

# Lecture 23

**Definition:** The *covariance* between the rvs  $X$  and  $Y$  is given by

$$\begin{aligned}\text{Cov}(X, Y) &= \text{E}( (X - \text{E}(X))(Y - \text{E}(Y)) ) \\ &= \text{E}(XY) - \text{E}(X)\text{E}(Y)\end{aligned}$$

**Interpretation:**

- **positive covariance**
  - large  $x$ 's occur with large  $y$ 's
  - small  $x$ 's occur with small  $y$ 's
- **negative covariance**
  - large  $x$ 's occur with small  $y$ 's
  - small  $x$ 's occur with large  $y$ 's

Correlation is the scaled and preferred version of covariance.

**Definition:** The *correlation* between the rvs  $X$  and  $Y$  is given by

$$\rho = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{V(X)}\sqrt{V(Y)}}$$

**Discussion points:**

- $-1 \leq \text{Corr}(X, Y) \leq 1$
- correlation is location/scale invariant
- $\rho$  is the population analogue of  $r$
- $\rho$  typically relevant to continuous rvs
- if  $a > 0$ , then  $\text{Corr}(X, aX + b) = 1$
- if  $a < 0$ , then  $\text{Corr}(X, aX + b) = -1$

**Example:** Obtain the correlation between  $X$  and  $Y$  where the joint pmf of  $X$  and  $Y$  is given in the following table.

	<b>X=1</b>	<b>X=2</b>	<b>X=3</b>
<b>Y=1</b>	0.1	0.2	0.3
<b>Y=2</b>	0.0	0.2	0.2

**Proposition:** If  $X$  and  $Y$  are independent, then

$$\text{Cov}(X, Y) = 0$$

**In addition,**  $\text{Corr}(X, Y) = 0$  provided  $V(X)$  and  $V(Y)$  are nonzero. **The converse is not true.**

**Also, recall that correlation does not imply causation.**

**Proposition:**  $V(X + Y) = V(X) + V(Y) + 2\text{Cov}(X, Y)$

**Proposition:** More generally,

$$V(aX + bY + c) = a^2V(X) + b^2V(Y) + 2ab\text{Cov}(X, Y)$$

**Proposition:** Even more generally,

$$V\left(\sum_{i=1}^n a_i X_i + c\right) = \sum_{i=1}^n a_i^2 V(X_i) + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$$
$$E\left(\sum_{i=1}^n a_i X_i + c\right) = c + \sum_{i=1}^n a_i E(X_i)$$

Lets put some of this stuff together to provide a useful result.

**Corollary:** Suppose that the rv's  $X_1, \dots, X_n$  are a *sample*. In other words, the  $X$ 's are independent and arise from a common distribution with mean  $\mu$  and variance  $\sigma^2$ . Then the sample mean has the following properties:

- $E(\bar{X}) = \mu$
- $V(\bar{X}) = \sigma^2/n$

Suprisingly, we have reached this point in our Statistics course and we have not yet defined the word *statistic*.

**Definition:** A *statistic* is a function of the data.

Some examples:

- $\bar{X} = \sum_{i=1}^n X_i/n$  is a statistic
- $S^2 = \sum_{i=1}^n (X_i - \bar{X})^2/(n - 1)$  is a statistic

Since data are variable, statistics are also variable. Sometimes we are interested in the distributions of statistics.

**Example:** Obtain the distribution of the statistic  $Q = X+Y$  where the joint pmf of  $X$  and  $Y$  is given in the following table.

	<b>X=1</b>	<b>X=2</b>	<b>X=3</b>
<b>Y=1</b>	0.1	0.1	0.2
<b>Y=2</b>	0.2	0.3	0.1



The previous example was simple. To generalize, we need to go a little crazy with notation.

Suppose that  $X_1, \dots, X_n$  are discrete with joint pmf  $p(x_1, \dots, x_n)$ . Then the pmf for the general statistic  $Q(X_1, \dots, X_n)$  is

$$p_Q(q) = \sum_A p(x_1, \dots, x_n)$$

where the sum is a multiple sum and  $A$  is the set of  $x_1, \dots, x_n$  such that  $Q(x_1, \dots, x_n) = q$ .

Suppose that  $X_1, \dots, X_n$  are continuous with joint pdf  $f(x_1, \dots, x_n)$ . Then the cdf for the general statistic  $Q(X_1, \dots, X_n)$  is

$$F_Q(q) = P(Q \leq q) = \int_A f(x_1, \dots, x_n) dx_1 \dots dx_n$$

where the integral is a multiple integral and  $A$  is the set of  $x_1, \dots, x_n$  such that  $Q(x_1, \dots, x_n) \leq q$ .

I have mentioned previously that statistical practice relies heavily on computation. Here is a simulation procedure that can be used to approximate distributions of statistics when the sums and integrals from the previous page are too difficult to obtain analytically.

- Repeat the following two steps  $M$  times where  $M$  is large and let  $i$  denote the  $i$ -th iteration
  - generate  $x_1, \dots, x_n$  according to  $p(x_1, \dots, x_n)$  or  $f(x_1, \dots, x_n)$  (depending whether the data are discrete or continuous)
  - calculate  $Q_i = Q(x_1, \dots, x_n)$  for the data
- approximate the distbn of  $Q$  with a histogram based on generated outcomes  $Q_1, \dots, Q_M$