

## Sec 2.4 - Correlation Coefficient.

Let  $X$  and  $Y$  have joint pdf  $f(x, y)$ .  
Assume all expected values exist (finite).

$$\text{Let } \mu_1 = E[X], \quad \sigma_1^2 = \text{Var}[X]$$

$$\text{Let } \mu_2 = E[Y], \quad \sigma_2^2 = \text{Var}[Y]$$

def Covariance of  $X$  and  $Y$

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - \mu_1)(Y - \mu_2)] \\ &= E[XY - \mu_2 X - \mu_1 Y + \mu_1 \mu_2] \\ &= E[XY] - \mu_2 E[X] - \mu_1 E[Y] \\ &\quad + \mu_1 \mu_2 \\ &= E[XY] - \mu_1 \mu_2 - \mu_1 \mu_2 \\ &\quad + \mu_1 \mu_2 \end{aligned}$$

$$\begin{aligned} &= E[XY] - \mu_1 \mu_2 \\ \text{OR } &E[XY] - E[X]E[Y] \end{aligned}$$

(2)

def) Correlation Coefficient of X and Y

If  $\sigma_1 > 0$  and  $\sigma_2 > 0$  then  
the correlation coefficient of X  
and Y is

$$\rho = \frac{E[(X - \mu_1)(Y - \mu_2)]}{\sigma_1 \sigma_2}$$

$$= \frac{\text{COV}(X, Y)}{\sigma_1 \sigma_2}$$

$$\begin{aligned} \text{Note: } E[XY] &= \mu_1 \mu_2 + \rho \sigma_1 \sigma_2 \\ &= \mu_1 \mu_2 + \text{COV}(X, Y) \end{aligned}$$

ex) Let the RV's X and Y have  
joint pdf

$$f(x, y) = \begin{cases} x+y, & 0 < x < 1, \quad 0 < y < 1 \\ 0, & \text{o.w.} \end{cases}$$

Compute  $\rho$ .

(3)

$$E[XY] = \int_0^1 \int_0^1 xy(x+y) dx dy$$

$$= \int_0^1 \int_0^1 (x^2 y + xy^2) dx dy$$

$$\vdots$$

$$= 1/3$$

$$\mu_1 = E[X] = \int_0^1 \int_0^1 x(x+y) dx dy$$

$$= \int_0^1 \int_0^1 x(x+y) dy dx$$

$$\mu_2 = E[Y] = 7/12$$

~~$$\int_0^1 \int_0^1 x dx$$~~

marginal pdf. of X

$$f_x(x) = \int_0^1 (x+y) dy$$

$$E[X] = \int_0^1 x \left( \int_0^1 (x+y) dy \right) dx$$

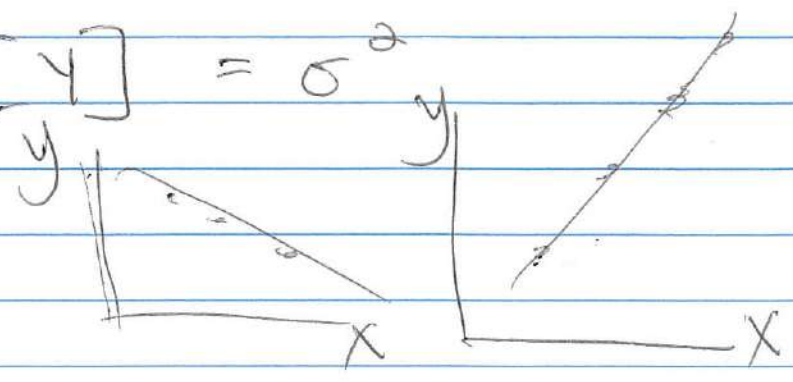
$$E[X^2] = \int_0^1 \int_0^1 x^2 (x+y) dx dy = \frac{5}{12}$$

$$= E[Y^2]$$

$$\text{Var}[X] = \frac{5}{12} - \left(\frac{7}{12}\right)^2 = \frac{11}{144}$$

$$= \text{Var}[Y] = \sigma^2$$

$$\sigma = \sqrt{\frac{11}{144}}$$



$$\text{Cov}(X, Y) = E[XY] - \mu_1 \mu_2$$

$$= \frac{1}{3} - \left(\frac{7}{12}\right)\left(\frac{7}{12}\right) = -\frac{1}{144}$$

Correlation Coefficient of X and Y.

$$\rho = \frac{-\frac{1}{144}}{\sqrt{\frac{11}{144}} \sqrt{\frac{11}{144}}} = -\frac{1}{11}$$

Properties

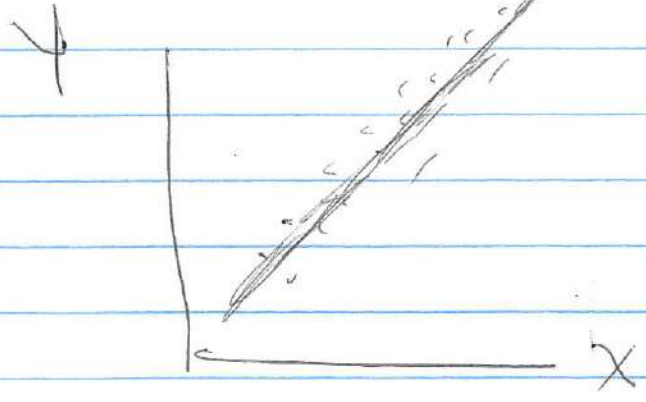
- ①  $-1 \leq \rho \leq 1$
- ② measure intensity of the concentration of the probability for X and Y about a line

③ If  $\rho = 1$ , then we can make a line  $y = a + bx$ ,  $b > 0$  such that it contains all of the distribution of  $X$  and  $Y$ .

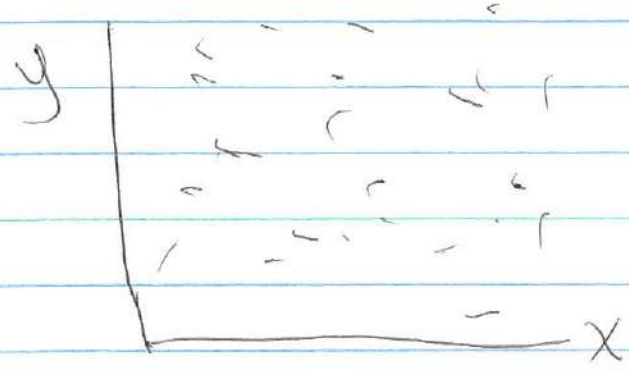
④ If  $\rho = -1$ , then we can make a line  $y = a + bx$ ,  $b < 0$  such that it contains all of the probability of the dist. of  $X$  and  $Y$

⑤ If  $X$  and  $Y$  are independent (sec 2.5), then  $\rho = 0$ .

Relationship via data.

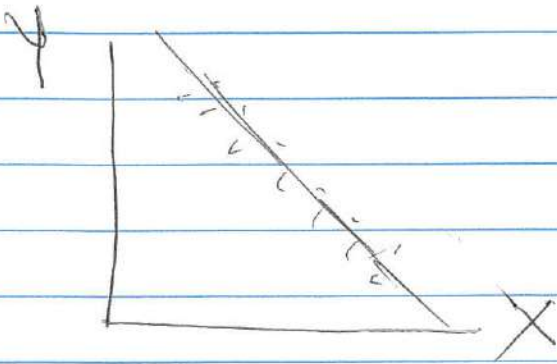


←  $r$  close to 1  
strong relationship.



vague cloud.  
 $r$  close to 0

6

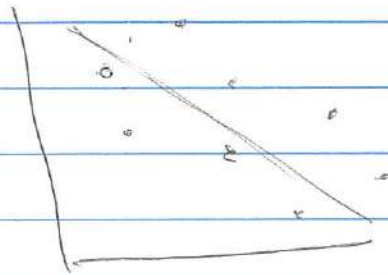


$r$  close to  $-1$

Strong relationship

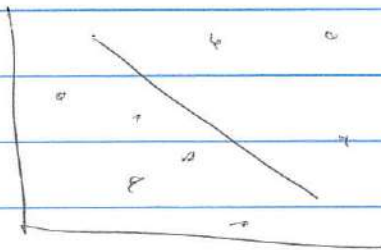
$$-1 \leq r \leq -0.7$$
$$-0.7 \leq r \leq -0.3$$

Strong neg.  
moderate neg.



$$-0.3 \leq r \leq -0.1$$

Weak neg.



$$-0.1 \leq r \leq 0.1$$

No rel.

$$0.1 \leq r \leq 0.3$$

Weak pos.

$$0.3 \leq r \leq 0.7$$

Moderate pos.

$$0.7 \leq r \leq 1$$

Strong pos.

7

Thm 1

Suppose  $(X, Y)$  have joint dist.  $f(x, y)$

$$\text{Let } \mu_1 = E[X]$$

$$\mu_2 = E[Y]$$

$$\sigma_1^2 = \text{Var}[X] > 0 \text{ and finite}$$

$$\sigma_2^2 = \text{Var}[Y] > 0 \text{ and finite.}$$

Let  $\rho$  be the correlation coeff.  
between  $X$  and  $Y$ .

If  $E[Y|X]$  is linear in  $X$

$$\text{then } E[Y|X] = \mu_2 + \rho \cdot \frac{\sigma_2}{\sigma_1} (X - \mu_1)$$

and

$$E[\text{Var}(Y|X)] = \sigma_2^2 (1 - \rho^2)$$

---

$$\text{If } E[Y|X=x] = a + bx$$

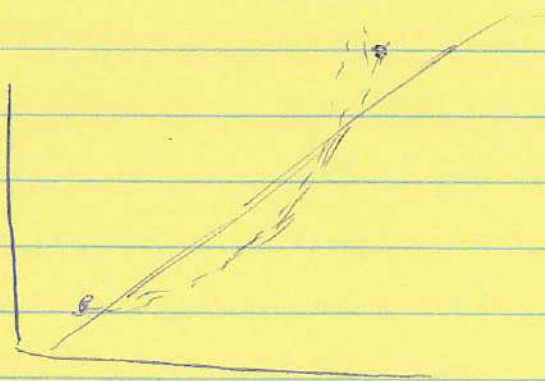
$$\text{slope } = b = \rho \frac{\sigma_2}{\sigma_1} \quad \text{and}$$

$$\text{y-int } a = \mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1$$

$$\mu_2 + \rho \frac{\sigma_2}{\sigma_1} (X - \mu_1)$$

$$= \mu_2 + \rho \frac{\sigma_2}{\sigma_1} X - \rho \frac{\sigma_2}{\sigma_1} \mu_1$$

$$= \underbrace{\mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1}_a + \underbrace{\rho \frac{\sigma_2}{\sigma_1}}_b X$$



For data =

$$\hat{y} = a + bX.$$

predicted value  $\hat{y}$        $a$        $b$        $X$        $\leftarrow$  std dev.

$$\bar{y} = r \cdot \frac{S_y}{S_x} \cdot \bar{X}$$

mean of  $y$  from data.       $\bar{y}$        $r$        $\frac{S_y}{S_x}$        $\bar{X}$       mean of  $X$  from data