

STAT 23400 Lecture 9: Expected Values, Covariance and Correlation

Review: Joint Continuous distributions

- **Joint pdf** for two continuous r.v.'s:

$$P[(X, Y) \in A] = \iint_A f(x, y) dx dy .$$

- **Marginal pdf:**

$$f_X(x) = \int_{\mathbb{R}} f(x, y) dy \quad \text{for all } x \in \mathbb{R}$$

$$f_Y(y) = \int_{\mathbb{R}} f(x, y) dx \quad \text{for all } y \in \mathbb{R}$$

- **Conditional probability distribution:**

$$f_{X|Y}(x | y) = \frac{f(x, y)}{f_Y(y)}, \quad f_{Y|X}(y | x) = \frac{f(x, y)}{f_X(x)}$$

- X and Y are **independent** if and only if

$$p(x, y) = p_X(x)p_Y(y) \quad \text{if } X \text{ and } Y \text{ are discrete,}$$

$$f(x, y) = f_X(x)f_Y(y) \quad \text{if } X \text{ and } Y \text{ are continuous.}$$

Outline for Today

- Expected Values, Covariance and Correlation (Section 5.2 in MMSA).
- Lab05: Central limit theorem (link posted on canvas).
- HW 5 is posted on canvas.

Expected Values

Expected Values of Functions of X, Y

For two r.v.'s X, Y with joint pmf/pdf $f(x, y)$, the expected value of a function $g(X, Y)$ of X and Y is defined as

$$E[g(X, Y)] = \begin{cases} \sum_{x,y} g(x, y)f(x, y) & \text{in discrete case,} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f(x, y) dx dy & \text{in continuous case.} \end{cases}$$

Halloween Game Example Revisited

Recall that X is the smaller number of the two dice, and Y is the bigger number. The joint pmf is

$f(x, y)$	1	2	$\frac{y}{3}$	4	5	6
1	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$
2	0	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$
x 3	0	0	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{2}{36}$	$\frac{2}{36}$
4	0	0	0	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{2}{36}$
5	0	0	0	0	$\frac{1}{36}$	$\frac{2}{36}$
6	0	0	0	0	0	$\frac{1}{36}$

Let us find $E(X)$ and $E(Y)$.

Recall the marginal distributions are

x	1	2	3	4	5	6
$p_X(x)$	11/36	9/36	7/36	5/36	3/36	1/36
y	1	2	3	4	5	6
$p_Y(y)$	1/36	3/36	5/36	7/36	9/36	11/36

What is $E(XY)$?

Thm. If $g(\mathbf{X}, \mathbf{Y}) = a\mathbf{X} + b\mathbf{Y}$ for any $a \in \mathbb{R}$ and any $b \in \mathbb{R}$, then $g(X, Y)$ is a **linear combination of X and Y** and

$$\mathbf{E}[g(\mathbf{X}, \mathbf{Y})] = a\mathbf{E}(\mathbf{X}) + b\mathbf{E}(\mathbf{Y})$$

Thm. If X and Y are **independent**, then for any functions g and h , we have

$$E[g(X)h(Y)] = E[g(X)]E[h(Y)].$$

Covariance

For any dependent random variables X and Y ,

How can we measure the relationship between X and Y ?

How “strong” is this relationship?

- Positive Association, e.g., SAT score v.s. hours spent on studying.
- Negative Association: e.g., SAT score v.s. sleeping hours per night.

Covariance of two random variables

Covariance of two random variables X, Y with means μ_X, μ_Y .

Take $g(X, Y) = (X - \mu_X) \cdot (Y - \mu_Y)$. The covariance of X and Y is defined as

$$\text{Cov}(X, Y) = \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)].$$

$\text{Cov}(X, Y)$ is a measure of how the variables 'covary'.

$\text{Cov}(X, Y) > 0 \rightarrow$ when X increases Y tends to increase.

$\text{Cov}(X, Y) < 0 \rightarrow$ when X increases Y tends to decrease.

A helpful formula: $\text{Cov}(X, Y) = E(XY) - \mu_X \mu_Y.$

Properties of Covariance

The covariance has the following properties.

For any rv's X , Y , Z and constant numbers a , b , we have

- $\text{Cov}(X, X) = \text{Var}(X)$
- Symmetry: $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- Homogeneity: $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$
- The covariance of a r.v. and a constant number is 0, i.e.,
 $\text{Cov}(a, X) = 0$
- Right-linearity: $\text{Cov}(X + Y, Z) = \text{Cov}(X, Z) + \text{Cov}(Y, Z)$
- Left-linearity: $\text{Cov}(X, Y + Z) = \text{Cov}(X, Y) + \text{Cov}(X, Z)$

Proofs for Properties of Covariance

The proofs of the above are all straightforward from definition. Here we prove the Right-linearity as an example.

$$\begin{aligned}\text{Cov}(X + Y, Z) &= E((X + Y)Z) - E(X + Y)E(Z) \\ &= E(XZ) + E(YZ) - [E(X) + E(Y)]E(Z) \\ &= \underbrace{E(XZ) - E(X)E(Z)}_{\text{Cov}(X,Z)} + \underbrace{E(YZ) - E(Y)E(Z)}_{\text{Cov}(Y,Z)} \\ &= \text{Cov}(X, Z) + \text{Cov}(Y, Z)\end{aligned}$$

Note in the proof above, we used the property of expected value that

$$\begin{aligned}E(X + Y) &= E(X) + E(Y) \\ E(XZ + YZ) &= E(XZ) + E(YZ)\end{aligned}$$

Example 1: Application of the Properties of Covariance

Let X_1, X_2, X_3 be uncorrelated r.v.'s each with variances $\sigma_1^2, \sigma_2^2, \sigma_3^2$ respectively. Find $\text{Cov}(X_1 - X_2, X_1 - X_3)$.

Example 2: Halloween Game Example Revisited

In the halloween game example, find the covariance of X and Y .

Correlation

Problems with Covariance

But, the “relationship” between X and Y should be the same no matter what units we choose to measure in.

Solution: Use **correlation coefficient** to measure linear relationship.

$$\rho_{XY} = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X\sigma_Y}.$$

- Correlation is unitless
- We always have $-1 \leq \rho_{XY} \leq 1$
 - $\rho_{XY} = +1 \implies$ perfect, positive, linear relationship
 - $\rho_{XY} = 0 \implies$ no *linear* relationship (other relationship?)
 - $\rho_{XY} = -1 \implies$ perfect, negative, linear relationship

Halloween Game Example Revisited

Recall that, previously, in the Halloween Game Example, we have calculated that

$$\text{Cov}(X, Y) = \frac{1225}{1296} \approx 0.95.$$

Can we say the relationship between X and Y are almost perfectly linear?

Halloween Game Example Revisited

Recall that, previously, in the Halloween Game Example, we have calculated that

$$\text{Cov}(X, Y) = \frac{1225}{1296} \approx 0.95.$$

Can we say the relationship between X and Y are almost perfectly linear?

No, because covariance is not scale free. Need to consider correlation.

We can further obtain $\sigma_X = 2.88$, $\sigma_Y = 1.40$ and

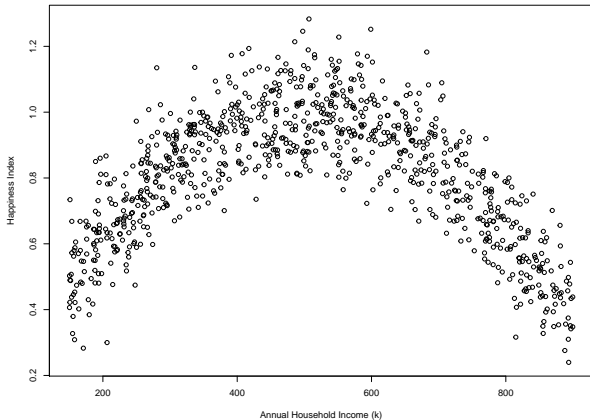
$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \approx 0.24,$$

indicating a less strong linear relationship.

What is the correlation if X and Y are independent?

Example: Zero Correlation Does NOT Imply Independence

Consider the relationship between the Annual Household Income (in thousand US dollars) and the Family Happiness Index.



Linear Combinations of Random Variables

- We can define **joint probability mass function** for more than two r.v.'s, $\mathbf{X} = (X_1, X_2, \dots, X_n)$, as

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \\ &= P(\{X_1 = x_1\} \cap \{X_2 = x_2\} \cap \dots \cap \{X_n = x_n\}). \end{aligned}$$

- Define the **joint probability density function** for continuous r.v.'s $\mathbf{X} = (X_1, \dots, X_n)$ such that

$$P[\mathbf{X} \in A] = \int \dots \int_A f(x_1, \dots, x_n) dx_1 dx_2 \dots dx_n.$$

Multivariate Random Variables: Independence

- More generally, a sequence of random variables X_1, X_2, \dots, X_n are **(mutually) independent** if and only if

$$P(X_1 \in A_1, \dots, X_n \in A_n) = P(X_1 \in A_1) \cdots P(X_n \in A_n).$$

for all sequence of events A_1, A_2, \dots

- Equivalently, the random variables X_1, X_2, \dots, X_n are **(mutually) independent** if and only if their joint distributions factors into the product of their marginal distributions.

$$f(x_1, x_2, \dots, x_n) = f_1(x_1)f_2(x_2) \cdots f_n(x_n)$$

for all x_1, x_2, \dots, x_n .

Linear Combinations of Random Variables

Definition. For any random variables X_1, X_2, \dots, X_n , a **linear combination** of X_1, X_2, \dots, X_n is

$$a_1X_1 + a_2X_2 + \cdots + a_nX_n,$$

where a_1, a_2, \dots, a_n are constant numbers. For example,

- The sum

$$X_1 + X_2 + \cdots + X_n$$

is a linear combination of X_1, X_2, \dots, X_n with all $a_i = 1$.

- The average

$$\frac{X_1 + X_2 + \cdots + X_n}{n}$$

is a linear combination of X_1, X_2, \dots, X_n with all $a_i = 1/n$.

- The difference $X - Y$ is a linear combination of X and Y with $a_1 = 1, a_2 = -1$.

Example: Linear Combinations

Suppose the bus fare is

\$2 for senior citizens, \$1 for children, and \$3 for all other people

Let

X = the number of senior citizens on the bus,

Y = the number of children on the bus,

Z = the number of all other passengers on the bus

The total amount of bus fares collected is then

$$2X + Y + 3Z$$

which is a linear combination of X , Y , Z .

Expected Values for Linear Combinations of RV's

$$E(a_1X_1 + a_2X_2 + \cdots + a_nX_n) = a_1E(X_1) + a_2E(X_2) + \cdots + a_nE(X_n)$$

Variance of a Linear Combination of RV's

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

Example: Expected Value and Variance of $\text{Binom}(n, p)$

Recall that, if $X \sim \text{Binom}(n, p)$, we have

$$E(X) = np, \quad \text{Var}(X) = np(1 - p).$$

We are now ready to prove these formulas using linear combinations.

Recall a Binomial random variable $X \sim \text{Binom}(n, p)$ is the total number of successes obtained in n independent Bernoulli trials. The expected value and variance of X are thus

$$\begin{aligned} E(X) &= \underbrace{E(X_1)}_{=p} + \underbrace{E(X_2)}_{=p} + \cdots + \underbrace{E(X_n)}_{=p} = np \\ \text{Var}(X) &= \underbrace{\text{Var}(X_1)}_{=p(1-p)} + \underbrace{\text{Var}(X_2)}_{=p(1-p)} + \cdots + \underbrace{\text{Var}(X_n)}_{=p(1-p)} = np(1-p) \end{aligned}$$