

Topic 2: Matrix Algebra, Random Vectors, and Multivariate Normal

Reference:

1. Searle, S. R. (1982). Matrix Algebra useful for statistics. Wiley, Inter Science.
2. Horn, R. A. and Johnson, C. R. (1985). Matrix Analysis. Cambridge, University Press.

Matrix: A $n \times p$ matrix X is a rectangular array of numbers, which can be expressed as

$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \cdot & \cdot & \cdot & x_{1p} \\ x_{21} & x_{22} & \cdot & \cdot & \cdot & x_{2p} \\ \cdot & \cdot & & & \cdot & \cdot \\ \cdot & \cdot & & & \cdot & \cdot \\ \cdot & \cdot & & & \cdot & \cdot \\ x_{n1} & x_{n2} & \cdot & \cdot & \cdot & x_{np} \end{bmatrix} = (x_{ij}) = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ \vdots \\ x_n^T \end{pmatrix} = (x_{(1)}, x_{(2)}, \dots, x_{(p)}),$$

where $x_i = (x_{i1}, \dots, x_{ip})^T$ and $x_{(j)} = (x_{1j}, \dots, x_{nj})^T$, $i = 1, \dots, n$, $j = 1, \dots, k$.

Square Matrix: A $n \times p$ matrix X is called a square matrix if $n = p$.

Unit Vector: A $p \times 1$ vector with all numbers being one is called a unit vector and is denoted by $\mathbf{1}_p$.

Diagonal Matrix: A $p \times p$ diagonal matrix, denoted by D , is a square matrix with off-diagonal elements being zero.

$$D = \text{diag}(a_{11}, \dots, a_{pp}) = \begin{pmatrix} a_{11} & & & & \\ & a_{22} & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & a_{pp} \end{pmatrix}$$

Identity Matrix: A $p \times p$ identity matrix, denoted by I_p , is a diagonal matrix with diagonal elements being one.

Symmetric Matrix: A square ($p \times p$) matrix \mathbf{A} is symmetric if $A = A^T$, where “ T ” denotes the transpose.

Triangular Matrix: Upper triangular matrix $U = \begin{pmatrix} & & & \\ & & \text{---} & \\ & \text{---} & \text{---} & \\ & & & \text{---} \end{pmatrix}$

Lower triangular matrix $L = \begin{pmatrix} & & & \\ & & & \\ & & \text{---} & \\ & & & \text{---} \end{pmatrix}$

Trace: Let $A = (a_{ij})$ be a $p \times p$ matrix. The trace of A is defined as $\text{tr}(A) = \sum_{i=1}^p a_{ii}$.

Idempotent Matrix: A square matrix A is called an idempotent matrix if $A^2 = A$, where $A^2 = AA$.

Determinant and Cofactor: The determinant of a square matrix A , denoted by $\det(A)$ or $|A|$, is defined as

$$\det(A) = \sum_{j=1}^p (-1)^{i+j} a_{ij} \det(A(i \mid j)), \text{ where } \det(A(i \mid j)) \text{ is the } (i, j) \text{ cofactor of } A.$$

Inverse: The inverse of a square matrix A , if exists, A^{-1} is unique and satisfies

$$A^{-1}A = AA^{-1} = I_p.$$

Nonsingular: A is nonsingular iff the inverse of A exists iff $\det(A) \neq 0$.

Orthogonal: A $p \times p$ matrix A is orthogonal iff $AA^T = I_p$, which implies that $A^T = A^{-1}$.

Rank: The rank of a $p \times q$ matrix A , denoted by $\text{rank}(A)$, is the maximum number of linearly independent rows (columns) in A .

Eigenvalues and Eigenvectors: Let A be a $p \times p$ matrix and y a nonzero $p \times 1$ vector such that $Ay = \lambda y$. Then, y is called the eigenvector of A and λ is the corresponding eigenvalue.

eigenvalue of A . Here, $P(\lambda) = \det(A - \lambda I_p)$ is the p th order polynomial in λ , which is called the characteristic polynomial of A .

Example 1. (See Appendix)

Partition Matrices:

Let $A_{m \times n} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ and $B_{n \times k} = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$, where A_{ij} and $B_{j\ell}$ are separately $m_i \times n_j$ and $n_j \times k_\ell$ matrices. Then, $AB = \begin{pmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{pmatrix}$.

Assume Moreover, assume that A is a $p \times p$ nonsingular matrix and $B = A^{-1}$. One has

$$B_{11} = (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} \quad (\triangleq A_{11,2}^{-1}), \quad B_{22} = (A_{22} - A_{21}A_{11}^{-1}A_{12})^{-1} \quad (\triangleq A_{22,1}^{-1})$$

$$B_{12} = -A_{11}^{-1}A_{12}A_{22,1}^{-1} = -A_{11,2}^{-1}A_{12}A_{22}^{-1}, \text{ and } B_{21} = -A_{22}^{-1}A_{21}A_{11,2}^{-1} = -A_{22,1}^{-1}A_{21}A_{11}^{-1}.$$

Example 2. (See Appendix)

Spectral Decomposition Theorem (Jordan Decomposition Theorem):

Any symmetric $p \times p$ matrix A can be expressed as $A = \Gamma \Lambda \Gamma^T = \sum_{i=1}^p \lambda_i r_{(i)} r_{(i)}^T$, where Λ is a diagonal matrix of eigenvalues of A and Γ is an orthogonal matrix whose columns are the corresponding standardized eigenvectors.

Singular Value Decomposition Theorem:

Let A be a $p \times q$ matrix with $\text{rank}(A) = r$. Then, A can be expressed as $A = U \Lambda V^T$, where $U_{p \times r}$ and $V_{q \times r}$ are orthonormal matrices and Λ is a diagonal matrix with positive elements such that $A^T A = V \Lambda V^T$, $AA^T = U \Lambda U^T$, and $\Lambda = \Lambda^{\frac{1}{2}}$.

Quadratic Forms and Definiteness:

Let A be a $p \times p$ matrix and y be a $p \times 1$ vector. $q = y^T A y$ is called a quadratic form in y and A is the matrix of the quadratic form.

$$\left\{ \begin{array}{l} A \text{ is positive definite (p.d.) if } y^T A y > 0 \text{ for all } y \neq 0. \\ A \text{ is positive semidefinite (p.s.d.) if } y^T A y \geq 0 \text{ and } y^T A y = 0 \text{ for some } y \neq 0. \end{array} \right.$$

Moments of Random Vectors and Matrices:

$$\text{Let } x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}, y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_q \end{pmatrix}, \text{ and } U = \begin{pmatrix} u_{11} & u_{12} & \dots & \dots & \dots & u_{1q} \\ u_{21} & u_{22} & \dots & \dots & \dots & u_{2q} \\ \vdots & & & & & \vdots \\ \vdots & & & & & \vdots \\ u_{p1} & u_{p2} & \dots & \dots & \dots & u_{pq} \end{pmatrix}.$$

$$E[x] = \begin{pmatrix} E[x_1] \\ E[x_2] \\ \vdots \\ E[x_p] \end{pmatrix}, V(x) = \begin{pmatrix} \text{var}[x_1] & \text{cov}(x_1, x_2) & \dots & \dots & \text{cov}(x_1, x_p) \\ \text{cov}(x_2, x_1) & \text{var}[x_2] & \dots & \dots & \text{cov}(x_2, x_p) \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ \text{cov}(x_p, x_1) & \text{cov}(x_p, x_2) & \dots & \dots & \text{var}[x_p] \end{pmatrix},$$

$$E[U] = \begin{pmatrix} E[u_{11}] & E[u_{12}] & \dots & \dots & E[u_{1q}] \\ E[u_{21}] & E[u_{22}] & \dots & \dots & E[u_{2q}] \\ \vdots & & & & \vdots \\ \vdots & & & & \vdots \\ E[u_{p1}] & E[u_{p2}] & \dots & \dots & E[u_{pq}] \end{pmatrix}, Cov(x, y) = E[(x - E[x])(y - E[y])^T].$$

When $x = y$, $Cov(x, y) \triangleq V(x)$.

Property 2.1.

- (a) $V(x) = E[xx^T] - E[x](E[x])^T$. $\text{Var}(x) = (\text{Cov}(x_i, x_j)) = (E[x_i x_j]) - (E[x_i] E[x_j])$
- (b) $E[Ax + b] = AE[x] + b$ and $V(Ax + b) = AV(x)A^T$, where A and b are separately a $p \times p$ constant matrix and a $p \times 1$ constant vector.

$$(c) Cov(Ax, By) = ACov(x, y)B^T.$$

$$(d) E[y^T Ay] = \text{tr}(AV(y)) + (E[y])^T A(E[y]).$$

$$= \sum \sum a_{ij} y_i y_j$$

Example 3. (See Appendix)

Example 4. (See Appendix)

Example 5. (See Appendix)

$$\begin{aligned} E[Ax + b] &= (E[a_i^T x] + b) \\ &= (a_i^T E[x]) + (b_i) \\ V(Ax + b) &= E[(Ax + b) - E(Ax + b)](Ax + b - E(Ax + b))^T \\ &= E[(A(x - E[x]))(x - E[x])^T A^T] \\ &= (E[a_i^T (x - E[x]) (x - E[x])^T a_j]) \\ &= (a_i^T V(x) a_j) \end{aligned}$$

Characteristic Function: The characteristic function of a $p \times 1$ random vector x is

defined to be $\varphi_x(t) = E[e^{it^T x}]$, where t is a $p \times 1$ constant vector in some neighborhood of $\mathbf{0}$ vector.

Vector and Matrix Differentiation:

Let $u = g(x)$.

$$\frac{\partial u}{\partial x} \triangleq \begin{bmatrix} \frac{\partial u}{\partial x_1} \\ \frac{\partial u}{\partial x_2} \\ \vdots \\ \frac{\partial u}{\partial x_p} \end{bmatrix} \quad \text{and} \quad \frac{\partial^2 u}{\partial x \partial x} \triangleq \begin{bmatrix} \frac{\partial}{\partial x_1} \left(\frac{\partial u}{\partial x_1} \right) & \frac{\partial}{\partial x_2} \left(\frac{\partial u}{\partial x_1} \right) & \cdots & \frac{\partial}{\partial x_p} \left(\frac{\partial u}{\partial x_1} \right) \\ \frac{\partial}{\partial x_1} \left(\frac{\partial u}{\partial x_2} \right) & \frac{\partial}{\partial x_2} \left(\frac{\partial u}{\partial x_2} \right) & \cdots & \frac{\partial}{\partial x_p} \left(\frac{\partial u}{\partial x_2} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} \left(\frac{\partial u}{\partial x_p} \right) & \frac{\partial}{\partial x_2} \left(\frac{\partial u}{\partial x_p} \right) & \cdots & \frac{\partial}{\partial x_p} \left(\frac{\partial u}{\partial x_p} \right) \end{bmatrix}.$$

$\frac{\partial x_j}{\partial x_i} (\frac{\partial x_i}{\partial x})$

Properties 2.2.

(a) Let $u = a^T y$, $\frac{\partial u}{\partial y} = a$.

(b) Let $u = y^T y$, $\frac{\partial u}{\partial y} = 2y$.

(c) Let $u = y^T A y$, $\frac{\partial u}{\partial y} = Ay + A^T y$. (When A is symmetric, $\frac{\partial u}{\partial y} = 2Ay$.)

$$u = y^T A y = \sum_{i,j} a_{ij} y_i y_j$$

$$\frac{\partial u}{\partial y_i} = \sum_j a_{ij} y_j + \sum_j a_{ji} y_j = a_i^T y + a_{ii} y_i.$$

Transformations of Random Vectors: Let x be a continuous $p \times 1$ random vector with p.d.f. $f_X(x)$, and let $y = u(x)$ be a 1-1 transformation. Then, the p.d.f. of y is

$$f_Y(y) = f_X(u^{-1}(y)) |\det(J)|, \text{ where } J = \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \cdots & \frac{\partial x_1}{\partial y_p} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} & \cdots & \frac{\partial x_2}{\partial y_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_p}{\partial y_1} & \frac{\partial x_p}{\partial y_2} & \cdots & \frac{\partial x_p}{\partial y_p} \end{pmatrix}.$$

$\left(\frac{\partial y_i}{\partial x_j} \right)$

Measures of Multivariate Scatter (Dispersion):

- (a) The generalized variance: $\det(V(x))$.
- (b) The total variation: $\text{tr}(V(x))$.

Linear Combinations: Let A be a $q \times p$ constant matrix; b be a $q \times 1$ constant vector; and x be a $p \times 1$ random vector with $x \sim (\mu, \Sigma)$, where Σ is a p.d. matrix. Then, $y = Ax + b \sim (A\mu + b, A\Sigma A^T)$.

Scaling transformation: $y = (\text{diag}(\Sigma))^{-1/2}(x - \mu) \sim (0, (\text{diag}(\Sigma))^{-1/2}\Sigma(\text{diag}(\Sigma))^{-1/2})$, where $P = (\text{diag}(\Sigma))^{-1/2}\Sigma(\text{diag}(\Sigma))^{-1/2}$ is a correlation matrix.

Mahalanobis transformation: $z = \Sigma^{-1/2}(x - \mu) \sim (0, I_p)$.

Distances or Dissimilarities:

Euclidean Distance: $\|x_1 - x_2\|_2 = ((x_1 - x_2)^T (x_1 - x_2))^{1/2}$.

Mahalanobis Distance: $\|z_1 - z_2\|_2 = ((x_1 - x_2)^T \Sigma^{-1} (x_1 - x_2))^{1/2}$.

Karl Pearson Distance: $\|y_1 - y_2\|_2 = ((x_1 - x_2)^T (\text{diag}(\Sigma))^{-1} (x_1 - x_2))^{1/2}$.

✳

Multivariate Normal Distribution:

Definition. The p -variate random vector X follows a multivariate normal distribution if for any $p \times 1$ constant vector b , $U = b^T x$ is either a constant or a univariate normal, ~~for any $p \times 1$ constant vector b~~ .

Theorem 2.1. If x is p -variate normal and $y = Ax + b$ for a $q \times p$ constant matrix A and a $q \times 1$ constant vector b , then y is q -variate normal.

Corollary 2.1.1. If Z is p -variate standard normal with p.d.f. $f_Z(z) = (2\pi)^{-p/2} e^{-\frac{1}{2}z^T z}$, and $Y = AZ + b$, then Y is q -variate normal with mean b and variance matrix AA^T .

Corollary 2.1.2. Let X be p -variate normal with $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_p(\mu, \Sigma)$, where X_1 and X_2 are separately $q \times 1$ and $(p - q) \times 1$ random vectors, respectively.

$$(Q_{p-q} \otimes I_{q \times (p-q)}) X$$

Let $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ and $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$. Then, $X_1 \sim N_q(\mu_1, \Sigma_{11})$ and $X_2 \sim N_{p-q}(\mu_2, \Sigma_{22})$.

$$(I_{q \times q} \otimes I_{2 \times (p-q)}) X$$

Theorem 2.2. For any $p \times 1$ constant vector μ and $p \times p$ constant matrix Σ , which is a p.s.d. matrix, there exists a unique multivariate normal distribution with mean μ and variance Σ .

Proof: (By the uniqueness of characteristic function $\varphi_x(t) = e^{it^T \mu - t^T \Sigma t/2}$ for $N_p(\mu, \Sigma)$.)

Lemma 2.1. If a $p \times p$ constant matrix Σ is p.s.d., there exists a $p \times p$ matrix M such that $\Sigma = M^2$ (M is often denoted by $\Sigma^{1/2}$).

Theorem 2.3. For any $p \times 1$ constant vector μ and $p \times p$ constant matrix Σ , which is a p.s.d. matrix, a random vector $Y \sim N_p(\mu, \Sigma)$ can be generated as $Y = \Sigma^{1/2} Z + \mu$.

Proof: (Corollary 2.1.1 and Theorem 2.2)

Theorem 2.4. Let $X \sim N_p(\mu, \Sigma)$, where Σ is p.d. The density function of X is

$$f_X(x) = (2\pi)^{-p/2} (\det(\Sigma))^{-1/2} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) I_R(x).$$

Theorem 2.5. Let X be p -variate normal with $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_p(\mu, \Sigma)$, where X_1 and X_2 are

separately $q \times 1$ and $(p - q) \times 1$ random vectors. Let $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ and $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$. The following two conditions are equivalent:

- (a) $\Sigma_{12} = 0$, i.e., x_1 and x_2 are uncorrelated.
- (b) x_1 and x_2 are independent.

Example 6. (Appendix)

Theorem 2.6. Let $y \sim N_n(X\beta, \sigma^2 I_n)$. Then, the least squares estimators $\hat{\beta}$ and the residual vector e have the following properties:

(a) $\begin{pmatrix} \hat{\beta} \\ e \end{pmatrix}$ is multivariate normal.

(b) $\hat{\beta}$ and e are independent.

$$(c) \begin{pmatrix} \hat{\beta} \\ e \end{pmatrix} \sim N_{p+n} \left(\begin{pmatrix} \beta \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 (X^T X)^{-1} & 0 \\ 0 & \sigma^2 (I_n - H) \end{pmatrix} \right).$$

Theorem 2.7. Let $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_p(\mu, \Sigma)$ with Σ being a positive definite matrix. The

conditional distribution of X_2 given $X_1 = x_1$ is $(p-q)$ -variate normal with mean ~~and variance being~~ $\mu_{2|1} = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1)$ and $\Sigma_{22|1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$, respectively.

Proof :

Let $Y \triangleq \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} I_q & 0 \\ -\Sigma_{21}\Sigma_{11}^{-1} & I_{p-q} \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$. By Theorem 2.5 and $Cov(Y_1, Y_2) = 0$, it implies
 $\neg\Sigma_{11}\Sigma_{11}^{-1}\Sigma_{12} + \Sigma_{12} = 0$

that Y_1 and Y_2 are independent. Thus, the distribution of Y_2 given $X_1 = x_1$ is same with the unconditional distribution of Y_2 , where $Y_2 \sim N_{p-q}(\mu_2 \cancel{+} \Sigma_{21}\Sigma_{11}^{-1}\mu_1, \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$.

Conditioning on $X_1 = x_1$, one has $\underline{X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1} \sim N_{p-q}(\mu_2 \cancel{+} \Sigma_{21}\Sigma_{11}^{-1}\mu_1, \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$,

i.e., $X_2 | X_1 = x_1 \sim N_{p-q}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1), \Sigma_{22|1})$.

Remark. From Theorem 2.6, one has $\det(\Sigma) = \det(\Sigma_{11})\det(\Sigma_{22|1})$, which is consistent

with the property $\det \begin{pmatrix} A_{11} & 0 \\ 0 & A_{22|1} \end{pmatrix} = \det(C_1 A C_2) = \det(A)$ with $C_1 = \begin{pmatrix} I_q & 0 \\ -A_{21}A_{11}^{-1} & I_{p-q} \end{pmatrix}$

and $C_2 = \begin{pmatrix} I_q & -A_{11}^{-1}A_{12} \\ 0 & I_{p-q} \end{pmatrix}$.

with a p.d. Σ .

Theorem 2.8. Assume that $\mathbf{X} \sim N_p(\mu, \Sigma)$, where Σ is p.d.. Let $\chi^2 = (\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu)$.

Then, $\chi^2 \sim \chi_p^2$.

Proof:

Let $Z = \Sigma^{-1/2}(\mathbf{X} - \mu)$. $Z \sim N_p(0, I_p)$, and, hence $\chi^2 = Z^T Z \sim \chi^2_p$.

Example 7. (See Appendix)

Theorem 2.9. (Multivariate Central Limit Theorem) Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent and identically distributed p -variate random vectors with mean μ and

variance Σ . Let $\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$. Then, $\sqrt{n}(\bar{\mathbf{X}}_n - \mu) \xrightarrow{d} N_p(0, \Sigma)$ as $n \rightarrow \infty$.

Spherical Distribution: The p -variate random vector z with density function $f_z(z)$ is spherical or spherically symmetric if $f_z(z)$ depends on z only through $z^T z$. (an equivalent statement: $\Gamma z \sim z$ for any orthogonal matrix Γ).

Remark: Let $u = \Gamma z$, where $\Gamma = \begin{pmatrix} ct^T \\ \Gamma_1 \end{pmatrix}$ with $\Gamma^T \Gamma = I_p$ and $c = (tt^T)^{-1/2}$. The c.h.f. of z is

$$\phi_z(t) = \iiint e^{u_1/c} f_z(\sum_{i=1}^p u_i^2) du_1 \dots du_p \triangleq g(t^T t).$$

Elliptical Distribution: Assume that z follows a p -variate spherical distribution, A and m are separately $p \times p$ constant matrix and $p \times 1$ fixed vector. Then, $x = Az + m$ is said to be elliptical or elliptically symmetric.

Note: Let $z \sim (0, cI_p)$ and $x \sim (\mu, \Sigma)$, one can derive that the characteristic function of

$$x \text{ is } \phi_x(t) = E[e^{it^T X}] = e^{it^T \mu} \cdot g(t^T \Sigma t).$$

Theorem 2.10. Let x follow a p -variate elliptical distribution with variance matrix Σ . If diagonality of Σ implies independence of all components of x , then x is multivariate normal.

Proof: (See Muirhead (1982)).

Without loss of generality, we assume that $E[x] = 0$. By the the diagonality of Σ and the independence of x_i 's, one can derive $g(\sum_{i=1}^p t_i^2 \sigma_i^2) = \prod_{i=1}^p g(t_i^2 \sigma_i^2)$. Thus, $g(z) = e^{kz}$ for

$k \leq 0$, which implies that $\phi_x(t) = e^{-t^T \Sigma t / 2}$ because it is the characteristic function. By the uniqueness of the c.h.f, x is shown to be multivariate normal.

Wishart Distribution ($W_p(\Sigma, n)$): Let $X_{n \times p}$ be a data matrix with x_i 's being a random sample from a $N_p(0, \Sigma)$. Then, $X^T X = \sum_{i=1}^n x_i x_i^T \sim W_p(\Sigma, n)$.

Theorem 2.11. Let x_1, \dots, x_n be a random sample from $N_p(0, \Sigma)$, $X = (x_1, \dots, x_n)^T$ and C be an $n \times n$ idempotent matrix with $\text{rank}(C) = r$. Then, $X^T C X \sim W_p(\Sigma, r)$.

Proof:

Let $C = \Gamma_n \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} \Gamma_n^T$ and $Y = \Gamma_n^T X$ with $\Gamma_n = (\gamma_{(1)}, \dots, \gamma_{(n)})$. $\Gamma_n^T \Gamma_n = I_n$. $y_j = \sum_{i=1}^n \gamma_{i(j)} x_i \sim N_p(0, \Sigma)$ and y_j 's are independent. One has $X^T C X = Y^T \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} Y = \sum_{j=1}^r y_j y_j^T$. Since y_j 's are i.i.d random vectors from a $N_p(0, \Sigma)$, it implies that $X^T C X \sim W_p(\Sigma, r)$.

Example 8. (Appendix)

Hotelling T^2 Distribution ($T^2(p, n)$): Let $d \sim N_p(0, I_p)$ and $M \sim W_p(I, n)$ be mutually independent. Then, $nd^T M^{-1} d \sim T^2(p, n)$.

Example 9. (Appendix)

Theorem 2.12. Let $M = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} \sim W_p(\Sigma, n)$ with $\dim(M_{11}) = a$. Then,

$$M_{22.1} = M_{22} - M_{21} M_{11}^{-1} M_{12} \sim W_{p-a}(\Sigma_{22.1}, n-a).$$

Proof:

Let $M = X^T X$ with $X = (X_1, X_2)$ being a data matrix from a $N_p(0, \Sigma)$ and $M_{ij} = X_i^T X_j$.

Then, $M_{22.1} = X_2^T P X_2$ with $P = I_n - X_1 (X_1^T X_1)^{-1} X_1$.

$$\underline{x_{2.1}} = \underline{x_2 - \Sigma_{21} \Sigma_{11}^{-1} x_1}$$

Since $P X_1 = 0$, one has $M_{22.1} = X_2^T P X_2 = X_{2.1}^T P X_{2.1}$, where $X_{2.1} = X_2 - X_1 \Sigma_{11}^{-1} \Sigma_{12}$.

From the proof of Theorem 2.7, we know that $X_{2.1}$ is a data matrix from $N_{p-a}(0, \Sigma_{22.1})$.

It is implied and, hence, by Theorem 2.11, it can be derived that $M_{22.1} \sim W_{p-a}(\Sigma_{22.1}, n-a)$.

Theorem 2.13. $T^2(p, n) \sim \frac{np}{n-p+1} F_{p, n-p+1}$.

Proof:

Let $d \sim N_p(0, I_p)$ be independent of $M \sim W_p(I_p, n)$, and $D = (D_1, d)$ be a nonsingular

matrix with $d^T D_1 = 0$. One has $\alpha = nd^T M^{-1} d = \frac{nd^T M^{-1} d}{d^T d} d^T d \sim T^2(p, n)$.

From $N = D^{-1} M (D^{-1})^T \sim W_p(D^{-1} (D^{-1})^T, n)$, one has $(N_{22,1})^{-1} = d^T M^{-1} d$ with
 $N_{22,1} \sim W_1((d^T d)^{-1}, n-p+1)$.

Thus, conditioning on d , $\beta = (d^T d / d^T M^{-1} d) \sim \chi^2_{n-p+1}$, which does not depend on d .

By using $d^T d \sim \chi^2_p$, we have $\alpha = n \frac{(d^T d / p)}{(\beta / (n-p+1))} \frac{p}{(n-p+1)} \sim \frac{np}{(n-p+1)} F_{p, n-p+1}$.

Remark: It implies from Theorem 2.13 that $\frac{n(n-p)}{(n-1)p} (\bar{x}_n - \mu)^T (S_n^2)^{-1} (\bar{x}_n - \mu) \sim F_{p, n-p}$.

Theorem 2.14. Let $d \sim N_p(0, I_p)$ be independent of $M \sim W_p(I_p, n)$. Then, $\frac{|M|}{|M + dd^T|} \sim Beta\left(\frac{n-p+1}{2}, \frac{p}{2}\right)$.

Proof:

$$\frac{|M|}{|M + dd^T|} = |I_p + M^{-1} dd^T|^{-1} = (1 + d^T M^{-1} d)^{-1} = \frac{\beta}{d^T d + \beta} \sim Beta\left(\frac{n-p+1}{2}, \frac{p}{2}\right).$$

Wilks' Lambda Distribution ($\Lambda(p, n, m)$): Let $A \sim W_p(I, n)$ and $B \sim W_p(I, m)$ be mutually independent with $n \geq p$. Then, $\Lambda = |A| / |A + B| = |I + A^{-1}B|^{-1} \sim \Lambda(p, n, m)$.

Theorem 2.15. $\Lambda(p, n, m) \sim \prod_{i=1}^m u_i$, where u_i 's are independent Beta random variables with parameters $(\frac{n+i-p}{2}, \frac{p}{2})$, $i = 1, \dots, m$.

Proof:

Let $B = X^T X$ with x_1, \dots, x_m being a sample from a $N_p(0, I_p)$, and $M_i = M_{i-1} + x_i x_i^T$, $i =$

$$= 1, \dots, m, \text{ with } M_0 = A. \text{ One has } \Lambda(p, n, m) = \prod_{i=1}^m u_i \text{ with } u_i = \frac{|M_{i-1}|}{|M_i|}, i = 1, \dots, m.$$

By Theorem 2.14, one has $u_i \sim Beta(\frac{n+i-p}{2}, \frac{p}{2})$, $i = 1, \dots, m$.

Using the independence of M_i and u_i , and the independence of u_i and x_{i+1}, \dots, x_m , it follows that u_i 's are independent.

Appendix

Example 1.

Let X be a $n \times p$ matrix with $\text{rank}(X) = p < n$, $H = X(X^T X)^{-1} X^T$.

(a) H is an idempotent matrix. ($H^2 = (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T) = H$)

(b) The eigenvalues of H are either 0 or 1.

Let λ and y be any corresponding eigenvalue and eigenvector of H . One has $Hy = \lambda y$ and hence, $H(Hy) = H(\lambda y)$.

Since H is an idempotent matrix, $H(Hy) = H(\lambda y)$.

Since the left-hand side $= H^2 y = Hy = \lambda y$ and the right-hand side $= \lambda Hy = \lambda^2 y$, we get have $\lambda(\lambda - 1)y = 0$. It follows that $\lambda = 0$ or 1.

~~$\lambda(\lambda - 1)y = 0$, which implies that $\lambda = 0$ or 1.~~

(c) $\text{rank}(H) = \text{trace}(H) = p$.

$$\text{trace}(H) = \text{trace}(X(X^T X)^{-1} X^T) = \text{trace}((X^T X)^{-1} X^T X) = \text{trace}(I_p) = p.$$

Let $y = (X^T X)^{-1/2} X^T$. One has

$$\begin{aligned} \text{rank}(H) &= \text{rank}(X(X^T X)^{-1} X^T) = \text{rank}(X Y Y^T) = \text{rank}(Y) \\ &= \text{rank}((X^T X)^{-1/2} X^T) = \text{rank}(X) = p. \end{aligned}$$

Example 2. (Least Squares Estimation)

Let $y = (y_1, \dots, y_n)^T \sim (X\beta, \sigma^2 I_n)$, where X is a $n \times p$ known matrix with $p < n$ and

$\text{rank}(X) = p$, β is a $p \times 1$ unknown parameter vector, and σ^2 is an unknown positive

parameter.

The least squares estimation criterion for β is to find a minimizer of the following sum

of squares: $S(\beta) = (y - X\beta)^T (y - X\beta)$. The minimizer, say $\hat{\beta}$ of $S(\beta)$ can be derived

to be $\hat{\beta} = (X^T X)^{-1} X^T y$ and is called the least squares estimator of β .

The estimated or predicted vector is $\hat{y} = X\hat{\beta} \triangleq Hy$, where $H = X(X^T X)^{-1} X^T$ is called

a hat matrix. The residual vector is defined to be $e = y - \hat{y} = (I_n - H)y$.

Let $X = (X_1, X_2)$, $\beta = (\beta^T_1, \beta^T_2)^T$. $H = X(X^T X)^{-1} X^T$ and $H_{22,1} = X_2^T(I - H_1)X_2$ with $H_1 = X_1(X_1^T X_1)^{-1} X_1^{-1}$.

One has $\hat{\beta}_2 = (-H_{22,1}^{-1}(X_2^T X_1)(X_1^T X_1)^{-1}, H_{22,1}^{-1}) \begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} Y$, which implies that

$$\hat{\beta}_2 = (-H_{22,1}^{-1} X_2^T H_1 + H_{22,1}^{-1} X_2^T) Y = (H_{22,1}^{-1} X_2^T (I - H_1)) Y = (e_2^T e_2)^{-1} e_2^T Y$$

(or $= (e_2^T e_2)^{-1} e_2^T ((I - H_1)Y)$),

where $e_2 = (I - H_1)X_2$.

Reference:

1. Rao, C. R. and Toutenburg, H. (1999). Linear Models: Least Squares and Alternatives. Second Edition. Springer-Verlag.
2. Rawlings, J. O., Pantula, S. G., and Dickey, D. A. (1998). Applied Regression Analysis: A Research Tool. Second Edition. Springer-Verlag.
3. Sen , A. K. and Srivastava , M. S. (1991) Regression Analysis: Theory, Methods, and Applications. Springer-Verlag.

Example 3.

Let Σ be a variance matrix of a $p \times 1$ random vector X .

(a) Σ is at least a p.s.d. matrix.

Let a be any $p \times 1$ constant vector, $a^T \Sigma a = a^T V(X) a = V(a^T X) \geq 0$.

(b) By the spectral decomposition theorem, Σ can be expressed as $\Sigma = \sum_{i=1}^p \lambda_i \gamma_{(i)} \gamma_{(i)}^T$.

One has $\lambda_i \geq 0$ ($\because \lambda_i = \gamma_{(i)}^T \Sigma \gamma_{(i)} \geq 0$).

Without loss of generality, let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. If $\lambda_p > 0$, then Σ is non-singular. (\because

$\text{rank}(\Sigma) = \text{rank}(\Gamma \Lambda \Gamma^T) = \text{rank}(\Lambda) = p$.)

(c) Find a vector η satisfying $\eta^T \eta = 1$ such that the variance of $U = \eta^T X$ is maximized.

Since η can be expressed as $\eta = \sum_{i=1}^p \eta_i \gamma_{(i)}$, it can be derived that

$\text{Var}(U) = \eta^T V(X) \eta = \eta^T \sum_{i=1}^p \lambda_i \gamma_{(i)} \gamma_{(i)}^T \eta = \sum_{i=1}^p \lambda_i \eta_i^2 \leq \lambda_1 \sum_{i=1}^p \eta_i^2 = \lambda_1$. Equality holds if

and only if $\eta = \gamma_{(1)}$.

Example 4. (Example 2 continued)

Let $e = Y - X\hat{\beta} = (I_n - H)Y$.

$$(a) E[\hat{\beta}] = E[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T E[Y] = (X^T X)^{-1} X^T X \beta = \beta \text{ and}$$

$$V(\hat{\beta}) = (X^T X)^{-1} X^T V(Y) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}.$$

$$(b) E[e] = (I - H)E[Y] = (I - H)X\beta = (X - X)\beta = 0 (\because HX = X) \text{ and}$$

$$V[e] = (I - H)(\sigma^2 I_n)(I - H) = \sigma^2 (I - H).$$

$$(c) Cov(\hat{\beta}, e) = Cov((X^T X)^{-1} X^T Y, (I_n - H)Y) = (X^T X)^{-1} X^T (\sigma^2 I_n)(I_n - H) = 0.$$

$$(d) E[e^T e] = E[Y^T (I - H)Y] = \text{tr}((I_n - H)(\sigma^2 I_n)) + \beta^T X^T (I_n - H)X\beta = (n - p)\sigma^2$$

($\because I_n - H$ is idempotent).

Example 5.

Let $\begin{pmatrix} x \\ y \end{pmatrix} \sim \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$ with Σ_{11} and Σ_{22} being non-singular. Moreover,

let $\eta = a^T x$ and $\phi = b^T y$.

Consider $\rho(\eta, \phi) = \frac{a^T \Sigma_{12} b}{(a^T \Sigma_{11} a)^{1/2} (b^T \Sigma_{22} b)^{1/2}}$ subject to $a^T \Sigma_{11} a = b^T \Sigma_{22} b = 1$, the objective

here is to find (a, b) such that $\rho(\eta, \phi)$ is maximized.

Let $a^{*T} = \Sigma_{11}^{-1/2} a$ and $b^{*T} = \Sigma_{22}^{-1/2} b$. $\rho(\eta, \phi)$ can be re-expressed as

$$\rho(\eta, \phi) = a^{*T} (\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}) b^{*T} \triangleq a^{*T} K b^{*T}, \text{ where } K = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}.$$

By the singular value decomposition theorem, $K = ULV^T$ with $U = (\alpha_1, \dots, \alpha_p)$,

$V = (\beta_1, \dots, \beta_q)$, and $L = \begin{pmatrix} \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_k}) & 0_{k \times (q-k)} \\ 0_{(p-k) \times k} & 0_{(p-k) \times (q-k)} \end{pmatrix}$ satisfying $KK^T = UL^2 U^T$

and $K^T K = VL^2 V^T$, where $\lambda_1 \geq \dots \geq \lambda_k > 0$.

Since a^* and b^* can be expressed as $a^* = \sum_{i=1}^p a_i^* \alpha_i$ and $b^* = \sum_{i=1}^q b_i^* \beta_i$, one has

$$\rho(\eta, \phi) = (\sum_{i=1}^k a_i^* b_i^* \sqrt{\lambda_i}) \leq \left[(\sum_{i=1}^k a_i^{*2} \sqrt{\lambda_i})(\sum_{i=1}^k b_i^{*2} \sqrt{\lambda_i}) \right]^{\frac{1}{2}}.$$

(Equality holds if $a_i^* = c b_i^*, i = 1, \dots, k, c > 0.$)

$$\leq \sqrt{\lambda_1}.$$

(Equality holds if $a_1^* = b_1^* = 1, a_2^* = \dots = a_p^* = 0$, and $b_2^* = \dots = b_q^* = 0$.)

Thus, $\rho(\eta, \phi)$ is maximized when $(a, b) = (\Sigma_{11}^{-\frac{1}{2}} \alpha_1, \Sigma_{22}^{-\frac{1}{2}} \beta_1)$.

Example 6. Let X_1, \dots, X_n be a random sample from $N_p(\mu, \sigma^2)$. Then, the sample mean

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ and the sample variance $S_n^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{(n-1)}$ are independent.

Proof:

Let $\mathbf{X} = (X_1, \dots, X_n)^T$ and $I_n - H = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$.

One has $\bar{X}_n = \frac{1}{n} \mathbf{1}_n^T \mathbf{X}$ and $S^2 = \frac{\mathbf{X}^T (I_n - H) \mathbf{X}}{(n-1)} = \frac{((I_n - H) \mathbf{X})^T ((I_n - H) \mathbf{X})}{(n-1)}$.

From $\begin{pmatrix} \mathbf{1}_n^T \\ I_n - H \end{pmatrix} \mathbf{X} \sim N_{n+1} \left(\begin{pmatrix} n\mu \\ 0 \end{pmatrix}, \begin{pmatrix} n\sigma^2 & 0 \\ 0 & (I_n - H)\sigma^2 \end{pmatrix} \right)$, it implies that \bar{X}_n and $(I_n - H) \mathbf{X}$

are independent. Since S_n^2 is a function of $(I_n - H) \mathbf{X}$, \bar{X}_n and S_n^2 are shown to be independent.

Example 7.

Let X_1, \dots, X_n be a random sample from a $N(\mu, \sigma^2)$. Show that $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$.

Proof :

Since $I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ is an idempotent matrix with rank $(n-1)$, by the spectral decomposition theorem, one has

an application of

enables us to have

$$\frac{\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}_n)^2}{\sigma^2} = \frac{\mathbf{X}^T \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \mathbf{X}}{\sigma^2} = \tilde{\mathbf{Z}}^T \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \tilde{\mathbf{Z}} = \tilde{\mathbf{Z}}^T \left(\Gamma_{n-1} \Gamma_{n-1}^T \right) \tilde{\mathbf{Z}},$$

where $\Gamma_{n-1} = (r_{(1)}, \dots, r_{(n-1)})$ with $\Gamma_{n-1} \mathbf{1}_n = 0$ and $\tilde{\mathbf{Z}} = \frac{\mathbf{X} - \mu \mathbf{1}_n}{\sigma}$ $\sim N_n(0, I_n)$.

Since $\tilde{\mathbf{Z}} = \Gamma_{n-1}^T \tilde{\mathbf{Z}} \sim N_{n-1}(0, \mathbf{I}_{n-1})$, it implies that \tilde{Z}_i 's are i.i.d. χ^2_1 random variables, and,

$$\text{hence, } \tilde{\mathbf{Z}}^T \left(\Gamma_{n-1} \Gamma_{n-1}^T \right) \tilde{\mathbf{Z}} = \sum_{i=1}^{n-1} \tilde{Z}_i^2 \sim \chi^2_{n-1}.$$

Example 8.

Let x_1, \dots, x_n be a random sample from a $N_p(\mu, \Sigma)$ and $X = (x_1, \dots, x_n)^T$ be a data

matrix. Then, $(n-1)S_n^2 = X^T (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) X \sim W_p(\Sigma, (n-1))$.

Proof:

Let $\tilde{\mathbf{X}} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ with $\tilde{x}_i = x_i - \mu$, where $\tilde{\mathbf{X}}$ is a data matrix from a $N_p(0, \Sigma)$. $\downarrow i=1, \dots, n$

One has $(n-1)S_n^2 = X^T (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) X = \tilde{\mathbf{X}}^T (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) \tilde{\mathbf{X}}$. $\leftarrow \tilde{\mathbf{X}} = X - \mathbf{1}_n \mu^T$

Since $(I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T)$ is an idempotent matrix with $\text{rank}(I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) = (n-1)$, it follows

from Theorem 2.11 that $(n-1)S_n^2 \sim W_p(\Sigma, (n-1))$.

Example 9. Let $\bar{x}_n = \frac{1}{n} \mathbf{X}^T \mathbf{1}_n$. Then, $n(\bar{x}_n - \mu)^T (S_n^2)^{-1} (\bar{x}_n - \mu) \sim T^2(p, (n-1))$.

$$\begin{aligned} & \sum_{i=1}^n (\sqrt{n}(\bar{x}_n - \mu)) \sim N_p(0, I_p) \\ & \text{&} \quad (n-1)S_n^2 \sim W_p(\Sigma, (n-1)), \end{aligned}$$

which implies that $(n-1) \sum_{i=1}^n S_n^{-2} \sum_{j=1}^n \sim W_p(I_p, n-1)$.

$$\Rightarrow (n-1) \left(\sum_{i=1}^n \sqrt{n}(\bar{x}_n - \mu) \right)^T \left((n-1) \sum_{i=1}^n S_n^{-2} \sum_{j=1}^n \right)^{-1} \left(\sum_{i=1}^n \sqrt{n}(\bar{x}_n - \mu) \right) \sim T^2(p, n-1)$$