



# Multiple Random Variables

## Lecture 5

Sept. 23, 25, 30, and Oct. 7 and 14, 2025

## Outline

- 1 Outline
- 2 Functions of Two Random Variables
- 3 Two Functions of Two Random Variables
- 4 Random Vectors
- 5 Mean and Covariance Matrix
- 6 Gaussian Random Vector
- 7 Mean-Square Estimation
- 8 Linear Mean-Square Estimation

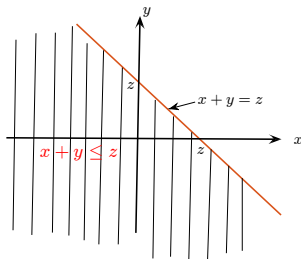
# Functions of Two Random Variables

- Suppose we are given two RVs  $X$  and  $Y$  with known joint cdf  $F_{XY}(x, y)$  (or pdf  $f_{XY}(x, y)$ ) and a function  $z = g(x, y)$ . What is the cdf (or pdf) of the random variable  $Z = g(X, Y)$ ?
- Use

$$F_Z(z) = P(Z \leq z) = P((x, y) : g(x, y) \leq z) \\ = \int \int_{\{(x, y) : g(x, y) \leq z\}} f_{XY}(x, y) \, dx \, dy$$

- Then  $f_Z(z) = \frac{dF_Z(z)}{dz}$ .

**Example: Sum of Two RVs.** Let  $Z = X + Y$  and  $(X, Y) \sim f_{XY}(x, y)$ . Find  $f_Z(z)$ .



$$F_Z(z) = \int \int_{\{(x,y): x+y \leq z\}} f_{XY}(x, y) dx dy = \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{z-y} f_{XY}(x, y) dx \right] dy$$

$$\Rightarrow f_Z(z) = \frac{dF_Z(z)}{dz} = \int_{-\infty}^{\infty} f_{XY}(z - y, y) dy$$

If  $X$  and  $Y$  are **independent**, then  $f_Z(z) = \int_{-\infty}^{\infty} f_X(z - y)f_Y(y) dy$ . This is **convolution** of  $f_X(\cdot)$  with  $f_Y(\cdot)$



## Two Functions of Two Random Variables

- Suppose we are given two RVs  $X$  and  $Y$  with known joint cdf  $F_{XY}(x, y)$  (or pdf  $f_{XY}(x, y)$ ) and two functions  $z = g(x, y)$  and  $w = h(x, y)$ . What is the joint cdf (or pdf) of random variables  $Z = g(X, Y)$  and  $W = h(X, Y)$ ?
- **Approach 1:** Use

$$\begin{aligned} F_{ZW}(z, w) &= P(Z \leq z, W \leq w) \\ &= P(\{(x, y) : g(x, y) \leq z, h(x, y) \leq w\}) \\ &= \int \int_{\{(x, y) : g(x, y) \leq z, h(x, y) \leq w\}} f_{XY}(x, y) \, dx \, dy \end{aligned}$$

- Then  $f_{ZW}(zw) = \frac{\partial^2 F_{ZW}(z,w)}{\partial z \partial w}$ .

- **Approach 2:** Suppose  $(X, Y) \sim f_{XY}(x, y)$ ,  $Z = g(X, Y)$ ,  $W = h(X, Y)$ , and  $g(x, y)$  and  $h(x, y)$  are differentiable.
- For a fixed  $(z, w)$ , simultaneously solve  $g(x, y) = z$  and  $h(x, y) = w$  for real-valued  $(x, y)$ . If there exists no real-valued solution, then  $f_{ZW}(z, w) = 0$ .
- Else, let there be  $n$  solutions  $(x_1, y_1), \dots, (x_n, y_n)$  satisfying  $g(x_i, y_i) = z$  and  $h(x_i, y_i) = w$ . Then

$$f_{ZW}(z, w) = \sum_{i=1}^n \frac{f_{XY}(x_i, y_i)}{|\det J(x_i, y_i)|}$$

if  $\det J(x_i, y_i) \neq 0$  for any  $i$ , where the Jacobian of transformation

$$J(x_i, y_i) = \left[ \begin{array}{cc} \frac{\partial g(x, y)}{\partial x} & \frac{\partial g(x, y)}{\partial y} \\ \frac{\partial h(x, y)}{\partial x} & \frac{\partial h(x, y)}{\partial y} \end{array} \right] \bigg|_{(x, y) = (x_i, y_i)}$$

- This method fails if  $\det J(x_i, y_i) = 0$

**Example:** Let  $Z = X + Y$  and  $(X, Y) \sim f_{XY}(x, y)$ . Find  $f_Z(z)$ .

We have  $Z = X + Y = g(X, Y)$ . We create a new RV:

$W = X = h(X, Y)$ , and then use the “expression” given earlier. For given  $(z, w)$ , we have a unique solution ( $n = 1$ ):  $(x_1, y_1) = (w, z - w)$ .

The Jacobian is

$$J(x, y) = \begin{bmatrix} \frac{\partial(x+y)}{\partial x} & \frac{\partial(x+y)}{\partial y} \\ \frac{\partial x}{\partial x} & \frac{\partial x}{\partial y} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

Then  $\det J(x, y) = -1$  and  $|\det J(x, y)| = 1$ . Thus

$$\begin{aligned}
 f_{ZW}(z, w) &= f_{XY}(w, z - w) \\
 \Rightarrow f_Z(z) &= \int_{-\infty}^{\infty} f_{ZW}(zw) dw = \int_{-\infty}^{\infty} f_{XY}(w, z - w) dw
 \end{aligned}$$

If  $X$  and  $Y$  are independent,  $f_Z(z) = \int_{-\infty}^{\infty} f_X(w)f_Y(z - w) dw$ :  
convolution!



**Example:** Let  $Z = \sqrt{X^2 + Y^2} = g(X, Y)$ ,  $(X, Y) \sim f_X(x)f_Y(y)$  and both  $X$  and  $Y \sim \mathcal{N}(0, \sigma^2)$ . Find  $f_Z(z)$ .

Approach 2: Given

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

We create a new RV:  $W = X = h(X, Y)$ . For given  $(z, w)$ , we have  $x = w$  and therefore,  $y^2 = z^2 - x^2 = z^2 - w^2 \Rightarrow y = \pm\sqrt{z^2 - w^2}$  if  $z^2 > w^2$ .

- If  $z < |w|$ , there is no real-valued solution, hence,  $f_{ZW}(z, w) = 0$  for  $z < |w|$ .
- For  $z > |w| \geq 0$ , there are two solutions:  $(x_1, y_1) = (w, \sqrt{z^2 - w^2})$ ,  $(x_2, y_2) = (w, -\sqrt{z^2 - w^2})$

The Jacobian is

$$J(x, y) = \begin{bmatrix} \frac{\partial \sqrt{x^2+y^2}}{\partial x} & \frac{\partial \sqrt{x^2+y^2}}{\partial y} \\ \frac{\partial x}{\partial x} & \frac{\partial x}{\partial y} \end{bmatrix} = \begin{bmatrix} \frac{x}{\sqrt{x^2+y^2}} & \frac{y}{\sqrt{x^2+y^2}} \\ 1 & 0 \end{bmatrix}$$

$$\Rightarrow |\det J(x, y)| = \frac{|y|}{|\sqrt{x^2 + y^2}|} = \frac{\sqrt{z^2 - w^2}}{z} \text{ for } (x, y) = (x_1, y_1), (x_2, y_2)$$

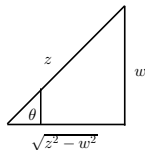
Thus, for  $z > |w| \geq 0$ ,

$$f_{ZW}(z, w) = \sum_{i=1}^2 \frac{f_{XY}(x_i, y_i)}{|\det J(x_i, y_i)|} = \frac{z}{\sqrt{z^2 - w^2}} \left[ \frac{1}{2\pi\sigma^2} e^{-\frac{z^2}{2\sigma^2}} + \frac{1}{2\pi\sigma^2} e^{-\frac{z^2}{2\sigma^2}} \right]$$

$$\Rightarrow f_{ZW}(z, w) = \begin{cases} \frac{z}{\pi \sigma^2 \sqrt{z^2 - w^2}} e^{-\frac{z^2}{2\sigma^2}} & z > |w| \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Noting that  $z > |w| \geq 0$  is equivalent to  $-z < w < z$ ,  $z > 0$ ,

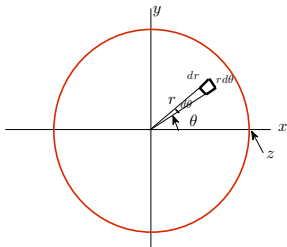
$$\begin{aligned} f_Z(z) &= \int_{-\infty}^{\infty} f_{ZW}(z, w) dw = \int_{-z}^z \frac{z}{\pi \sigma^2 \sqrt{z^2 - w^2}} e^{-\frac{z^2}{2\sigma^2}} dw \\ &= \frac{ze^{-\frac{z^2}{2\sigma^2}}}{\pi \sigma^2} \int_{-z}^z \frac{1}{\sqrt{z^2 - w^2}} dw = \frac{2ze^{-\frac{z^2}{2\sigma^2}}}{\pi \sigma^2} \int_0^z \frac{1}{\sqrt{z^2 - w^2}} dw \end{aligned}$$



Set  $w = z \sin(\theta) \Rightarrow dw = z \cos(\theta) d\theta$  and  $\sqrt{z^2 - w^2} = z \cos(\theta)$ . Then  $\int_0^z \frac{1}{\sqrt{z^2 - w^2}} dw = \int_0^{\pi/2} d\theta = \frac{\pi}{2}$ . Thus, we have the Rayleigh pdf

$$f_Z(z) = \frac{z}{\sigma^2} e^{-\frac{z^2}{2\sigma^2}} u(z)$$

**Approach 1:** Clearly  $F_Z(z) = 0$  for  $z \leq 0$ . Set  $x^2 + y^2 = r^2$ ,  $dx dy = r dr d\theta$ :



$$\begin{aligned} F_Z(z) &= \int \int_{\{(x,y): \sqrt{x^2+y^2} \leq z\}} \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}} dx dy \\ &= \frac{1}{2\pi\sigma^2} \int_0^z r e^{-\frac{r^2}{2\sigma^2}} \left[ \int_0^{2\pi} d\theta \right] dr = \frac{1}{\sigma^2} \int_0^z r e^{-\frac{r^2}{2\sigma^2}} dr \\ \Rightarrow f_Z(z) &= \frac{dF_Z(z)}{dz} = \frac{z}{\sigma^2} e^{-\frac{z^2}{2\sigma^2}} u(z) \end{aligned}$$

## Random Vectors

- Let  $X_1, X_2, \dots, X_n$  be random variables defined on the same probability space. We define a random vector (RV) as

$$\mathbf{x} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

- $\mathbf{X}$  is completely specified by its joint CDF for  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ :

$$F_{\mathbf{X}}(\mathbf{x}) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n), \quad \mathbf{x} \in \mathbb{R}^n$$

- If  $\mathbf{X}$  is continuous, i.e.,  $F_{\mathbf{X}}(\mathbf{x})$  is a continuous function of  $\mathbf{x}$ , then  $\mathbf{X}$  can be specified by its joint pdf

$$f_{\mathbf{X}}(\mathbf{x}) = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \frac{\partial^n F_{\mathbf{X}}(\mathbf{x})}{\partial x_1 \dots \partial x_n}$$



- Conditional cdf (pdf, pmf) can also be defined in the usual way.

E.g., the conditional pdf of  $\mathbf{X}_{k+1}^n \stackrel{\text{def}}{=} (X_{k+1}, X_{k+2}, \dots, X_n)$  given  $\mathbf{X}^k \stackrel{\text{def}}{=} (X_1, X_2, \dots, X_k)$  is

$$f_{\mathbf{X}_{k+1}^n | \mathbf{X}^k}(\mathbf{x}_{k+1}^n | \mathbf{x}^k) = \frac{f_{\mathbf{X}}(x_1, x_2, \dots, x_n)}{f_{\mathbf{X}^k}(x_1, x_2, \dots, x_k)} = \frac{f_{\mathbf{X}}(\mathbf{x})}{f_{\mathbf{X}^k}(\mathbf{x}^k)}$$

- Chain rule:** We can write

$$f_{\mathbf{X}}(\mathbf{x}) = f_{X_1}(x_1) f_{X_2 | X_1}(x_2 | x_1) f_{X_3 | X_2, X_1}(x_3 | x_2, x_1) \cdots f_{X_n | \mathbf{X}^{n-1}}(x_n | \mathbf{x}^{n-1})$$

**Proof:** By definition of conditional pdf, the rule holds for  $n = 2$ . Now use induction: suppose it holds for  $n - 1$ . Then

$$\begin{aligned}
 f_{\mathbf{X}}(\mathbf{x}) &= \underbrace{f_{\mathbf{X}^{n-1}}(\mathbf{x}^{n-1})}_{\text{iterate}} f_{X_n | \mathbf{X}^{n-1}}(x_n | \mathbf{x}^{n-1}) \\
 &= f_{\mathbf{X}^{n-2}}(\mathbf{x}^{n-2}) f_{X_{n-1} | \mathbf{X}^{n-2}}(x_{n-1} | \mathbf{x}^{n-2}) f_{X_n | \mathbf{X}^{n-1}}(x_n | \mathbf{x}^{n-1}) \\
 &\vdots \\
 &= f_{X_1}(x_1) f_{X_2 | X_1}(x_2 | x_1) f_{X_3 | X_2, X_1}(x_3 | x_2, x_1) \cdots f_{X_n | \mathbf{X}^{n-1}}(x_n | \mathbf{x}^{n-1})
 \end{aligned}$$

# Independence and Conditional Independence

- Independence is defined in the usual way, e.g.,  $X_1, X_2, \dots, X_n$  are independent if

$$f_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n f_{X_i}(x_i) \text{ for all } (x_1, x_2, \dots, x_n)$$

- Important special case: **i.i.d. RVs**:  $X_1, X_2, \dots, X_n$  are said to be independent and identically distributed (i.i.d.) if they are independent and have the same marginals.  
Example: if we flip a coin  $n$  times independently, we generate i.i.d.  $\text{Bern}(p)$  RVs  $X_1, X_2, \dots, X_n$
- RVs  $X_1$  and  $X_3$  are said to be **conditionally independent** given  $X_2$  if

$$f_{X_1 X_3 | X_2}(x_1, x_3 | x_2) = f_{X_1 | X_2}(x_1 | x_2) f_{X_3 | X_2}(x_3 | x_2) \text{ for all } (x_1, x_2, x_3)$$

- Conditional independence neither implies nor is implied by independence;  $X_1$  and  $X_3$  independent given  $X_2$  does not mean that  $X_1$  and  $X_3$  are independent (or vice versa).



- Example: **Coin with random bias**. Given a coin with random bias  $P \sim f_P(p)$ , flip it  $n$  times independently to generate the RVs  $X_1, X_2, \dots, X_n$ , where  $X_i = 1$  if the  $i$ -th flip is heads,  $= 0$  otherwise. Example: if we flip a coin  $n$  times independently, we generate i.i.d.  $\text{Bern}(p)$  RVs  $X_1, X_2, \dots, X_n$ 
  - $X_1, X_2, \dots, X_n$  are **not** independent.
  - However,  $X_1, X_2, \dots, X_n$  are conditionally independent given  $P$ ; in fact, they are i.i.d.  $\text{Bern}(p)$  for every  $P = p$ .
- Example: **Additive noise channel**. Consider an additive noise channel with signal  $X$ , noise  $Z$ , and observation  $Y = X + Z$ , where  $X$  and  $Z$  are independent random variables.
  - Although  $X$  and  $Z$  are independent, they are not in general conditionally independent given  $Y$ .

## Mean and Covariance Matrix

- The mean of the random vector  $\mathbf{X}$  is defined componentwise as

$$E\{\mathbf{X}\} = \begin{bmatrix} E\{X_1\} & E\{X_2\} & \cdots & E\{X_n\} \end{bmatrix}^\top$$

- Denote the covariance between  $X_i$  and  $X_j$ ,  $\text{Cov}(X_i, X_j)$ , by  $\sigma_{ij}$  (so the variance of  $X_i$  is denoted by  $\sigma_{ii}$ ,  $\text{Var}(X_i)$ , or  $\sigma_{X_i}^2$ )
- The **covariance matrix** of  $\mathbf{X}$  is defined as

$$\Sigma_{\mathbf{X}} = E\{[\mathbf{X} - E\{\mathbf{X}\}][\mathbf{X} - E\{\mathbf{X}\}]^{\top}\} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn} \end{bmatrix}$$

- For  $n = 2$ , use the definition of correlation coefficient to obtain

$$\Sigma_{\mathbf{X}} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_{X_1}^2 & \rho_{X_1 X_2} \sigma_{X_1} \sigma_{X_2} \\ \rho_{X_1 X_2} \sigma_{X_1} \sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix}$$

## Properties of Covariance Matrix $\Sigma_X$

- $\Sigma_{\mathbf{X}}$  is **real**, and **symmetric** (i.e.,  $\Sigma_{\mathbf{X}} = \Sigma_{\mathbf{X}}^{\top}$ ) since  $\sigma_{ij} = \sigma_{ji}$ .
- $\Sigma_{\mathbf{X}}$  is **positive semidefinite**, i.e., the **quadratic form**

$$\mathbf{z}^\top \Sigma_X \mathbf{z} \geq 0 \text{ for every real vector } \mathbf{z}$$

Equivalently, all the **eigenvalues** of  $\Sigma_{\mathbf{X}}$  are nonnegative.

- To show that  $\Sigma_X$  is **positive semidefinite**, consider

$$\begin{aligned} \mathbf{z}^\top \Sigma_{\mathbf{X}} \mathbf{z} &= \mathbf{z}^\top E\{[\mathbf{X} - E\{\mathbf{X}\}][\mathbf{X} - E\{\mathbf{X}\}]^\top\} \mathbf{z} \\ &= E\{\mathbf{z}^\top [\mathbf{X} - E\{\mathbf{X}\}][\mathbf{X} - E\{\mathbf{X}\}]^\top \mathbf{z}\} \\ &= E\{(\mathbf{z}^\top [\mathbf{X} - E\{\mathbf{X}\}])^2\} \geq 0 \end{aligned}$$

## Gaussian Random Vector

- A random vector  $\mathbf{X} = (X_1, X_2, \dots, X_n)$  is a Gaussian random vector (or  $X_1, X_2, \dots, X_n$  are jointly Gaussian RVs) if the joint pdf is of the form

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\top} \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

where  $\boldsymbol{\mu} = E\{\mathbf{X}\}$ ,  $\Sigma$  is the covariance matrix of  $\mathbf{X}$  and  $|\Sigma| > 0$ , i.e.,  $\Sigma$  is positive definite. We use the notation  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$

- The above definition requires that  $|\Sigma| > 0$ . An alternative definition that does not require  $|\Sigma| > 0$  is as follows. Random vector  $\mathbf{X}$  is Gaussian, or  $X_1, X_2, \dots, X_n$  are jointly Gaussian RVs, if

$$\mathbf{c}^\top \mathbf{X} = \sum_{i=1}^n c_i X_i \sim \mathcal{N}(\mu, \sigma^2)$$

for any nonzero vector  $\mathbf{c}$ . That is, any linear combination of the components of  $\mathbf{X}$  is a scalar Gaussian random variable.

## Characteristic Function

- The **characteristic function**  $\phi_{\mathbf{X}}(\boldsymbol{\nu})$  of a random vector  $\mathbf{X} \in \mathbb{R}^n$  is defined as (the integral is  $n$ -dimensional and  $j = \sqrt{-1}$ )

$$\phi_{\mathbf{X}}(\boldsymbol{\nu}) = E\{e^{j\boldsymbol{\nu}^\top \mathbf{X}}\} = E\{e^{j\sum_{i=1}^n \nu_i X_i}\} = \int_{-\infty}^{\infty} e^{j\boldsymbol{\nu}^\top \mathbf{x}} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}$$

- In terms of multidimensional integration

$$\phi_{\mathbf{X}}(\boldsymbol{\nu}) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{j\sum_{i=1}^n \nu_i X_i} f_{\mathbf{X}}(x_1, \dots, x_n) dx_1 \cdots dx_n$$

It is  $n$ -dimensional Fourier transform of  $f_{\mathbf{X}}(\mathbf{x})$ .

- Inverse Fourier transform of  $\phi_{\mathbf{X}}(\boldsymbol{\nu})$  yields  $f_{\mathbf{X}}(\mathbf{x})$ :

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^n} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{-j\sum_{i=1}^n \nu_i X_i} \phi_{\mathbf{X}}(\nu_1, \dots, \nu_n) d\nu_1 \cdots d\nu_n$$

- The characteristic function of  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is given by

$$\phi_{\mathbf{X}}(\boldsymbol{\nu}) = e^{j\boldsymbol{\nu}^\top \boldsymbol{\mu} - \frac{1}{2}\boldsymbol{\nu}^\top \boldsymbol{\Sigma} \boldsymbol{\nu}}$$

This expression is valid whether or not  $|\boldsymbol{\Sigma}| > 0$ . Recall that  $|\boldsymbol{\Sigma}| \geq 0$

# Properties of Gaussian Random Vectors (GRV)

- **Property 1:** Uncorrelation implies independence.

This can be verified by substituting  $\sigma_{ij} = 0$  for all  $i \neq j$  in the joint pdf. Then  $\Sigma$  becomes diagonal and so does  $\Sigma^{-1}$ , and the joint pdf reduces to the product of the marginals  $X_i \sim \mathcal{N}(\mu_i, \sigma_{ii})$ .

- **Property 2:** Linear transformation of a GRV yields a GRV.

Suppose  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ ,  $\mathbf{X} \in \mathbb{R}^n$ . Define  $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b} \in \mathbb{R}^m$ . The characteristic function of  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  is given by

$$\phi_{\mathbf{X}}(\boldsymbol{\nu}) = e^{j\boldsymbol{\nu}^T \boldsymbol{\mu} - \frac{1}{2}\boldsymbol{\nu}^T \Sigma \boldsymbol{\nu}}$$

The characteristic function of  $\mathbf{Y}$  is

$$\begin{aligned} \phi_{\mathbf{Y}}(\boldsymbol{\nu}) &= E\{e^{j\boldsymbol{\nu}^T \mathbf{Y}}\} = E\{e^{j\boldsymbol{\nu}^T (\mathbf{A}\mathbf{X} + \mathbf{b})}\} \\ &= E\{e^{j(\mathbf{A}^T \boldsymbol{\nu})^T \mathbf{X}}\} e^{j\boldsymbol{\nu}^T \mathbf{b}} = e^{j(\mathbf{A}^T \boldsymbol{\nu})^T \boldsymbol{\mu} - \frac{1}{2}(\mathbf{A}^T \boldsymbol{\nu})^T \Sigma \mathbf{A}^T \boldsymbol{\nu}} e^{j\boldsymbol{\nu}^T \mathbf{b}} \\ &= e^{j\boldsymbol{\nu}^T (\mathbf{A}\boldsymbol{\mu} + \mathbf{b}) - \frac{1}{2}\boldsymbol{\nu}^T (\mathbf{A}\Sigma\mathbf{A}^T) \boldsymbol{\nu}}. \end{aligned}$$

Thus,  $\mathbf{Y} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^T)$ .

- **Property 3:** Marginals of a GRV are Gaussian, i.e., if  $\mathbf{X}$  is GRV then for any subset  $\{i_1, i_2, \dots, i_k\} \subset \{1, 2, \dots, n\}$  of indexes, the RV

$$\mathbf{Y} = [X_{i_1} \ X_{i_2} \ \dots \ X_{i_k}]^\top$$

is a GRV.

This follows from Property 2.

- The converse of Property 3 does not hold in general, i.e., Gaussian marginals do not necessarily mean that the RVs are jointly Gaussian.
- **Property 4:** Conditionals of a GRV are Gaussian, more specifically, if

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \dots \\ \mathbf{X}_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_1 \\ \dots \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \vdots & \boldsymbol{\Sigma}_{12} \\ \dots & \vdots & \dots \\ \boldsymbol{\Sigma}_{21} & \vdots & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right)$$

where  $\mathbf{X}_1$  is a  $k$ -dim RV and  $\mathbf{X}_2$  is an  $n - k$ -dim RV, then

$$\mathbf{X}_2 \Big| \{\mathbf{X}_1 = \mathbf{x}_1\} \sim \mathcal{N}(\boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})$$

## Mean-Square Estimation

- Consider two random vectors  $\mathbf{X} \in \mathbb{R}^n$  and  $\mathbf{Y} \in \mathbb{R}^m$ , where  $\mathbf{Y}$  is observed but  $\mathbf{X}$  is not. (For example,  $\mathbf{Y}$  could be noisy measurements of some function of  $\mathbf{X}$ .) Given  $\mathbf{Y}$ , we wish to estimate  $\mathbf{X}$  as  $\hat{\mathbf{X}}(\mathbf{Y}) = \mathbf{g}(\mathbf{Y})$  to minimize the mean-square error (MSE)

$$\mathbf{g}(\mathbf{Y}) = \arg \min E\{\|\mathbf{X} - \mathbf{g}(\mathbf{Y})\|^2\}$$

where

$$\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = \sum_{i=1}^n x_i^2$$

- The solution is given by

$$\mathbf{g}(\mathbf{Y}) = E\{\mathbf{X} | \mathbf{Y}\} \Rightarrow \mathbf{g}(\mathbf{y}) = E\{\mathbf{X} | \mathbf{Y} = \mathbf{y}\}$$

We have

$$E\{\|\mathbf{X} - \mathbf{g}(\mathbf{Y})\|^2\} = E_Y\{E_X\{\|\mathbf{X} - \mathbf{g}(\mathbf{Y})\|^2 | \mathbf{Y} = \mathbf{y}\}\}$$

Now for each  $\mathbf{Y} = \mathbf{y}$ ,  $E_X\{\|\mathbf{X} - \mathbf{g}(\mathbf{y})\|^2 | \mathbf{Y} = \mathbf{y}\}$  is minimized if  $\mathbf{g}(\mathbf{y}) = E\{\mathbf{X} | \mathbf{Y} = \mathbf{y}\}$



- To establish the claim, consider the first-order optimality condition:

$$E_X\{\|\mathbf{X} - \mathbf{g}(\mathbf{y})\|^2 \mid \mathbf{Y} = \mathbf{y}\} = \int \sum_{i=1}^n [x_i - g_i(\mathbf{y})]^2 f_{X|Y}(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

$$0 = \frac{\partial ()}{\partial g_\ell(\mathbf{y})} = -2 \int [\mathbf{x}_\ell - g_\ell(\mathbf{y})] f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x}, \quad \ell = 1, 2, \dots, n$$

$$\Rightarrow g_\ell(\mathbf{y}) = E\{X_\ell \mid \mathbf{Y} = \mathbf{y}\} \Rightarrow \mathbf{g}(\mathbf{y}) = E\{\mathbf{X} \mid \mathbf{Y} = \mathbf{y}\}$$

- Thus,

$$E\{\|\mathbf{X} - \mathbf{g}(\mathbf{Y})\|^2\} = \int E_X\{\|\mathbf{X} - \mathbf{g}(\mathbf{y})\|^2 \mid \mathbf{Y} = \mathbf{y}\} f_Y(\mathbf{y}) d\mathbf{y}$$

is minimized for  $\mathbf{g}(\mathbf{Y}) = E\{\mathbf{X} | \mathbf{Y}\}$

- Thus  $E\{\mathbf{X} | \mathbf{Y}\}$  minimizes the MSE conditioned on every  $\mathbf{Y} = \mathbf{y}$  and not just its average over  $\mathbf{Y}$  !

## Jointly Gaussian $\mathbf{X}$ and $\mathbf{Y}$

- Suppose  $(\mathbf{X} \in \mathbb{R}^n$  and  $\mathbf{Y} \in \mathbb{R}^m)$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Y} \\ \dots \\ \mathbf{X} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_Y \\ \dots \\ \boldsymbol{\mu}_X \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \vdots & \boldsymbol{\Sigma}_{YX} \\ \dots & \vdots & \dots \\ \boldsymbol{\Sigma}_{XY} & \vdots & \boldsymbol{\Sigma}_{XX} \end{bmatrix} \right)$$

- Recall that

$$\mathbf{X} \Big| \{ \mathbf{Y} = \mathbf{y} \} \sim \mathcal{N} \left( \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} (\mathbf{y} - \boldsymbol{\mu}_Y) + \boldsymbol{\mu}_X, \boldsymbol{\Sigma}_{XX} - \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{YX} \right)$$

- Therefore, the optimal estimator is

$$E\{\mathbf{X} \mid \mathbf{Y}\} = \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} (\mathbf{Y} - \boldsymbol{\mu}_Y) + \boldsymbol{\mu}_X$$

and the corresponding minimum MSE (MMSE)

$$E\{\|\mathbf{X} - E\{\mathbf{X} \mid \mathbf{Y}\}\|^2\}$$

$$\text{tr}(\boldsymbol{\Sigma}_{XX} - \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{YX})$$

**Example 8.21 (Gubner)** If scalar signal  $X \sim \mathcal{N}(0, 1)$  and noise  $W \sim \mathcal{N}(0, \sigma^2)$ , find the MMSE estimator of  $X$  given noisy observation  $Y = X + W$ . The signal and noise are independent.

We have  $\Sigma_{XY} = \Sigma_{XX} = 1$  and  $\Sigma_{YY} = \Sigma_{XX} + \Sigma_{WW} = 1 + \sigma^2$  since  $\Sigma_{XW} = 0$ . Also  $\mu_X = \mu_Y = 0$ . Therefore,

$$\begin{aligned} E\{X | Y = y\} &= \Sigma_{XY} \Sigma_{YY}^{-1} (y - \mu_Y) + \mu_X \\ &= \frac{1}{1 + \sigma^2} (y - 0) + 0 = \frac{y}{1 + \sigma^2} \end{aligned}$$

The corresponding MMSE is

$$\Sigma_{XX} - \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX} = 1 - \frac{1}{1 + \sigma^2} = \frac{\sigma^2}{1 + \sigma^2}$$

## Linear Mean-Square Estimation

- Consider two random vectors  $\mathbf{X} \in \mathbb{R}^n$  and  $\mathbf{Y} \in \mathbb{R}^m$ , where  $\mathbf{Y}$  is observed but  $\mathbf{X}$  is not. Given  $\mathbf{Y}$ , we wish to estimate  $\mathbf{X}$  as  $\hat{\mathbf{X}}(\mathbf{Y}) = \mathbf{A}\mathbf{Y} + \mathbf{b}$  to minimize the mean-square error (MSE)

$$\{\hat{\mathbf{A}}, \hat{\mathbf{b}}\} = \arg \min_{\mathbf{A}, \mathbf{b}} E\{\|\mathbf{X} - \hat{\mathbf{X}}(\mathbf{Y})\|^2\}$$

where  $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = \sum_{i=1}^n x_i^2$ .

- Let

$$E \left\{ \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \right\} = \begin{bmatrix} \mathbf{m}_X \\ \mathbf{m}_Y \end{bmatrix}, \text{cov} \left( \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \right) = \begin{bmatrix} \text{cov}(\mathbf{X}, \mathbf{X}) & \text{cov}(\mathbf{X}, \mathbf{Y}) \\ \text{cov}(\mathbf{Y}, \mathbf{X}) & \text{cov}(\mathbf{Y}, \mathbf{Y}) \end{bmatrix}$$

- Define  $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{m}_{\mathbf{X}}$  and  $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{m}_{\mathbf{Y}}$ . Then  $E\{\tilde{\mathbf{X}}\} = \mathbf{0} = E\{\tilde{\mathbf{Y}}\}$ ,  $\text{cov}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) = \text{cov}(\mathbf{X}, \mathbf{X})$ ,  $\text{cov}(\tilde{\mathbf{Y}}, \tilde{\mathbf{Y}}) = \text{cov}(\mathbf{Y}, \mathbf{Y})$ , and  $\text{cov}(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}) = \text{cov}(\mathbf{X}, \mathbf{Y}) = E\{[\mathbf{X} - \mathbf{m}_{\mathbf{X}}][\mathbf{Y} - \mathbf{m}_{\mathbf{Y}}]^{\top}\} = (\text{cov}(\mathbf{Y}, \mathbf{X}))^{\top}$ .

- Set  $\mathbf{X} = \tilde{\mathbf{X}} + \mathbf{m}_X$  etc, and rewrite

$$\begin{aligned} \text{MSE} &= E\{\|\mathbf{X} - (\mathbf{A}\mathbf{Y} + \mathbf{b})\|^2\} = E\{[\mathbf{X} - (\mathbf{A}\mathbf{Y} + \mathbf{b})]^\top [\mathbf{X} - (\mathbf{A}\mathbf{Y} + \mathbf{b})]\} \\ &= E\{\|\tilde{\mathbf{X}} + \mathbf{m}_X - \mathbf{A}(\tilde{\mathbf{Y}} + \mathbf{m}_Y) - \mathbf{b}\|^2\} \\ &= E\{\|\tilde{\mathbf{X}} - \mathbf{A}\tilde{\mathbf{Y}}\|^2\} + \|\mathbf{m}_X - \mathbf{A}\mathbf{m}_Y - \mathbf{b}\|^2 \\ &\quad + 2 \underbrace{E\{(\tilde{\mathbf{X}} - \mathbf{A}\tilde{\mathbf{Y}})^\top\}}_{=0} (\mathbf{m}_X - \mathbf{A}\mathbf{m}_Y - \mathbf{b}) \end{aligned}$$

- Therefore, to minimize MSE, we must estimate  $\mathbf{b}$  as

$$\hat{\mathbf{b}} = \mathbf{m}_X - \mathbf{A}\mathbf{m}_Y = \mathbf{m}_X - \hat{\mathbf{A}}\mathbf{m}_Y$$

which minimizes the second term  $\|\mathbf{m}_X - \mathbf{A}\mathbf{m}_Y - \mathbf{b}\|^2$  whatever the choice of  $\hat{\mathbf{A}}$ , and it does not affect  $E\{\|\tilde{\mathbf{X}} - \mathbf{A}\tilde{\mathbf{Y}}\|^2\}$ .

- Now consider minimization of  $E\{\|\tilde{\mathbf{X}} - \mathbf{A}\tilde{\mathbf{Y}}\|^2\}$  w.r.t.  $\mathbf{A}$ . Using  $\Sigma_{XY} = \text{cov}(\mathbf{X}, \mathbf{Y})$ , rewrite (explained in next slide)

$$E\{\|\tilde{\mathbf{X}} - \mathbf{A}\tilde{\mathbf{Y}}\|^2\} = \text{tr}(\Sigma_{XX}) + \text{tr}(\Sigma_{YY}\mathbf{A}^\top\mathbf{A}) - \text{tr}(\Sigma_{YX}\mathbf{A}) - \text{tr}(\Sigma_{XY}\mathbf{A}^\top)$$



# Matrix Calculus

- The **gradient**  $\nabla f(\mathbf{x})$  of the differentiable scalar function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f(\mathbf{x})$ , is a column vector  $\in \mathbb{R}^n$ , given by

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix}. \quad (2)$$

- $f(\mathbf{x}) = \mathbf{b}^\top \mathbf{x} = \sum_{i=1}^n b_i x_i$ . Then  $\nabla f(\mathbf{x}) = \mathbf{b}$ . Since  $\mathbf{b}^\top \mathbf{x} = \mathbf{x}^\top \mathbf{b}$ ,  $\nabla \mathbf{x}^\top \mathbf{b} = \mathbf{b}$ .
- $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n x_i A_{ij} x_j$ . Then  $\nabla f(\mathbf{x}) = \mathbf{A} \mathbf{x} + \mathbf{A}^\top \mathbf{x}$ . This follows from

$$\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial x_\ell} = \sum_{j=1}^n A_{\ell j} x_j + \sum_{i=1}^n x_i A_{i\ell} = \sum_{j=1}^n A_{\ell j} x_j + \sum_{j=1}^n A_{\ell i}^\top x_j = [\mathbf{A} \mathbf{x} + \mathbf{A}^\top \mathbf{x}]_\ell.$$

- Similarly, define **gradient**  $\nabla f(X)$  of the differentiable scalar function  $f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ ,  $f(X)$ , as a matrix  $\in \mathbb{R}^{n \times m}$ , whose  $(i, j)$ th element is given by

$$[\nabla f(X)]_{ij} = \frac{\partial f(X)}{\partial X_{ij}}. \quad (3)$$

We will also write  $\nabla f(X)$  as  $\nabla_X f(X)$  and  $\frac{\partial f(X)}{\partial X}$

- **Some useful results:**  $A$ ,  $X$  and  $B$  below are all matrices.
  - $\frac{\partial \text{tr}(AXB)}{\partial X} = A^\top B^\top$ .
  - $\frac{\partial \text{tr}(AX^\top B)}{\partial X} = BA$
  - $\frac{\partial \text{tr}(X^\top AXB)}{\partial X} = 2AXB$  if  $A = A^\top$  and  $B = B^\top$ .
  - $\frac{\partial \text{tr}(XBX^\top A)}{\partial X} = 2AXB$  if  $A = A^\top$  and  $B = B^\top$ .



○○○○

- 

$$C = \text{tr}(\Sigma_{XX}) + \text{tr}(\Sigma_{YY} \mathbf{A}^\top \mathbf{A}) - \text{tr}(\Sigma_{YX} \mathbf{A}) - \text{tr}(\Sigma_{XY} \mathbf{A}^\top)$$

w.r.t.  $\mathbf{A}$ , we set

$$\begin{aligned} \mathbf{0} &= \frac{\partial C}{\partial \mathbf{A}} \\ &= \mathbf{0} + 2\mathbf{A}\Sigma_{YY} - \Sigma_{YX}^\top - \Sigma_{XY} \\ &= 2\mathbf{A}\Sigma_{YY} - 2\Sigma_{XY} \end{aligned}$$

- 

$$\hat{\mathbf{A}}\Sigma_{YY} = \Sigma_{XY} \Rightarrow \hat{\mathbf{A}} = \Sigma_{XY}\Sigma_{YY}^{-1} \text{ if the inverse exists.}$$

- Thus, the solution to  $\{\hat{\mathbf{A}}, \hat{\mathbf{b}}\} = \arg \min_{\mathbf{A}, \mathbf{b}} E\{\|\mathbf{X} - (\mathbf{A}\mathbf{Y} + \mathbf{b})\|^2\}$  is given by  $\hat{\mathbf{A}} = \Sigma_{XY}\Sigma_{YY}^{-1}$  and  $\hat{\mathbf{b}} = \mathbf{m}_X - \hat{\mathbf{A}}\mathbf{m}_Y$ , leading to

$$\hat{\mathbf{X}} = \hat{\mathbf{A}}\mathbf{Y} + \hat{\mathbf{b}} = \hat{\mathbf{A}}(\mathbf{Y} - \mathbf{m}_Y) + \mathbf{m}_X$$

- The optimal MSE is

$$\begin{aligned} \text{MSE}_o &= E\{\|\mathbf{X} - (\hat{\mathbf{A}}\mathbf{Y} + \hat{\mathbf{b}})\|^2\} = E\{\|\tilde{\mathbf{X}} - \hat{\mathbf{A}}\tilde{\mathbf{Y}}\|^2\} \\ &= \text{tr}(\Sigma_{XX}) + \text{tr}(\Sigma_{YY}\hat{\mathbf{A}}^\top\hat{\mathbf{A}}) - \text{tr}(\Sigma_{YX}\hat{\mathbf{A}}) - \text{tr}(\Sigma_{XY}\hat{\mathbf{A}}^\top) \\ &= \text{tr}(\Sigma_{XX}) - \text{tr}(\Sigma_{YX}\hat{\mathbf{A}}) = \text{tr}(\Sigma_{XX}) - \text{tr}(\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX}) \end{aligned}$$

where we have used

$$\text{tr}(\Sigma_{YY}\hat{\mathbf{A}}^\top\hat{\mathbf{A}}) - \text{tr}(\Sigma_{XY}\hat{\mathbf{A}}^\top) = \text{tr}((\hat{\mathbf{A}}\Sigma_{YY} - \Sigma_{XY})\hat{\mathbf{A}}^\top) = 0$$

## Example

Let  $X$  be the random variable representing a signal with mean  $m_X$  and variance  $P$ . The observations are  $Y_i = X + Z_i$ , for  $i = 1, 2, \dots, n$ , where the  $Z_i$ s are zero-mean, uncorrelated noise with variance  $\sigma^2$ , and  $X$  and  $Z_i$ s are also uncorrelated. Find the MMSE linear estimate of  $X$  given  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_n]^\top$ , and its MSE.

**Solution:** We have  $\text{cov}(X, Y_i) = \text{cov}(X, X) = P$ , and

$$\begin{aligned} \text{cov}(Y_i, Y_j) &= \text{cov}(X + Z_i, X + Z_j) = \text{cov}(X, X) + \text{cov}(Z_i, Z_j) + 0 \\ &= \begin{cases} P + \sigma^2 & \text{if } i = j \\ P & \text{if } i \neq j \end{cases} \end{aligned}$$

Thus,  $\Sigma_{XX} = P$ ,  $\Sigma_{XY} = P\mathbf{1}_n^\top$ , and  $\Sigma_{YY} = P\mathbf{1}_n\mathbf{1}_n^\top + \sigma^2\mathbf{I}_n$  where  $\mathbf{I}_n$  denotes the  $n \times n$  identity matrix and  $\mathbf{1}_n$  denotes the  $n$ -dimensional column vector of all ones. The matrix inversion lemma states

$$\begin{aligned} (A + BCD)^{-1} &= A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \\ \Rightarrow \Sigma_{YY}^{-1} &= (\sigma^2\mathbf{I}_n + \mathbf{1}_nP\mathbf{1}_n^\top)^{-1} = \frac{1}{\sigma^2} \left( \mathbf{I}_n - \frac{P}{\sigma^2 + nP} \mathbf{1}_n\mathbf{1}_n^\top \right) \end{aligned}$$

We have  $\hat{\mathbf{A}} = \Sigma_{XY} \Sigma_{YY}^{-1}$  leading to

$$\begin{aligned}\hat{\mathbf{A}}\mathbf{Y} &= \frac{P}{\sigma^2} \mathbf{1}_n^\top \left( \mathbf{I}_n - \frac{P}{\sigma^2 + nP} \mathbf{1}_n \mathbf{1}_n^\top \right) \mathbf{Y} \\ &= \frac{P}{\sigma^2} \left( \mathbf{1}_n^\top - \frac{P}{\sigma^2 + nP} n \mathbf{1}_n^\top \right) \mathbf{Y} \\ &= \frac{P}{\sigma^2} \left( 1 - \frac{nP}{\sigma^2 + nP} \right) \sum_{i=1}^n Y_i = \frac{P}{\sigma^2 + nP} \sum_{i=1}^n Y_i\end{aligned}$$

$$\begin{aligned}\Rightarrow \hat{X} &= \hat{\mathbf{A}}\tilde{\mathbf{Y}} + m_X = \frac{P}{\sigma^2 + nP} \sum_{i=1}^n (Y_i - m_X) + m_X \\ &= \frac{P}{\sigma^2 + nP} \sum_{i=1}^n Y_i + \frac{\sigma^2}{\sigma^2 + nP} m_X\end{aligned}$$

Optimal MSE is (you work the details!)

$$\text{MSE}_o = \frac{P\sigma^2}{\sigma^2 + nP}$$