

LECTURE 34 : CORRELATION AND COVARIANCE

• TOPICS TO COVER (BASED ON CH 4.8)

→ CORRELATION AND COVARIANCE

→ CORRELATION AND COVARIANCE

LET X AND Y BE TWO RANDOM VARIABLES WITH FINITE SECOND MOMENTS, I.E.

$$E X^2 < \infty \text{ AND } E Y^2 < \infty.$$

DEFINE THREE MEASURES OF A LINEAR RELATIONSHIP BETWEEN X AND Y :

CORRELATION BETWEEN X AND Y = $E(XY)$ JOINT EXPECTATION OF X AND Y

COVARIANCE BETWEEN X AND Y : $\text{Cov}(X, Y) = E(X - E(X))(Y - E(Y))$

CORRELATION COEFFICIENT BETWEEN X AND Y = $\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$

OBSERVE THE FOLLOWING :

$$\rightarrow E(XY) = \begin{cases} \sum_i \sum_j x_i y_j p_{X,Y}(x_i, y_j), & \text{DISCRETE CASE,} \\ \int_{\mathbb{R}^2} uv f_{X,Y}(u, v) dv du, & \text{CONT. CASE.} \end{cases}$$

$$\rightarrow \text{Cov}(X, X) = E(X - E(X))(X - E(X)) = E(X - E(X))^2 = \text{Var}(X)$$

$$\begin{aligned}
 \rightarrow \quad \text{COV}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\
 &\stackrel{\text{LINEARITY OF EXPECTATION}}{=} E[XY] - E[XE(Y)] - E[E(X)Y] + E[E(X)E(Y)] \\
 &\stackrel{E(X) \text{ AND } E(Y) \text{ ARE CONSTANTS}}{=} E[XY] - E(X)E(Y) - E(X)E(Y) + E(X)E(Y) \\
 &= E[XY] - E(X)E(Y)
 \end{aligned}$$

$$\rightarrow \quad \text{IF } E(X) = 0 \text{ OR } E(Y) = 0 : \text{COV}(X, Y) = E(XY)$$

DEFINITION : IF $\text{COV}(X, Y) = 0$: X AND Y ARE CALLED UNCORRELATED.

IF $\text{COV}(X, Y) > 0$: X AND Y ARE CALLED +VELY UNCORRELATED.

IF $\text{COV}(X, Y) < 0$: X AND Y ARE CALLED -VELY UNCORRELATED.

IF $\sigma_X^2 = \text{Var}(X) > 0$ AND $\sigma_Y^2 = \text{Var}(Y) > 0$, I.E., $\rho_{X,Y}$ IS WELL-DEFINED,

$$\begin{aligned}
 \text{COV}(X, Y) & > 0 & \iff & \rho_{X,Y} > 0 \\
 & < 0 & & < 0
 \end{aligned}$$

INDEPENDENCE AND UNCORRELATEDNESS

THEOREM: X AND Y ARE INDEPENDENT \Rightarrow X AND Y ARE UNCORRELATED.

HOWEVER THE CONVERSE IS NOT TRUE.

PROOF: LET X AND Y ARE INDEPENDENT, I.E.

$$f_{X,Y}(u,v) = f_X(u) \cdot f_Y(v) \quad \forall (u,v) \in \mathbb{R}^2 \quad (1)$$

$$\Rightarrow \int_{\mathbb{R}^2} uv f_{X,Y}(u,v) dv du = \int_{\mathbb{R}^2} uv f_X(u) \cdot f_Y(v) dv du \quad (2)$$

$$\Rightarrow E(XY) = \int_{\mathbb{R}} u f_X(u) du \int_{\mathbb{R}} v f_Y(v) dv$$

$$\Rightarrow E(XY) = E(X) E(Y) \Rightarrow X \text{ AND } Y \text{ ARE UNCORRELATED}$$

THE CONVERSE DOESN'T HOLD TRUE AS (2) \nRightarrow (1), I.E., IF TWO INTEGRALS

ARE THE SAME \nRightarrow INTEGRANDS / FUNCTIONS INTEGRATED ARE THE SAME. E.G.

$$\int_0^1 x dx = \int_0^1 (1-x) dx = \frac{1}{2} \quad \text{BUT} \quad x \neq (1-x) \quad x \in [0,1]$$

ANOTHER WAY TO REALIZE THAT THE CONVERSE DOESN'T HOLD IS AS FOLLOWS.

TAKE $X \sim \text{UNIFORM}(-1, 1)$ AND DEFINE $Y = X^2$. CLEARLY X AND Y

ARE NOT INDEPENDENT.

$$\text{CONSIDER } E(X) = \int_{-1}^1 u \cdot \frac{1}{2} du = \left. \frac{u^2}{4} \right|_{-1}^1 = \frac{1^2 - (-1)^2}{4} = 0$$

$$E(Y) = E(X^2) = \int_{-1}^1 u^2 \cdot \frac{1}{2} du = \left. \frac{u^3}{6} \right|_{-1}^1 = \frac{1^3 - (-1)^3}{6} = \frac{2}{6}$$

$$E(XY) = E(X \cdot X^2) = E(X^3)$$

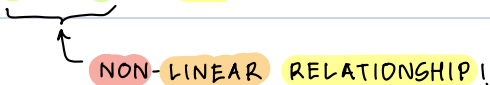
$$= \int_{-1}^1 u^3 \cdot \frac{1}{2} du = \left. \frac{u^4}{8} \right|_{-1}^1 = \frac{1^4 - (-1)^4}{8} = 0$$

$$\Rightarrow \text{Cov}(X, Y) = E(XY) - E(X) \cdot E(Y)$$

$$= 0 - 0 \cdot \frac{1}{3}$$

$$= 0$$

\Rightarrow X AND Y ARE UNCORRELATED.

BUT $Y = X^2$ ARE CLEARLY DEPENDENT.
 NON-LINEAR RELATIONSHIP!

CORRELATION CAN ONLY DETECT LINEAR RELATIONSHIPS! \leftarrow YOU SEE IT NOW! 😊

EXERCISE. PROVE THE ABOVE THEOREM FOR THE DISCRETE CASE.

UNCORRELATEDNESS OF A SET (POSSIBLY OF SIZE ≥ 2) OF RANDOM VARIABLES.

DEFINITION: A SET (POSSIBLY OF SIZE ≥ 2) OF RANDOM VARIABLES IS CALLED
 UNCORRELATED IF THEY ARE PAIRWISE UNCORRELATED, I.E. LET

$S = \{X_1, \dots, X_n\}$ BE A SET OF n RVs, S IS CALLED

UNCORRELATED IF $\text{Cov}(X_i, X_j) = 0 \quad \forall i \neq j = 1, \dots, n.$

REMARK: RECALL THAT UNLIKE UNCORRELATEDNESS, PAIRWISE INDEPENDENCE IS NOT SUFFICIENT FOR MUTUAL INDEPENDENCE.

OBSERVE THE FOLLOWING:

$$\rightarrow \text{Cov}(X+Y, U+V) = \text{Cov}(X, U) + \text{Cov}(X, V) + \text{Cov}(Y, U) + \text{Cov}(Y, V)$$

PROOF: CONSIDER:

$$\begin{aligned} \text{Cov}(X+Y, U+V) &= E[(X+Y)(U+V)] - E(X+Y)E(U+V) \\ &= E(XU) + E(XV) + E(YU) + E(YV) \\ &\quad - (E(X)E(U) + E(X)E(V) + E(Y)E(U) + E(Y)E(V)) \\ &= E(XU) - E(X)E(U) + E(XV) - E(X)E(V) \\ &\quad + E(YU) - E(Y)E(U) + E(YV) - E(Y)E(V) \\ &= \text{Cov}(X, U) + \text{Cov}(X, V) + \text{Cov}(Y, U) + \text{Cov}(Y, V) \end{aligned}$$

$$\rightarrow \text{Cov}(aX+b, cY+d) = ac \text{Cov}(X, Y) \quad a, b, c, d \text{ ARE CONSTANTS}$$

PROOF: EXERCISE. HINT: VERY SIMILAR TO THE PROOF OF THE VARIANCE OF A RANDOM VARIABLE IS AFFECTED BY THE CHANGE OF SCALE BUT NOT OF THE ORIGIN.

$$\rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X,Y)$$

PROOF: CONSIDER:

$$\begin{aligned} \text{Var}(X+Y) &= \text{Cov}(X+Y, X+Y) \\ &= \text{Cov}(X,X) + \text{Cov}(X,Y) + \text{Cov}(Y,X) + \text{Cov}(Y,Y) \\ &= \text{Var}(X) + \text{Cov}(X,Y) + \text{Cov}(X,Y) + \text{Var}(Y) \\ &= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X,Y) \end{aligned}$$

\rightarrow IF X AND Y ARE UNCORRELATED THEN: $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$

PROOF: EXERCISE. HINT: FOLLOWS DIRECTLY FROM THE RESULT ABOVE.

$$\rightarrow \rho_{aX+b, cY+d} = \rho_{X,Y} \quad \text{FOR } a, c > 0$$

PROOF: CONSIDER:

$$\begin{aligned} \rho_{aX+b, cY+d} &= \frac{\text{Cov}(aX+b, cY+d)}{\sqrt{\text{Var}(aX+b) \text{Var}(cY+d)}} = \frac{ac \text{Cov}(X,Y)}{\sqrt{a^2 \text{Var}(X) c^2 \text{Var}(Y)}} \\ &= \frac{ac \text{Cov}(X,Y)}{ac \sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} \\ &= \rho_{X,Y} \end{aligned}$$

REPRESENTING THE CHANGE OF ORIGIN AND SCALE

STANDARDIZED VERSION OF X

STANDARDIZED VERSION OF Y

$$\rightarrow \text{Cov}\left(\frac{X - E(X)}{\sigma_X}, \frac{Y - E(Y)}{\sigma_Y}\right) = \rho_{X,Y}$$

CAN BE REGARDED AS THE STANDARDIZED VERSION OF $\text{Cov}(X,Y)$

MEAN = 0, VAR = 1

PROOF: CONSIDER:

$$\begin{aligned} \text{COV} \left(\frac{X - E(X)}{\sigma_X}, \frac{Y - E(Y)}{\sigma_Y} \right) & \stackrel{\text{CHANGE OF ORIGIN}}{=} \text{COV} \left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y} \right) \\ & \stackrel{\text{CHANGE OF SCALE}}{=} \frac{\text{COV}(X, Y)}{\sigma_X \sigma_Y} = \rho_{X,Y} \end{aligned}$$

RANGE OF THE CORRELATION COEFFICIENT

THEOREM: $-1 \leq \rho_{X,Y} \leq 1$

PROOF: THE PROOF FOLLOWS FROM THE

CAUCHY-SCHWARZ INEQUALITY.

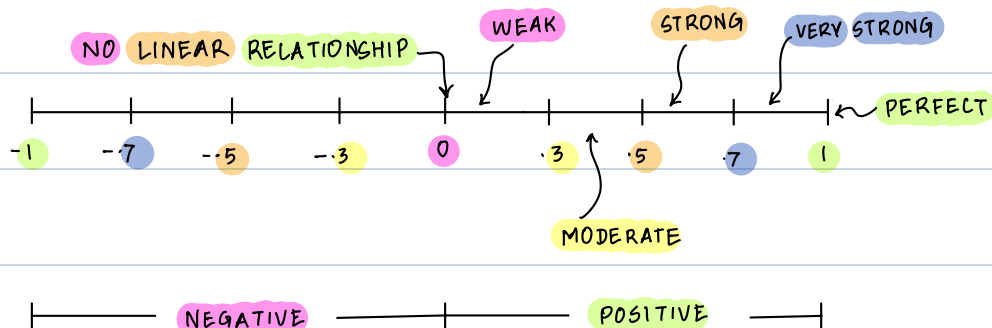
REFER TO THE TEXTBOOK!

RELATED TO THE TRIANGLE INEQUALITY.

$-1 \leq \rho_{X,Y} \leq 1$

PERFECT NEGATIVE LINEAR RELATIONSHIP
 $Y = aX + b$ FOR SOME $a < 0$

PERFECT POSITIVE LINEAR RELATIONSHIP
 $Y = aX + b$ FOR SOME $a > 0$



EXERCISE: CALCULATE THE CORRELATION COEFFICIENT $\rho_{X,Y}$ BETWEEN THE HW

SCORE (X) AND EXAM SCORE (Y) FOR THE FOLLOWING DATA:

• HW SCORE (X): 9 10 7 5 8 5

• EXAM SCORE (Y): 8 9 6 4 7 3

CLASSIFY $\rho_{X,Y}$ INTO ONE OF THE ABOVE RANGES.