

Chapter 2. Vectors and Vector Spaces

Section 2.3. Centered Vectors and Variances and Covariances of Vectors

Note. In this section, we introduce some operations on vectors which relate to treating them as data sets. In particular, we define variance, covariance, and correlation.

Note. Recall that if x_1, x_2, \dots, x_n are data points from a population of size n , then

- the mean of the population is $\mu = \frac{\sum_{i=1}^n x_i}{n}$,
- the variance of the population is $\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$, and
- the standard deviation of the population is $\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$.

Note. If x_1, x_2, \dots, x_n are data points *sampled* from a population, then

- the sample mean is $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$,
- the sample variance is $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n - 1)}$, and
- the sample standard deviation is $s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n - 1)}}$.

Definition. For a given n -vector x , its *centered counterpart*, denoted x_c , is $x_c = x - \bar{x}$ where \bar{x} is the mean vector of x , $\bar{x} = [\bar{x}, \bar{x}, \dots, \bar{x}]$. (Notice $\|\bar{x}\|^2 = n\bar{x}^2$.) Any n -vector with entries that sum to 0 is a *centered vector*.

Note. In Exercise 2.14, it is shown that for any $x, y \in \mathbb{R}^n$, $(x + y)_c = x_c + y_c$.

Note. We can interpret vector \bar{x} as

$$\bar{x} = \text{proj}_{1_n}(x) = \frac{\langle 1_n, x \rangle 1_n}{\|1_n\|^2} = \left(\frac{1_n^T x}{n} \right) 1_n$$

(recall that $\|1_n\|^2 = (\sqrt{n})^2 = n$). As shown in the class notes for Section 2.2 (see the Note after the definition of “projection”), $\text{proj}_x(y) \perp (y - \text{proj}_x(y))$ and so $\bar{x} \perp x - \bar{x}$ or $\bar{x} \perp x_c$. So by the Pythagorean Theorem (see the same note), $\|x\|^2 = \|\bar{x}\|^2 + \|x_c\|^2$. Notice that in terms of vector entries (i.e., scalars) this implies

$$\sum_{i=1}^n x_i^2 = n\bar{x}^2 + \sum_{i=1}^n (x_i - \bar{x})^2 \text{ or } \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2.$$

With $\{x_1, x_2, \dots, x_n\}$ as a data set for a population, we have the variance $\sigma^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / n$ and mean $\mu = \bar{x}$ so that the familiar formula $\sigma^2 = \sum_{i=1}^n x_i^2 / n - \mu^2$ (or $n\sigma^2 = \sum_{i=1}^n x_i^2 - n\mu^2$) is equivalent to the vector equation $\|x_c\|^2 = \|x\|^2 - \|\bar{x}\|^2$.

Definition. For $x \in \mathbb{R}^n$, the *scaled vector*, denoted x_s , is $x_s = \sqrt{n-1}x / \|x_c\|$. The *centered and scaled vector* is $x_{cs} = \sqrt{n-1}x_c / \|x_c\|$. The process of creation of x_{cs} from x is called *standardizing*. The *standard deviation* of x is the scalar quantity $s_x = \|x_c\| / \sqrt{n-1}$ and the *variance* of x is $V(x) = s_x^2 = \|x_c\|^2 / (n-1)$.

Note. In terms of the entries of n -vector x we have

$$s_x = \frac{\|x_c\|}{\sqrt{n-1}} = \frac{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}{\sqrt{n-1}},$$

as would be expected for the *sample* standard deviation of a data set. Also, $V(x) = \sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)$, as would be expected for the *sample* variance.

Note. For $x \in \mathbb{R}^n$ we have the variance of the scaled vector x_s is

$$\begin{aligned} V(x_s) &= V\left(\sqrt{n-1} \frac{x}{\|x_c\|}\right) = \frac{\left\|\left(\sqrt{n-1} \frac{x}{\|x_c\|}\right)_c\right\|^2}{n-1} \\ &= \frac{\left\|\sqrt{n-1} \frac{x}{\|x_c\|} - \sqrt{n-1} \frac{\bar{x}}{\|x_c\|}\right\|^2}{n-1} \text{ since } \overline{a\bar{x}} = a\bar{x} \\ &= \frac{\|x - \bar{x}\|^2}{\|x_c\|^2} = \frac{\|x - \bar{x}\|^2}{\|x - \bar{x}\|^2} \text{ since } x_c = x - \bar{x} \\ &= 1. \end{aligned}$$

thus the term “scaled” vector x_s is justified; its sample variance and sample standard deviation are both 1.

Note. Recall for sample (x_i, y_i) for $i = 1, 2, \dots, n$ of discrete random variable pair (X, Y) , the sample covariance is

$$\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

The sample correlation is

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{s_x s_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}.$$

Definition. Let x and y be n -vectors. The (sample) *covariance* between x and y is

$$\text{Cov}(x, y) = \frac{\langle x - \bar{x}, y - \bar{y} \rangle}{n - 1} = \frac{\langle x_c, y_c \rangle}{n - 1}.$$

Note. With 1_n as the summing vector, we have $x - \bar{x} = x - \bar{x}1_n$ and $y - \bar{y} = y - \bar{y}1_n$ (with the careful interpretation of \bar{x} and \bar{y} as scalars *or* vectors), so

$$\begin{aligned} \text{Cov}(x, y) &= \frac{\langle x - \bar{x}, y - \bar{y} \rangle}{n - 1} = \frac{\langle x - \bar{x}1_n, y - \bar{y}1_n \rangle}{n - 1} \\ &= \frac{\langle x, y \rangle - \langle \bar{x}1_n, y \rangle - \langle x, \bar{y}1_n \rangle + \langle \bar{x}1_n, \bar{y}1_n \rangle}{n - 1} \\ &= \frac{\langle x, y \rangle - \sum_{i=1}^n \bar{x}y_i - \sum_{i=1}^n x_i\bar{y} + n\bar{x}\bar{y}}{n - 1} \\ &= \frac{\langle x, y \rangle - n\bar{x} \sum_{i=1}^n y_i/n - n\bar{y} \sum_{i=1}^n x_i/n + n\bar{x}\bar{y}}{n - 1} \\ &= \frac{\langle x, y \rangle - n\bar{x}\bar{y} - n\bar{x}\bar{y} + n\bar{x}\bar{y}}{n - 1} \\ &= \frac{\langle x, y \rangle - n\bar{x}\bar{y}}{n - 1}. \end{aligned}$$

Also,

$$\begin{aligned} \text{Cov}(x, x) &= \frac{\langle x - \bar{x}, x - \bar{x} \rangle}{n - 1} = \frac{\langle x, x \rangle - 2\langle x, \bar{x} \rangle + \langle \bar{x}, \bar{x} \rangle}{n - 1} \\ &= \frac{\sum_{i=1}^n x_i^2 - 2\sum_{i=1}^n x_i\bar{x} + \sum_{i=1}^n \bar{x}^2}{n - 1} = \frac{\sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2)}{n - 1} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = V(x). \end{aligned}$$

Theorem 2.3.1. Properties of Covariance.

Let x, y, z be n -vectors and let $a \in \mathbb{R}$. Then:

1. $\text{Cov}(a1_n, y) = 0$,
2. $\text{Cov}(ax, y) = a\text{Cov}(x, y)$,
3. $\text{Cov}(y, y) = V(y)$,
4. $\text{Cov}(x + z, y) = \text{Cov}(x, y) + \text{Cov}(z, y)$, in particular $\text{Cov}(x + y, y) = \text{Cov}(x, y) + V(y)$ and $\text{Cov}(x + a1_n, y) = \text{Cov}(x, y)$.

Definition. Let x and y be n -vectors. The *correlation* between x and y is

$$\text{Corr}(x, y) = \text{Cov}(x_{cs}, y_{cs}) = \left\langle \frac{x_c}{\|x_c\|}, \frac{y_c}{\|y_c\|} \right\rangle = \frac{\langle x_c, y_c \rangle}{\|x_c\| \|y_c\|}.$$

Note. Gentle describes correlation as “a measure of the extent to which the vectors point in the same direction.” This is justified with the observation that $\text{Corr}(x, y) = \frac{\langle x_c, y_c \rangle}{\|x_c\| \|y_c\|} = \cos \theta$ where θ is the angle between x_c and y_c (though not the angle between x and y ; see Exercise 2.15). It then follows that (as we expect) $\text{Corr}(x, y) \in [-1, 1]$.

Note. We can also express $\text{Corr}(x, y)$ as

$$\begin{aligned} \text{Corr}(x, y) &= \frac{\langle x_c, y_c \rangle}{\|x_c\| \|y_c\|} = \frac{\langle x_c, y_c \rangle}{(n-1) \sqrt{\frac{\|x_c\|^2}{n-1} \frac{\|y_c\|^2}{n-1}}} \\ &= \frac{\langle x_c, y_c \rangle}{(n-1) \sqrt{V(x)V(y)}} = \frac{\langle x - \bar{x}, y - \bar{y} \rangle}{(n-1) \sqrt{V(x)V(y)}} = \frac{\text{Cov}(x, y)}{\sqrt{V(x)V(y)}}. \end{aligned}$$

Note. As might be expected, scaling x by a scalar $a \in \mathbb{R}$ affects the correlation only in sign:

$$\begin{aligned}\text{Corr}(ax, y) &= \frac{\langle (ax)_c, y_c \rangle}{\|(ax)_c\| \|y_c\|} = \frac{\langle ax - a\bar{x}, y_c \rangle}{\|ax - a\bar{x}\| \|y_c\|} \\ &= \frac{a \langle x - \bar{x}, y_c \rangle}{|a| \|x - \bar{x}\| \|y_c\|} = \frac{a}{|a|} \frac{\langle x_c, y_c \rangle}{\|x_c\| \|y_c\|} = \text{sign}(a) \text{Corr}(x, y).\end{aligned}$$

Revised: 5/27/2020