

# MAT 4375

## Multivariate Statistical Methods

### Study Guide



Fall 2024

## Vectors and Matrices

- A matrix is a rectangular array of scalar elements

$$A_{p \times q} = (a_{ij}) = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1q} \\ a_{21} & a_{22} & \dots & a_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pq} \end{pmatrix}$$

$p \times q$  dimension of matrix

P: number of rows

q: number of columns

Row vector:  $a_{1j} = (a_{11}, a_{12}, a_{13}, \dots, a_{1q})$

Column vector:

$$a_{1i} = \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \\ \vdots \\ a_{p1} \end{pmatrix}$$

## Elementary Matrix Operations

Sum and difference:  $A+B$  or  $A-B$ , done by doing element-wise addition and subtraction

Multiplication (written  $AB$ )

$$AB = \left( \sum_k a_{ik} b_{kj} \right)$$

- For multiplication to work, we need the matrices to be conformable

Inner product: Product of row vector and column vector, leads to scalar

Outer product: product of column vector and row vector, leads to matrix

Transpose: interchanging rows and columns of  $A$ , denoted  $A'$

A square matrix has same rows and columns. A square matrix  $A$  is symmetric if  $A=A'$

Properties of transpose:

- $(A')' = A$
- $(A+B)' = A'+B'$
- $(AB)' = B'A'$
- $(cA)' = cA'$

A diagonal matrix is a square matrix with zero off-diagonal elements

Upper triangular matrices have zero elements below the diagonal and lower triangular matrices have zero elements above the diagonal

Kronecker Product:

$$A = (a_{ij}) : p \times q$$

$$B = (b_{ij}) : m \times n$$

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1q}B \\ a_{21}B & a_{22}B & \dots & a_{2q}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1}B & a_{p2}B & \dots & a_{pq}B \end{pmatrix}, Pm \times qn$$

Properties of Kronecker Product:

$$(A+B) \otimes C = (A \otimes C) + (B \otimes C)$$

$$A \otimes (B+C) = (A \otimes B) + (A \otimes C)$$

$$- A \otimes B \neq B \otimes A$$

$$(A \otimes B)(C \otimes D) = AC \otimes BD$$

$$-(A \otimes B)' = A' \otimes B'$$

Vectorization, denote by  $\text{vec}(A)$ , is done by stacking the columns of  $A$  and we form a column vector

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$\text{vec}(A) = \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix}$$

## Trace of Matrix

- Sum of diagonals of matrix  
 $\text{trace}(A) = \sum_i a_{ii}$
- Properties of trace
  - $\text{trace}(A') = \text{trace}(A)$
  - $\text{trace}(A+B) = \text{trace}(A) + \text{trace}(B)$
  - $\text{trace}(AB) = \text{trace}(BA)$
  - $\text{trace}(AA') = \text{trace}(A'A)$
  - $\text{trace}(ABC) = \text{trace}(CAB) = \text{trace}(BCA)$
  - $\text{trace}(A \otimes B) = \text{trace}(A)\text{trace}(B)$

## Determinant of Matrix

- For 1x1 matrix, or scalar, the determinant is the value itself
- For 2x2 matrix  $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$ , the determinant, denoted by  $|A|$  is given by

$$|A| = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}$$

For 3x3 matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

The determinant is calculated as

$$a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

General formula

$$\det(A) = \sum_{i=1}^n (-1)^{i+j} a_{ij} |A_{ij}| \text{ for all } j$$

or

$$\det(A) = \sum_{j=1}^n (-1)^{j+i} a_{ij} |A_{ij}| \text{ for all } i$$

## Properties

- $-|A'| = |A|$
- $-|AB| = |A||B| = |BA|$
- $-k|A| = k^p |A|, A \text{ p} \times p$
- $-|AA'| \geq 0$

Positive semi-definite matrix: a symmetric matrix  $A$  is called positive semi-definite if  $\bar{x}'A\bar{x} \geq 0$  for all  $\bar{x} \in \mathbb{R}^p$  and  $\bar{x}'A\bar{x} = 0$  for some nonzero  $\bar{x}$ . A positive semi-definite matrix can be written as  $X\bar{X}'$  for some matrix  $X$ , hence  $|A| \geq 0$

Positive definite matrix: a symmetric matrix  $A$  is called positive definite if  $\bar{x}'A\bar{x} > 0$  for any non-zero  $\bar{x} \in \mathbb{R}^p$  vector. Again,  $A$  can be written as  $X\bar{X}'$ , but in this case  $X$  is nonsingular,  $|A| > 0$

## Rank of Matrix

Denoted  $P(A)$ , is the number of linearly independent rows of the matrix

Can use elementary row operations to get the row-echelon form (or reduced row-echelon form). The rank of the matrix is the number of non-zero rows of the reduced row-echelon form

A square matrix  $A_{p \times p}$  is said to be of full rank if  $|A|=p$ . That is all rows are linearly independent

Full rank  $\Leftrightarrow |A| > 0 \Leftrightarrow$  positive-definite

## Inverse of Matrix

- A square matrix is invertible (or non-singular) if it is of full rank. Otherwise, it is singular
- For a non-singular matrix  $A_{p \times p}$ , there exists a unique matrix (inverse), often denoted by  $A^{-1}$  s.t.  $AA^{-1} = I_p$
- Can obtain inverse of  $A$  using reduced row echelon form  $(|A|I_p) \rightarrow [I_p | A^{-1}]$
- Properties of inverse:
  - $(A')^{-1} = (A^{-1})'$ ,  $(A^{-1})^{-1} = A$
  - $(AB)^{-1} = B^{-1}A^{-1}$ ,  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$
  - $(kA)^{-1} = \frac{1}{k} A^{-1}$  for  $k \neq 0$ ,  $|A^{-1}| = \frac{1}{|A|}$

## Generalized Inverses

- Denoted by  $A^-$  s.t.  $AA^-A = A$
- Not usually unique
- Moore-Penrose inverse is unique:
  - Moore-Penrose Conditions
    - $AA^-$  and  $A^-A$  are symmetric
    - $AA^-A = A$  and  $A^-AA^- = A^-$

## Null Matrix

- Denoted by  $A^0$  s.t.  $AA^0 = 0$
- $I - A(A^0)^T A'$  is a null matrix for  $A$

## Eigenvalues and Eigenvectors

- They play important roles in multivariate analysis
- Eigenvalues of a square matrix  $A_{p \times p}$  are roots of  $|A - \lambda I_p| = 0$
- Eigenvector  $\bar{v}$  corresponding to eigenvalue  $\lambda$  of  $A$  can be found by solving  $(A - \lambda I_p)\bar{v} = 0$
- Condition number:  $\frac{\lambda_{\max}}{\lambda_{\min}}$
- Condition index:  $\frac{\lambda_{\max}}{\lambda_{\min}}$
- A full rank matrix has all non-zero eigenvalues
- Large condition number means we have an ill-conditioned problem
- Determinant and Trace
  - $\text{trace}(A) = \sum_{i=1}^p \lambda_i$
  - $|A| = \prod_{i=1}^p \lambda_i$
- $\lambda_i$ : eigenvalues of  $A_{p \times p}$

## Matrix Decompositions

- Singular Value Decomposition: For any matrix  $A$ , there exist two orthogonal matrices  $U$  and  $V$  and a diagonal matrix  $D$  such that  $A = UDV'$ .
  - Diagonal values of  $D$  are singular values of  $A$
  - Columns of  $U$  and  $V$  are called left and right singular vectors
- For a symmetric matrix  $A$ , there exists an orthogonal matrix  $T$  and a diagonal matrix  $\Lambda$  such that  $A = T\Lambda T'$ . The diagonal elements of  $\Lambda$  are the eigenvalues of  $A$  and the columns of  $T$  are the corresponding eigenvectors. This is known as eigen or spectral decomposition
- Cholesky Decomposition: For a symmetric positive definite matrix  $A$ , we have  $A = TT'$ ,  $T$  lower triangular

## Quadratic forms

If we have a vector  $\bar{x} \in \mathbb{R}^n$ , then we have  $\bar{x}'A\bar{x} \in \mathbb{R}$  is a quadratic form for some  $A \in \mathbb{R}^{n \times n}$

If we want to write  $\sum_{i=1}^n (x_i - \bar{x}_i)^2$  in matrix form, we use the following:

- $\bar{x}$  is a vector of size  $n$  where all elements are equal to 1
- $J_{nn}$  is an  $n \times n$  matrix where all elements are equal to 1 (write as  $J$  for simplicity)
- $J^2 = nJ$
- If we let  $J' = \frac{1}{n}J$ , we get  $J'^2 = J'$
- Define  $C = I - \frac{1}{n}J = I - J'$
- We get  $\bar{x}'C\bar{x} = \sum_{i=1}^n (x_i - \bar{x}_i)^2$

## Some common Matrix Derivatives

If we have a column vector

$$\bar{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

Then we have

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = \begin{pmatrix} \frac{\partial f(\bar{a})}{\partial a_1} \\ \frac{\partial f(\bar{a})}{\partial a_2} \\ \vdots \\ \frac{\partial f(\bar{a})}{\partial a_n} \end{pmatrix}$$

For  $f(\bar{a}) = \bar{c}'\bar{a}$ , where  $\bar{c}$  is a column vector of the same length as  $\bar{a}$ . We then get

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = \bar{c}$$

For  $f(\bar{a}) = \bar{a}'B\bar{a}$ ,  $B = (b_{ij}) = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{pmatrix}$  a symmetric matrix. We have

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = 2B'\bar{a} = 2Ba$$

## Probability Distributions

- If  $X$  is a random variable, its cumulative distribution function (cdf) is  $F(x) = P(X \leq x)$  for all  $x \in (-\infty, \infty)$

## Discrete Random Variables

- A random variable  $X$  is called a discrete random variable if it assumes only finite or denumerably infinite number of values
- $F_X(x)$  is a step function
- The probability mass function of a discrete random variable  $X$  that takes values  $x_i$ , where  $i=1, 2, 3, \dots$ , is  $P(X=x_i) = P(X=x_i)$  for all  $i$
- Common Discrete Distributions: Discrete Uniform, Bernoulli, Binomial, Poisson, Geometric, Negative Binomial

## Continuous Random Variables

- A random variable is continuous if it takes on any real numbers within its support
- A random variable is continuous if  $F(x)$  is continuous
- The probability density function  $f_X(x)$  for a continuous random variable is a function such that
  - $F(x) = \int_{-\infty}^x f_X(t)dt$  for  $x \in (-\infty, \infty)$
  - or  $\frac{d}{dx} F_X(x) = f_X(x)$
  - $f_X(x) \geq 0$  for all  $x$
  - $\int_{-\infty}^{\infty} f_X(x)dx = 1$
- Common Continuous Distributions: Normal, Chi-Square, Exponential, Gamma, Beta, Student t, F

## Expectation and Moments

- For a continuous random variable  $X$ , the expected value is

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx$$

- For a discrete random variable  $X$ , the expected value is

$$E[X] = \sum x_i P(X=x_i) \text{ for all } x_i$$

- Also known as the mean of distribution
- Is a measure of central tendency
- Other measures of central tendency include median and mode
- For a random variable  $X$ , the  $r$ th non-central moment is defined as  $E[X^r]$
- The expected value is therefore the first non-central moment

The  $r$ th central moment is defined as  $E[(X-EX)]^r$

The variance of a random variable  $X$  is defined as its second central moment, or  $E[(X-EX)^2]$

The square root of the variance is referred to as the standard deviation

Other measures of variability include interquartile range (IQR) and range

### Normal Distribution

- Symmetric, unimodal

- Closed under linear combinations, conditioning, and marginalizing

- Most commonly used distributions are transformations of the normal distribution

- CLT

- Density function

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

$-\infty < \mu < \infty, \sigma > 0$

- $E[X] = \mu$  and  $\text{Var}[X] = \sigma^2$

Distribution characterized by its mean and variance, written  $X \sim N(\mu, \sigma^2)$

- Density usually denoted  $\phi_{\mu, \sigma^2}(x)$  and CDF  $\Phi_{\mu, \sigma^2}(x)$

- $\frac{X-\mu}{\sigma} \sim N(0, 1)$  standard normal, pdf  $\phi(x)$  and cdf  $\Phi(x)$

### Relationship of Normal with Other Distributions

- If  $Z \sim N(0, 1)$ , then  $Z^2 \sim \chi^2(1)$

- If  $Z_1, Z_2, \dots, Z_n \sim N(0, 1)$ , then  $\sum_{i=1}^n Z_i^2 \sim \chi^2(n)$

- If  $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$ , then we have  $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$

- If  $Z \sim N(0, 1)$  and  $V \sim \chi^2(n)$ , then

$$\frac{Z}{\sqrt{V/n}} \sim t(n)$$

- If  $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$ , then

$$\frac{\bar{X}-\mu}{\frac{s}{\sqrt{n}}} \sim t(n-1)$$

- If  $V \sim \chi^2(r)$  and  $W \sim \chi^2(s)$ , then

$$\frac{V}{r} \sim F(r, s)$$

### Random Samples and Statistics

- A random sample is a sequence of independent and identically distributed random variables from a given distribution (i.e.  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f(x)$ )

- A statistic is a function of a random sample that is observable

$$T = f(X_1, X_2, \dots, X_n)$$

$T = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  is a statistic

$T = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)$  is not a statistic

- The sampling distribution of a statistic

i.e.  $\bar{X} \sim N(\mu, \sigma^2/n)$

### Estimation Methods

- Method of moments: equate sample and population moments of the same order, then solve for the estimators of parameters

- Maximum likelihood estimate: The likelihood function is

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

Then we solve for  $\hat{\theta}$  that maximizes  $L(\theta)$

- Unbiased estimator:  $E[\hat{\theta} - \theta] = 0$

- Confidence interval: We say that  $[L, U]$  is a  $100(1-\alpha)\%$  confidence interval for a parameter  $\theta$  if  $P(L < \theta < U) = 1 - \alpha$

### Hypothesis Testing

- We usually formulate hypotheses

$$H_0: \theta = \theta_0$$

$$H_1: \theta \neq \theta_0 \text{ or } \theta > \theta_0 \text{ or } \theta < \theta_0$$

- Compute a test statistic  $T$  and the p-value is  $P(T > t_c)$ , where  $t_c$  is the critical value. Reject  $H_0$  if the p-value is less than a significance level  $\alpha$

### Errors

	Accept	Reject
$H_0$ True	✓	Type I Error
$H_0$ False	Type II Error	✓

### Random Vectors

- Random vectors are defined as a vector of random variables  $\vec{X} = (X_1, X_2, \dots, X_p)$  and are associated with joint distributions, which is multivariate in nature

- The vector of random variables might consist of different outcomes taken from the same individuals, or a single variable/outcome measured at multiple time points

- If  $\vec{X} = (X_1, X_2, \dots, X_p)$  is a random vector, its joint distribution is given by

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_p \leq x_p)$$

for  $-\infty < x_i < \infty$  and  $i=1, 2, \dots, p$

- The expected value of a random vector is given by

$$E[\vec{X}] = (E[X_1], E[X_2], \dots, E[X_p]) \\ = (\mu_1, \mu_2, \dots, \mu_p) = \vec{\mu}$$

- For a random vector  $\vec{X} = (X_1, X_2, \dots, X_p)$  with expected value  $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_p)$ , the 2nd central moment, also known as the dispersion matrix, is given by

$$E[(\vec{X} - \vec{\mu})(\vec{X} - \vec{\mu})']$$

Note: if  $\vec{X}$  is a column vector, then it is given by

$$E[(\vec{X} - \vec{\mu})(\vec{X} - \vec{\mu})']$$

- Often denoted by

$$\Sigma_{pp} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} & \dots & \sigma_{2p} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} & \dots & \sigma_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \sigma_{p3} & \dots & \sigma_{pp} \end{pmatrix}$$

where  $\sigma_{ii} = E[(X_i - \mu_i)^2] = \sigma_i^2$  and  $\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)]$

### Correlation matrix

$$\rho = (\rho_{ij}) = \frac{E[(X_i - \mu_i)(X_j - \mu_j)]}{\sqrt{\sigma_i^2 \sigma_j^2}}$$

- Both  $\Sigma$  and  $\rho$  are symmetric and positive semi-definite matrices, in most practical cases they're positive definite

- The covariance matrix is often estimated by the sample moments corresponding to each of its elements

- The unbiased estimate of  $\Sigma$  is referred to as the sample covariance matrix

### Moment Generating Functions

- For a univariate random variable

$$M_x(t) = E[e^{tX}]$$

- The mgf uniquely identifies distributions

- We also have

$$\frac{d^r}{dt^r} M_x(t)|_{t=0} = E[X^r]$$

- The mgf of a random vector of length  $p$  (row vector) is defined as

$$M_{\vec{X}}(t) = M_{\vec{X}}(t_1, t_2, \dots, t_p) = E[e^{t^T \vec{X}}] = E[e^{t_1 X_1 + t_2 X_2 + \dots + t_p X_p}]$$

- We can obtain univariate mgf's by

$$M_{X_i}(t_i) = M_{\vec{X}}(0, 0, \dots, t_i, 0, \dots, 0)$$

### Random Matrices

- Matrix random variables are defined as

$$X_{pp} = \begin{pmatrix} X_{11} & X_{12} & X_{13} & \dots & X_{1p} \\ X_{21} & X_{22} & X_{23} & \dots & X_{2p} \\ X_{31} & X_{32} & X_{33} & \dots & X_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ X_{p1} & X_{p2} & X_{p3} & \dots & X_{pp} \end{pmatrix}$$

where the elements  $X_{ij}$ 's are themselves random variables

- The expectation for matrix random variables is

$$E[\vec{X}] = \begin{pmatrix} E[X_{11}] & E[X_{12}] & E[X_{13}] & \dots & E[X_{1p}] \\ E[X_{21}] & E[X_{22}] & E[X_{23}] & \dots & E[X_{2p}] \\ E[X_{31}] & E[X_{32}] & E[X_{33}] & \dots & E[X_{3p}] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ E[X_{p1}] & E[X_{p2}] & E[X_{p3}] & \dots & E[X_{pp}] \end{pmatrix}$$

However, we run into problem when we try to consider higher order moments for matrix random variables

- The most common approach is to use the vectorized form of the matrix  $\vec{X}$ , where the covariance matrix is represented as  $\Sigma \otimes \Psi$ , where  $\Sigma$  represents the covariance for column vectors and  $\Psi$  represents the covariance for row vectors

### Multivariate Bias

- Suppose  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_p)$  is an estimator for a vector parameter  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$ . Multivariate bias is defined as

$$E[\hat{\theta} - \theta] = (E[\hat{\theta}_1] - \theta_1, E[\hat{\theta}_2] - \theta_2, \dots, E[\hat{\theta}_p] - \theta_p)$$

- An estimator is referred to as an unbiased estimator if  $E[\hat{\theta} - \theta] = \vec{0}$

- Although the definition of bias is quite natural, we face a challenge when we attempt to compare two vector estimators with respect to their respective bias vectors

- For each vector estimator, we can calculate

$$D = \sqrt{(E[\hat{\theta}] - \theta)^T (E[\hat{\theta}] - \theta)'} = \sqrt{\sum_{i=1}^p (E[\hat{\theta}_i] - \theta_i)^2}$$

Compare the values for both estimators

- Can also use absolute bias across elements of the vector

### Multivariate MSE

- For univariate estimators, it is defined as  $E[(\hat{\theta} - \theta)^2]$  and can be rewritten as

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta} - \theta])^2 = \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2$$

- For vector estimators, the MSE is defined as

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)'(\hat{\theta} - \theta)]$$

This is an outer product, hence  $\text{MSE}(\hat{\theta})$  is a matrix.

- We can rewrite the MSE as

$$\text{MSE}(\hat{\theta}) = \Sigma_{\hat{\theta}} + \text{Bias}(\hat{\theta})' \text{Bias}(\hat{\theta})$$

$\Sigma_{\hat{\theta}}$ : variance covariance matrix of  $\hat{\theta}$

- If we want to compare 2 vector estimators, we vectorize the matrix and calculate the Euclidean distance to compare MSE for two vector estimators

- $\text{MSE}(\hat{\theta})$  involves covariances in addition to variances, hence difficult to interpret

- To overcome difficulty in interpretation, we can use another definition of MSE

$$\text{MSE}^*(\hat{\theta}) = (\text{Var}(\hat{\theta}_1) + \text{Bias}(\hat{\theta}_1)^2, \dots, \text{Var}(\hat{\theta}_p) + \text{Bias}(\hat{\theta}_p)^2)$$

This corresponds to the sum of the diagonal elements of  $\Sigma_{\hat{\theta}}$  and the diagonal elements of  $\text{Bias}(\hat{\theta})' \text{Bias}(\hat{\theta})$ . Also equivalent to vector of element-wise MSEs

- Also easier to extend to matrix parameters

$$\text{MSE} = \begin{pmatrix} \text{Var}(\hat{\theta}_1) + \text{Bias}(\hat{\theta}_1)^2 & \dots & \text{Var}(\hat{\theta}_p) + \text{Bias}(\hat{\theta}_p)^2 \\ \vdots & \ddots & \vdots \\ \text{Var}(\hat{\theta}_{p1}) + \text{Bias}(\hat{\theta}_{p1})^2 & \dots & \text{Var}(\hat{\theta}_{pp}) + \text{Bias}(\hat{\theta}_{pp})^2 \end{pmatrix}$$

Multivariate bias for matrix estimators is straightforward

$$\text{Bias} = \begin{pmatrix} E[\hat{\theta}_{11} - \theta_{11}] & \dots & E[\hat{\theta}_{1p} - \theta_{1p}] \\ \vdots & \ddots & \vdots \\ E[\hat{\theta}_{p1} - \theta_{p1}] & \dots & E[\hat{\theta}_{pp} - \theta_{pp}] \end{pmatrix}$$

- MSE for vector estimators are vectors. Use distance or average of element-wise MSEs to compare estimators

- Both bias and MSE for matrix estimators are matrices. We vectorize the matrices first and do the same as with vector estimators to compare estimators

### Joint and Marginal Distributions

- If we have a random vector  $\vec{X} = (X_1, X_2, \dots, X_p)$ , the joint CDF is given by

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_p \leq x_p)$$

- The joint density function  $f_{\vec{X}}(x_1, x_2, \dots, x_p)$  is defined in a way such that

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \dots \int_{-\infty}^{x_p} f_{\vec{X}}(x_1, x_2, \dots, x_p) dx_1 dx_2 \dots dx_p$$

Similarly, the joint density function can be found from the CDF by differentiating w.r.t. all elements of the p-vector

- The marginal distribution of  $X_i$  in a p-variate random vector is obtained by taking p-1 integrals w.r.t. all the remaining random variables

### Conditional Distributions and Independence

- Consider a p-variate vector  $\vec{X} = (X_1, X_2, \dots, X_p)$ , the conditional distribution of a sub-vector  $\vec{X}' = (X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_p)$  given  $X_i$  is

$$f_{\vec{X}'}|_{X_i}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_p | x_i) = \frac{f_{\vec{X}}(x_1, x_2, \dots, x_p)}{f_{X_i}(x_i)}$$

Conditional distribution of a sub-vector of random variables conditioned on another sub-vector of random variables can also be defined in a similar way.

A vector of random variables  $\vec{X} = (X_1, \dots, X_p)$  are independent if and only if the joint distribution factorizes. That is

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = \prod_{i=1}^p F_{X_i}(x_i)$$

The joint density also factorizes

$$f_{\vec{X}}(x_1, x_2, \dots, x_p) = \prod_{i=1}^p f_{X_i}(x_i)$$

Alternatively, random variables are independent iff the joint MGF factorizes

$$M_{\vec{X}}(\vec{t}) = \prod_{i=1}^p M_{X_i}(t_i)$$

We also have if a vector of random variables  $\vec{X} = (Y_1, Y_2, \dots, Y_p)$  are independent, then

$$E[g(X_1)g(X_2)\dots g(X_p)] = \prod_{i=1}^p E[g(X_i)]$$

Independence also means that conditional distributions are the same as unconditional distributions

### Properties of Expectations and Covariances

Consider a random vector  $\vec{X} = (X_1, X_2, \dots, X_p)$  and let  $E[\vec{X}] = \vec{\mu}$ , then

$$E[A\vec{X}'] = AE[\vec{X}] = A\vec{\mu}'$$

for any non-random matrix  $A: q \times p$ ,  $q$  any constant

Similarly

$$E[\vec{X}B] = E[\vec{X}]B = \vec{\mu}B$$

for any non-random matrix  $B:p \times r$ ,  $r$  any constant

$$E[\vec{X}'] = E[\vec{X}]'$$

Let  $Cov(\vec{X}) = \Sigma$ , then

$$\text{Cov}(A\vec{X}') = ACov(\vec{X})A' = A\Sigma A'$$

for any non-random matrix  $A: q \times p$

Similarly, we can show that

$$\text{Cov}(\vec{X}B) = B'\text{Cov}(\vec{X})B = B'\Sigma B$$

for any non-random matrix  $B:p \times r$

$$\text{Cov}(\vec{X}') = \text{Cov}(\vec{X})$$

Covariance of a non-random vector is 0

Covariance matrix is always positive semi-definite

### Multivariate Normal Distribution

Univariate normal distribution density

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

A random vector  $\vec{X} = (X_1, X_2, \dots, X_p)$  has a multivariate normal distribution iff every linear combination of  $\vec{X}$  has a univariate normal distribution

Multivariate normal distribution can also be defined with the aid of the univariate standard normal distribution, and by taking advantage of the fact that the normal distribution belongs to the location-scale family

We first consider a random vector  $\vec{U} = (U_1, U_2, \dots, U_p)$ ,  $U_i \sim N(0, 1)$ . We have

$$E[\vec{U}] = \vec{0} \text{ and } \text{Cov}(\vec{U}) = I_p$$

The distribution of  $\vec{U}'$  is the same as the joint distribution of the  $U_i$ 's and is given by

$$f_{\vec{U}}(u_1, u_2, \dots, u_p) = \prod_{i=1}^p f_{U_i}(u_i) = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\sum_{i=1}^p u_i^2}$$

This can also be written as

$$\begin{aligned} & \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}\vec{u}')} \\ & = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}'\vec{u})} \\ & = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}'\vec{u})} \end{aligned}$$

This distribution is referred to as the standard  $p$ -variate normal distribution, often denoted by  $\vec{U} \sim N_p(\vec{0}, I_p)$ .

A  $p$ -dimensional random vector  $\vec{X} = (X_1, X_2, \dots, X_p)$  is said to have a multivariate normal distribution if it has the same distribution as  $\vec{U} + \vec{U}'B$ , where  $\vec{U} \sim N_p(\vec{0}, I_p)$ ,  $B$  is a non-singular and symmetric matrix of dimension  $p \times p$  and  $\Sigma = B'B = BB$ .

The multivariate normal distribution is often denoted as  $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$  and its density function is given as

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\vec{X}-\vec{\mu})'\Sigma^{-1}(\vec{X}-\vec{\mu})'}$$

This can be derived using change-of-variables

$B$  is referred to as the square root of  $\Sigma$ , denote it by  $\Sigma^{\frac{1}{2}}$

The density function can also be written as

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2}\text{trace}(\Sigma(\vec{X}-\vec{\mu})')^2}$$

Using the trace property, we can also write

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2}\text{tr}(\Sigma(\vec{X}-\vec{\mu})')^2}$$

When  $p=1$ , the above reduces to univariate normal density

### Properties of Multivariate Normal

Consider  $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$ , we have

$$\begin{aligned} E[\vec{X}] &= \vec{\mu} \\ \text{Cov}(\vec{X}) &= \Sigma \end{aligned}$$

We can write  $\Sigma$  as

$$\Sigma_{p,p} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$$

The multivariate normal distribution can also be defined using MGFs

MGF of univariate standard normal:  $M_U(t) = e^{\frac{1}{2}t^2}$

If we have a random vector  $\vec{U} = (U_1, U_2, \dots, U_p) \sim N_p(0, I_p)$ , its mgf is  $M_{\vec{U}}(t) = e^{\frac{1}{2}\vec{t}'\vec{t}}$

$$\vec{t} = (t_1, t_2, \dots, t_n)$$

If we have  $\vec{X} = \vec{\mu} + \vec{U}\Sigma$ , we get

$$M_{\vec{X}}(t) = e^{\vec{\mu}'\vec{t} + \frac{1}{2}\vec{t}'\Sigma\vec{t}}$$

This is the mgf for  $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$

Normal distribution is closed under linear combination

Theorem: If  $\vec{X} = (X_1, X_2, \dots, X_p) \sim N_p(\vec{\mu}, \Sigma)$  and  $C:p \times m$ , then

$$\vec{X}C \sim N_p(\vec{\mu}C, C'\Sigma C)$$

Can be used to prove normality of individual elements

If we have  $k$  independent  $p$ -variate normal random vectors  $\vec{X}_i \sim N_p(\vec{\mu}_i, \Sigma_i)$  for  $i=1, 2, \dots, k$  and  $A_i:p \times p$ . Then

$$\sum_{i=1}^k \vec{X}_i A_i \sim N_p\left(\sum_{i=1}^k \vec{\mu}_i A_i, \sum_{i=1}^k \vec{\Sigma}_i A_i A_i'\right)$$

### Limit Theorems

$X_1, X_2, \dots, X_n$  any random sample

$$E[X_i] = \mu \quad \text{Var}[X_i] = \sigma^2$$

Univariate Weak Law of Large Numbers:

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

Univariate Central Limit Theorem:

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

Now, consider iid  $p$ -variate random vectors  $\vec{X}_i$  with mean  $\vec{\mu}$  and covariance matrix  $\Sigma$

Multivariate Weak Law of Large Numbers

$$\bar{X} = (X_1, X_2, \dots, X_p) \xrightarrow{P} \vec{\mu}$$

Multivariate CLT

$$\sqrt{n}(\bar{X} - \vec{\mu}) \xrightarrow{d} N_p(0, \Sigma)$$

### Wishart Distribution

For univariate random variables, we have

$$X \sim N(0, 1) \Rightarrow X^2 \sim \chi^2(1)$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(0, 1) \Rightarrow \sum_{i=1}^n X_i^2 \sim \chi^2(n)$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \Rightarrow \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

The Wishart distribution is a multivariate extension of chi-squared distribution. Let

$$\vec{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} \sim N_p(0, I_p)$$

Then we get that  $W = \vec{X}'\vec{X}$  has the Wishart distribution, written as  $w_p(I, I_p)$

If  $\vec{X} \sim N_p(0, \Sigma)$ , then  $W = \vec{X}'\vec{X} \sim w_p(I, \Sigma)$

If  $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$ , then  $W = \vec{X}'\vec{X}$  has the non-central Wishart distribution

If we let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$  be the columns of  $\vec{X}$ , we can write  $\vec{X} \sim N_p(\vec{\mu}, \Sigma, I_n)$  where  $\vec{\mu}' = (\mu_1, \mu_2, \dots, \mu_p)$

A random matrix  $W:p \times n$  is said to have the Wishart distribution if and only if it can be written as  $W = \vec{X}'\vec{X}$ , where  $\vec{X} \sim N_p(\vec{\mu}, \Sigma, I_n)$ ,  $\Sigma$  positive definite nonsingular matrix. This is written as  $w_p(n, \Sigma, \Delta)$ ,  $\Delta = \vec{\mu}'\vec{\mu}$  non-centrality parameter. When  $\vec{\mu} = 0$ , we have the central wishart distribution  $w_p(n, \Sigma)$

We can write  $W = XX' = \sum_{i=1}^n X_i X_i'$   $X_i \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$

Elements of  $W$  are sums of independent random variables. Diagonal elements have  $X_i'^2(n)$  distribution

Wishart density:

$$f_W(w) = c |w|^{(n-p)/2} |w|^{-\frac{1}{2}\text{trace}(\Sigma^{-1}w)}$$

$$c = \left(\frac{1}{2}\right)^{np} T_p\left(\frac{n}{2}\right)^{-1} T_p(\cdot) \text{ multivariate Gamma}$$

Properties:

If  $W_1 \sim w_p(n, \Sigma, \Delta_1)$  and  $W_2 \sim w_p(m, \Sigma, \Delta_2)$  and independent, then we have  $W_1 + W_2 \sim w_p(n+m, \Sigma, \Delta_1 + \Delta_2)$

If  $W \sim w_p(n, \Sigma, \Delta)$  and  $A:r \times p$ , then  $AWA' \sim w_r(n, A\Sigma A', \Delta A A')$

### Sampling from Multivariate Normal

If we let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$ , we get  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N_p(\vec{\mu}, \frac{1}{n} I_p)$

We can also get if  $S = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$ , we get  $(n-1)S \sim w_p(n-1, \Sigma)$

We also have that  $\bar{X}$  and  $S$  are jointly sufficient for  $\vec{\mu}$  and  $\Sigma$

### Estimation of $\vec{\mu}$ and $\Sigma$

We can obtain that the method of moments estimators for  $\vec{\mu}$  and  $\Sigma$  are  $\hat{\vec{\mu}} = \bar{X}$  and  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$  by solving the system

$$E[X_i] = \bar{X} \text{ and } E[X_i X_i'] = \frac{1}{n} \sum_{i=1}^n X_i X_i'$$

If we want to do MLE, the likelihood is

$$L(\vec{\mu}, \Sigma | X) = \left(\frac{1}{2\pi}\right)^p |\Sigma|^{\frac{1}{2}} \exp\left(-\frac{1}{2}\sum_{i=1}^n (X_i - \bar{X})'\Sigma^{-1}(X_i - \bar{X})\right)$$

The log-likelihood is given by

$$\ell(\vec{\mu}, \Sigma | X) = -n\log(\sqrt{2\pi}) - \frac{1}{2} \log|\Sigma| - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})'\Sigma^{-1}(X_i - \bar{X})$$

If we take derivatives wr.t.  $\vec{\mu}$ , we get

$$\frac{\partial \ell}{\partial \vec{\mu}} = -\frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})$$

If we set this to 0, we get  $\hat{\vec{\mu}} = \bar{X}$

To get the MLE for  $\Sigma$ , we need the following matrix derivatives

$$\frac{\partial |A|}{\partial A} = (A')^{-1}|A| \text{ and } \frac{\partial \text{trace}(AB)}{\partial B} = A'$$

Rewrite the log-likelihood as

$$\ell(\vec{\mu}, \Sigma | X) = -n\log(\sqrt{2\pi}) + \frac{1}{2} \log|\Sigma| - \frac{1}{2} \text{trace}((X - \bar{X})(X - \bar{X})')$$

We can get

$$\frac{\partial \ell}{\partial \Sigma^{-1}} = \frac{n}{2} \frac{1}{|I_p|} \cdot \Sigma \cdot |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

If we set this to 0, we get

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

Note:  $\hat{\Sigma}$  is not an unbiased estimator for  $\Sigma$ .

### Hypothesis Testing for $\vec{\mu}$

For univariate case, if we have a sample from  $N(\mu, \sigma^2)$ , where  $\sigma^2$  is fixed but unknown. If we want to test  $H_0: \mu = \mu_0$  vs.  $H_1: \mu \neq \mu_0$

We reject  $H_0$  when

$$T = \left| \frac{\bar{X} - \mu_0}{\frac{1}{\sqrt{n}}} \right| > t_{\alpha/2}(n-1)$$

We also have  $T^2 \sim F_{1, n-1}$ , so we can also reject  $H_0$  when  $T^2 > F_{1, n-1}(1, n-1)$

For multivariate normal, if we want to test

$$H_0: \vec{\mu} = \vec{\mu}_0 \text{ vs. } H_1: \vec{\mu} \neq \vec{\mu}_0$$

we use the Hotelling's  $T^2$  statistic

$$T^2 = n(\bar{X} - \vec{\mu}_0)' S^{-1}(\bar{X} - \vec{\mu}_0)$$

- The Hotelling's  $T^2$  statistic is also the likelihood ratio test statistic
- The distribution for  $T^2$  is not a standard distribution, so we need to scale it
- Roy's Union-Intersection Test statistic  

$$RF = \left( \frac{n-p}{p(n-p)} \right) T^2 \sim F(p, n-p)$$
- We reject  $H_0$  when  $RF > F_{1-\alpha}(p, n-p)$

### Testing Multivariate Normality

- One way is to vectorize the random sample  $X = (X_1, X_2, \dots, X_n)$ .
- We have if  $X_i \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ , then  $\text{Vec}(X) \sim N_{pn}(\text{Vec}(\mu), \Sigma \otimes I_n)$ . Then we have

$$\text{Vec}(X) \sim N_{pn}(\text{Vec}(\mu), \Sigma \otimes I_n)$$

- Then we can use univariate tests and traditional graphical methods like q-q plot to test for normality

- Limitation: Vectorized data does not necessarily consist of independent samples
- Another way is to use Small's graphical method. Let  $c_i$  be defined as

$$c_i = \frac{n}{(n-1)^2} (x_i - \bar{x})' S^{-1} (x_i - \bar{x})$$

- When  $X_i \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ , it can be shown that  $c_i$ 's are approximately distributed as  $\text{Beta}(\frac{n}{2}, \frac{n-1-p}{2})$ .
- We can use Kolmogorov-Smirnov's test as well as graphical approaches to test

### Two Sample Test

- Let  $X = (X_1, X_2, \dots, X_m)$ ,  $Y = (Y_1, Y_2, \dots, Y_n)$  random samples from two independent  $p$ -variate normal distributions with common covariance matrix  $\Sigma$ . Say  $N_p(\mu_1, \Sigma)$  and  $N_p(\mu_2, \Sigma)$ , respectively.
- We want to test the following:

$$H_0: \mu_1 = \mu_2 \quad \text{vs.} \quad H_1: \mu_1 \neq \mu_2$$

- For univariate case, we reject  $H_0$  when

$$T = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} \left| \frac{\bar{y} - \bar{x}}{s_p} \right| > t_{1-\alpha/2}(n_1 + n_2 - 2)$$

$$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$$

$s_1^2, s_2^2$  sample variance for groups

- We can also reject  $H_0$  when  

$$T^2 > F_{1-\alpha}(1, n_1 + n_2 - 2)$$

- For multivariate case, we have the two sample Hotelling's  $T^2$  statistic is

$$T^2 = \left( \frac{n_1 n_2}{n_1 + n_2} \right) (\bar{X} - \bar{Y})' S_p^{-1} (\bar{X} - \bar{Y})$$

$$S_p = \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1 + n_2 - 2}$$

$S_1, S_2$  sample covariance matrices for groups

- Two-sample Roy's UIT

$$RF = \left( \frac{f-p+1}{f-p} \right) T^2 \sim F(p, f-p+1)$$

$$f = n_1 + n_2 - 2$$

Reject  $H_0$  when  $RF > F_{1-\alpha}(p, f-p+1)$

### Testing for Normality

- Shapiro-Wilk's test: Consider a univariate random sample  $X_1, X_2, \dots, X_n$  and suppose  $Y_1, Y_2, \dots, Y_n$  are the corresponding order statistics. The Shapiro-Wilk's test for normality is given by

$$W = \frac{\sum_{i=1}^n a_i Y_i}{S}$$

$S$  sample standard deviation  
 $a_i$  function of order statistic of standard normal distribution

- Kolmogorov-Smirnov's test: uses the difference between the empirical distribution of two random samples, and reject when the difference is large

$$D = \sup_x |F_U(x) - F_W(x)|$$

reject when  $D$  is large

- To test for normality, compare the empirical distribution corresponding to the data to standard normal

- In practice, we usually use graphical methods, by themselves or combined with formal tests