

# Multivariate Analysis Lecture 4: A Random Sample from A Multivariate Distribution

Zhaoxia Yu  
Professor, Department of Statistics

2025-04-10

## Section 1

### Outline

# Outline

- Review of Lecture 03
- Inference of Means
- Inference of a Linear Combination of Mean
- A Simulation Study
- Generalized Variance
- Normal and Multivariate Normal

## Section 2

### Lecture 3

## Subsection 1

### Random Samples

# Random Samples

- Univariate distribution: If  $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2)$ , then
  - $\bar{X} \sim (\mu, \sigma^2/n)$
  - $E[s^2] = \sigma^2$  where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ .
- Multivariate distribution: If  $X_1, \dots, X_n \stackrel{iid}{\sim} (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then
  - $\bar{X} \sim (\boldsymbol{\mu}, \frac{1}{n} \boldsymbol{\Sigma})$
  - $E[\mathbf{S}] = \boldsymbol{\Sigma}$  where  $\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T$ .

## Subsection 2

### Linear Combinations of a Random Vector

Outline **Lecture 3** Inference of Means Inference of a Linear C

- Outline **Lecture 3** Inference of Means Inference of a Linear C

Outline **Lecture 3** Inference of Means Inference of a Linear C

- Outline **Lecture 3** Inference of Means Inference of a Linear C

Outline **Lecture 3** Inference of Means Inference of a Linear C

- Outline **Lecture 3** Inference of Means Inference of a Linear C



- $$Y = a^T X = \frac{1}{3}(X_1 + X_2 + X_3)$$

- 9 / 50

## Linear Combinations: Example

- Example 1: Assume we have a random sample from a distribution with mean  $\mu$  and variance  $\sigma^2$ , i.e.,  $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2)$ .
- We often stack the random variables vertically:

$$\mathbf{X}_{n \times 1} = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}.$$

- An equivalent expression,  $\mathbf{X} = (X_1, \dots, X_n)^T$ .
- Note that  $\mathbf{X}$  is a **random vector** with mean vector and covariance matrix

$$E[\mathbf{X}] = \mu \mathbf{1}_n = \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix}, \text{Cov}(\mathbf{X}) = \sigma^2 \mathbf{I}_n = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

# Linear Combinations: Example

- We can express the sample mean as a linear combination of the random vector  $\mathbf{X}$ :

$$\bar{X} = \frac{1}{n} \mathbf{1}^T \mathbf{X},$$

where  $\mathbf{1} = (1, \dots, 1)^T$  is a  $n \times 1$  vector.

- By the linear combination results, we have

$$E[\bar{X}] = \frac{1}{n} \mathbf{1}^T E[\mathbf{X}] = \frac{1}{n} \mathbf{1}^T \mu \mathbf{1} = \mu$$

$$\text{Var}[\bar{X}] = \frac{1}{n^2} \mathbf{1}^T \text{Cov}(\mathbf{X}) \mathbf{1} = \frac{1}{n^2} \mathbf{1}^T \sigma^2 \mathbf{I}_n \mathbf{1} = \frac{1}{n} \sigma^2.$$

## Section 3

### Inference of Means

## Subsection 1

### Univariate

# Univariate

- A random sample  $X_1, \dots, X_n$  from a univariate distribution with mean  $\mu$  and variance  $\sigma^2$ .
- We are interested in making inference about the population mean  $\mu$ .
- The sample mean  $\bar{X}$  is an unbiased estimator of  $\mu$ , i.e.,  $E[\bar{X}] = \mu$ . We often use  $\hat{\mu} = \bar{X}$  to estimate  $\mu$ .
- How to quantify the uncertainty of  $\hat{\mu}$ ? Recall that  $Var(\bar{X}) = \frac{\sigma^2}{n}$ ..
- $\sigma^2$  is unknown. But the sample variance  $s^2$  is an unbiased estimator of  $\sigma^2$ , i.e.,  $E[s^2] = \sigma^2$ . We often use  $\hat{\sigma}^2 = s^2$  to estimate  $\sigma^2$ .
- The sample variance  $s^2$  is defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

# Standard Error of $\bar{X}$

- The standard error of  $\bar{X}$  is defined as  $se(\bar{X}) = \sqrt{Var(\bar{X})} = \frac{\sigma}{\sqrt{n}}$ .
- We can estimate it by  $se(\bar{X}) = \frac{s}{\sqrt{n}}$ .
- A “large-sample” (approximate) confidence interval for  $\mu$  is given by

$$\bar{X} \pm z_{\alpha/2} se(\bar{X})$$

where  $z_{\alpha/2}$  is the upper  $\alpha/2$  quantile of the standard normal distribution.

- A “small-sample” (approximate) confidence interval for  $\mu$  is given by

$$\bar{X} \pm t_{\alpha/2} se(\bar{X})$$

where  $t_{\alpha/2}$  is the upper  $\alpha/2$  quantile of the t-distribution with  $n - 1$  degrees of freedom.

## Subsection 2

### Multivariate



# A Random Sample From a Multivariate Distribution

- Consider a random sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  from a multivariate distribution with mean vector  $\boldsymbol{\mu}_{p \times 1}$  and covariance matrix  $\boldsymbol{\Sigma}_{p \times p}$ .
- We often stack the random vectors to form an  $n \times p$  matrix:

$$\mathbf{X}_{n \times p} = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix}$$



# A Random Sample From a Multivariate Distribution: Sample Covariance Matrix

- The sample covariance matrix is

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

- It is unbiased for  $\Sigma$ , i.e.,  $E[\mathbf{S}] = \Sigma$ .
- We showed that

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^T \mathbf{C} \mathbf{X}$$

where  $\mathbf{C}_{n \times n} = \mathbf{I} - \frac{1}{n} \mathbf{J} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T$

- This expression is helpful when we derive the distribution of  $\mathbf{S}$ .

## Section 4

### Inference of a Linear Combination of Means

## Subsection 1

### Linear Combinations of Means

# Linear Combinations of Means

- In many situations, the parameter of interest is a function of the means.
- For example, we may be interested in the mean of a linear combination of the means, i.e.,  $\mu^T \mathbf{a}$ , where  $\mathbf{a} = (a_1, \dots, a_p)^T$  is a  $p \times 1$  vector.
- In the following simulated study, we will show how to construct a large-sample confidence interval for  $\mu^T \mathbf{a}$ .

## Subsection 2

### A Simulated Study

# Daily Intake of Protein

- This is a simulated data set
- For adults, the recommended range of daily protein intake is between 0.8 g/kg and 1.8 g/kg of body weight
- 60 observations
- 4 sources of proteins
  - meat
  - dairy
  - vegetables / nuts / tofu
  - other



# Choose Mean Vector and Covariance Matrix

- The multivariate distribution has

- mean vector

$$\mu = (24, 16, 8, 8)^T$$

- covariance matrix

$$\Sigma = 4 * \begin{pmatrix} 1.3 & 0.3 & 0.3 & 0.3 \\ 0.3 & 1.3 & 0.3 & 0.3 \\ 0.3 & 0.3 & 1.3 & 0.3 \\ 0.3 & 0.3 & 0.3 & 1.3 \end{pmatrix}$$

# Define Mean Vector and Covariance Matrix in R

```
#the library "MASS" is required
library(MASS)
my.cov=4*(diag(4) + 0.3* rep(1,4)%o%rep(1,4))
eigen(my.cov)#to check whether the cov matrix is p.d.
```

```
## eigen() decomposition
## $values
## [1] 8.8 4.0 4.0 4.0
##
## $vectors
##      [,1]      [,2]      [,3]      [,4]
## [1,] -0.5  0.8660254  0.0000000  0.0000000
## [2,] -0.5 -0.2886751 -0.5773503 -0.5773503
## [3,] -0.5 -0.2886751 -0.2113249  0.7886751
## [4,] -0.5 -0.2886751  0.7886751 -0.2113249
```

```
my.mean=8*c(3,2,1,1)
n=60
```

# Simulate A Random Sample

```

set.seed(1)
x=mvrnorm(n, mu=my.mean, Sigma=my.cov)
dim(x)

```

```
## [1] 60 4
```

```

protein=as.matrix(data.frame(meat=x[,1],dairy=x[,2],
                             veg=x[,3], other=x[,4]))

```

# The simulated data

protein

##		meat	dairy	veg	other
##	[1,]	29.08891	17.54865	5.814221	7.264953
##	[2,]	23.65965	13.06336	8.734581	9.452868
##	[3,]	26.43410	16.83504	9.278807	8.409798
##	[4,]	21.68232	15.51922	3.379171	5.954558
##	[5,]	22.22387	15.45446	8.804571	7.562144
##	[6,]	25.54395	16.46835	8.556332	10.299174
##	[7,]	20.15075	14.71290	10.660378	7.584075
##	[8,]	25.44330	14.98680	4.866275	6.323171
##	[9,]	23.41142	16.34138	6.667006	6.164109
##	[10,]	28.21604	16.64242	5.874860	7.078538
##	[11,]	22.58127	13.61817	5.178349	5.652878
##	[12,]	22.19211	16.04745	8.714666	6.732854
##	[13,]	25.97926	16.80008	7.189986	9.716474
##	[14,]	25.66703	20.61869	13.775770	9.078236
##	[15,]	20.16010	16.09623	5.020107	8.049388
##	[16,]	24.57145	18.88862	6.884708	5.917308

# Sample Mean and Sample Covariance

```
xbar=matrix(colMeans(protein), 4, 1)
t(xbar)
```

```
##           [,1]      [,2]      [,3]      [,4]
## [1,] 24.03403 15.92836 7.66049 7.738634
```

```
S=cov(protein)
S
```

```
##           meat      dairy      veg      other
## meat  4.2956426 0.8150757 1.1294478 0.5532420
## dairy 0.8150757 4.4052993 0.3497889 0.2337300
## veg   1.1294478 0.3497889 5.1705794 0.5897121
## other 0.5532420 0.2337300 0.5897121 4.5287293
```

# Estimation

- An unbiased estimator of  $\mu$  is the sample mean vector, i.e.,  $\hat{\mu} = \bar{\mathbf{X}}$ .
- An unbiased estimator of  $\Sigma$  is the sample covariance matrix  $\mathbf{S}$ , i.e.,  $\hat{\Sigma} = \mathbf{S}$
- We have shown that  $Var(\bar{\mathbf{X}}) = \frac{1}{n}\Sigma$ , where  $n = 60$ .
- We can estimate it by

$$\hat{Var}(\bar{\mathbf{X}}) = \frac{1}{60}\mathbf{S}$$

## Linear Functions/Combinations: Three Questions

- Suppose we only have a random sample and we would like to make inference of the following:
- Q1: Construct a large-sample (approximate) C.I. for protein from meat. In other words, the parameter of interest is  $\mu_1$ .
- Q2: Construct a large-sample C.I. for the total protein intake
- Q3: Construct a large-sample C.I. for the difference of protein intake between from meat and from vegetable

# Linear Functions/Combinations: Question 1

- Q1: Construct a large-sample (approximate) C.I. for protein from meat. In other words, the parameter of interest is  $\mu_1$ .
- Estimate  $\bar{X}_{(1)} = 24.0$ .
- We need compute the standard error (s.e.) of  $\bar{X}_1$ , which is defined as  $se(\bar{X}_{(1)}) = \sqrt{\hat{var}(\bar{X}_{(1)})}$
- Two ways to compute the s.e.,
  - 1  $se(\bar{X}_{(1)}) = \sqrt{4.2956/60} = 0.27$
  - 2 The calculation can also be done by noticing that  $\bar{X}_1$  is a linear combination of  $\bar{\mathbf{X}}$ :  $\bar{X}_{(1)} = \mathbf{a}^T \bar{\mathbf{X}}$ , where  $\mathbf{a}^T = (1, 0, 0, 0)$ . Thus,

$$\hat{Var}(\bar{X}_{(1)}) = \mathbf{a}^T \frac{\mathbf{S}}{60} \mathbf{a}$$



## Linear Functions/Combinations: Question 2

- Q2: Construct a large-sample C.I. for the total protein intake
- The parameter of interest is  $\mu_1 + \mu_2 + \mu_3 + \mu_4 = \mathbf{a}^T \boldsymbol{\mu}$ , where  $\mathbf{a} = (1, 1, 1, 1)^T$ .
- Estimate:  $\mathbf{a}^T \bar{\mathbf{X}}$
- Standard error:  $\sqrt{\mathbf{a}^T \frac{\mathbf{S}}{n} \mathbf{a}}$

## Linear Functions/Combinations: Question 3

- Q3: Construct a large-sample C.I. for the difference of protein intake between from meat and from vegetable
- The parameter of interest is  $\mu_1 - \mu_3 = \mathbf{a}^T \boldsymbol{\mu}$ , where  $\mathbf{a} = (1, 0, -1, 0)^T$ .
- Estimate:  $\mathbf{a}^T \bar{\mathbf{X}}$
- Standard error:  $\sqrt{\mathbf{a}^T \frac{\mathbf{S}}{n} \mathbf{a}}$

## Linear Functions/Combinations: Question 1 (continued)

- R code to compute using the above two ways

```
sqrt(S[1,1]/60) # Method 1
```

```
## [1] 0.2675706
```

```
# Method 2
```

```
a=matrix(c(1,0,0,0),4,1)
```

```
sqrt(t(a)%*%S%*%a/60)
```

```
## [1]
```

```
## [1,] 0.2675706
```

- Both methods give

$$s.e.(\bar{X}_{(1)}) = 0.27$$

## Linear Functions/Combinations: Question 2 (continued)

```
a=matrix(1,4,1)
t(a)%*% xbar #estimate
```

```
##           [,1]
## [1,] 55.36152
```

```
sqrt(t(a)%*%S%*%a/60) #standard error
```

```
##           [,1]
## [1,] 0.6550095
```

```
#a large-sample 95% C.I.
c(t(a)%*% xbar- 1.96*sqrt(t(a)%*%S%*%a/60),
  t(a)%*% xbar+ 1.96*sqrt(t(a)%*%S%*%a/60))
```

```
## [1] 54.07770 56.64534
```

## Linear Functions/Combinations: Question 3 (continued)

```
a=matrix(c(1,0,-1,0),4,1)
t(a)%*% xbar #estimate
```

```
##           [,1]
## [1,] 16.37354
```

```
sqrt(t(a)%*%S%*%a/60) #standard error
```

```
##           [,1]
## [1,] 0.3465864
```

```
#a large-sample 95% C.I.
c(t(a)%*% xbar- 1.96*sqrt(t(a)%*%S%*%a/60),
  t(a)%*% xbar+ 1.96*sqrt(t(a)%*%S%*%a/60))
```

```
## [1] 15.69423 17.05285
```

## Section 5

### Generalized Variance

# Why Do We Need Generalized Variance?

- For a random variable (i.e., univariate), we quantify dispersion using variance and standard deviation.
- For a random vector (i.e., multivariate), we use its covariance matrix to quantify the dispersion as well as the relationships between different variables / features.
  - The dispersion information is represented by a matrix, which has  $p(p + 1)/2$  unique parameters

# What is Generalized Variance?

- It is attempting to have a scalar summary (i.e., a single number) to quantify the “total” amount of dispersion for a multivariate distribution
- Generalized variance
  - Provides an overall measure of dispersion of the multivariate distribution
  - One choice is the determinant:  $|\Sigma|$ .
  - A larger determinant indicates a greater degree of dispersion



## Generalized Variance: An Example

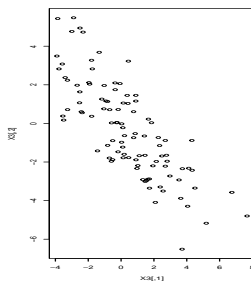
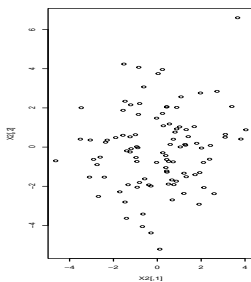
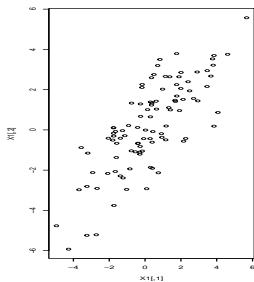
$$\Sigma_1 = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 3 & 0 \\ 0 & 3 \end{pmatrix}, \Sigma_3 = \begin{pmatrix} 5 & -4 \\ -4 & 5 \end{pmatrix}$$

```

Sigma1=matrix(c(5,4,4,5), 2,2)
X1=mvrnorm(100, mu=c(0,0), Sigma=Sigma1)
Sigma2=matrix(c(3,0,0,3), 2,2)
X2=mvrnorm(100, mu=c(0,0), Sigma=Sigma2)
Sigma3=matrix(c(5,-4,-4,5), 2,2)
X3=mvrnorm(100, mu=c(0,0), Sigma=Sigma3)
    
```

# Generalized Variance: An Example

```
par(mfrow=c(1,3))
plot(X1); plot(X2); plot(X3)
```



## Generalized Variance might (NOT) be useful

- In the example above,  $|\Sigma_1| = |\Sigma_2| = |\Sigma_3| = 9!$
- $|\Sigma|$  does not tell the orientations.
- $|\Sigma|$  is useful to compare two patterns when they have nearly the same orientations.
- The generalized variance does not capture all the information contained in the covariance matrix.
- The eigenvalues provide more information than the determinant - Principal Component Analysis!

## Section 6

Normal Distributions: univariate, multivariate,  
matrix normal distributions

# The Big Picture: Univariate vs Multivariate

- Review:** A random sample, denoted by  $X_1, \dots, X_n$ , from a (univariate) normal distribution  $N(\mu, \sigma^2)$ 
  - What are the distributions of  $\bar{X}, s^2$ ? What useful statistics can be constructed?
- New material:** A random sample, denoted by  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , from a multivariate normal distribution  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 
  - What are the distributions of  $\bar{\mathbf{X}}, \mathbf{S}$ ? What useful statistics can be constructed?

## Subsection 1

### Derivation of t-test

# Derivation of t-test

- A random sample  $X_1, \dots, X_n$  from a univariate Normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma)^2$$

- The sample mean  $\bar{X} \sim N(\mu, \sigma^2/n)$ . Standardized the sample mean:

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

- The sample variance  $s^2 \sim \chi_{n-1}^2$ , i.e.,  $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$ .
- The sample variance  $s^2$  is independent of  $\bar{X}$ .

# Derivation of t-test

- The t-statistic is defined as

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}} = \frac{\sqrt{n}(\bar{X} - \mu)}{s}$$

- It follows a t-distribution with  $n - 1$  degrees of freedom, denoted by  $t_{n-1}$ .
- The t-distribution is a family of distributions that are symmetric and bell-shaped, like the standard normal distribution, but have heavier tails.
- The t-distribution is used in hypothesis testing and confidence intervals for small sample sizes, particularly when the population standard deviation is unknown.
- The t-distribution approaches the standard normal distribution as the sample size increases.





# The Big Picture: Multivariate

- A random sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  from a multivariate normal distribution  $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
- Let

$$\mathbf{X}_{n \times p} = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix}$$

$\mathbf{X}$  follows a matrix normal distribution.

- 1 Sample mean vector follows a multivariate normal, i.e.,  $\bar{\mathbf{X}} \sim \mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}/n)$
- 2 Sample covariance matrix  $(n-1)\mathbf{S}$  follows a Wishart distribution, i.e.,  $(n-1)\mathbf{S} \sim \text{Wishart}_p(n-1, \boldsymbol{\Sigma})$
- 3 Independence between  $\bar{\mathbf{X}}$  and  $\mathbf{S}$ .
- 4 Hoetelling's  $T^2$ :  $T^2 = (\bar{\mathbf{X}} - \boldsymbol{\mu})^T \left(\frac{\mathbf{S}}{n}\right)^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu})$