

Corollary. Suppose that  $X_1, \dots, X_n$  are uncorrelated and have same mean  $\mu$  and variance  $\sigma^2$ . Let

Function of  $X_1, \dots, X_n$  →  $E(\bar{X}_n) = \mu$  (c.f.)

relax i.i.d condition

$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n-1}$  (a r.v.)

definition of variance =  $E(\underline{X} - \underline{\mu})^2 = \sigma^2$  (c.f.)

then  $E(S^2) = \sigma^2 \Rightarrow S^2 \xrightarrow{p} \sigma^2$

Proof.  $(n-1)S^2 = \sum_{i=1}^n (X_i - \bar{X}_n)^2$

It can be shown that when  $n \rightarrow \infty$   $\text{Var}(S^2) \rightarrow 0$

Note.  $\bar{X}_n \xrightarrow{p} \mu$   
 $E(\bar{X}_n) = \mu$   
 $\text{Var}(\bar{X}_n) = \sigma^2/n$

$(n-1)S^2 = \sum_{i=1}^n [(X_i - \mu) - (\bar{X}_n - \mu)]^2$

$= \sum_{i=1}^n [(X_i - \mu)^2] + [\sum_{i=1}^n (\bar{X}_n - \mu)^2] - 2(\bar{X}_n - \mu) [\sum_{i=1}^n (X_i - \mu)]$

$= [\sum_{i=1}^n (X_i - \mu)^2] + n(\bar{X}_n - \mu)^2 - 2n(\bar{X}_n - \mu)^2$

$= [\sum_{i=1}^n (X_i - \mu)^2] - n(\bar{X}_n - \mu)^2$

Therefore,

$(n-1)E(S^2) = \left\{ \sum_{i=1}^n E[(X_i - \mu)^2] \right\} - nE[(\bar{X}_n - \mu)^2]$

$\because$  uncorrelated  $\Rightarrow \text{Var}(X_i) = \sigma^2$   $\Rightarrow n\sigma^2 - n\text{Var}(\bar{X}_n) = (n-1)\sigma^2$

$= \sigma^2/n$  (LN p.8-13)

Note. The previous three corollaries also hold if  $X_1, \dots, X_n$  are "uncorrelated" is replaced by "independent."

$\because$  "independence" implies "uncorrelated"

Theorem ( $\rho$  of linear transformation).

Recall. 2nd Corollary in LN p.8-10

$\text{Cor}(a_0 + a_1 X, b_0 + b_1 Y) = \text{sign}(a_1 b_1) \times \text{Cor}(X, Y)$

and  $| \text{Cor}(a_0 + a_1 X, b_0 + b_1 Y) | = | \text{Cor}(X, Y) |$

Why?  $\rho_{XY}$  is invariant under location and scale changes.

standardization

Proof. Let  $S = a_0 + a_1 X$  and  $T = b_0 + b_1 Y$ , then

$\text{Cov}(S, T) = \text{Cov}(a_0 + a_1 X, b_0 + b_1 Y) = a_1 b_1 \text{Cov}(X, Y)$

$\text{Var}(S) = a_1^2 \text{Var}(X)$ , and  $\text{Var}(T) = b_1^2 \text{Var}(Y)$

Therefore,

$\rho_{ST} = \frac{\text{Cov}(S, T)}{\sigma_S \sigma_T} = \frac{a_1 b_1 \text{Cov}(X, Y)}{|a_1| |b_1| \sigma_X \sigma_Y} = \frac{a_1 b_1}{|a_1 b_1|} \rho_{XY} = \text{sign}(a_1 b_1) \rho_{XY}$

Thm in LN p.8-12

$\text{Var}=1$

$\frac{Cov(X,Y)}{\sigma_X \sigma_Y} \rightarrow$  Theorem (some properties of  $\rho$ ).  $\hookrightarrow$  Cauchy-Schwarz inequality

$0 \leq |\rho_{XY}| \leq 1 \iff (1) -1 \leq \rho_{XY} \leq 1. (\iff |Cov(X,Y)| \leq \sigma_X \sigma_Y)$

$\rho$  is unit-free  $(2) \rho_{XY} = \pm 1$  if and only if there exist  $a, b \in \mathbb{R}$

such that  $P(Y=aX+b)=1$ .

$Y=aX+b$   
almost surely

(3) Furthermore,  $\rho_{XY}=1$ , if  $a>0$  and  $\rho_{XY}=-1$ , if  $a<0$ .

$P(E_A)=1$

$\rho(E_A^c)=0$

Proof of (1).

柯西不等式  $\rightarrow$  Cauchy-Schwarz inequality

Cauchy-Schwarz inequality

$\underline{u}=(u_1, \dots, u_n), \underline{v}=(v_1, \dots, v_n)$

$|\sum_i u_i v_i| = |\langle \underline{u}, \underline{v} \rangle| \leq \|\underline{u}\| \cdot \|\underline{v}\| = \sqrt{\sum_i u_i^2} \cdot \sqrt{\sum_i v_i^2}$

$\uparrow \begin{matrix} Y(\omega) \cdot \frac{1}{\sigma_Y} \\ X(\omega) \cdot \frac{1}{\sigma_X} \end{matrix} \uparrow \begin{matrix} \sum_i \rightarrow \sum_{\omega \in \Omega} \end{matrix}$

Thm in  
LNp. 8-13

$$0 \leq Var\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right)$$

$$= Var\left(\frac{X}{\sigma_X}\right) + Var\left(\frac{Y}{\sigma_Y}\right) + 2Cov\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right)$$

$$= \frac{Var(X)}{\sigma_X^2} + \frac{Var(Y)}{\sigma_Y^2} + 2 \frac{Cov(X,Y)}{\sigma_X \sigma_Y}$$

$$= 1 + 1 + 2\rho_{XY} \Rightarrow \rho_{XY} \geq -1.$$

$$|E(XY)| \leq \sqrt{E(X^2)} \sqrt{E(Y^2)}$$

Similarly,

$$|Cov(X,Y)| = |E[(X-\mu_X)(Y-\mu_Y)]| \leq \sqrt{E[(X-\mu_X)^2]} \sqrt{E[(Y-\mu_Y)^2]} = \sigma_X \sigma_Y$$

$$0 \leq Var\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = 1 + 1 - 2\rho_{XY} \Rightarrow \rho_{XY} \leq 1.$$

Proof of (2) and (3). We see from the proof of (1),



$$\rho_{XY} = 1 \iff Var\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = 0,$$

dist.  
unknown

• different transformation  $g$   
 $\Rightarrow$  different information  
• These expectations are called  
parameters in statistics  
• parameters can be estimated  
by r.v.'s (transformation  
of data), e.g.,

$$\bar{X}_n \xrightarrow{P} \mu$$

$$S_n^2 \xrightarrow{P} \sigma^2$$

$Var(Z)=0 \iff$   
 $Z=c$  almost  
surely, for  
a constant  $c$

$$\iff P\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} = c\right) = 1,$$

where  $c$  is a constant.

$$\iff P\left(Y = \frac{\sigma_Y}{\sigma_X} X + c\sigma_Y\right) = 1.$$

fixed

$$\text{Similarly, } \rho_{XY} = -1 \iff P\left(Y = -\frac{\sigma_Y}{\sigma_X} X + c\sigma_Y\right) = 1.$$

- **Q:** How to use expectations to (roughly) characterize the distribution of random variables  $X_1, \dots, X_n$ ?

$\triangleright g(X_1, \dots, X_n) = X_i \Rightarrow E[g(\mathbf{X})] = \mu_{X_i}$ : mean of  $X_i$ .

$g$ : 1st order  
polynomials  
of  $X_1, \dots, X_n$

$\triangleright g(X_1, \dots, X_n) = (X_i - \mu_{X_i})^2 \Rightarrow E[g(\mathbf{X})] = \sigma_{X_i}^2$ : variance of  $X_i$ .

$\triangleright g(X_1, \dots, X_n) = (X_i - \mu_{X_i})(X_j - \mu_{X_j})$  for  $i \neq j$   
 $\Rightarrow E[g(\mathbf{X})] = \sigma_{X_i X_j}$ : covariance of  $X_i$  and  $X_j$ .

$g$ : 2nd order  
polynomials  
of  $X_1, \dots, X_n$

$\triangleright g(X_1, \dots, X_n) = [(X_i - \mu_{X_i})/\sigma_{X_i}][(X_j - \mu_{X_j})/\sigma_{X_j}]$  for  $i \neq j$   
 $\Rightarrow E[g(\mathbf{X})] = \rho_{X_i X_j}$ : correlation coefficient of  $X_i$  and  $X_j$ .

$\triangleright$  Notes.  $\mu_{X_i}, \sigma_{X_i}^2, \sigma_{X_i X_j}, \rho_{X_i X_j}$  are constants, not random



Recall. Conditional dist.  
LNp. 7-51~59

# Conditional Expectation $\mathbf{X} \in \mathbb{R}^n, \mathbf{Y} \in \mathbb{R}^m$ <sup>p. 8-19</sup>

- Recall.  $p_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})$  or  $f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})$  is a pmf/pdf for  $\mathbf{y}$  ( $\mathbf{y}$ : random,  $\mathbf{x}$ : fixed).
- Definition. For random vectors  $\mathbf{X}$  and  $\mathbf{Y}$ , the conditional expectation of  $\mathbf{Z} = h(\mathbf{Y})$  given  $\mathbf{X} = \mathbf{x}$ , where  $h: \mathbb{R}^m \rightarrow \mathbb{R}^1$ , is a function of  $\mathbf{x}$

平均:  $h(\mathbf{Y})$   
權重:  $P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})$

$$E_{\mathbf{Y}|\mathbf{X}}(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} h(\mathbf{y}) p_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) \quad (1)=(2) \text{ in LNp. 8-1}$$

in the discrete case, or,  $\mathbf{Z}$

平均:  $h(\mathbf{Y})$   
權重:  $f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x})$

$$E_{\mathbf{Y}|\mathbf{X}}(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}) = \int_{\mathbb{R}^m} h(\mathbf{y}) f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mathbf{y} \quad (3)=(4) \text{ in LNp. 8-1}$$

in the continuous case,  $\mathbf{Z}$

$$= \int_{\mathbb{R}^m} \mathbf{z} f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

provided that the sum or integral converges absolutely.

Some Notes.

▪  $E_{\mathbf{Y}|\mathbf{X}}(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x})$ : a function of  $\mathbf{x}$  and free of  $\mathbf{Y}$ .

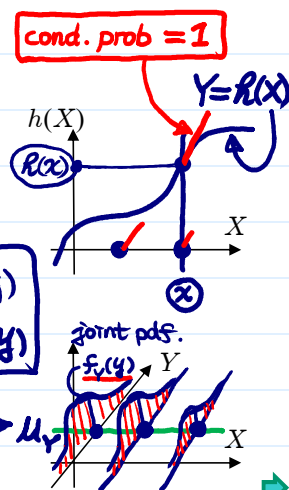
▪  $E_{\mathbf{Y}|\mathbf{X}}[h(\mathbf{X}) | \mathbf{X} = \mathbf{x}] = h(\mathbf{x})$ .

$$E_{\mathbf{Y}|\mathbf{X}}[h(\mathbf{X}, \mathbf{Y}) | \mathbf{X} = \mathbf{x}] = \int h(\mathbf{x}, \mathbf{y}) f_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mathbf{y}$$

▪ If  $\mathbf{X}$  and  $\mathbf{Y}$  are independent, then

$$E_{\mathbf{Y}|\mathbf{X}}(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}) = E_{\mathbf{Y}}[h(\mathbf{Y})]$$

a constant line of  $\mathbf{x}$



Let  $g(\mathbf{x}) = E_{\mathbf{Y}|\mathbf{X}}[h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}]$ , where  $g: \mathbb{R}^n \rightarrow \mathbb{R}^1$ , then we write

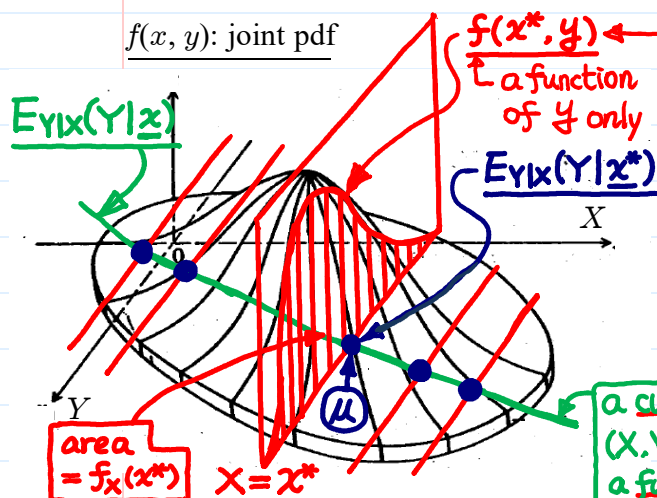
after cf. before

$$E_{\mathbf{Y}|\mathbf{X}}(h(\mathbf{Y}) | \mathbf{X})$$

when  $\mathbf{x}$  in  $g$  is replaced by  $\mathbf{X}$  (a fixed value replaced by a r.v.).

Notice that  $g(\mathbf{X})$  is a random variable.

$f(x, y)$ : joint pdf



▪  $f(x, y)$ : a joint pdf.

▪ Fix  $x^*$ , is  $f(x^*, y)$  a pdf of  $y$ ? i.e.,

$$f_X(x^*) = \int_{-\infty}^{\infty} f(x^*, y) dy \stackrel{?}{=} 1.$$

▪  $f_{Y|\mathbf{X}}(y|x^*) = f(x^*, y) / f_X(x^*)$  is a pdf of  $y$  since

$$\frac{\int_{-\infty}^{\infty} f(x^*, y) dy}{f_X(x^*)} = 1.$$

a curve on (X, Y) plane: a function of  $\mathbf{x}$

$$h(\mathbf{Y}) = \mathbf{Y}$$

center of gravity

▪  $E_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y} | x^*)$ : mean of  $f_{Y|\mathbf{X}}(y|x^*)$ .

▪ Do it for any  $x = x^*$ , and get a function of  $x$   $\Rightarrow E_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y} | x)$

$$\begin{aligned} \int (y - \mu) f(x^*, y) dy &= 0 \\ \Rightarrow \int y f(x^*, y) dy &= \mu \int f(x^*, y) dy = \mu f_X(x^*) \\ \Rightarrow \mu &= \int y f(x^*, y) / f_X(x^*) dy = E_{Y|\mathbf{X}}(Y | x^*) \\ \Rightarrow f_{Y|\mathbf{X}}(y|x^*) &\text{ \& } f(x^*, y) \text{ have same center of gravity} \end{aligned}$$



➤ Example. Sample a student from an elementary school. Let

$X = \text{age}$  (unit: year),  $Y = \text{height}$  (unit: cm)

of the student. **Population**: all students of the school.

Q: What's the source of their randomness?

•  $Y|X=x$ : a random variable (unit: cm) that represents the height distribution of students with age=x.

•  $g(x) = E_{Y|X}(Y|X=x)$  or  $E_{Y|X}(Y|x)$ : a function maps from age (unit: year) to average height (unit: cm) of students with age=x.

Note.  $E_{Y|X}(Y|x)$  is not a random variable.

•  $g(X) = E_{Y|X}(Y|X)$ : a random variable because it is a function of age  $X$ , where  $X$  is a random variable.

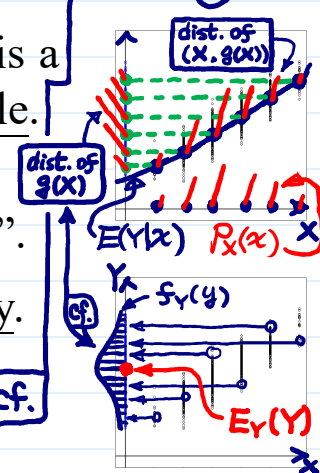
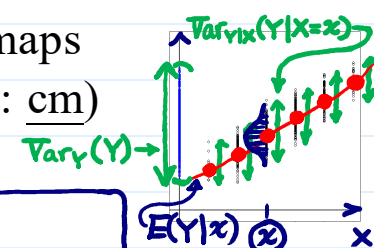
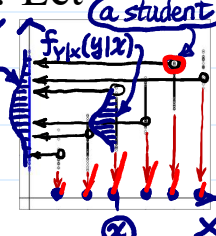
$g(x)$ : 1-to-1 in this case  $\Rightarrow P(E_{Y|X}(Y|x) = E_{Y|X}(Y|z)) = P(X=z)$

Note.  $g(X) = E_{Y|X}(Y|X)$  is height, its unit is "cm".

•  $Var_{Y|X}(Y|X=x)$  &  $Var_{Y|X}(Y|X)$  defined similarly.

•  $E_Y(Y)$ : average height of all students;

$Var_Y(Y)$ : variation of height of all students.



• Theorem (Law of Total Expectation). For two random vectors

$\underline{X} (\in \mathbb{R}^m)$  and  $\underline{Y} (\in \mathbb{R}^n)$ ,

$$E_{\underline{X}, \underline{Y}}[R(\underline{Y})] = E_{\underline{Y}} E_{\underline{X}|\underline{Y}}$$

$$E_{\underline{X}}\{E_{\underline{Y}|\underline{X}}[h(\underline{Y})|\underline{X}]\} = E_{\underline{Y}}[h(\underline{Y})]$$

use the example in LNp.8-21 to realize the meaning of these terms.

In particular, let  $h(\underline{Y}) = Y_i$ , we have

$$E_{\underline{X}}[E_{\underline{Y}|\underline{X}}(Y_i|\underline{X})] = E_{\underline{Y}}(Y_i)$$

$Y_i$  &  $E_{Y|X}(Y_i|X)$  have same mean

Proof.

$$g(x)$$

$$E_{Y_i}(Y_i) = E_{\underline{X}, \underline{Y}}(Y_i)$$

(only prove it for the continuous case)

$$E_{\underline{X}}\{E_{\underline{Y}|\underline{X}}[h(\underline{Y})|\underline{X}]\}$$

$$= \int_{\mathbb{R}^m} E_{\underline{Y}|\underline{X}}(h(\underline{Y})|\underline{x}) f_{\underline{X}}(\underline{x}) d\underline{x}$$

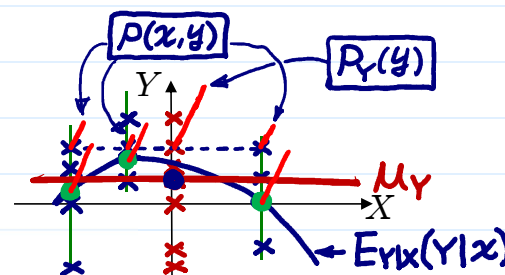
$$= \int_{\mathbb{R}^m} \left[ \int_{\mathbb{R}^n} h(\underline{y}) f_{\underline{Y}|\underline{X}}(\underline{y}|\underline{x}) d\underline{y} \right] f_{\underline{X}}(\underline{x}) d\underline{x}$$

$$= \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} h(\underline{y}) \frac{f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y})}{f_{\underline{X}}(\underline{x})} f_{\underline{X}}(\underline{x}) d\underline{x} d\underline{y}$$

$$= \int_{\mathbb{R}^n} h(\underline{y}) \left[ \int_{\mathbb{R}^m} f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) d\underline{x} \right] d\underline{y}$$

$$= \int_{\mathbb{R}^n} h(\underline{y}) f_{\underline{Y}}(\underline{y}) d\underline{y}$$

$$= E_{\underline{Y}}[h(\underline{Y})]$$



Interchange dy & dx  
dydx → dxdy

$$E_{\underline{X}, \underline{Y}}[R(\underline{Y})]$$

$$= \int_{\mathbb{R}^n} h(\underline{y}) f_{\underline{Y}}(\underline{y}) d\underline{y}$$

$$= E_{\underline{Y}}[h(\underline{Y})]$$

$$= E_{\underline{Y}}[h(\underline{Y})]$$

$$= E_{\underline{Y}}[h(\underline{Y})]$$

multiplication law (LNp. 7-55)

$$= \int_{\mathbb{R}^n} \int_{\mathbb{R}^m} h(\underline{y}) f_{\underline{X}|\underline{Y}}(\underline{x}|\underline{y}) f_{\underline{Y}}(\underline{y}) d\underline{x} d\underline{y}$$

$$\Rightarrow E_{\underline{X}, \underline{Y}} = E_{\underline{Y}} E_{\underline{X}|\underline{Y}}$$

generalization:  
 $E_{\underline{X}, \underline{Y}}[R(\underline{X}, \underline{Y})]$   
 $= E_{\underline{X}} E_{\underline{Y}|\underline{X}}[R(\underline{X}, \underline{Y})|\underline{X}]$   
 $= E_{\underline{X}} E_{\underline{X}|\underline{Y}}[R(\underline{X}, \underline{Y})|\underline{Y}]$