

### 3.6 Correlation

Def. a)  $\text{Cor}(X, Y) = \bar{E}[(X - \mu_1)(Y - \mu_2)]$ ,  $\mu_1 = \bar{E}(X)$ ,  $\mu_2 = \bar{E}(Y)$ .

b)  $\rho = \rho_{X,Y} = \frac{\text{Cor}(X, Y)}{\sigma_1 \sigma_2}$ ,  $\sigma_1 = \sqrt{\text{Var}(X)}$ ,  $\sigma_2 = \sqrt{\text{Var}(Y)}$ .

Note  $\text{Cor}(X, Y) = \rho \sigma_1 \sigma_2$

Ex 1. Let  $X, Y$  be r.v. with  $\sigma_2^2 = \text{Var}(Y)$ ,  $\sigma_1^2 = \text{Var}(X)$ , and correlation  $\rho = \rho(X, Y)$ . Confirm that

(1) a)  $Y = \rho \frac{\sigma_2}{\sigma_1} X + Z$ , and  $X, Z$  are uncorrelated:  $\text{Cor}(X, Z) = 0$

b)  $\text{Var}(Z) = \sigma_2^2 (1 - \rho^2)$

Answer. a) Since  $Z = Y - \rho \frac{\sigma_2}{\sigma_1} X$ ,

$$\begin{aligned} \text{Cor}(X, Z) &= \text{Cor}\left(X, Y - \rho \frac{\sigma_2}{\sigma_1} X\right) = \text{Cor}(X, Y) - \rho \frac{\sigma_2}{\sigma_1} \text{Var}(X) = \\ &= \text{Cor}(X, Y) - \rho \sigma_2 \sigma_1 = 0 \end{aligned}$$

$$\begin{aligned} \text{b) } \text{Var}(Z) &= \text{Var}(Y) + \rho^2 \frac{\sigma_2^2}{\sigma_1^2} \text{Var}(X) - 2\rho \frac{\sigma_2}{\sigma_1} \text{Cor}(X, Y) = \\ &= \sigma_2^2 + \rho^2 \frac{\sigma_2^2}{\sigma_1^2} \cdot \cancel{\sigma_1^2} - 2\rho \frac{\sigma_2}{\sigma_1} \cdot \cancel{\sigma_1 \sigma_2} = \sigma_2^2 (1 - \rho^2) \end{aligned}$$

Comments on Ex 1. ① Since  $\text{Var}(Z) \geq 0$ , we have  $1 - \rho^2 \geq 0$ , equivalently,  $|\rho| \leq 1$  ( $-1 \leq \rho \leq 1$ ).

② "Extreme" correlation is when  $\rho^2 = 1$  or  $\rho = \pm 1$ : in this case,  $Z = c$  (constant) with prob. 1, and  $Y = \pm \frac{\sigma_2}{\sigma_1} X + c$

③ The roles of  $X$  and  $Y$  can be switched:

$X = \rho \frac{\sigma_1}{\sigma_2} Y + \tilde{Z}$ , where  $Y, \tilde{Z}$  are uncorrelated, and

$$\text{Var}(\tilde{Z}) = \sigma_1^2(1-\rho^2)$$

Remark 1. Let  $E(X) = \mu_1$ ,  $E(Y) = \mu_2$ .  $E\rho(1)$  can be rewritten as

$$(1) \quad Y - \mu_2 = \rho \frac{\sigma_2}{\sigma_1} (X - \mu_1) + V, \text{ where } X, V \text{ are uncorrelated,}$$

$$\underline{E(V) = 0}, \quad \text{Var}(V) = \sigma_2^2(1-\rho^2)$$

on average, positive and negative values of  $V$  are "balanced".

$|\rho| \leq 1$  by Cauchy-Schwarz inequality:

$$|\rho| = \frac{|\text{Cov}(X, Y)|}{\sigma_1 \sigma_2} \leq 1 \text{ means}$$

$$(1) \quad |E[(X - \mu_1)(Y - \mu_2)]| \leq \sqrt{E[(X - \mu_1)^2]} \sqrt{E[(Y - \mu_2)^2]};$$

(2) is a consequence of Cauchy-Schwarz inequality

$$(*) \quad |E(UV)| \leq \sqrt{E(U^2)} \sqrt{E(V^2)}$$

$$\begin{aligned} U &= X - \mu_1 \\ V &= Y - \mu_2 \end{aligned}$$

Proof of (\*). For any  $t \in \mathbb{R}$ ,

$$0 \leq E[(U + tV)^2] = E(U^2) + 2E(UV)t + E(V^2)t^2$$

$$= c + 2bt + at^2, \text{ where } c = E(U^2), b = E(UV), a = E(V^2).$$

So, discriminant  $D = (2b)^2 - 4ac \leq 0$ , equivalently,  $b^2 \leq ac$ ,

$$|b| \leq \sqrt{a} \sqrt{c}.$$

Some inequalities

1. If  $X \geq 0$ , then  $E(X) \geq 0$ .

1. If  $X \geq Y$ , then  $E(X) \geq E(Y)$ .

2.  $|E(X)| \leq E(|X|)$ , because  $-|X| \leq X \leq |X|$  implies  $-E(|X|) \leq E(X) \leq E(|X|)$ .

3.  $|E(X)| \leq E(|X|) \leq \sqrt{E(X^2)}$  because by Cauchy-Schwarz

$$E(|X|) = E(|X| \cdot 1) \leq \sqrt{E(X^2)} \cdot \sqrt{E(1)} = \sqrt{E(X^2)}.$$

### 3.4 Indicator method

Recall

(a)  $I_A = \begin{cases} 1 & \text{if } A \\ 0 & \text{if not } A \end{cases}$  is Bernoulli( $p$ ) with  $p = P(A)$

$$I_{A^c} = 1 - I_A, \quad I_A I_B = I_{A \cap B}$$

(b)  $E(I_A) = P(A) = p$ ,  $\text{Var}(I_A) = P(A) - P(A)^2$

$$\begin{aligned} \text{Cor}(I_A, I_B) &= E(I_A I_B) - E(I_A)E(I_B) \\ &= P(A \cap B) - P(A)P(B) \end{aligned}$$

$A, B$  independent iff  $\text{Cor}(I_A, I_B) = 0$ .

Ex 1. Hats of  $n$  people are mixed up. Everyone randomly one by one picks a hat. Let  $X = \#$  of matches.

Find  $E(X) = ?$ , and  $\text{Var}(X) = ?$ .

Answer.  $A_i = \text{"}i\text{th person gets right hat"}$ ,  $i = 1, \dots, n$ .

$$\text{Then } X = \sum_{i=1}^n I_{A_i}, \quad E(X) = \sum_{i=1}^n P(A_i) = n \cdot \frac{1}{n} = 1.$$

$$P(A_i) = \frac{(n-1)!}{n!} = \frac{1}{n}$$

$$V_{\text{var}}(X) = E(X^2) - E(X)^2 = 2 - 1 = 1.$$

$$E(X^2) = E\left(\sum_{i,j=1}^n 1_{A_i} 1_{A_j}\right) = E\left(\sum_{i=1}^n 1_{A_i} + 2 \sum_{i < j} 1_{A_i} 1_{A_j}\right) = \sum_{i=1}^n P(A_i) + 2 \sum_{i < j} P(A_i A_j)$$

$$= n \cdot \frac{1}{n} + 2 \cdot \frac{n(n-1)}{2} \cdot \frac{1}{n(n-1)} = 2, \text{ because}$$

$$P(A_i A_j) = \frac{(n-2)!}{n!} = \frac{1}{n(n-1)}.$$

If  $X = \sum_{i=1}^n 1_{A_i}$ , then  $X^2 = \left(\sum_{i=1}^n 1_{A_i}\right)^2 = \sum_{i,j=1}^n 1_{A_i} 1_{A_j} = \sum_{i=1}^n 1_{A_i} + 2 \sum_{i < j} 1_{A_i} 1_{A_j}$ ,  
 and  $E(X^2) = \sum_{i=1}^n P(A_i) + 2 \sum_{i < j} P(A_i A_j)$

$E \times 2$ . There are  $m$  red and  $w$  white balls in the box,  
 $K = m + w$ .  $n$  balls are removed one by one. Assume,  
 $n < m < K$ . Let  $X = \#$  of red balls among  $n$ .  $X$  is  
 hypergeometric.

a) Find  $E(X)$ .

Answer.  $A_i = \text{"}i\text{th removed ball is red"}$ . Then

$$X = 1_{A_1} + \dots + 1_{A_n}, \quad E(X) = \sum_{i=1}^n P(A_i) = np = n \cdot \frac{m}{K}.$$

$$P(A_i) = \frac{m \cdot (K-1) \dots (K-1-(n-1)+1)}{K(K-1) \dots (K-n+1)} = \frac{m}{K} =: p$$

b) Find  $\text{Var}(X) = E(X^2) - E(X)^2 = E(X^2) - n^2 p^2$ .

Answer.  $X^2 = \sum_{i,j} 1_{A_i} 1_{A_j} = \sum_{i=1}^n 1_{A_i} + 2 \sum_{i < j} 1_{A_i} 1_{A_j}$

$$E(X^2) = \sum_{i=1}^n P(A_i) + 2 \sum_{i < j} P(A_i A_j) =$$

$$P(A_i) = p, \quad P(A_i A_j) = \frac{m(m-1)(K-2) \dots (K-2-(n-2)+1)}{K(K-1) \dots (K-n+1)} = \frac{m}{K} \frac{m-1}{K-1} (\approx p^2)$$

$$= np + 2 \cdot \frac{n(n-1)}{2} \cdot \frac{m}{K} \cdot \frac{m-1}{K-1} = np + n(n-1) \frac{m}{K} \cdot \frac{m-1}{K-1}$$

$$= np + \cancel{2} \frac{n(n-1)}{2} \cdot \frac{m}{K} \cdot \frac{m-1}{K-1} = np + n(n-1) \frac{m}{K} \cdot \frac{m-1}{K-1}$$

$$\text{Var}(X) = np + n(n-1) \frac{m}{K} \frac{m-1}{K-1} - n^2 p^2.$$