

## 2. Covariance and Correlation (Rice Section 4.3, Blitzstein 7.3)

- a) Define the covariance  $\sigma_{xy}$  and the correlation  $\rho_{xy}$  for random variables  $X$  and  $Y$  and derive the computation formula  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$ . Compare to formulas for  $\text{Var}(X)$ .

### DEFINITION

If  $X$  and  $Y$  are jointly distributed random variables with expectations  $\mu_X$  and  $\mu_Y$ , respectively, the covariance of  $X$  and  $Y$  is

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

provided that the expectation exists.

### DEFINITION

If  $X$  and  $Y$  are jointly distributed random variables and the variances and covariances of both  $X$  and  $Y$  exist and the variances are nonzero, then the correlation of  $X$  and  $Y$ , denoted by  $\rho$ , is

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

using linearity of Expectation

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

$$\begin{aligned} \text{Cov}(X, X) &= E[(X - \mu_X)^2] = E[X^2] - (E[X])^2 \leftarrow \begin{array}{l} \text{can be} \\ \text{generalized.} \end{array} \\ &= E[XY] - E[X]E[Y] \\ &= \text{Var}(X) \quad \text{or } E[XY] - \mu_X \mu_Y. \end{aligned}$$

- b) As an example, consider a fat coin that has some probability of landing on its edge. Suppose the probability of it landing heads is  $p$  and the probability of it landing tails is also  $p$  ( $0 < p \leq 0.5$ ). If you flip the coin  $n = 2$  times and let  $X$  be the number of heads and  $Y$  the number of tails, find the covariance and correlation of  $X$  and  $Y$  as a function of  $p$ , and give the values when  $p = 0.4$ .

  $P(0 < p \leq 0.5)$  flip this coin  $n=2$  times,  $X: \# \text{heads}$   
 $Y: \# \text{tails}$

		X			marginality:
		0	1	2	
Y	2	$p^2$	$\emptyset$	$\emptyset$	$p^2$
	1	$2p(1-p)$	$2p^2$	$\emptyset$	$2p(1-p)$
0	$(1-p)^2$	$2p(1-p)$	$p^2$	$(1-p)^2$	

	X	Y
HH	$p^2$	H T
HT	$p^2$	2 0
HE	$p(1-p)$	1 1
TT	$p^2$	0 2
TH	$p^2$	1 1
TE	$p(1-p)$	0 1
EE	$(1-p)^2$	0 0
ET	$p(1-p)$	0 1
EH	$p(1-p)$	1 0

$X$  and  $Y$  are both binomial marginally!

- $X \sim \text{Binom}(2, p)$
- $Y \sim \text{Binom}(2, p)$

$$\rightarrow E[X] = np = 2p = E[Y]$$

$$\text{Var}[X] = npq = 2p(1-p) = \text{Var}[Y]$$

$$p=0.4$$

$$\therefore \text{Cov}(X, Y) = E[XY] - E[X]E[Y] \\ = 2p^2 - 4p^2 = -2p^2 = -0.32.$$

bc  $E[XY] \neq 0$  only when  $X=1$  and  $Y=1$ .

$$\therefore \rho(X, Y) = \frac{-2p^2}{2p(1-p)} = -\frac{p}{(1-p)} = -\frac{0.4}{0.6} = -\frac{2}{3}.$$

- c) Derive the formulas for  $\text{Var}(X + Y)$  and  $\text{Var}(X - Y)$ . Use the fact that  $\text{Var}(X + Y)$  and  $\text{Var}(X - Y)$  are both non-negative to prove that  $-1 \leq \rho_{xy} \leq 1$  (proving either the upper or lower bound will be sufficient for the presentation).

$$\begin{aligned}\text{• } \text{Var}(X+Y) &= E[((X+Y) - \mu_{X+Y})^2] = E[(X-\mu_X) + (Y-\mu_Y)]^2 \\ &= \underbrace{E[(X-\mu_X)^2]}_{\text{Var}(X)} + \underbrace{E[(Y-\mu_Y)^2]}_{\text{Var}(Y)} + \underbrace{E[2(X-\mu_X)(Y-\mu_Y)]}_{2\text{Cov}(X,Y)} \\ &= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X,Y)\end{aligned}$$

$$\begin{aligned}\text{• } \text{Var}(X-Y) &= E[((X-Y) - \mu_{X-Y})^2] = E[(X-\mu_X) - (Y-\mu_Y)]^2 \\ &= \underbrace{E[(X-\mu_X)^2]}_{\text{Var}(X)} + \underbrace{E[(Y-\mu_Y)^2]}_{\text{Var}(Y)} - 2E[(X-\mu_X)(Y-\mu_Y)] \\ &= \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X,Y)\end{aligned}$$

fact:  $\text{Cov}(a+bX, c+dY) = bd \text{Cov}(X,Y)$

Show  $-1 \leq \rho \leq 1$ . Suppose  $X, Y$  have mean 0.

$$\begin{aligned}\textcircled{1} \quad \text{Var}\left(\frac{X}{\sigma_x} + \frac{Y}{\sigma_y}\right) &= \frac{\text{Var}(X)}{\sigma_x^2} + \frac{\text{Var}(Y)}{\sigma_y^2} + \frac{2\text{Cov}(X,Y)}{\sigma_x \sigma_y} \\ &= 1 + 1 + 2\rho \geq 0 \rightarrow \rho \geq -1\end{aligned}$$

$$\begin{aligned}\textcircled{2} \quad \text{Var}\left(\frac{X}{\sigma_x} - \frac{Y}{\sigma_y}\right) &= \frac{\text{Var}(X)}{\sigma_x^2} + \frac{\text{Var}(Y)}{\sigma_y^2} - \frac{2\text{Cov}(X,Y)}{\sigma_x \sigma_y} \\ &= 1 + 1 - 2\rho \geq 0 \rightarrow \rho \leq 1\end{aligned}$$

$$\therefore -1 \leq \rho \leq 1.$$

- d) Show that independent implies uncorrelated, and explain why uncorrelated does not always imply independent. As an example, consider  $X$  and  $Y$  uniformly distributed over the triangular region  $0 < |Y| < X < 1$ .

Independent  $\Rightarrow$  uncorrelated

$(\Rightarrow)$  pf) \* uncorrelated iff  $E[XY] = E[X]E[Y]$

\* independent iff  $f_{XY}(x,y) = f_X(x)f_Y(y)$

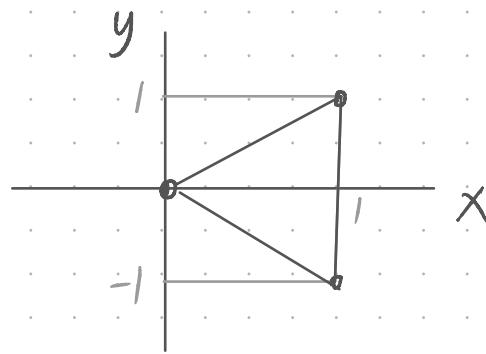
• MV LOTVS:  $E(XY) = \iint xy f_{XY}(x,y) dx dy$

$$\text{If indep.} \rightarrow = \iint xy \cdot f_X(x)f_Y(y) dx dy$$

$$= \int x f_X(x) dx \cdot \int y f_Y(y) dy$$

$$= E(X)E(Y). \Rightarrow \text{thus, uncorrelated}$$

( $\Leftarrow$ ) pf)



marginal pdf for  $X$ :  $f_x(x) = 2x$ ,  $0 < x < 1$

$$\therefore E(X) = \int x f(x) dx = \frac{2}{3}$$

marginal pdf for  $Y$  is  $f_y(y) = 1 - |y|$ ,  $-1 \leq y \leq 1$ .

$$\therefore E(Y) = 0$$

$$\begin{aligned} \text{Cov}(X, Y) &= \underline{E(XY)} - \underline{E(X)E(Y)} \\ &\quad \downarrow \quad \frac{2}{3} \cdot 0 = 0 \quad \underbrace{f_{xy}(x,y) = I(0 < |y| < x < 1)}_{\substack{\text{inner integral} \\ \text{outer integral}}} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{xy}(x,y) dx dy = \int_{-1}^{1} \int_{|y|}^{1} xy dx dy = 0. \end{aligned}$$

outer integral:

$$\frac{1}{2} \int_{-1}^1 y dy - \frac{1}{2} \int_{-1}^1 y^3 dy$$

$$= \left[ \frac{1}{2} y^2 \right]_{-1}^1 - \left[ \frac{1}{4} y^4 \right]_{-1}^1$$

$$= 0 - 0 = 0$$

$$\begin{aligned} \int_{|y|}^1 xy dx &= \left[ \frac{x^2}{2} y \right]_{|y|}^1 \\ &= \frac{1}{2} y - \frac{|y|^2}{2} y \\ &= \frac{1}{2} y - \frac{y^3}{2} \end{aligned}$$

Thus, uncorrelated.

However,  $Y|X=x \sim \text{Unif}(-x, x)$  depends on  $x$  (not indep.)

$\therefore$  Uncorrelated does not always imply independent.