

# Chapter 3: Population mean & covariance and Multivariate normal distribution

Zhixiang ZHANG

University of Macau

zhixzhang@um.edu.mo

January 16, 2025

# Table of Contents

1 Population mean and covariance matrices

2 Multivariate Normal distribution

# Random samples from a distribution

- In order to study the sampling variability of statistics like  $\bar{\mathbf{x}}$  and  $\mathbf{S}_n$  with the ultimate aim of making inferences, we need to make assumptions about the variables whose observed values constitute the data set  $\mathbf{X}$ .
- Suppose, then, that the data have not yet been observed, but we intend to collect  $n$  sets of measurements on  $p$  variables. Consequently, we treat them as random variables.
- A random sample can now be defined. If the row vectors  $\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_n^\top$  in  $\mathbf{X} \in \mathbb{R}^{n \times p}$  represent independent observations drawn from a common joint distribution with density function  $f(\mathbf{x}) = f(x_1, x_2, \dots, x_p)$ , then  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  are said to form a *random sample* from  $f(\mathbf{x})$ .

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix}$$

$$x_i \sim f(\vec{x})$$

$$(f(\mathbf{x}))$$

# Expectation of random matrices(vectors)

- A random matrix/vector is a matrix/vector whose elements are random variables.
- Let  $\mathbf{X} = \{X_{ij}\}$  be an  $n \times p$  random matrix. Then the expected value of  $\mathbf{X}$ , denoted by  $E(\mathbf{X})$ , is the  $n \times p$  matrix of numbers (if they exist)

$$E(\mathbf{X}) = \begin{bmatrix} E(X_{11}) & E(X_{12}) & \cdots & E(X_{1p}) \\ E(X_{21}) & E(X_{22}) & \cdots & E(X_{2p}) \\ \vdots & \vdots & \ddots & \vdots \\ E(X_{n1}) & E(X_{n2}) & \cdots & E(X_{np}) \end{bmatrix}.$$

- Let  $\mathbf{X}$  and  $\mathbf{Y}$  be random matrices of the same dimension, and let  $\mathbf{A}$  and  $\mathbf{B}$  be conformable matrices of constants. Then

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}) \quad \begin{matrix} E[\text{row } A \times X \times \text{column } B] \\ = E[a_i^T X b_j] \end{matrix}$$

*A is fixed matrix*

$$E(\mathbf{AXB}) = \mathbf{A}E(\mathbf{X})\mathbf{B}$$

$$E(\mathbf{AXX}^T\mathbf{A}^T) = \mathbf{A}E(\mathbf{XX}^T)\mathbf{A}^T \neq \mathbf{A}E(\mathbf{X})E(\mathbf{X}^T)\mathbf{A}^T \quad \text{because } \mathbf{X} \text{ and } \mathbf{X}^T \text{ is correlated}$$

$$E(Z^2) = (E(Z))^2 \Leftrightarrow \text{Var}(Z) = 0$$

# Population mean and variance

The marginal mean and variance of a random vector  $\mathbf{X}^\top = [X_1, X_2, \dots, X_p]$ :

$$\mu_i = \begin{cases} \int_{-\infty}^{\infty} x_i f_i(x_i) dx_i & \text{if } X_i \text{ is a continuous random variable with probability} \\ & \text{density function } f_i(x_i) \\ \sum_{\text{all } x_i} x_i p_i(x_i) & \text{if } X_i \text{ is a discrete random variable with probability} \\ & \text{function } p_i(x_i) \end{cases}$$

$$\sigma_i^2 = \begin{cases} \int_{-\infty}^{\infty} (x_i - \mu_i)^2 f_i(x_i) dx_i & \text{if } X_i \text{ is a continuous random variable} \\ & \text{with probability density function } f_i(x_i) \\ \sum_{\text{all } x_i} (x_i - \mu_i)^2 p_i(x_i) & \text{if } X_i \text{ is a discrete random variable} \\ & \text{with probability function } p_i(x_i) \end{cases}$$

It will be convenient later to denote the marginal variances by  $\sigma_{ii}$ .

# Population covariance

The behavior of any pair of random variables, such as  $X_i$  and  $X_k$  is described by their joint probability function, and a measure of the linear association between them is provided by the covariance:

$$\sigma_{ik} = E(X_i - \mu_i)(X_k - \mu_k)$$

$$= \begin{cases} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_i - \mu_i)(x_k - \mu_k) f_{ik}(x_i, x_k) dx_i dx_k & \text{if } X_i, X_k \text{ are continuous} \\ & \text{random variables with} \\ & \text{the joint density} \\ & \text{function } f_{ik}(x_i, x_k) \\ \sum_{\text{all } x_i} \sum_{\text{all } x_k} (x_i - \mu_i)(x_k - \mu_k) p_{ik}(x_i, x_k) & \text{if } X_i, X_k \text{ are discrete} \\ & \text{random variable with} \\ & \text{joint probability} \\ & \text{function } p_{ik}(x_i, x_k) \end{cases}$$

# Independence

- More generally, the collective behavior of the  $p$  random variables  $X_1, X_2, \dots, X_p$  is described by a joint probability density function  $f(x_1, x_2, \dots, x_p) = f(\mathbf{x})$ . ( $f(\mathbf{x})$  will often be the multivariate normal density function in this course)
- If

$$P[X_i \leq x_i \text{ and } X_k \leq x_k] = P[X_i \leq x_i] P[X_k \leq x_k]$$

for all pairs of values  $x_i, x_k$ , then  $X_i$  and  $X_k$  are said to be *independent*. When  $X_i$  and  $X_k$  are continuous random variables with joint density  $f_{ik}(x_i, x_k)$  and marginal densities  $f_i(x_i)$  and  $f_k(x_k)$ , the independence condition becomes

$$f_{ik}(x_i, x_k) = f_i(x_i) f_k(x_k)$$

for all pairs  $(x_i, x_k)$ .

# Independence

$X_i \in \mathbb{R}^{4 \times 1}$   
 $i=1 \text{ and } 2.$

two variables  
 $X_1 \quad X_2$   
 $\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{pmatrix}$

- The  $p$  continuous random variables  $X_1, X_2, \dots, X_p$  are *mutually independent* if their joint density can be factored as

$$f_{12\dots p}(x_1, x_2, \dots, x_p) = f_1(x_1) f_2(x_2) \cdots f_p(x_p)$$

for all  $p$ -tuples  $(x_1, x_2, \dots, x_p)$ .

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

- Independence of  $X_i$  and  $X_k$  implies that  $\text{Cov}(X_i, X_k) = 0$ .

However, the converse of is not true in general.

$$X_1 = Z \sim N(0,1) \\ X_2 = Z^2 \\ \text{Cov}(X_1, X_2) = 0 \Rightarrow X_1 \perp\!\!\!\perp X_2$$

only linear independence

more comprehensive independence

$$\text{if } X \perp\!\!\!\perp Y \Rightarrow E(g(X)) = E(g(X)|Y=y) = E(g(X))$$

$$E(g(X)h(Y)) = E(g(X))E(h(Y))$$

$$E(g(X)h(Y)) = E(g(X)|Y=y)h(y) = E(g(X))h(y) = E(g(X))E(h(Y))$$



# Population mean vec and population covariance mat

We shall refer to  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  below as the *population mean* (vector) and *population covariance* (matrix), respectively:

$$\boldsymbol{\mu} = E(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}$$

$$\begin{aligned} \boldsymbol{\Sigma} &= E(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top \\ &= E \left( \begin{bmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \\ \vdots \\ X_p - \mu_p \end{bmatrix} [X_1 - \mu_1, X_2 - \mu_2, \dots, X_p - \mu_p] \right) \\ &= (\sigma_{ij})_{1 \leq i, j \leq p}. \end{aligned}$$

Here and below in this Chapter  $\mathbf{X}$  refers to a random vector.

# Unbiasedness

Handwritten derivations showing the unbiasedness of sample mean and covariance estimators.

Left side (Sample Mean):

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

$$= \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu + \mu - \bar{x})(x_i - \mu + \mu - \bar{x})^T$$

$$= \sum$$

Right side (Sample Covariance):

Let  $X_1, X_2, \dots, X_p$  be the components of  $x_i$ .

The sample covariance matrix is defined as:

$$\text{Cov}(X_1) \in \mathbb{R}^{p \times p}$$

$$= \text{Cov} \begin{pmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1p} \end{pmatrix}$$

It is noted that:

$$\text{Cov}(X_1) = \text{Cov}(X_2) = \dots = \text{Cov}(X_n)$$

If  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are random samples drawn from a distribution with mean  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , sample mean  $\bar{\mathbf{x}}$  and sample covariance

$$\frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

are unbiased estimators of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , respectively.

# Correlation matrix

The *population correlation matrix*

$$\boldsymbol{\rho} = \begin{bmatrix} \frac{\sigma_{11}}{\sqrt{\sigma_{11}}\sqrt{\sigma_{11}}} & \frac{\sigma_{12}}{\sqrt{\sigma_{11}}\sqrt{\sigma_{22}}} & \cdots & \frac{\sigma_{1p}}{\sqrt{\sigma_{11}}\sqrt{\sigma_{pp}}} \\ \frac{\sigma_{12}}{\sqrt{\sigma_{11}}\sqrt{\sigma_{22}}} & \frac{\sigma_{22}}{\sqrt{\sigma_{22}}\sqrt{\sigma_{22}}} & \cdots & \frac{\sigma_{2p}}{\sqrt{\sigma_{22}}\sqrt{\sigma_{pp}}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_{1p}}{\sqrt{\sigma_{11}}\sqrt{\sigma_{pp}}} & \frac{\sigma_{2p}}{\sqrt{\sigma_{22}}\sqrt{\sigma_{pp}}} & \cdots & \frac{\sigma_{pp}}{\sqrt{\sigma_{pp}}\sqrt{\sigma_{pp}}} \end{bmatrix} = \mathbf{D}_{\sigma}^{-1/2} \mathbf{\Sigma} \mathbf{D}_{\sigma}^{-1/2},$$

where  $\mathbf{D}_{\sigma} = \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$ .

Correspondingly, the sample correlation matrix can be represented by  $\mathbf{D}_s^{-1/2} \mathbf{S} \mathbf{D}_s^{-1/2}$  where  $\mathbf{D}_s$  is the diagonal matrix of  $\mathbf{S}$ .

# Linear combination of random vector

random vector  
 $\vec{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^p$   
 $E\vec{X} = \mu, \text{Cov}(\vec{X}) = \Sigma$   
 Assume  $x_1, \dots, x_n$  are independent  
 and have same distributions as  
 that of  $X$   
 for  $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i, S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$   
 we have  $E\bar{X} = \mu, ES = \Sigma$

- The linear combination  $\mathbf{c}^\top \mathbf{X} = c_1 X_1 + \dots + c_p X_p$  has

$$\text{mean} = E(\mathbf{c}^\top \mathbf{X}) = \mathbf{c}^\top \boldsymbol{\mu}$$

$$\text{variance} = \text{Var}(\mathbf{c}^\top \mathbf{X}) = \mathbf{c}^\top \boldsymbol{\Sigma} \mathbf{c}$$

where  $\boldsymbol{\mu} = E(\mathbf{X})$  and  $\boldsymbol{\Sigma} = \text{Cov}(\mathbf{X})$ .

- Given  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and the sample mean  $\bar{\mathbf{x}}$  and sample covariance matrices  $\mathbf{S}$ . What are the sample mean and sample variance of  $\mathbf{c}^\top \mathbf{x}_1, \dots, \mathbf{c}^\top \mathbf{x}_n$ ?

•  $X \in \mathbb{R}^p$   $c \in \mathbb{R}^p$  fixed find  $E c^T X$   $\text{Var}(c^T X)$  Assume  $E(X) = \mu$  and  $\text{Cov}(X) = \Sigma$

↓  
r.v.

$$E c^T X = c^T E(X) = c^T \mu$$

↓  
 $c$  is not random

$$\text{Var}(c^T X) = E(c^T X - E c^T X)(c^T X - E c^T X) = E(c^T X - E c^T X)(X^T c - E X^T c) = c^T \{E(X - EX)(X - EX)^T\} c = c^T \Sigma c = \|c\|_2^2 \Sigma$$

by above,  $E(c^T X_1) = c^T E(X_1) = c^T \bar{x}$

$$\text{Var}(c^T X_1) = c^T \text{Var}(X_1) c = c^T S c = \|c\|_2^2 S$$

for  $i=1$ , similar for  $i > 1$

$$c^T X = \sum_{i=1}^n c_i x_i$$

↓ several combination

$$CX = \begin{pmatrix} c_1^T X \\ \vdots \\ c_n^T X \end{pmatrix}$$

$$E CX = C E X = C \mu$$

$$\text{Cov}(CX) = C \Sigma C^T$$

result

# Linear combination of random vector

In general, consider the  $q$  linear combinations of the  $p$  random variables  $X_1, \dots, X_p$  :

$$\mathbf{W} = \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_q \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1p} \\ c_{21} & c_{22} & \cdots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{q1} & c_{q2} & \cdots & c_{qp} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = \mathbf{C}\mathbf{X}.$$

We have

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{W}} &= E(\mathbf{W}) = E(\mathbf{C}\mathbf{X}) = \mathbf{C}\boldsymbol{\mu}_{\mathbf{X}}, \\ \boldsymbol{\Sigma}_{\mathbf{W}} &= \text{Cov}(\mathbf{W}) = \text{Cov}(\mathbf{C}\mathbf{X}) = \mathbf{C}\boldsymbol{\Sigma}_{\mathbf{X}}\mathbf{C}^{\top}, \end{aligned} \tag{1}$$

where  $\boldsymbol{\mu}_{\mathbf{X}}$  and  $\boldsymbol{\Sigma}_{\mathbf{X}}$  are the mean vector and variance-covariance matrix of  $\mathbf{X}$ .

# Example

Let  $\mathbf{X} = [X_1, X_2]^\top$  be a random vector with mean vector  $\boldsymbol{\mu}_{\mathbf{X}} = [\mu_1, \mu_2]^\top$  and variance-covariance matrix

$$\boldsymbol{\Sigma}_{\mathbf{X}} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix}$$

Find the covariance matrix for  $\mathbf{Z} = [Z_1, Z_2]^\top$  where

$$Z_1 = X_1 - X_2, \quad Z_2 = X_1 + X_2.$$

Solution: We can write

$$\mathbf{Z} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} := \mathbf{C}\mathbf{X}.$$

Therefore, the covariance matrix is given by

$$\begin{aligned} \text{Cov}(\mathbf{Z}) &= \mathbf{C}\mathbf{\Sigma}_\mathbf{X}\mathbf{C}^\top = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} \sigma_{11} - 2\sigma_{12} + \sigma_{22} & \sigma_{11} - \sigma_{22} \\ \sigma_{11} - \sigma_{22} & \sigma_{11} + 2\sigma_{12} + \sigma_{22} \end{bmatrix} \end{aligned}$$



# Cross-covariance

where  $E(E(Y)^T) = E(Y)^T$

$$E\left[\left(\mathbf{X}\mathbf{Y}^T - \mathbf{X}E(\mathbf{Y})^T - E(\mathbf{X})\mathbf{Y}^T + E(\mathbf{X})E(\mathbf{Y})^T\right)\right]$$

$$= E(\mathbf{X}\mathbf{Y}^T) - E(\mathbf{X})E(\mathbf{Y})^T$$

- Cross-covariance of random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  equals

$$\begin{aligned}\text{Cov}(\mathbf{X}, \mathbf{Y}) &= E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T] \\ &= E[\mathbf{X}\mathbf{Y}^T] - E[\mathbf{X}]E[\mathbf{Y}]^T.\end{aligned}$$

- Let  $\mathbf{A}$ ,  $\mathbf{B}$  be conformable matrices of constants, then

$$\text{Cov}(\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{Y}) = \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}^T.$$

- $\text{Cov}(\mathbf{C}\mathbf{X}) = \mathbf{C}\Sigma_X\mathbf{C}^T$  is a special case of the above.

# Table of Contents

- 1 Population mean and covariance matrices
- 2 Multivariate Normal distribution

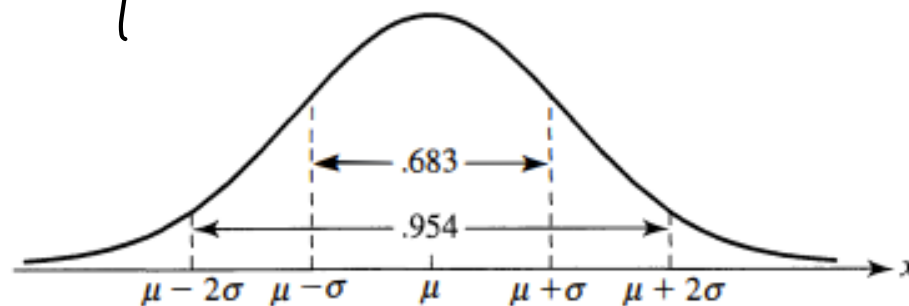
# (Univariate) normal distribution

Univariate Gaussian (normal) density:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}} (\sigma^2)^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu) \frac{1}{\sigma^2} (x-\mu)}$$

$p(X=a) = 0$  in  $f_X(x)$  is continuous or  $X \sim \mathcal{N}(\mu, \sigma^2)$



but here  $P(\mu - \sigma \leq X \leq \mu + \sigma) = 0.683$

**Figure 1:** A normal density with mean  $\mu$  and variance  $\sigma^2$  and selected areas under the curve  
 which means the cumulative  $P$  is higher around  $\mu$

# Multivariate normal density

## Definition 1

A random vector  $\mathbf{X}$  has a Multivariate Normal (MVN) distribution if it has a joint probability density function (pdf) of the form

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

$\downarrow$   
 determinant

mean and covariance

We denote this  $p$ -dimensional normal density by  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , which is analogous to the normal density in the univariate case.

## Special cases: diagonal $\Sigma$

For p-variate case, if  $\Sigma$  is diagonal, then

$$\Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma_{11}} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \frac{1}{\sigma_{pp}} \end{bmatrix} \quad \text{and} \quad \underline{|\Sigma| = (\sigma_{11}) (\sigma_{22}) \cdots (\sigma_{pp})}$$

$$\begin{aligned} \text{so } f_{\mathbf{x}}(\mathbf{x}) &= \frac{1}{(2\pi)^{\frac{p}{2}} \sqrt{(\sigma_{11}) \cdots (\sigma_{pp})}} \\ &\quad \times \exp \left\{ -\frac{1}{2} \frac{(x_1 - \mu_1)^2}{\sigma_{11}} - \cdots - \frac{1}{2} \frac{(x_p - \mu_p)^2}{\sigma_{pp}} \right\} \\ &= f_{x_1}(x_1) \cdot f_{x_2}(x_2) \cdots f_{x_p}(x_p). \end{aligned}$$

*they are independent to each other*

# Bivariate case

- For bivariate case,
 
$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N_2 \left[ \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11}^2 & \rho_{12} \sigma_{11} \sigma_{22} \\ \rho_{12} \sigma_{11} \sigma_{22} & \sigma_{22}^2 \end{pmatrix} \right].$$

- If

$$\rho_{12} = 0, f_{\mathbf{x}} \left( \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) = f_{x_1}(x_1) \cdot f_{x_2}(x_2).$$

$$B \sim N(\mu, \sigma^2)$$

# Shape of the density

Two bivariate distributions are shown below.

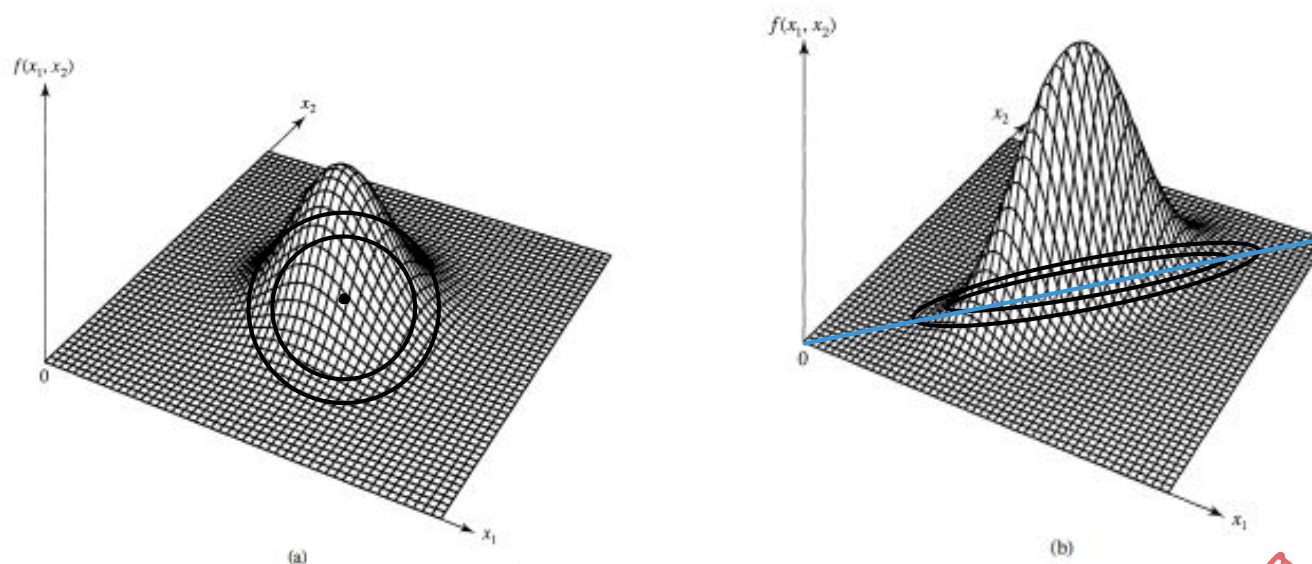


Figure 2: Two bivariate normal distributions. (a)  $\sigma_{11} = \sigma_{22}$  and  $\rho_{12} = 0$ . (b)  $\sigma_{11} = \sigma_{22}$  and  $\rho_{12} = .75$ .

Handwritten notes:

$$|\Sigma_1|^{\frac{1}{2}} \neq |\Sigma_2|^{\frac{1}{2}}$$

$$\Sigma_1 = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$$

Annotations:

- peak height (pointing to the peak of the distribution in plot b)
- $\therefore f_2(x)_{\max} > f_1(x)_{\max}$
- $\therefore |\Sigma_2| = \sigma_{11}\sigma_{22} - \sigma_{12}^2 < |\Sigma_1|$

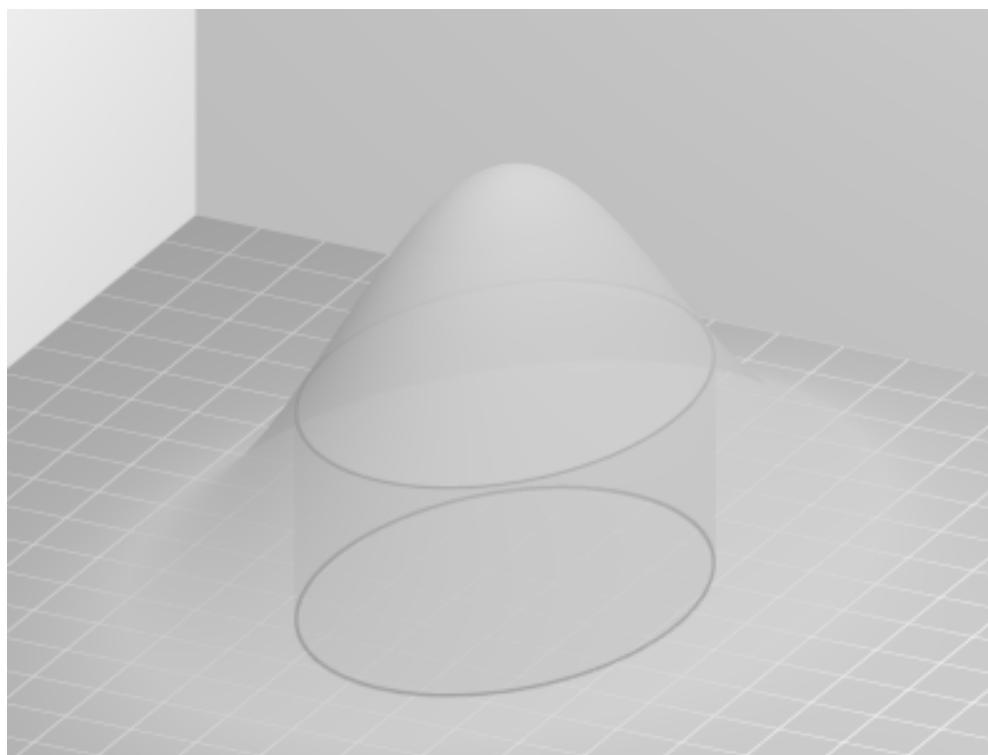
# Constant density contour

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{(\mathbf{x}-\mu)^\top \Sigma^{-1} (\mathbf{x}-\mu)}{2}} = d$$

Handwritten notes:  $(\mathbf{x}-\mu)^\top \Sigma^{-1} (\mathbf{x}-\mu) = 2 \log \left( \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \right) = \text{Constant}$

Values of  $\mathbf{x}$  yielding constant height for density are ellipsoids.

$$\text{Constant density contour} = \{\mathbf{x} : (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) = c^2\}$$



$$P(x_1 \in [a, b]) = \text{Area under curve}$$

in univariate case  
Area represent  $P$

in multivariate case  
Volume represent  $P$

$$P((x_1, x_2) \in O) = \text{Volume under surface}$$

Figure 3: Constant density contour for bivariate normal

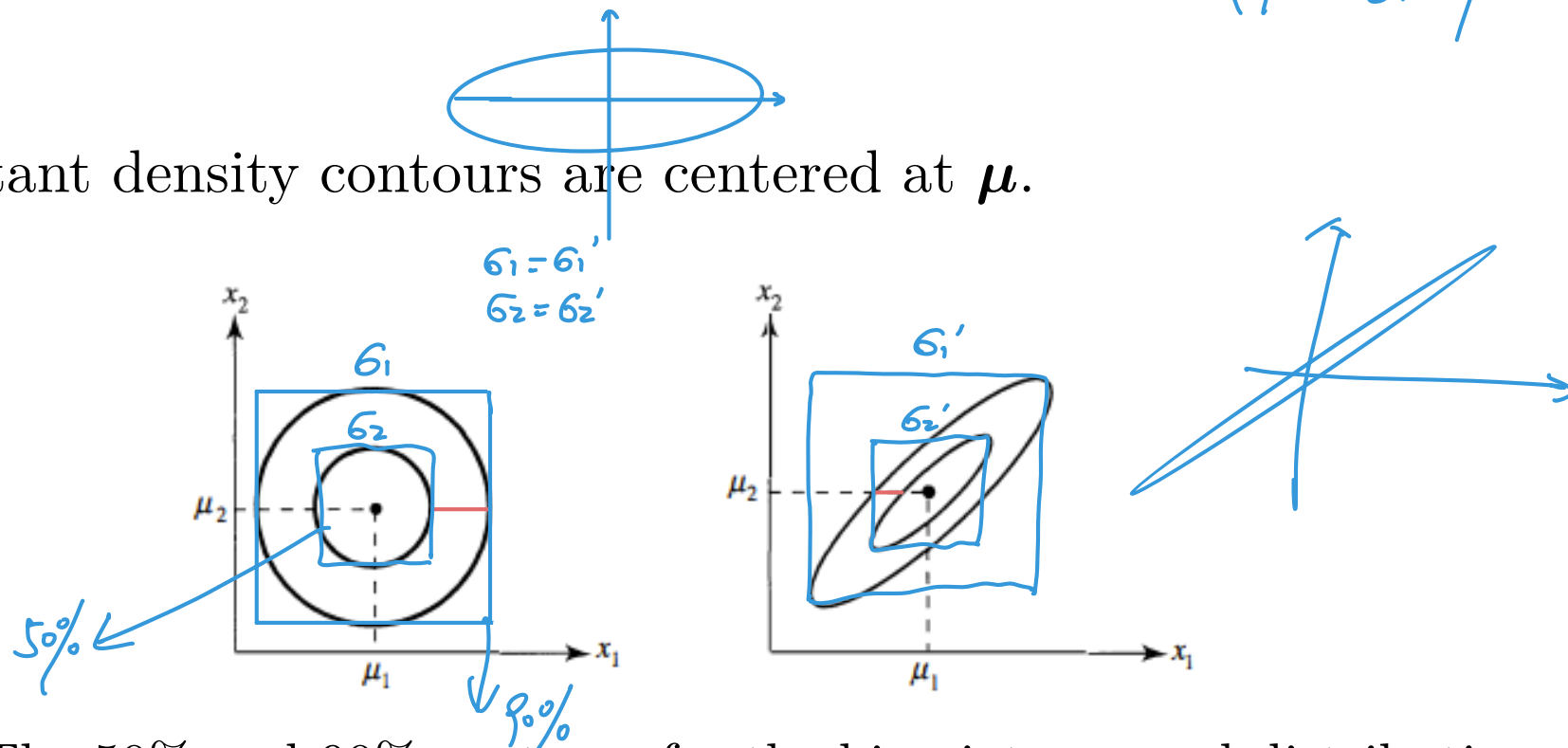


# Constant density contour

$$\sigma_{11} > \sigma_{22} \text{ and } \sigma_{12} = 0$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{pmatrix}$$

The constant density contours are centered at  $\mu$ .



**Figure 4:** The 50% and 90% contours for the bivariate normal distributions in Figure 2

# Some equivalent definitions of MVN

Besides specifying the density to define the multivariate normal distribution, we have some equivalent definitions:

## Definition 2

$$Z_i \sim \mathcal{N}(0, 1)$$

linear combination of  $Z$  will be multivariate normal

A random vector  $\mathbf{X}$  has a multivariate normal distribution if there exist  $\boldsymbol{\mu} \in \mathbb{R}^k$ ,  $\mathbf{A} \in \mathbb{R}^{k \times \ell}$  such that  $\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}$  where  $\mathbf{Z}$  is a random vector whose components are independent  $\mathcal{N}(0, 1)$  random variables.

## Definition 3

If for every  $\mathbf{a}$ ,  $\mathbf{a}^T \mathbf{X}$  is normal then  $\mathbf{C}\mathbf{X}$  is normal because  $\mathbf{a}^T \mathbf{C}\mathbf{X} = (\mathbf{C}^T \mathbf{a})^T \mathbf{X}$  is normal by df. 3 and  $\mathbf{X} \sim \text{MVN}$

A random vector  $\mathbf{X}$  has a multivariate normal distribution if for every real vector  $\mathbf{a}$ , the random variable  $\mathbf{a}^T \mathbf{X}$  is normal.

If  $\mathbf{X} \sim \text{MVN}$ , then  $\mathbf{C}\mathbf{X}$  is MVN  
 we can find  $\mathbf{A}, \boldsymbol{\mu}$ , s.t.  $\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}$   $\mathbf{Z} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_p \end{pmatrix}$   $Z_i \text{ iid } \mathcal{N}(0, 1)$

$$CX = CAZ + C\mu \quad \text{by df. 2 it is MVN}$$

# Some important properties for multivariate normal vectors

- ① Linear combinations of the components of  $\mathbf{X}$  are normally distributed.
- ② All subsets of the components of  $\mathbf{X}$  have a (multivariate) normal distribution. *linear combination of original vectors*
- ③ Zero covariance implies that the corresponding components are independently distributed.
- ④ The conditional distributions of the components are (multivariate) normal.

# Property 1: Linear combinations of a normal vector is normal

If  $X \sim N_p(\mu, \Sigma)$   
 then  $CX + d \sim N_p(C\mu + d, C\Sigma C^T)$

If  $\mathbf{X}$  is distributed as  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , the  $q$  linear combinations  
 if  $C \in \mathbb{R}^{q \times p}$

$$\underset{(q \times p)}{\mathbf{A}} \underset{(p \times 1)}{\mathbf{X}} = \begin{bmatrix} a_{11}X_1 + \cdots + a_{1p}X_p \\ a_{21}X_1 + \cdots + a_{2p}X_p \\ \vdots \\ a_{q1}X_1 + \cdots + a_{qp}X_p \end{bmatrix}$$

are distributed as  $N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . Also,  $\underset{(p \times 1)}{\mathbf{X}} + \underset{(p \times 1)}{\mathbf{d}}$ , where  $\mathbf{d}$  is a vector of constants, is distributed as  $N_p(\boldsymbol{\mu} + \mathbf{d}, \boldsymbol{\Sigma})$ .

# Example

For  $\mathbf{X} = (X_1, X_2, X_3)^\top$  distributed as  $N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , find the distribution of  $(X_1 - X_2, X_2 - X_3)^\top$ .

# Example

For  $\mathbf{X} = (X_1, X_2, X_3)^\top$  distributed as  $N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , find the distribution of  $(X_1 - X_2, X_2 - X_3)^\top$ .

We write

$$\begin{bmatrix} X_1 - X_2 \\ X_2 - X_3 \end{bmatrix} = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \mathbf{A}\mathbf{X}$$

# Example

Then the distribution of  $\mathbf{AX}$  is multivariate normal with mean

$$\mathbf{A}\boldsymbol{\mu} = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_1 - \mu_2 \\ \mu_2 - \mu_3 \end{bmatrix}$$

and covariance matrix

$$\begin{aligned} \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top &= \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_{22} & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_{33} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ 0 & -1 \end{bmatrix} \\ &= \begin{bmatrix} \sigma_{11} - 2\sigma_{12} + \sigma_{22} & \sigma_{12} + \sigma_{23} - \sigma_{22} - \sigma_{13} \\ \sigma_{12} + \sigma_{23} - \sigma_{22} - \sigma_{13} & \sigma_{22} - 2\sigma_{23} + \sigma_{33} \end{bmatrix} \end{aligned}$$

## Property 2: Subsets of components of normal are normal

$\mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$  is MVN, then  $\begin{pmatrix} x_3 \\ x_5 \end{pmatrix}$  is MVN.  
 because  $\begin{pmatrix} x_3 \\ x_5 \end{pmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 0 & 1 & \cdots & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$

All subsets of  $\mathbf{X}$  are normally distributed. If we respectively partition  $\mathbf{X}$ , its mean vector  $\boldsymbol{\mu}$ , and its covariance matrix  $\boldsymbol{\Sigma}$  as

$\text{Cov} \begin{pmatrix} x_3 \\ x_5 \end{pmatrix} = \begin{bmatrix} 6_{33} & 6_{35} \\ 6_{55} & 6_{55} \end{bmatrix}$

$$\underset{(p \times 1)}{\mathbf{X}} = \begin{bmatrix} \underset{(q \times 1)}{\mathbf{X}_1} \\ \underset{((p-q) \times 1)}{\mathbf{X}_2} \end{bmatrix} \quad \underset{(p \times 1)}{\boldsymbol{\mu}} = \begin{bmatrix} \underset{(q \times 1)}{\boldsymbol{\mu}_1} \\ \underset{((p-q) \times 1)}{\boldsymbol{\mu}_2} \end{bmatrix}$$

and

$$\underset{(p \times p)}{\boldsymbol{\Sigma}} = \begin{bmatrix} \underset{(q \times q)}{\boldsymbol{\Sigma}_{11}} & \underset{(q \times (p-q))}{\boldsymbol{\Sigma}_{12}} \\ \underset{((p-q) \times q)}{\boldsymbol{\Sigma}_{21}} & \underset{((p-q) \times (p-q))}{\boldsymbol{\Sigma}_{22}} \end{bmatrix}$$

then  $\mathbf{X}_1$  is distributed as  $N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ .



## Property 3

If  $\begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$  is MVN,  $\text{Cov}(X_1, X_2) = 0$ , then  $X_1 \perp X_2$

because  $\begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$    
  $\text{Cov}(X_1)$    
  $\text{Cov}(X_2)$

$$\textcircled{2} (X-\mu)^T \Sigma^{-1} (X-\mu)$$

$$= (X_1 - \mu_1, X_2 - \mu_2) \begin{pmatrix} \Sigma_1^{-1} & 0 \\ 0 & \Sigma_2^{-1} \end{pmatrix} \begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{pmatrix}$$

$$\text{joint density} = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{(X-\mu)^T \Sigma^{-1} (X-\mu)}{2}} = \frac{1}{(2\pi)^{\frac{k}{2}} |\Sigma_1|^{\frac{1}{2}}} e^{-\frac{(X_1-\mu_1)^T \Sigma_1^{-1} (X_1-\mu_1)}{2}} \times \frac{1}{(2\pi)^{\frac{p-k}{2}} |\Sigma_2|^{\frac{1}{2}}} e^{-\frac{(X_2-\mu_2)^T \Sigma_2^{-1} (X_2-\mu_2)}{2}}$$

here  $|\Sigma| = |\Sigma_1| \cdot |\Sigma_2|$

$\Leftarrow$  generally not true.

- If  $\mathbf{X}_1 \in \mathbb{R}^{q_1}$  and  $\mathbf{X}_2 \in \mathbb{R}^{q_2}$  are independent, then  $\text{Cov}(\mathbf{X}_1, \mathbf{X}_2) = \mathbf{0}$ , a  $q_1 \times q_2$  matrix of zeros.

- If  $\begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$  is  $N_{q_1+q_2} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$ , then  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are independent if and only if  $\Sigma_{12} = \mathbf{0}$ .

- Note that if two normal vectors have a covariance of zero, it does not imply that they are independent. A counterexample is by taking  $\mathbf{X}_2 = Z\mathbf{X}_1$  for some  $Z$  independent of  $\mathbf{X}_1$  and has distribution  $P(Z = \pm 1) = 1/2$ .

If  $X_1$  is MVN and  $X_2$  is MVN and  $\text{Cov}(X_1, X_2) = 0$  then  $X_1 \perp X_2$   $\times$

$$\text{Cov} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

$\in \mathbb{R}^k$    
  $\in \mathbb{R}^{p-k}$

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} X_1 \end{bmatrix}^T + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} X_2 \end{bmatrix}^T$$

# Example

Let  $\underset{(3 \times 1)}{\mathbf{X}}$  be  $N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with  $\vec{X} \in \mathbb{R}^{3 \times 1} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$   $\text{Cov}(X) = \boldsymbol{\Sigma}$

$$\boldsymbol{\Sigma} = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

$\text{Cov}(X_1, X_2) = \text{Cov}(X_2, X_1) = 1 \neq 0$   
 Are  $X_1$  and  $X_2$  independent? What about  $(X_1, X_2)$  and  $X_3$ ?

$$\text{Cov}\left(\begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, X_3\right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} \text{Cov}(X_1, X_3) \\ \text{Cov}(X_2, X_3) \end{pmatrix}$$

and  $(X_1, X_2)$  and  $X_3$  are jointly normal  
 because  $\vec{X}$  is  $N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

# Example

Let  $\underset{(3 \times 1)}{\mathbf{X}}$  be  $N_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with

$$\boldsymbol{\Sigma} = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Are  $X_1$  and  $X_2$  independent? What about  $(X_1, X_2)$  and  $X_3$  ?

$X_1$  and  $X_2$  are not independent;  $(X_1, X_2)$  and  $X_3$  are independent.

# Property 4: Conditional distributions of components are normal

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \text{ is MVN}$$

$$\text{then } X_1 | X_2 = x_2 \text{ is MVA}$$

$$X_2 | X_1 = x_1 \text{ is MVN}$$

Let  $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$  be distributed as  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with  $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ ,

$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$ , and  $|\boldsymbol{\Sigma}_{22}| > 0$ . Then the conditional distribution of  $\mathbf{X}_1$ , given that  $\mathbf{X}_2 = \mathbf{x}_2$ , is normal and has

$$\text{Mean} = \mu_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \mu_2)$$

and

$$X_1 | X_2 = x_2 \quad E(\dots)$$

$$EX = \begin{pmatrix} EX_1 \\ EX_2 \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \text{Cov}X = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Let  $Y = AX$  where  $A = \begin{bmatrix} I_p & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I_{p-q} \end{bmatrix}$ ,  $Y = \begin{pmatrix} X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 \\ X_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}\right)$

$X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 \checkmark$   
 $X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 | X_2 = x_2$   
 then we know  
 $X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 | X_2 = x_2$

Covariance =  $\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$

$E[X_1 | X_2 = x_2] = \mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2 + \Sigma_{12}\Sigma_{22}^{-1}x_2$

So  $X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 \perp X_2$

# Distribution of Sample mean

For any  $\mathbf{a} \in \mathbb{R}^p$ ,  $\mathbf{a}^T \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{a}^T \mathbf{x}_i$  is normal

So  $\bar{\mathbf{x}}$  is normal

$$E \bar{\mathbf{x}} = \boldsymbol{\mu}$$

- If  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then for any positive integer  $n$ ,

$$\text{Cov}(\bar{\mathbf{x}}) = \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\right) = \frac{1}{n} \boldsymbol{\Sigma} \quad \uparrow \downarrow \quad \text{remain same}$$

$$\sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \sim N_p(0, \boldsymbol{\Sigma}).$$

$$\bar{\mathbf{x}} - \boldsymbol{\mu} \sim N_p(0, \frac{1}{n} \boldsymbol{\Sigma}) \quad \therefore \sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \sim N_p(0, \boldsymbol{\Sigma})$$

(Proof hint: Use Definition 3, show that  $\mathbf{c}^T \bar{\mathbf{x}}$  is normal first.)

- For real data, the multivariate normal assumption may be too strong. Fortunately, the Central Limit Theorem implies that the above result still holds asymptotically as the sample size  $n$  approaches infinity.

# Central Limit Theorem

## (Multivariate) Central Limit Theorem

Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  be independent observations from any population with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ . Then

$$\sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \xrightarrow{d} N_p(0, \boldsymbol{\Sigma}) \text{ as } n \rightarrow \infty,$$

where convergence in distribution  $\xrightarrow{d}$  means that for any set  $\Omega \in \mathbb{R}^p$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathbf{S}_n \in \Omega) = \mathbb{P}(\mathbf{Y} \in \Omega)$$

for a random vector  $\mathbf{Y}$  with normal distribution  $N_p(0, \boldsymbol{\Sigma})$ .

# Why MVN

## Why emphasis the MVN?

- Genuinely good population model for some natural phenomena
- Even for nonnormal data, MVN is often useful approximation - such as for inferences involving sample mean vectors, which are asymptotically normal due to CLT *last page*
- Theoretical elegance, such as
  - MVN is completely characterized by its mean and covariance
  - Linear transformations of MVN are normal
  - Uncorrelated components of MVN are independent

# Maximum Likelihood Estimation

- Maximum likelihood estimation (MLE) selects parameter values that maximize the joint density of the observed data.
- In the next few pages, we derive the MLE for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  for multivariate normal data  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{i.i.d.}{\sim} N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
- This optimization leads to the MLEs:  $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$  and  $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$ .



# Maximum likelihood estimation of $\mu$ and $\Sigma$

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be a random sample from  $N_p(\mu, \Sigma)$ .

Joint density:

$$\begin{aligned}
 f(\mathbf{x}_1, \dots, \mathbf{x}_n) &= \prod_{i=1}^n f(\mathbf{x}_i) \\
 &= \prod_{i=1}^n \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \mu)^\top \Sigma^{-1} (\mathbf{x}_i - \mu) \right\} \\
 &= \frac{1}{(2\pi)^{\frac{np}{2}} |\Sigma|^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^\top \Sigma^{-1} (\mathbf{x}_i - \mu) \right\} \\
 &= L(\mu, \Sigma) = L(\mathbf{x}_1, \dots, \mathbf{x}_n)
 \end{aligned}$$

Goal: Find values of  $\mu$  and  $\Sigma$  that maximize the likelihood of observing  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

$$\begin{aligned}
 \hat{\mu}, \hat{\Sigma} &= \underset{\mu, \Sigma}{\operatorname{argmax}} L(\mathbf{x}_1, \dots, \mathbf{x}_n) \\
 \hat{\mu} &= \underset{\mu}{\operatorname{argmin}} \sum_{i=1}^n (\mathbf{x}_i - \mu)^\top \Sigma^{-1} (\mathbf{x}_i - \mu)
 \end{aligned}$$

$$\begin{aligned} \text{Consider } \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) & \quad (\text{tr } AB = \text{tr } BA) \\ &= \text{tr} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \text{tr} \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \\ &= \text{tr} \left[ \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \boldsymbol{\mu}) (\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \boldsymbol{\mu})^\top \right] \\ &= \text{tr} \left[ \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top + n \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) (\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \right] \end{aligned}$$

Use

$$\begin{aligned} (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) &= \text{tr} \left\{ (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &= \text{tr} \left\{ \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \right\}, \end{aligned}$$

rewrite

$$\begin{aligned} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) &= \text{tr} \left\{ \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &= \text{tr} \left\{ \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \right\} \\ &= \text{tr} \left\{ \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n [(\mathbf{x}_i - \bar{\mathbf{x}}) + (\bar{\mathbf{x}} - \boldsymbol{\mu})] [(\mathbf{x}_i - \bar{\mathbf{x}}) + (\bar{\mathbf{x}} - \boldsymbol{\mu})]^\top \right\} \\ &= \text{tr} \left\{ \boldsymbol{\Sigma}^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top + \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\bar{\mathbf{x}} - \boldsymbol{\mu})^\top + \right. \right. \\ &\quad \left. \left. \sum_{i=1}^n (\bar{\mathbf{x}} - \boldsymbol{\mu}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top + n(\bar{\mathbf{x}} - \boldsymbol{\mu}) (\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \right] \right\} \end{aligned}$$

So,

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{np}{2}} |\boldsymbol{\Sigma}|^{\frac{n}{2}}} \times \exp \left\{ -\frac{1}{2} \operatorname{tr} \left( \boldsymbol{\Sigma}^{-1} \left[ \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \right] \right) \right\}$$

Note that the value of  $\boldsymbol{\mu}$  maximizing  $L(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the value minimizing  $\operatorname{tr} \{ n\boldsymbol{\Sigma}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \}$ . Since  $\boldsymbol{\Sigma}^{-1}$  is positive definite,

$$\operatorname{tr} \{ n\boldsymbol{\Sigma}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \} = n(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}) > 0$$

unless

$$\boldsymbol{\mu} = \bar{\mathbf{x}}.$$

Thus, the likelihood is maximized with respect to  $\boldsymbol{\mu}$  at  $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$ .

It remains to maximize

$$L(\hat{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} e^{-\text{tr}[\boldsymbol{\Sigma}^{-1} (\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})^\top)]/2}$$

over  $\boldsymbol{\Sigma}$ .

Result: Given a  $p \times p$  symmetric positive definite <sup>but important</sup> matrix  $\mathbf{B}$  and a scalar  $b > 0$ , it follows that

$$\frac{1}{|\boldsymbol{\Sigma}|^b} e^{-\text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{B})/2} \leq \frac{1}{|\mathbf{B}|^b} (2b)^{pb} e^{-bp}$$

for all positive definite  $\boldsymbol{\Sigma}$ , with equality holding only for  $\boldsymbol{\Sigma} = (1/2b)\mathbf{B}$ .

By the above result with  $b = n/2$  and  $\mathbf{B} = \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})^\top$ , the maximum occurs at  $\hat{\boldsymbol{\Sigma}} = (1/n) \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})^\top = (n-1)\mathbf{S}/n$ .