.
- $\text{cov}[(X_i - X_j)^2, (Y_k - Y_l)^2] = \delta_2 - \sigma^2 \tau^2$  if  $i \neq j, k \neq l$ , and  $\#\{i, j\} \cap \{k, l\} = 1$ , and there are  $4n(n-1)(n-2)$  such terms.

Substituting and simplifying gives the result in (a). For (b), we use the definition of correlation and the formulas for  $\text{var}[S^2(\mathbf{X})]$  and  $\text{var}[S^2(\mathbf{Y})]$  from the section on the sample variance.

Asymptotically, the correlation between the sample variances is the same as the correlation between the special sample variances given above:

$$\text{cor}[S^2(\mathbf{X}), S^2(\mathbf{Y})] \rightarrow \frac{\delta_2 - \sigma^2 \tau^2}{\sqrt{(\sigma_4 - \sigma^4)(\tau_4 - \tau^4)}} \text{ as } n \rightarrow \infty \quad (6.7.52)$$

### Sample Covariance

Suppose first that the distribution means  $\mu$  and  $\nu$  are known. As noted earlier, this is almost always an unrealistic assumption, but is still a good place to start because the analysis is very simple and the results we obtain will be useful below. A natural estimator of the distribution covariance  $\delta = \text{cov}(X, Y)$  in this case is the *special sample covariance*

$$W(\mathbf{X}, \mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)(Y_i - \nu) \quad (6.7.53)$$

Note that the special sample covariance generalizes the special sample variance:  $W(\mathbf{X}, \mathbf{X}) = W^2(\mathbf{X})$ .

$W(\mathbf{X}, \mathbf{Y})$  is the sample mean for a random sample of size  $n$  from the distribution of  $(X - \mu)(Y - \nu)$  and satisfies the following properties:

1.  $\mathbb{E}[W(\mathbf{X}, \mathbf{Y})] = \delta$
2.  $\text{var}[W(\mathbf{X}, \mathbf{Y})] = \frac{1}{n}(\delta_2 - \delta^2)$
3.  $W(\mathbf{X}, \mathbf{Y}) \rightarrow \delta$  as  $n \rightarrow \infty$  with probability 1

Proof

These results follow directly from the section on the Law of Large Numbers. For part (b), note that

$$\text{var}[(X - \mu)(Y - \nu)] = \mathbb{E}[(X - \mu)^2(Y - \nu)^2] - (\mathbb{E}[(X - \mu)(Y - \nu)])^2 = \delta_2 - \delta^2 \quad (6.7.54)$$

As an estimator of  $\delta$ , part (a) means that  $W(\mathbf{X}, \mathbf{Y})$  is *unbiased* and part (b) means that  $W(\mathbf{X}, \mathbf{Y})$  is *consistent*.

Consider now the more realistic assumption that the distribution means  $\mu$  and  $\nu$  are unknown. A natural approach in this case is to average  $[(X_i - M(\mathbf{X}))][Y_i - M(\mathbf{Y})]$  over  $i \in \{1, 2, \dots, n\}$ . But rather than dividing by  $n$  in our average, we should divide by whatever constant gives an unbiased estimator of  $\delta$ . As shown in the next theorem, this constant turns out to be  $n-1$ , leading to the standard *sample covariance*:

$$S(\mathbf{X}, \mathbf{Y}) = \frac{1}{n-1} \sum_{i=1}^n [X_i - M(\mathbf{X})][Y_i - M(\mathbf{Y})] \quad (6.7.55)$$

$\mathbb{E}[S(\mathbf{X}, \mathbf{Y})] = \delta$ .

Proof

Expanding as above we have,

$$\sum_{i=1}^n [X_i - M(\mathbf{X})][Y_i - M(\mathbf{Y})] = \sum_{i=1}^n X_i Y_i - nM(\mathbf{X})M(\mathbf{Y}) \quad (6.7.56)$$

But  $\mathbb{E}[X_i Y_i] = \text{cov}(X_i, Y_i) + \mathbb{E}(X_i)\mathbb{E}(Y_i) = \delta + \mu\nu$ . Similarly, from the covariance of the sample means and the unbiased property,  $\mathbb{E}[M(\mathbf{X})M(\mathbf{Y})] = \text{cov}[M(\mathbf{X}), M(\mathbf{Y})] + \mathbb{E}[M(\mathbf{X})]\mathbb{E}[M(\mathbf{Y})] = \delta/n + \mu\nu$ . So taking expected values in the displayed equation above gives

$$\mathbb{E}\left(\sum_{i=1}^n [X_i - M(\mathbf{X})][Y_i - M(\mathbf{Y})]\right) = n(\delta + \mu\nu) - n(\delta/n + \mu\nu) = (n-1)\delta \quad (6.7.57)$$

$S(\mathbf{X}, \mathbf{Y}) \rightarrow \delta$  as  $n \rightarrow \infty$  with probability 1.

Proof

Once again, we have

$$S(\mathbf{X}, \mathbf{Y}) = \frac{n}{n-1} [M(\mathbf{XY}) - M(\mathbf{X})M(\mathbf{Y})] \quad (6.7.58)$$

where  $M(\mathbf{XY})$  denotes the sample mean for the sample of the products  $(X_1Y_1, X_2Y_2, \dots, X_nY_n)$ . By the strong law of large numbers,  $M(\mathbf{X}) \rightarrow \mu$  as  $n \rightarrow \infty$ ,  $M(\mathbf{Y}) \rightarrow \nu$  as  $n \rightarrow \infty$ , and  $M(\mathbf{XY}) \rightarrow \mathbb{E}(XY) = \delta + \mu\nu$  as  $n \rightarrow \infty$ , each with probability 1. So the result follows by letting  $n \rightarrow \infty$  in the displayed equation.

Of course, the *sample correlation* is

$$R(\mathbf{X}, \mathbf{Y}) = \frac{S(\mathbf{X}, \mathbf{Y})}{S(\mathbf{X})S(\mathbf{Y})} \quad (6.7.59)$$

Since the sample correlation  $R(\mathbf{X}, \mathbf{Y})$  is a nonlinear function of the sample covariance and sample standard deviations, it will not in general be an unbiased estimator of the distribution correlation  $\rho$ . In most cases, it would be difficult to even compute the mean and variance of  $R(\mathbf{X}, \mathbf{Y})$ . Nonetheless, we can show convergence of the sample correlation to the distribution correlation.

$R(\mathbf{X}, \mathbf{Y}) \rightarrow \rho$  as  $n \rightarrow \infty$  with probability 1.

Proof

This follows immediately from the strong law of large numbers and previous results. From the result [above](#)  $S(\mathbf{X}, \mathbf{Y}) \rightarrow \delta$  as  $n \rightarrow \infty$ , and from the section on the sample variance,  $S(\mathbf{X}) \rightarrow \sigma$  as  $n \rightarrow \infty$  and  $S(\mathbf{Y}) \rightarrow \tau$  as  $n \rightarrow \infty$ , each with probability 1. Hence  $R(\mathbf{X}, \mathbf{Y}) \rightarrow \delta/\sigma\tau = \rho$  as  $n \rightarrow \infty$  with probability 1.

Our next theorem gives a formula for the variance of the sample covariance, not to be confused with the covariance of the sample variances given [above](#)!

The variance of the sample covariance is

$$\text{var}[S(\mathbf{X}, \mathbf{Y})] = \frac{1}{n} \left( \delta_2 + \frac{1}{n-1} \sigma^2 \tau^2 - \frac{n-2}{n-1} \delta^2 \right) \quad (6.7.60)$$

Proof

Recall first that

$$S(\mathbf{X}, \mathbf{Y}) = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)(Y_i - Y_j) \quad (6.7.61)$$

Hence using the bilinearity of the covariance operator we have

$$\text{var}[S(\mathbf{X}, \mathbf{Y})] = \frac{1}{4n^2(n-1)^2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \text{cov}[(X_i - X_j)(Y_i - Y_j), (X_k - X_l)(Y_k - Y_l)] \quad (6.7.62)$$

We compute the covariances in this sum by considering disjoint cases:

- $\text{cov}[(X_i - X_j)(Y_i - Y_j), (X_k - X_l)(Y_k - Y_l)] = 0$  if  $i = j$  or if  $k = l$ , and there are  $2n^3 - n^2$  such terms.
- $\text{cov}[(X_i - X_j)(Y_i - Y_j), (X_k - X_l)(Y_k - Y_l)] = 0$  if  $i, j, k, l$  are distinct, and there are  $n(n-1)(n-2)(n-3)$  such terms.
- $\text{cov}[(X_i - X_j)(Y_i - Y_j), (X_k - X_l)(Y_k - Y_l)] = 2\delta_2 + 2\sigma^2\tau^2$  if  $i \neq j$  and  $\{k, l\} = \{i, j\}$ , and there are  $2n(n-1)$  such terms.
- $\text{cov}[(X_i - X_j)(Y_i - Y_j), (X_k - X_l)(Y_k - Y_l)] = \delta_2 - \delta^2$  if  $i \neq j, k \neq l$ , and  $\#\{i, j\} \cap \{k, l\} = 1$ , and there are  $4n(n-1)(n-2)$  such terms.

Substituting and simplifying gives the result

It's not surprising that the variance of the standard sample covariance (where we don't know the distribution means) is greater than the variance of the special sample covariance (where we do know the distribution means).

$\text{var}[S(\mathbf{X}, \mathbf{Y})] > \text{var}[W(\mathbf{X}, \mathbf{Y})]$

Proof

From results above, and some simple algebra,

$$\text{var}[S(\mathbf{X}, \mathbf{Y})] - \text{var}[W(\mathbf{X}, \mathbf{Y})] = \frac{1}{n(n-1)} (\delta^2 + \sigma^2\tau^2) > 0 \quad (6.7.63)$$

But note that the difference goes to 0 as  $n \rightarrow \infty$ .

$\text{var}[S(\mathbf{X}, \mathbf{Y})] \rightarrow 0$  as  $n \rightarrow \infty$ . Thus, the sample covariance is a *consistent* estimator of the distribution covariance.

## Regression

In our first discussion [above](#), we studied regression from a deterministic, descriptive point of view. The results obtained applied only to the sample. Statistically more interesting and deeper questions arise when the data come from a random experiment, and we try to draw inferences about the underlying distribution from the sample regression. There are two models that commonly arise. One is where the response variable is random, but the predictor variable is deterministic. The other is the model we consider here, where the predictor variable and the response variable are both random, so that the data form a random sample from a bivariate distribution.

Thus, suppose again that we have a basic random vector  $(X, Y)$  for an experiment. Recall that in the section on (distribution) correlation and regression, we showed that the best linear predictor of  $Y$  given  $X$ , in the sense of minimizing mean square error, is the random variable

$$L(Y | X) = \mathbb{E}(Y) + \frac{\text{cov}(X, Y)}{\text{var}(X)} [X - \mathbb{E}(X)] = \nu + \frac{\delta}{\sigma^2} (X - \mu) \quad (6.7.64)$$

so that the *distribution regression line* is given by

$$y = L(Y | X = x) = \nu + \frac{\delta}{\sigma^2}(x - \mu) \quad (6.7.65)$$

Moreover, the (minimum) value of the mean square error is  $\mathbb{E}\{[Y - L(Y | X)]\} = \text{var}(Y)[1 - \text{cor}^2(X, Y)] = r^2(1 - \rho^2)$ .

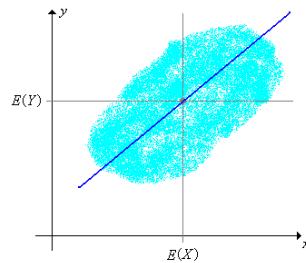


Figure 6.7.5: The distribution regression line

Of course, in real applications, we are unlikely to know the distribution parameters  $\mu$ ,  $\nu$ ,  $\sigma^2$ , and  $\delta$ . If we want to estimate the distribution regression line, a natural approach would be to consider a random sample  $((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$  from the distribution of  $(X, Y)$  and compute the sample regression line. Of course, the results are exactly the same as in the discussion [above](#), except that all of the relevant quantities are random variables. The sample regression line is

$$y = M(\mathbf{Y}) + \frac{S(\mathbf{X}, \mathbf{Y})}{S^2(\mathbf{X})}[x - M(\mathbf{X})] \quad (6.7.66)$$

The mean square error is  $S^2(\mathbf{Y})[1 - R^2(\mathbf{X}, \mathbf{Y})]$  and the coefficient of determination is  $R^2(\mathbf{X}, \mathbf{Y})$ .

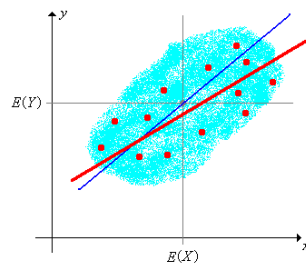


Figure 6.7.6: The distribution and sample regression lines

The fact that the sample regression line and mean square error are completely analogous to the distribution regression line and mean square error is mathematically elegant and reassuring. Again, the coefficients of the sample regression line can be viewed as estimators of the respective coefficients in the distribution regression line.

The coefficients of the sample regression line converge to the coefficients of the distribution regression line with probability 1.

1.  $\frac{S(\mathbf{X}, \mathbf{Y})}{S^2(\mathbf{X})} \rightarrow \frac{\delta}{\sigma^2}$  as  $n \rightarrow \infty$
2.  $M(\mathbf{Y}) - \frac{S(\mathbf{X}, \mathbf{Y})}{S^2(\mathbf{X})}M(\mathbf{X}) \rightarrow \nu - \frac{\delta}{\sigma^2}\mu$  as  $n \rightarrow \infty$

Proof

This follows from the strong law of large numbers and previous results. with probability 1,  $S(\mathbf{X}, \mathbf{Y}) \rightarrow \delta$  as  $n \rightarrow \infty$ ,  $S^2(\mathbf{X}) \rightarrow \sigma^2$  as  $n \rightarrow \infty$ ,  $M(\mathbf{X}) \rightarrow \mu$  as  $n \rightarrow \infty$ , and  $M(\mathbf{Y}) \rightarrow \nu$  as  $n \rightarrow \infty$ .

Of course, if the linear relationship between  $X$  and  $Y$  is not strong, as measured by the sample correlation, then transformation applied to one or both variables may help. Again, some typical transformations are explored in the exercises [below](#).

## Exercises

### Basic Properties

Suppose that  $x$  and  $y$  are population variables, and  $\mathbf{x}$  and  $\mathbf{y}$  samples of size  $n$  from  $x$  and  $y$  respectively. Suppose also that  $m(\mathbf{x}) = 3$ ,  $m(\mathbf{y}) = -1$ ,  $s^2(\mathbf{x}) = 4$ ,  $s^2(\mathbf{y}) = 9$ , and  $s(\mathbf{x}, \mathbf{y}) = 5$ . Find each of the following:

1.  $r(\mathbf{x}, \mathbf{y})$
2.  $m(2\mathbf{x} + 3\mathbf{y})$
3.  $s^2(2\mathbf{x} + 3\mathbf{y})$
4.  $s(2\mathbf{x} + 3\mathbf{y} - 1, 4\mathbf{x} + 2\mathbf{y} - 3)$

Suppose that  $x$  is the temperature (in degrees Fahrenheit) and  $y$  the resistance (in ohms) for a certain type of electronic component after 10 hours of operation. For a sample of 30 components,  $m(\mathbf{x}) = 113$ ,  $s(\mathbf{x}) = 18$ ,  $m(\mathbf{y}) = 100$ ,  $s(\mathbf{y}) = 10$ ,  $r(\mathbf{x}, \mathbf{y}) = 0.6$ .

1. Classify  $x$  and  $y$  by type and level of measurement.
2. Find the sample covariance.
3. Find the equation of the regression line.

Suppose now that temperature is converted to degrees Celsius (the transformation is  $\frac{5}{9}(x - 32)$ ).

4. Find the sample means.

5. Find the sample standard deviations.
6. Find the sample covariance and correlation.
7. Find the equation of the regression line.

Answer

1. continuous, interval
2.  $m = 45^\circ$ ,  $s = 10^\circ$

Suppose that  $x$  is the length and  $y$  the width (in inches) of a leaf in a certain type of plant. For a sample of 50 leaves  $m(\mathbf{x}) = 10$ ,  $s(\mathbf{x}) = 2$ ,  $m(\mathbf{y}) = 4$ ,  $s(\mathbf{y}) = 1$ , and  $r(\mathbf{x}, \mathbf{y}) = 0.8$ .

1. Classify  $x$  and  $y$  by type and level of measurement.
2. Find the sample covariance.
3. Find the equation of the regression line with  $x$  as the predictor variable and  $y$  as the response variable.

Suppose now that  $x$  and  $y$  are converted to inches (0.3937 inches per centimeter).

4. Find the sample means.
5. Find the sample standard deviations.
6. Find the sample covariance and correlation.
7. Find the equation of the regression line.

Answer

1. continuous, ratio
2.  $m = 25.4$ ,  $s = 5.08$

### Scatterplot Exercises

Click in the interactive scatterplot, in various places, and watch how the means, standard deviations, correlation, and regression line change.

Click in the interactive scatterplot to define 20 points and try to come as close as possible to each of the following sample correlations:

1. 0
2. 0.5
3. -0.5
4. 0.7
5. -0.7
6. 0.9
7. -0.9.

Click in the interactive scatterplot to define 20 points. Try to generate a scatterplot in which the regression line has

1. slope 1, intercept 1
2. slope 3, intercept 0
3. slope -2, intercept 1

### Simulation Exercises

Run the bivariate uniform experiment 2000 times in each of the following cases. Compare the sample means to the distribution means, the sample standard deviations to the distribution standard deviations, the sample correlation to the distribution correlation, and the sample regression line to the distribution regression line.

1. The uniform distribution on the square
2. The uniform distribution on the triangle.
3. The uniform distribution on the circle.

Run the bivariate normal experiment 2000 times for various values of the distribution standard deviations and the distribution correlation. Compare the sample means to the distribution means, the sample standard deviations to the distribution standard deviations, the sample correlation to the distribution correlation, and the sample regression line to the distribution regression line.

### Transformations

Consider the function  $y = a + bx^2$ .

1. Sketch the graph for some representative values of  $a$  and  $b$ .
2. Note that  $y$  is a linear function of  $x^2$ , with intercept  $a$  and slope  $b$ .
3. Hence, to fit this curve to sample data, simply apply the standard regression procedure to the data from the variables  $x^2$  and  $y$ .

Consider the function  $y = \frac{1}{a+bx}$ .

1. Sketch the graph for some representative values of  $a$  and  $b$ .
2. Note that  $\frac{1}{y}$  is a linear function of  $x$ , with intercept  $a$  and slope  $b$ .
3. Hence, to fit this curve to our sample data, simply apply the standard regression procedure to the data from the variables  $x$  and  $\frac{1}{y}$ .

Consider the function  $y = \frac{x}{a+bx}$ .

1. Sketch the graph for some representative values of  $a$  and  $b$ .
2. Note that  $\frac{1}{y}$  is a linear function of  $\frac{1}{x}$ , with intercept  $b$  and slope  $a$ .

- Hence, to fit this curve to sample data, simply apply the standard regression procedure to the data from the variables  $\frac{1}{x}$  and  $\frac{1}{y}$ .
- Note again that the names of the intercept and slope are reversed from the standard formulas.

Consider the function  $y = ae^{bx}$ .

- Sketch the graph for some representative values of  $a$  and  $b$ .
- Note that  $\ln(y)$  is a linear function of  $x$ , with intercept  $\ln(a)$  and slope  $b$ .
- Hence, to fit this curve to sample data, simply apply the standard regression procedure to the data from the variables  $x$  and  $\ln(y)$ .
- After solving for the intercept  $\ln(a)$ , recover the statistic  $a = e^{\ln(a)}$ .

Consider the function  $y = ax^b$ .

- Sketch the graph for some representative values of  $a$  and  $b$ .
- Note that  $\ln(y)$  is a linear function of  $\ln(x)$ , with intercept  $\ln(a)$  and slope  $b$ .
- Hence, to fit this curve to sample data, simply apply the standard regression procedure to the data from the variables  $\ln(x)$  and  $\ln(y)$ .
- After solving for the intercept  $\ln(a)$ , recover the statistic  $a = e^{\ln(a)}$ .

### Computational Exercises

All statistical software packages will perform regression analysis. In addition to the regression line, most packages will typically report the coefficient of determination  $r^2(\mathbf{x}, \mathbf{y})$ , the sums of squares  $\text{sst}(\mathbf{y})$ ,  $\text{ssr}(\mathbf{x}, \mathbf{y})$ ,  $\text{sse}(\mathbf{x}, \mathbf{y})$ , and the standard error of estimate  $\text{se}(\mathbf{x}, \mathbf{y})$ . Most packages will also draw the scatterplot, with the regression line superimposed, and will draw the various graphs of residuals discussed above. Many packages also provide easy ways to transform the data. Thus, there is very little reason to perform the computations by hand, except with a small data set to master the definitions and formulas. In the following problem, do the computations and draw the graphs with minimal technological aids.

Suppose that  $x$  is the number of math courses completed and  $y$  the number of science courses completed for a student at Enormous State University (ESU). A sample of 10 ESU students gives the following data:  $((1, 1), (3, 3), (6, 4), (2, 1), (8, 5), (2, 2), (4, 3), (6, 4), (4, 3), (4, 4))$

- Classify  $x$  and  $y$  by type and level of measurement.
- Sketch the scatterplot.

Construct a table with rows corresponding to cases and columns corresponding to  $i$ ,  $x_i$ ,  $y_i$ ,  $x_i - m(\mathbf{x})$ ,  $y_i - m(\mathbf{y})$ ,  $[x_i - m(\mathbf{x})]^2$ ,  $[y_i - m(\mathbf{y})]^2$ ,  $[x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$ ,  $\hat{y}_i$ ,  $\hat{y}_i - m(\mathbf{y})$ ,  $[\hat{y}_i - m(\mathbf{y})]^2$ ,  $y_i - \hat{y}_i$ , and  $(y_i - \hat{y}_i)^2$ . Add a row at the bottom for totals and means. Use precision arithmetic.

- Complete the first 8 columns.
- Find the sample correlation and the coefficient of determination.
- Find the sample regression equation.
- Complete the table.
- Verify the identities for the sums of squares.

Answer

$i$	$x_i$	$y_i$	$x_i - m(\mathbf{x})$	$y_i - m(\mathbf{y})$	$[x_i - m(\mathbf{x})]^2$	$[y_i - m(\mathbf{y})]^2$	$[x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$	$\hat{y}_i$	$\hat{y}_i - m(\mathbf{y})$	$[\hat{y}_i - m(\mathbf{y})]^2$	$y_i - \hat{y}_i$	$(y_i - \hat{y}_i)^2$
1	1	1	-3	-2	9	4	6	9/7	-12/7	144/49	-2/7	4/49
2	3	3	-1	0	1	0	0	17/7	-4/7	16/49	4/7	16/49
3	6	4	2	1	4	1	2	29/7	8/7	64/49	-1/7	1/49
4	2	1	-2	-2	4	4	4	13/7	-8/7	64/49	-6/7	36/49
5	8	5	4	2	16	4	8	37/7	16/7	256/49	-2/7	4/49
6	2	2	-2	-1	4	1	2	13/7	-8/7	64/49	1/7	1/49
7	4	3	0	0	0	0	0	3	0	0	0	0
8	6	4	2	1	4	1	2	29/7	8/7	64/49	-1/7	1/49
9	4	3	0	0	0	0	0	3	0	0	0	0
10	4	4	0	1	0	1	0	3	0	0	1	1
Total	40	30	0	0	42	16	24	30	0	96/7	0	16/7
Mean	4	3	0	0	14/3	16/9	8/3	3	0	96/7	0	2/7

- discrete, ratio
- $r = 2\sqrt{3/14} \approx 0.926$ ,  $r^2 = 6/7$
- $y = 3 + \frac{4}{7}(x - 4)$
- $16 = 96/7 + 16/7$

The following two exercise should help you review some of the probability topics in this section.

Suppose that  $(X, Y)$  has a continuous distribution with probability density function  $f(x, y) = 15x^2y$  for  $0 \leq x \leq y \leq 1$ . Find each of the following:

- $\mu = \mathbb{E}(X)$  and  $\nu = \mathbb{E}(Y)$
- $\sigma^2 = \text{var}(X)$  and  $\tau^2 = \text{var}(Y)$
- $\sigma_3 = \mathbb{E}[(X - \mu)^3]$  and  $\tau_3 = \mathbb{E}[(Y - \nu)^3]$
- $\sigma_4 = \mathbb{E}[(X - \mu)^4]$  and  $\tau_4 = \mathbb{E}[(Y - \nu)^4]$
- $\delta = \text{cov}(X, Y)$ ,  $\rho = \text{cor}(X, Y)$ , and  $\delta_2 = \mathbb{E}[(X - \mu)^2(Y - \nu)^2]$

6.  $L(Y | X)$  and  $L(X | Y)$

Answer

1.  $5/8$   $5/6$
2.  $17/448$   $5/252$
3.  $-5/1792$   $-5/1512$
4.  $305/86$   $0165/3024$
5.  $5/336$   $\sqrt{5/17}$   $1/768$
6.  $L(Y | X) = \frac{10}{17} + \frac{20}{51}X$ ,  $L(X | Y) = \frac{3}{4}Y$

Suppose now that  $((X_1, Y_1), (X_2, Y_2), \dots, (X_9, Y_9))$  is a random sample of size 9 from the distribution in the previous exercise. Find each of the following:

1.  $\mathbb{E}[M(\mathbf{X})]$  and  $\text{var}[M(\mathbf{X})]$
2.  $\mathbb{E}[M(\mathbf{Y})]$  and  $\text{var}[M(\mathbf{Y})]$
3.  $\text{cov}[M(\mathbf{X}), M(\mathbf{Y})]$  and  $\text{cor}[M(\mathbf{X}), M(\mathbf{Y})]$
4.  $\mathbb{E}[W^2(\mathbf{X})]$  and  $\text{var}[W^2(\mathbf{X})]$
5.  $\mathbb{E}[W^2(\mathbf{Y})]$  and  $\text{var}[W^2(\mathbf{Y})]$
6.  $\mathbb{E}[S^2(\mathbf{X})]$  and  $\text{var}[S^2(\mathbf{X})]$
7.  $\mathbb{E}[S^2(\mathbf{Y})]$  and  $\text{var}[S^2(\mathbf{Y})]$
8.  $\mathbb{E}[W(\mathbf{X}, \mathbf{Y})]$  and  $\text{var}[W(\mathbf{X}, \mathbf{Y})]$
9.  $\mathbb{E}[S(\mathbf{X}, \mathbf{Y})]$  and  $\text{var}[S(\mathbf{X}, \mathbf{Y})]$

Answer

1.  $5/8$   $17/4032$
2.  $5/6$   $5/2268$
3.  $5/3024$   $\sqrt{5/17}$
4.  $17/448$   $317/1354752$
5.  $5/252$   $5/35721$
6.  $17/448$   $5935/21676032$
7.  $5/252$   $115/762048$
8.  $5/336$   $61/508032$
9.  $5/336$   $181/1354752$

### Data Analysis Exercises

Use statistical software for the following problems.

Consider the height variables in Pearson's height data.

1. Classify the variables by type and level of measurement.
2. Compute the correlation coefficient and the coefficient of determination.
3. Compute the least squares regression line, with the height of the father as the predictor variable and the height of the son as the response variable.
4. Draw the scatterplot and the regression line together.
5. Predict the height of a son whose father is 68 inches tall.
6. Compute the regression line if the heights are converted to centimeters (there are 2.54 centimeters per inch).

Answer

1. Continuous, ratio
2.  $r = 0.501$ ,  $r^2 = 0.251$
3.  $y = 33.893 + 0.514x$
5. 68.85
6.  $y = 86.088 + 0.514x$

Consider the petal length, petal width, and species variables in Fisher's iris data.

1. Classify the variables by type and level of measurement.
2. Compute the correlation between petal length and petal width.
3. Compute the correlation between petal length and petal width by species.

Answer

1. Species: discrete, nominal; petal length and width: continuous ratio
2. 0.9559
3. Setosa: 0.3316, Verginica: 0.3496, Versicolor: 0.6162

Consider the number of candies and net weight variables in the M&M data.

1. Classify the variable by type and level of measurement.
2. Compute the correlation coefficient and the coefficient of determination.
3. Compute the least squares regression line with number of candies as the predictor variable and net weight as the response variable.
4. Draw the scatterplot and the regression line in part (b) together.
5. Predict the net weight of a bag of M&Ms with 56 candies.
6. Naively, one might expect a much stronger correlation between the number of candies and the net weight in a bag of M&Ms. What is another source of variability in net weight?

Answer

1. Number of candies: discrete, ratio; net weight: continuous, ratio
2.  $r = 0.793$ ,  $r^2 = 0.629$
3.  $y = 20.278 + 0.507x$
5. 48.657
6. Variability in the weight of individual candies.

Consider the response rate and total SAT score variables in the SAT by state data set.

1. Classify the variables by type and level of measurement.
2. Compute the correlation coefficient and the coefficient of determination.
3. Compute the least squares regression line with response rate as the predictor variable and SAT score as the response variable.
4. Draw the scatterplot and regression line together.
5. Give a possible explanation for the negative correlation.

Answer

1. Response rate: continuous, ratio. SAT score could probably be considered either discrete or continuous, but is only at the interval level of measurement, since the smallest possible scores is 400 (200 each on the verbal and math portions).
2.  $r = -0.849$ ,  $r^2 = 0.721$
3.  $y = 1141.5 - 2.1x$
5. States with low response rate may be states for which the SAT is optional. In that case, the students who take the test are the better, college-bound students. Conversely, states with high response rates may be states for which the SAT is mandatory. In that case, all students including the weaker, non-college-bound students take the test.

Consider the verbal and math SAT scores (for all students) in the SAT by year data set.

1. Classify the variables by type and level of measurement.
2. Compute the correlation coefficient and the coefficient of determination.
3. Compute the least squares regression line.
4. Draw the scatterplot and regression line together.

Answer

1. Continuous perhaps, but only at the interval level of measurement because the smallest possible score on each part is 200.
2.  $r = 0.614$ ,  $r^2 = 0.377$
3.  $y = 321.5 + 0.3x$

Consider the temperature and erosion variables in the first data set in the Challenger data.

1. Classify the variables by type and level of measurement.
2. Compute the correlation coefficient and the coefficient of determination.
3. Compute the least squares regression line.
4. Draw the scatter plot and the regression line together.
5. Predict the O-ring erosion with a temperature of 31° F.
6. Is the prediction in part (c) meaningful? Explain.
7. Find the regression line if temperature is converted to degrees Celsius. Recall that the conversion is  $\frac{5}{9}(x - 32)$ .

Answer

1. temperature: continuous, interval; erosion: continuous ratio
2.  $r = -0.555$ ,  $r^2 = 0.308$
3.  $y = 106.8 - 1.414x$
5. 62.9.
6. This estimate is problematic, because 31° is far outside of the range of the sample data.
7.  $y = 61.54 - 2.545x$

This page titled [6.7: Sample Correlation and Regression](#) is shared under a [CC BY 2.0](#) license and was authored, remixed, and/or curated by [Kyle Siegrist \(Random Services\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.