# Covariance

Covariance provides a measure of the strength of the correlation between two or more sets of random variates. The covariance for two random variates $X$ and $Y$, each with sample size $N$, is defined by the expectation value

$$\text{cov}(X,Y) = \langle (X - \mu_X)(Y - \mu_Y) \rangle \tag{1}$$
$$= \langle X\,Y \rangle - \mu_X\,\mu_y \tag{2}$$

where $\mu_x = \langle X \rangle$ and $\mu_y = \langle Y \rangle$ are the respective means, which can be written out explicitly as

$$\text{cov}(X,Y) = \sum_{i=1}^{N} \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}. \tag{3}$$

For uncorrelated variates,

$$\text{cov}(X,Y) = \langle X\,Y \rangle - \mu_X\,\mu_Y = \langle X \rangle \langle Y \rangle - \mu_X\,\mu_Y = 0, \tag{4}$$

so the covariance is zero. However, if the variables are correlated in some way, then their covariance will be nonzero. In fact, if $\text{cov}(X,Y) > 0$, then $Y$ tends to increase as $X$ increases, and if $\text{cov}(X,Y) < 0$, then $Y$ tends to decrease as $X$ increases. Note that while statistically independent variables are always uncorrelated, the converse is not necessarily true.

In the special case of $Y = X$,

$$\text{cov}(X,X) = \langle X^2 \rangle - \langle X \rangle^2 \tag{5}$$
$$= \sigma_X^2, \tag{6}$$

so the covariance reduces to the usual variance $\sigma_X^2 = \text{var}(X)$. This motivates the use of the symbol $\sigma_{XY} = \text{cov}(X,Y)$, which then provides a consistent way of denoting the variance as $\sigma_{XX} = \sigma_X^2$, where $\sigma_X$ is the standard deviation.

The derived quantity

$$\text{cor}(X,Y) = \frac{\text{cov}(X,Y)}{\sigma_X\,\sigma_Y} \tag{7}$$
$$= \frac{\sigma_{XY}}{\sqrt{\sigma_{XX}\,\sigma_{YY}}}, \tag{8}$$

is called **statistical correlation (Pearson)** of $X$ and $Y$.

The covariance is especially useful when looking at the variance of the sum of two random variates, since

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\,\text{cov}(X,Y). \tag{9}$$

The covariance is symmetric by definition since

$$\operatorname{cov}(X, Y) = \operatorname{cov}(Y, X). \tag{10}$$

Given $n$ random variates denoted $X_1$, ..., $X_n$, the covariance $\sigma_{ij} \equiv \operatorname{cov}(X_i, X_j)$ of $X_i$ and $X_j$ is defined by

$$\operatorname{cov}(X_i, X_j) = \langle (X_i - \mu_i)(X_j - \mu_j) \rangle \tag{11}$$
$$= \langle X_i X_j \rangle - \mu_i \mu_j, \tag{12}$$

where $\mu_i = \langle X_i \rangle$ and $\mu_j = \langle X_j \rangle$ are the means of $X_i$ and $X_j$, respectively. The matrix $(V_{ij})$ of the quantities $V_{ij} = \operatorname{cov}(X_i, X_j)$ is called the covariance matrix.

The covariance obeys the identities

$$\operatorname{cov}(X + Z, Y) = \langle (X + Z) Y \rangle - \langle X + Z \rangle \langle Y \rangle \tag{13}$$
$$= \langle X Y \rangle + \langle Z Y \rangle - (\langle X \rangle + \langle Z \rangle) \langle Y \rangle \tag{14}$$
$$= \langle X Y \rangle - \langle X \rangle \langle Y \rangle + \langle Z Y \rangle - \langle Z \rangle \langle Y \rangle \tag{15}$$
$$= \operatorname{cov}(X, Y) + \operatorname{cov}(Z, Y). \tag{16}$$

By induction, it therefore follows that

$$\operatorname{cov}\left( \sum_{i=1}^{n} X_i, Y \right) = \sum_{i=1}^{n} \operatorname{cov}(X_i, Y) \tag{17}$$

$$\operatorname{cov}\left( \sum_{i=1}^{n} X_i, \sum_{j=1}^{m} Y_j \right) = \sum_{i=1}^{n} \operatorname{cov}\left( X_i, \sum_{j=1}^{m} Y_j \right) \tag{18}$$

$$= \sum_{i=1}^{n} \operatorname{cov}\left( \sum_{j=1}^{m} Y_j, X_i \right) \tag{19}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{m} \operatorname{cov}(Y_j, X_i) \tag{20}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{m} \operatorname{cov}(X_i, Y_j). \tag{21}$$