

The multivariate normal distribution

Multivariate normal distribution has many advantageous and revealing properties.

There are several ways of introducing multivariate normal. One way is to directly define the joint density of multivariate normal, such as in the text by Johnson and Wichern.

In the following we introduce multivariate normal via linear combinations of its univariate components.

First let's review the univariate normal distribution.

1 The univariate normal distribution

Definition of univariate standard normal random variable

Z is a standard normal random variable, denoted as $Z \sim N(0, 1)$, if its probability density function is

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad z \in \mathbb{R}.$$

(Notation: $\mathbb{R} = (-\infty, \infty)$. $z \in \mathbb{R}$ means z is a real number.)

By convention, ϕ or φ is used for the standard normal density instead of the generic notation f for probability density functions. A common notation for the cumulative distribution function of the standard normal is

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt, \quad z \in \mathbb{R}$$

Definition of univariate normal random variable

For any real constants a, b , if $Y = aZ + b$, that is, if Y is a (affine) linear transformation of the standard normal Z , then by the properties of expectation (via the properties of integration), the expectation and variance of Y are

$$\mathbb{E}(Y) = a\mathbb{E}(Z) + b = b, \quad \text{var}(Y) = a^2 \text{var}(Z) = a^2.$$

By variable substitution, for $a \neq 0$, Y has the density function

$$f_Y(y) = \frac{1}{a} \phi\left(\frac{y-b}{a}\right) = \frac{1}{\sqrt{2\pi}a} e^{-\frac{(y-b)^2}{2a^2}}, \quad x \in \mathbb{R}.$$

Then Y is a normal random variable with mean b and variance a^2 , denoted as $Y \sim N(b, a^2)$.

Remarks on univariate normal

- The degenerate case $a = 0$ yields $Y \sim N(b, 0)$, which represents the point mass distribution at b .
- A univariate normal distribution $N(\mu, \sigma^2)$ is uniquely determined by its mean μ and variance σ^2 . Therefore the distribution of $Y = aZ + b$ with $Z \sim N(0, 1)$ is completely and uniquely determined by the values of b and a^2 .
- The cumulative distribution function (CDF) of normal variable does not have closed form, which means the normal CDF has to be expressed as integrals or infinite series of known functions (e.g., polynomials).
- If X, Y are independent normal random variables, then $X + Y$ is of normal distribution. (Exercise)
 - Consequently, if $X_i, i = 1, \dots, k$ are independent normal, then their sum $X_1 + \dots + X_k$ is normal. This important facts will be used frequently.

– Furthermore, the linear combination of independent normal $a_1X_1 + \dots + a_kX_k$ is normal, since each a_iX_i is normal. In general a_i 's and other constants used in this course are real (rather than complex).

- However if X, Y are normal but not independent, then $X + Y$ is not necessarily normal.

An example

Suppose $X \sim N(0, 1)$. Let random variable $S = \pm 1$ with probability $\frac{1}{2}$ each, and S be independent of X . Such an S is a random sign. Let $Y = SX$, so

$$Y = \begin{cases} X, & \text{if } S = 1, \\ -X, & \text{if } S = -1. \end{cases}$$

Then

$$\begin{aligned} \mathbb{P}(Y \leq y) &= \mathbb{P}(Y \leq y, S = 1) + \mathbb{P}(Y \leq y, S = -1) \\ &= \mathbb{P}(Y \leq y|S = 1)\mathbb{P}(S = 1) + \mathbb{P}(Y \leq y|S = -1)\mathbb{P}(S = -1) \\ &= \mathbb{P}(X \leq y|S = 1)\frac{1}{2} + \mathbb{P}(-X \leq y|S = -1)\frac{1}{2} \\ &= \mathbb{P}(X \leq y)\frac{1}{2} + \mathbb{P}(-X \leq y)\frac{1}{2} && \text{since } X \perp\!\!\!\perp S \\ &= \mathbb{P}(X \leq y)\frac{1}{2} + \mathbb{P}(X \leq y)\frac{1}{2} && \text{since } -X \sim N(0, 1) \\ &= \mathbb{P}(X \leq y) \end{aligned}$$

Hence $Y \sim N(0, 1)$. Therefore the bivariate vector (X, Y) is marginally normal, that is, each component on its own is normal. However the joint distribution of (X, Y) is not bivariate normal. In fact,

$$X + Y = \begin{cases} 0, & \text{with probability } \frac{1}{2}, \\ 2X, & \text{with probability } \frac{1}{2}. \end{cases}$$

So $0 < \mathbb{P}(X + Y = 0) < 1$, which means $X + Y$ is not a univariate normal. Consequently (X, Y) is not bivariate normal, because a linear combination of the components $(X + Y)$ is not of normal distribution.

Remarks

In this example, $X + Y$ is a mixture of a discrete (Bernoulli) distribution and a continuous ($2X$) distribution which is normal. Furthermore, X and Y are uncorrelated but not independent.

2 Multivariate normal distribution via linear combinations

In the following we give a definition of multivariate normal in terms of linear combinations of its individual components, which are the familiar univariate normal random variables.

Definition

A p -variate random vector $\mathbf{X} = (X_1, \dots, X_p)'$ is **multivariate normal** if for any $\mathbf{b} \in \mathbb{R}^p$, the linear combination $\mathbf{b}'\mathbf{X} = b_1X_1 + \dots + b_pX_p$ is univariate normal (including the degenerate univariate normal, the point mass). (*)

This definition provides a convenient tool to construct multivariate normal vectors, as shown in the following examples.

Examples

• **p -variate standard normal distribution.**

Let Z_1, \dots, Z_p be independent univariate standard normal random variables, so $Z_i \sim N(0, 1)$.

By the independence, the p -variate variable $\mathbf{Z} = (Z_1, \dots, Z_p)$ has joint density

$$\begin{aligned} f_{\mathbf{Z}}(\mathbf{z}) &= f_{Z_1}(z_1)f_{Z_2}(z_2) \cdots f_{Z_p}(z_p) \\ &= \prod_{j=1}^p \frac{1}{\sqrt{2\pi}} e^{-z_j^2/2} = \frac{1}{\sqrt{2\pi}^p} e^{-\sum_{j=1}^p z_j^2/2} = \frac{1}{(2\pi)^{p/2}} e^{-\mathbf{z}'\mathbf{z}/2} = \frac{1}{(2\pi)^{p/2}} e^{-\|\mathbf{z}\|^2/2} \end{aligned}$$

where

$$\|\mathbf{z}\| = \mathbf{z}'\mathbf{z} = \sqrt{z_1^2 + \cdots + z_p^2}$$

denotes the Euclidean norm of vector \mathbf{z} in \mathbb{R}^p .

\mathbf{Z} is of multivariate normal distribution, and is termed as the p -variate **standard normal**.

Proof. We show that the p vector \mathbf{Z} with independent standard normal components is of p -variate normal distribution.

By definition (*), it is sufficient to show that any linear combination of the components of the p vector \mathbf{Z} is a univariate normal random variable.

For any linear combination $\mathbf{b}'\mathbf{Z} = b_1Z_1 + \cdots + b_pZ_p$, since $b_jZ_j \sim N(0, b_j^2)$, b_jZ_j 's are independent for $j = 1, \dots, p$. By the fact that sum of independent univariate normal variables is still a univariate normal variable, we can conclude that the linear combination $\mathbf{b}'\mathbf{Z}$ is of univariate normal. Since this is true for any $\mathbf{b} \in \mathbb{R}^p$, this proves that \mathbf{Z} is of p -variate normal. □

• **Correlated bivariate normal** constructed from linear combinations of independent standard normal.

Let Z_1, Z_2 be independent standard normal random variables. Then the random vector

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = \begin{bmatrix} Z_1 + 2Z_2 \\ 3Z_1 + 4Z_2 \end{bmatrix}$$

is of **bivariate normal**, because any linear combination of the components X, Y is a linear combination of Z_1, Z_2 :

$$b_1X + b_2Y = b_1(Z_1 + 2Z_2) + b_2(3Z_1 + 4Z_2) = (b_1 + 3b_2)Z_1 + (2b_1 + 4b_2)Z_2,$$

which is univariate normal, again from the normality of sums of independent univariate normal variables.

Furthermore, because $cov(Z_1, Z_2) = 0$, the covariance

$$cov(X, Y) = cov(Z_1 + 2Z_2, 3Z_1 + 4Z_2) = 3V(Z_1) + 8V(Z_2) = 11.$$

Hence (X, Y) is of correlated bivariate normal distribution constructed from independent, uncorrelated standard normal variables.

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim N_2 \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 5 & 11 \\ 11 & 25 \end{bmatrix} \right)$$

• **Uncorrelated bivariate normal** from linear combinations of independent variables.

If a pair of random variables X, Y are of equal variance, then $X + Y$ and $X - Y$ are uncorrelated:

$$cov(X + Y, X - Y) = var(X) + cov(X, Y) - cov(Y, X) - var(Y) = var(X) - var(Y) = 0$$

If in addition to equal-variance, X and Y are independent and of normal distributions, then any linear combination $b_1(X + Y) + b_2(X - Y) = (b_1 + b_2)X + (b_1 - b_2)Y$ is of univariate normal. Consequently $(X + Y, X - Y)$ is bivariate normal, its components not only uncorrelated but also independent.

3 Properties and density of multivariate normal distribution

Theorem (from p -variate normal to k -variate normal)

If $\mathbf{X} \in \mathbb{R}^p$ is p -variate normal, and $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{c}$, where \mathbf{A} is a matrix of dimension $k \times p$ and \mathbf{c} is a k -vector of constants, then $\mathbf{Y} \in \mathbb{R}^k$ is of k -variate normal distribution.

Proof. For any k -vector \mathbf{b} , we need to show that $\mathbf{b}'\mathbf{Y}$ is univariate normal. Write

$$\mathbf{b}'\mathbf{Y} = \mathbf{b}'\mathbf{A}\mathbf{X} + \mathbf{b}'\mathbf{c} = \mathbf{b}^*\mathbf{X} + \mathbf{c}^*$$

where $\mathbf{b}^* = (\mathbf{b}'\mathbf{A})' = \mathbf{A}'_{p \times k} \mathbf{b}_{k \times 1}$ is a vector in \mathbb{R}^p , and $\mathbf{c}^* = (\mathbf{b}'\mathbf{c})' = \mathbf{c}'\mathbf{b}$ is a constant in \mathbb{R} . \mathbf{X} is p -variate normal, so any linear combination $\mathbf{b}^*\mathbf{X}$ of the components of \mathbf{X} is of univariate normal. Since any linear combination of the k component of \mathbf{Y} is a linear combination of the components of \mathbf{X} plus a constant, thus univariate normal, we conclude that \mathbf{Y} is of k -variate normal. □

Remarks

- By the mean and variance formulas we proved in the previous notes on "Multivariate random sample matrices",

$$\mathbb{E}(\mathbf{Y}) = \mathbf{A}\mathbb{E}(\mathbf{X}) + \mathbf{c} = \mathbf{A}\boldsymbol{\mu}_x + \mathbf{c}, \quad Cov(\mathbf{Y}) = \mathbf{A}Cov(\mathbf{X})\mathbf{A}' = \mathbf{A}\Sigma_x\mathbf{A}'$$

Thus,

$$\mathbf{Y} \sim N_k(\mathbf{A}\boldsymbol{\mu}_x + \mathbf{c}, \mathbf{A}\Sigma_x\mathbf{A}') \tag{1}$$

- In writing up (1) the uniqueness of multivariate normal with given mean and variance is implicitly assumed. The theorem (uniqueness of multivariate normal) is stated near the end in this section with a brief outline of the proof.
- Each component of the k -variate $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{c}$ is a linear combination of the p -components of \mathbf{X} plus a constant.
- In applications, usually $rank(\mathbf{A}) = k \leq p$ to avoid degenerate cases.

Corollary (Construct k -variate normal with any mean and any legit covariance)

If \mathbf{Z} is p -variate standard normal, and $\mathbf{Y} = \mathbf{A}\mathbf{Z} + \mathbf{c}$, where matrix \mathbf{A} is $k \times p$ and $\mathbf{c} \in \mathbb{R}^k$, then \mathbf{Y} is k -variate normal with

$$\mathbb{E}(\mathbf{Y}) = \mathbf{c}, \quad Cov(\mathbf{Y}) = \mathbf{A}\mathbf{A}', \quad \text{that is, } \mathbf{Y} \sim N(\mathbf{c}, \mathbf{A}\mathbf{A}')$$

Corollary (Marginals of multivariate normal are multivariate normal)

Suppose $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$. Write $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$, where \mathbf{X}_1 consists of the first $q < p$ components of \mathbf{X} , so \mathbf{X}_1 is q -variate, consequently \mathbf{X}_2 is $(p - q)$ -variate. Partition $\boldsymbol{\mu}$ and Σ accordingly,

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \tag{2}$$

Then

$$\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \Sigma_{11}), \quad \mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \Sigma_{22}).$$

Proof. We are to derive that $\mathbf{X}_1, \mathbf{X}_2$ are normal with the stated means and covariance. $E(\mathbf{X}_i) = \boldsymbol{\mu}_i$ and $Cov(\mathbf{X}_i) = \Sigma_{ii}$ ($i = 1, 2$) can be verified directly. To show \mathbf{X}_i is normal, we express each \mathbf{X}_i as a linear transformation of \mathbf{X} .

$$\mathbf{X}_1 = [\mathbf{I}_q \ \mathbf{O}_{p-q}] \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} = \mathbf{A}_1 \mathbf{X}, \quad \mathbf{X}_2 = [\mathbf{O}_p \ \mathbf{I}_{p-q}] \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} = \mathbf{A}_2 \mathbf{X}$$

where \mathbf{I} and \mathbf{O} are the identity matrix and the components = 0 matrix with the corresponding dimensions. Recall $\mathbf{A}\mathbf{X} + \mathbf{c}$ is of multivariate normal for any \mathbf{A} and \mathbf{c} , as formulated in (1). Therefore $\mathbf{X}_i = \mathbf{A}_i \mathbf{X}$ ($i = 1, 2$) are of multivariate normal. □

Remarks

- An immediate conclusion from the second corollary is that each component of \mathbf{X} is univariate normal,

$$X_j \sim N_1(\mu_j, \sigma_{jj})$$

where σ_{jj} is the (j, j) th element or the i th diagonal entry of Σ .

- The converse is not necessarily true. If each component X_j 's is normal, $\mathbf{X} = (X_1, \dots, X_p)$ may not be p -variate normal, as illustrated by the example given earlier.
- For normal distributions, \mathbf{X}_1 and \mathbf{X}_2 are independent if and only if $Cov(\mathbf{X}_1, \mathbf{X}_2) = \mathbf{0}$, which can be derived from the form of the joint density function. (exercise)

Corollary (Conditional distributions of multivariate normal are multivariate normal)

Assume that $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \Sigma_{11})$, $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \Sigma_{22})$, and $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)'$ is p -variate normal with mean $\boldsymbol{\mu}$ and variance Σ of the partitioned form in (2). Then the conditional distribution of \mathbf{X}_1 given $\mathbf{X}_2 = \mathbf{x}_2$ is normal, with conditional mean and covariance as in the following.

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

The above result can be derived from the form of the conditional density.

Example: For $p = 2$, $q = 1$,

$$\mathbf{X} \sim N_2(\boldsymbol{\mu}, \Sigma), \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$$

The notations imply

$$X_1 \sim N(\mu_1, \sigma_1^2), \quad X_2 \sim N(\mu_2, \sigma_2^2), \quad \text{corr}(X_1, X_2) = \rho = \frac{\sigma_{12}}{\sigma_1 \sigma_2}$$

The conditional distribution of X_1 given that $X_2 = x_2$ has the more familiar expression

$$X_1 | X_2 = x_2 \sim N\left(\mu_1 + \frac{\sigma_{12}}{\sigma_{22}}(x_2 - \mu_2), \sigma_{11} - \frac{\sigma_{12}^2}{\sigma_{22}}\right) = N\left(\mu_1 + \rho \frac{\sigma_1}{\sigma_2}(x_2 - \mu_2), \sigma_1^2(1 - \rho^2)\right)$$

Theorem (Joint density of multivariate normal)

If $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$ is p -variate normal with mean $\boldsymbol{\mu}$ and covariance matrix Σ , where Σ is symmetric and positive definite, then \mathbf{X} has density function

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

Proof.

We need to show that there is a normal random vector with the given mean and covariance matrix, its density is of the desired form.

Since Σ is positive definite, by matrix theory, there is a $p \times p$ invertible matrix \mathbf{A} such that

$$\Sigma = \mathbf{A}\mathbf{A}'$$

with

$$\Sigma^{-1} = \mathbf{A}'^{-1} \mathbf{A}^{-1}, \quad |\Sigma| = \det(\Sigma) = (\det(\mathbf{A}))^2$$

Let $\mathbf{Z}_p \sim N(0, \mathbf{I}_p)$ be the p -variate standard normal, and define

$$\mathbf{X}^* = \mathbf{A}\mathbf{Z}_p + \boldsymbol{\mu}$$

Then

$$\mathbf{X}^* \sim N_p(\boldsymbol{\mu}, \mathbf{A}\mathbf{A}') = N_p(\boldsymbol{\mu}, \Sigma)$$

By the uniqueness theorem (stated below), \mathbf{X} and $\mathbf{X}^* = \mathbf{A}\mathbf{Z}_p + \boldsymbol{\mu}$ are of the same mean, same covariance matrix, thus they are of the same multivariate normal distribution.

By variable transformation, we may derive the common density function of \mathbf{X} and \mathbf{X}^* from the known density of

$$\mathbf{Z} = \mathbf{A}^{-1}(\mathbf{X} - \boldsymbol{\mu})$$

by the method of variable substitution for continuous functions. Note that, by matrix calculus, the variable transformation matrix (the Jacobian) is

$$\mathbf{J} = \left[\frac{\partial z_i}{\partial x_j} \right]_{i,j=1,\dots,p} = \mathbf{A}^{-1}$$

Recall the variable transformation formula for density functions

$$f_{\mathbf{X}}(\mathbf{x}) = |\mathbf{J}| f_{\mathbf{Z}}(\mathbf{z})$$

and the density of standard p -variate normal \mathbf{Z} ,

$$f_{\mathbf{Z}}(\mathbf{z}) = \frac{1}{(2\pi)^{p/2}} e^{-\|\mathbf{z}\|^2/2} = \frac{1}{(2\pi)^{p/2}} e^{-\mathbf{z}'\mathbf{z}/2}$$

we can then derive the density of the p -variate normal vector \mathbf{X} ,

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}) &= |\det(\mathbf{J})| f_{\mathbf{Z}}(\mathbf{A}^{-1}(\mathbf{x} - \boldsymbol{\mu})) \\ &= |\det(\mathbf{A}^{-1})| \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2}(\mathbf{A}^{-1}(\mathbf{x} - \boldsymbol{\mu}))' (\mathbf{A}^{-1}(\mathbf{x} - \boldsymbol{\mu}))} \\ &= \frac{1}{|\det(\mathbf{A})|} \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \mathbf{A}'^{-1} \mathbf{A}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \\ &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})} \end{aligned}$$

□

Remarks

The above theorem on joint density formula provides a necessary and sufficient condition for multivariate normal random vectors, thus this theorem itself often serves as the definition.

Theorem* (uniqueness of multivariate normal)*

For and vector $\mu \in \mathbb{R}^p$ and symmetric positive definite matrix Σ of dimension $p \times p$, there exists a unique multivariate normal distribution with mean vector μ and covariance matrix Σ .

- The uniqueness is a subtle key step in introducing multivariate normal via (affine) linear transformation.
- To prove the uniqueness, we would need to show that if \mathbf{X}, \mathbf{Y} are p -variate normal with the same mean and covariance matrix, then \mathbf{X} and \mathbf{Y} are of the same distribution.
 - The proof may start with the fact that linear transformation transforms multivariate normal vectors from one dimension to another.
 - So we may write $\mathbf{X} = \mathbf{AZ} + \mathbf{c}$, $\mathbf{Y} = \mathbf{BW} + \mathbf{d}$, where \mathbf{Z}, \mathbf{W} are multivariate standard normal.
 - $\mathbf{c} = E(\mathbf{X}) = E(\mathbf{Y}) = \mathbf{d}$. WLOG we may only consider the case $\mathbf{c} = \mathbf{d} = 0$.
 - We only consider the case the A and B both are $p \times k$ (otherwise add zero columns to one of them).
 - $\mathbf{AA}^T = Cov(\mathbf{X}) = Cov(\mathbf{Y}) = \mathbf{BB}^T$.
 - Show that then $\mathbf{B}^T = \mathbf{QA}^T$ for some $\mathbf{QQ}^T = \mathbf{I}_k$.
 - Then $\mathbf{Y} = \mathbf{BW} = \mathbf{A}(\mathbf{Q}^T\mathbf{W})$
 - $\mathbf{Q}^T\mathbf{W} = \mathbf{Z}^*$, where \mathbf{Z}^* are multivariate standard normal.
 - Then $\mathbf{Y} = \mathbf{AZ}^*$, $\mathbf{X} = \mathbf{AZ}$, same transformation of multivariate standard normal, are of the some distribution.
- A detailed proof of the uniqueness theorem was provided by Muirhead (ref. Muirhead 1982).

Corollary (Uncorrelatedness = Independence for multivariate normal)

Suppose $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \sim N_{p+q}(\mu, \Sigma)$, where \mathbf{X}_1 consists of the first p components of \mathbf{X} . Then

$$Cov(\mathbf{X}_1, \mathbf{X}_2) = 0_{p \times q} \implies \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \quad (\text{independence})$$

Proof.

By the marginal distribution property of multivariate normal, \mathbf{X}_1 is of p -variate normal, \mathbf{X}_2 is of q -variate normal. We can write

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & 0_{p \times q} \\ 0_{q \times p} & \Sigma_{22} \end{bmatrix}$$

Then the determinant and the inverse of Σ can be expressed as

$$|\Sigma| = |\Sigma_{11}| \cdot |\Sigma_{22}|, \quad \Sigma^{-1} = \begin{bmatrix} \Sigma_{11}^{-1} & 0_{p \times q} \\ 0_{q \times p} & \Sigma_{22}^{-1} \end{bmatrix}$$

assuming the non-degenerate case that all inverse matrices exist.

Bring the above expressions into the density of \mathbf{X} ,

$$\begin{aligned} f_{X_1, X_2}(\mathbf{x}_1, \mathbf{x}_2) &= f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{(p+q)/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu)} \\ &= \frac{1}{(2\pi)^{(p+q)/2} |\Sigma_{11}|^{1/2} |\Sigma_{22}|^{1/2}} e^{-\frac{1}{2}[\mathbf{x}_1 - \mu_1 \quad \mathbf{x}_2 - \mu_2] \begin{bmatrix} \Sigma_{11}^{-1} & 0_{p \times q} \\ 0_{q \times p} & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 - \mu_1 \\ \mathbf{x}_2 - \mu_2 \end{bmatrix}} \\ &= \frac{1}{(2\pi)^{p/2} |\Sigma_{11}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_1 - \mu_1)' \Sigma_{11}^{-1} (\mathbf{x}_1 - \mu_1)} \cdot \frac{1}{(2\pi)^{q/2} |\Sigma_{22}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_2 - \mu_2)' \Sigma_{22}^{-1} (\mathbf{x}_2 - \mu_2)} \\ &= f_{X_1}(\mathbf{x}_1) f_{X_2}(\mathbf{x}_2) \end{aligned}$$

The second from the last equality results from the rules of matrix multiplication (exercise). We have shown $f_{X_1, X_2}(\mathbf{x}_1, \mathbf{x}_2) = f_{X_1}(\mathbf{x}_1) f_{X_2}(\mathbf{x}_2)$. By the definition of the independence, we have proved $\mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2$. □

Remark

Since independence always implies uncorrelatedness, the above implies that, the equivalent relation between uncorrelatedness and Independence

$$Cov(\mathbf{X}_1, \mathbf{X}_2) = 0_{q \times (p-q)} \iff \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2$$

holds for multivariate normal (only).

4 Moment generation functions

The moment generating function for a univariate random variable X is

$$M_X(t) = \mathbb{E}(e^{tX})$$

for t values at which the expectation exists. The most important property is that the k th derivative of moment generation function at the origin is

$$M^{(k)}(0) = \mathbb{E}(X^k),$$

the k th moment of X (thus the name). For $X \sim N(\mu, \sigma^2)$, the moment generating function is

$$M_X(t) = \mathbb{E}(e^{tX}) = e^{t\mu + \frac{1}{2}\sigma^2 t^2}$$

The moment-generating function of a p -variate random vector $\mathbf{X} = (X_1, \dots, X_p)'$ has the form

$$M_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}(e^{\mathbf{t}'\mathbf{X}}) = \mathbb{E}(e^{t_1 X_1 + \dots + t_p X_p})$$

which is defined for $\mathbf{t} = (t_1, \dots, t_p)'$ wherever the expectation exists. The moment-generating function of p -variate normal $\mathbf{X} \sim N_p(\mu, \Sigma)$ has the simple expression

$$M_{\mathbf{X}}(\mathbf{t}) = e^{\mathbf{t}'\mu + \frac{1}{2}\mathbf{t}'\Sigma\mathbf{t}} = e^{t_1\mu_1 + \dots + t_p\mu_p + \frac{1}{2}Var(t_1 X_1 + \dots + t_p X_p)},$$

which can be derived from the moment generating function for univariate normal,

$$M_X(t) = M_{t\mathbf{X}}(1)$$

Remarks

In this course we rarely use moment generating functions directly. However many theoretical results we assume or use in this course, such as the central limit theorem and its various versions, are relatively easier to derive using moment generating functions or characteristic functions.

5 Maximum likelihood estimation

There are many methods to estimate the population parameters $\boldsymbol{\mu}$ and Σ from sample data. Maximum likelihood method is a popular and powerful method.

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample of independent p -variate random vectors of distribution $N_p(\boldsymbol{\mu}, \Sigma)$. Then each p -variate sample vector \mathbf{X}_j has density function

$$f(\mathbf{x}_j) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu})}$$

By independence, $\mathbf{X}_1, \dots, \mathbf{X}_n$ has joint density function

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\mu}, \Sigma) = \prod_{j=1}^n f(\mathbf{x}_j) = \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu})}$$

Given sample data $\mathbf{x}_1, \dots, \mathbf{x}_n$, the above joint density becomes a function with unknown parameters μ_i 's and σ_{ik} 's ($i, k = 1, \dots, p$) in $\boldsymbol{\mu}$ and Σ .

Express the joint density as a function of the unknown parameters given the data, we write

$$L(\boldsymbol{\mu}, \Sigma) = L(\boldsymbol{\mu}, \Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{j=1}^n f(\mathbf{x}_j) = \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu})}$$

This function is called the likelihood function of the sample data of size n .

The values of μ_i 's and σ_{ik} 's, thus the values of $\boldsymbol{\mu}$ and Σ that maximize the likelihood function are the MLE's (maximum likelihood estimates) for observed data $\mathbf{x}_1, \dots, \mathbf{x}_n$.

For p -variate normal, the MLE of the population mean $\boldsymbol{\mu}$ is the sample mean

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

and the MLE of the population covariance Σ is the sample variance matrix with denominator n (note: not $n-1$)

$$\hat{\Sigma} = \mathbf{S}_n = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$$

Claim: The MLE (maximum likelihood estimators) $\bar{\mathbf{x}}$ and \mathbf{S}_n maximize the likelihood function, that is,

$$L(\bar{\mathbf{x}}, \mathbf{S}_n) = L(\hat{\boldsymbol{\mu}}, \hat{\Sigma}) = \max_{\boldsymbol{\mu}, \Sigma} L(\boldsymbol{\mu}, \Sigma)$$

where

$$\begin{aligned} L(\boldsymbol{\mu}, \Sigma) &= L(\boldsymbol{\mu}, \Sigma | \mathbf{x}_1, \dots, \mathbf{x}_n) = f(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\mu}, \Sigma) = \prod_{j=1}^n f(\mathbf{x}_j) \\ f(\mathbf{x}_j) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu})} \end{aligned}$$

Proof. The derivation of MLE $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$ and $\hat{\Sigma} = \mathbf{S}_n$ consists of two parts:

The $\boldsymbol{\mu}$ part

Note that $\boldsymbol{\mu}$ only occurs in the exponent. We will show that, for any Σ ,

$$\max_{\boldsymbol{\mu}, \Sigma} L(\boldsymbol{\mu}, \Sigma) = \max_{\boldsymbol{\mu}} \left(\max_{\Sigma} L(\boldsymbol{\mu}, \Sigma) \right)$$

Then we show that, given any Σ , the exponent (thus the likelihood function) is maximized when the estimator of $\boldsymbol{\mu}$ is $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$.

$$L(\bar{\mathbf{x}}, \Sigma) = L(\hat{\boldsymbol{\mu}}, \Sigma) = \max_{\boldsymbol{\mu}} L(\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} (n\mathbf{S}_n))} \quad (3)$$

The Σ part

We will show that, among all positive definite Σ , $L(\hat{\boldsymbol{\mu}}, \Sigma)$ is maximized when Σ is $\hat{\Sigma} = \mathbf{S}_n$, which yields

$$L(\hat{\boldsymbol{\mu}}, \hat{\Sigma}) = \max_{\Sigma} L(\hat{\boldsymbol{\mu}}, \Sigma) = \max_{\boldsymbol{\mu}, \Sigma} L(\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{np/2} |\mathbf{S}_n|^{n/2}} e^{-\frac{np}{2}} \quad (4)$$

The following are the steps of the proof.

Proof of the $\boldsymbol{\mu}$ part

First we need to rewrite the exponent in terms of matrix trace, which is the sum of the diagonal elements of a square matrix.

- Use the property $\text{tr}(AB) = \text{tr}(BA)$ and $\text{tr}(c) = c$ for any number c , the scalar

$$(\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) = \text{tr} \{ (\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) \} = \text{tr} \{ \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu})(\mathbf{x}_j - \boldsymbol{\mu})' \}$$

- Using the property $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$,

$$\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \sum_{i=1}^n \text{tr} \{ \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})' \} = \text{tr} \left\{ \Sigma^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})' \right\}$$

- Regrouping, and using the fact that $\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) = \mathbf{0}_p$,

$$\begin{aligned} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})' &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \boldsymbol{\mu})(\mathbf{x}_i - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \boldsymbol{\mu})' \\ &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' + \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}} - \boldsymbol{\mu})' + \sum_{i=1}^n (\bar{\mathbf{x}} - \boldsymbol{\mu})(\mathbf{x}_i - \bar{\mathbf{x}})' + \sum_{i=1}^n (\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' \\ &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' \\ &= n\mathbf{S}_n + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' \end{aligned}$$

- The exponent of the likelihood function becomes

$$-\frac{1}{2}tr \left\{ \Sigma^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})' \right\} = -\frac{1}{2}tr [\Sigma^{-1}(n\mathbf{S}_n)] - \frac{1}{2}tr \left\{ \Sigma^{-1} [n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})'] \right\} \quad (5)$$

By the positive definiteness of Σ^{-1} , the last term in (5) is non-negative:

$$\frac{1}{2}tr \left\{ \Sigma^{-1} [n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})'] \right\} = \frac{n}{2}tr \left\{ (\bar{\mathbf{x}} - \boldsymbol{\mu})' \Sigma^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \right\} = \frac{n}{2}(\bar{\mathbf{x}} - \boldsymbol{\mu})' \Sigma^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \begin{cases} = 0, & \boldsymbol{\mu} = \bar{\mathbf{x}}, \\ > 0, & \boldsymbol{\mu} \neq \bar{\mathbf{x}}. \end{cases}$$

Thus the exponent of the likelihood function

$$-\frac{1}{2}tr \left(\Sigma^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})' \right) \begin{cases} = -\frac{1}{2}tr(\Sigma^{-1}n\mathbf{S}_n), & \boldsymbol{\mu} = \bar{\mathbf{x}}, \\ < -\frac{1}{2}tr(\Sigma^{-1}n\mathbf{S}_n), & \boldsymbol{\mu} \neq \bar{\mathbf{x}}. \end{cases}$$

That is, the exponent is maximized when $\boldsymbol{\mu} = \bar{\mathbf{x}}$, for any Σ . So (3) is proved.

Proof of the Σ part

For any symmetric positive definite Σ , we may write $\Sigma = AA = \Sigma^{1/2}\Sigma^{1/2}$, where A , denoted as $\Sigma^{1/2}$, can be chosen as symmetric positive definite. Assuming \mathbf{S}_n is positive definite, then $\Sigma^{-1/2}n\mathbf{S}_n\Sigma^{-1/2}$ is also positive definite since $(\Sigma^{-1/2}v)'n\mathbf{S}_n(\Sigma^{-1/2}v) > 0, \forall v \in \mathbb{R}^p \setminus 0_p$. Therefore, we may denote the eigenvalues of $\Sigma^{-1/2}n\mathbf{S}_n\Sigma^{-1/2}$ as $\lambda_1 \geq \dots \geq \lambda_p > 0$.

By the relationship of matrix trace and matrix eigenvalues,

$$tr \{ \Sigma^{-1}(n\mathbf{S}_n) \} = tr \left\{ \Sigma^{-1/2}(n\mathbf{S}_n)\Sigma^{-1/2} \right\} = \sum_{k=1}^p \lambda_k$$

By the relationship of matrix determinants and matrix eigenvalues,

$$\frac{\det(n\mathbf{S}_n)}{\det(\Sigma)} = \det \{ \Sigma^{-1}(n\mathbf{S}_n) \} = \det \left\{ \Sigma^{-1/2}(n\mathbf{S}_n)\Sigma^{-1/2} \right\} = \prod_{k=1}^p \lambda_k$$

Thus

$$\frac{1}{|\Sigma|^{n/2}} e^{-\frac{1}{2}tr \{ \Sigma^{-1}(n\mathbf{S}_n) \}} = \left(\frac{\prod_{k=1}^p \lambda_k}{|n\mathbf{S}_n|} \right)^{n/2} e^{-\frac{1}{2} \sum_{k=1}^p \lambda_k} = \frac{1}{n^{np/2} |\mathbf{S}_n|^{n/2}} \prod_{k=1}^p \lambda_k^{n/2} e^{-\frac{1}{2} \lambda_k}$$

Since the function $t^{n/2}e^{-t/2}$ achieves unique maximum $n^{n/2}e^{-n/2}$ at $t = n$ for given n (calculus exercise),

$$t^{n/2}e^{-t/2} \leq n^{n/2}e^{-n/2} \implies \lambda_k^{n/2} e^{-\frac{1}{2}\lambda_k} \leq n^{n/2}e^{-n/2}$$

We obtained an upper bound

$$\frac{1}{|\Sigma|^{n/2}} e^{-\frac{1}{2}tr \{ \Sigma^{-1}(n\mathbf{S}_n) \}} \leq \frac{1}{n^{np/2} |\mathbf{S}_n|^{n/2}} \prod_{k=1}^p n^{n/2} e^{-n/2} = \frac{1}{|\mathbf{S}_n|^{n/2}} e^{-\frac{np}{2}}$$

The right hand side upper bound can be achieved when and only when $\lambda_k = n$ for all $k = 1, \dots, p$, that is, if and only if

$$\Sigma^{-1/2}(n\mathbf{S}_n)\Sigma^{-1/2} = nI_p \iff \Sigma = \mathbf{S}_n$$

Thus the MLE for Σ is

$$\hat{\Sigma} = \mathbf{S}_n$$

Therefore

$$\max_{\boldsymbol{\mu}, \Sigma} L(\boldsymbol{\mu}, \Sigma) = \max_{\Sigma} L(\hat{\boldsymbol{\mu}}, \Sigma) = \max_{\Sigma} \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2}tr \{ \Sigma^{-1}(n\mathbf{S}_n) \}} = \frac{1}{(2\pi)^{np/2} |\mathbf{S}_n|^{n/2}} e^{-\frac{np}{2}} = L(\hat{\boldsymbol{\mu}}, \hat{\Sigma}),$$

which is (4), as desired.

We have proved that the maximum likelihood is achieved at $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$ and $\hat{\Sigma} = \mathbf{S}_n$. □

Remarks on MLE

- Recall the sample generalized variance is $|\mathbf{S}| = \frac{n}{n-1} |\mathbf{S}_n|$, the maximum likelihood can be written as

$$\frac{1}{(2\pi)^{np/2} |\mathbf{S}_n|^{n/2}} e^{-\frac{np}{2}} = C_1 |\mathbf{S}_n|^{-n/2} = C |\mathbf{S}|^{-n/2} = \text{constant} \times (\text{sample generalized variance})^{-n/2}$$

This expression will be useful in deriving likelihood based tests.

- Maximum likelihood estimates have several nice properties, such as functional invariance, consistency, and asymptotic normality.

- The estimators match the univariate maximum likelihood estimators, where

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 = \frac{n-1}{n} s^2$$

That is, while the MLE of $\boldsymbol{\mu}$ is an unbiased estimator, the MLE of σ^2 is not unbiased.

- The separability of the two steps in the maximization process is unique for normal distributions.

6 Distribution properties of multivariate sample mean and covariance

We already know that the sample mean of an *i.i.d.* sample from $N_p(\boldsymbol{\mu}, \Sigma)$ is p -variate normal:

$$\bar{\mathbf{X}} \sim N_p\left(\boldsymbol{\mu}, \frac{1}{n}\Sigma\right)$$

the distribution of the sample variance matrix \mathbf{S} is related to the Wishart distribution, which is a lot more involved.

The following properties related to sample variance matrix \mathbf{S} are commonly used in multivariate applications as facts. The proofs can be found in Anderson (2003) or Murihead (1982), and are omitted in this course (other than the first one, to be proved later).

- $\bar{\mathbf{X}}$ and \mathbf{S} are independent.
- For *i.i.d.* p -vectors $\mathbf{X}_i \sim (\boldsymbol{\mu}, \Sigma)$ (not necessarily normal), as $n \rightarrow \infty$, their mean vector

$$\bar{\mathbf{X}} \rightarrow \boldsymbol{\mu} \quad \text{in probability,}$$

and their sample covariance matrix

$$\mathbf{S} \rightarrow \Sigma \quad \text{in probability.}$$

- Under the same assumption, as $n \rightarrow \infty$, by the multivariate central limit theorem,

$$\sqrt{n}(\bar{\mathbf{X}} - \boldsymbol{\mu}) \rightarrow N_p(0, \Sigma)$$

where the convergence is in distribution.

- Another central limit theorem type result:

$$n(\bar{\mathbf{X}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \rightarrow \chi_p^2$$

where the convergence is also in distribution.

- For the sample covariance matrix \mathbf{S} , the probability distribution of $(n-1)\mathbf{S}$ is called **Wishart distribution** with dimension p , degrees of freedom $n-1$, and parametric matrix Σ .

$$(n-1)\mathbf{S} \sim W_p(n-1, \Sigma)$$

Properties of Wishart distribution* (not required for this course)

- Density function

For $p < n$, a $p \times p$ symmetric positive definite random matrix M of $W_p(n, \Sigma)$ distribution has the density function

$$f(M) = \frac{|M|^{n-p-1}}{2^{pn/2} \Gamma_p(\frac{1}{2}n) |\Sigma|^{n/2}} e^{-\frac{1}{2}\text{tr}(\Sigma^{-1}M)}$$

where Γ_m is the generalized Gamma function, for real $\alpha > \frac{1}{2}(p-1)$,

$$\Gamma_p(\alpha) = \int_{A \text{ positive definite}} e^{\text{tr}(A)} |A|^{\alpha-(p+1)/2} (dA) = \pi^{p(p-1)/4} \prod_{i=1}^p \Gamma\left[\alpha - \frac{1}{2}(i-1)\right]$$

- Due to the symmetry of Wishart matrix M , the density function actually is a joint density on the $p(p+1)/2$ unique components of M .

- Wishart distribution is a multivariate generalization of the chi-squared distribution. For $p=1$,

$$(n-1)s^2 = \sum_{j=1}^n (X_j - \bar{X})^2 \sim \sigma^2 \chi_{n-1}^2$$

- An alternative definition

Suppose $\mathbf{X}_j, j=1, \dots, n$ are *i.i.d.* $\sim N_p(0, \Sigma)$, then

$$\sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j' \sim W_p(n, \Sigma)$$

Note that, for $p=1$, write $\Sigma = \sigma^2$, we obtain the univariate result:

Suppose $X_j, j=1, \dots, n$ are *i.i.d.* $\sim N_1(0, \sigma^2)$, then

$$\sum_{j=1}^n X_j X_j' = \sum_{j=1}^n X_j^2 \sim \sigma^2 \chi_n^2$$

where χ_n^2 is the chi-square distribution of n degrees of freedom.

- Summability

If $M_1 \sim W_p(n_1, \Sigma)$, $M_2 \sim W_p(n_2, \Sigma)$, and M_1 and M_2 are independent, then

$$M_1 + M_2 \sim W_p(n_1 + n_2, \Sigma)$$

This property is particularly useful when we compare multiple samples with the same covariance structure, where we need to estimate their common covariance matrix.

7 Checking normal assumption

In practice, how do we check if the normality assumption is reasonable for our data, even approximately?

Analogous to the univariate case, we can check if the data violate the normality assumption too much, and we may use component-wise variable transformations to make the data closer to normal distribution assumption.

Chi-square Q-Q plot

Univariate case: Check the normal Q-Q plot

Multivariate case: Check the chi-square Q-Q plot

When $x_j, j = 1, \dots, n$ are an i.i.d. sample from $N_p(\boldsymbol{\mu}, \Sigma)$,

$$d_j^2 = (\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) \quad j = 1, \dots, n$$

form an i.i.d. random sample from χ_p^2 . Thus the sample estimates

$$\hat{d}_j^2 = (\mathbf{x}_j - \bar{\mathbf{x}})' \Sigma^{-1} (\mathbf{x}_j - \bar{\mathbf{x}}) \quad j = 1, \dots, n$$

should be approximately χ_p^2 observations (not completely i.i.d.). A quantile to quantile comparison should approximately follow a line, analogous to the univariate normal Q-Q plot.

Variable transformations

The component variables of a multivariate vector can be transformed individually or jointly to be closer to normal distribution, in order to use many theoretical results that are based on the normality assumption.

The following are common univariate transformations based on data types.

- Counts: logarithm or square root transformation
- Proportions: logistic $\ln \frac{p}{1-p}$
- Correlations: Fisher transformation $\frac{1}{2} \ln \frac{1+r}{1-r}$
- Continuous scale: Power transformation or Box-Cox transformation

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda}, & \lambda \neq 0, \\ \ln(x), & \lambda = 0 \end{cases}$$

where λ is chosen to maximize

$$\ell(\lambda) = -\frac{n}{2} \ln \left(\frac{1}{n} \sum_{i=1}^n (x_i^{(\lambda)} - \bar{x}^{(\lambda)})^2 \right) + (\lambda - 1) \sum_{j=1}^n \ln(x_j), \quad \bar{x}^{(\lambda)} = \frac{1}{n} \sum_{i=1}^n x_i^{(\lambda)}$$

8 Spherical and elliptic distributions* (not required)

In many multivariate normal theory based applications in statistics, as we can see in model output of statistical software (e.g. mixture models), the applications often extend to include spherical and elliptic distributions, two close cousins of normal distributions. In the following we briefly introduce the two types of distributions. The derivation of their properties are not covered in the course.

Most results discussion in this course under the assumption of multivariate normal distribution can be generalized to spherical and elliptic distributions which include multivariate normal. Here we only outline basic and interesting properties of elliptic distributions.

The classical theory is treated thoroughly in Muirhead (1982). In the large data or high dimensionally case, multivariate normal assumption is often too restrictive. Spherical and elliptic distributions are used to generalize multivariate normal to high dimensional case (e.g. Barber and Kolar, 2016).

A p -variate random vector \mathbf{X} is a spherical distribution if $Q\mathbf{X}$ is of the same distribution as \mathbf{X} for any $p \times p$ orthogonal matrix Q . A spherical distribution is spherically symmetric and does not change under rotations and reflections of the coordinate system. The density function $f_{\mathbf{X}}(\mathbf{x})$ depends on \mathbf{x} via $\|\mathbf{x}\| = \mathbf{x}'\mathbf{x}$ only, so that $f_{\mathbf{X}}(\mathbf{x}) = g(\mathbf{x}'\mathbf{x})$ for some function g .

Examples of spherical distributions:

- $p = 2$, $f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{2\pi} e^{-(\mathbf{x}'\mathbf{x})^{1/2}}$, a bivariate generalization of univariate double-exponential distribution.
- $p = 2$, $f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{2\pi} (1 + (\mathbf{x}'\mathbf{x})^{-3/2})$, a bivariate Cauchy distribution.
- For $a_1, a_2 > 0$,

$$f_{\mathbf{X}}(\mathbf{x}) \sim \begin{cases} N_p(\mathbf{0}, a_1 \mathbf{I}), & \text{with probability } q \\ N_p(\mathbf{0}, a_2 \mathbf{I}), & \text{with probability } 1 - q \end{cases}$$

is a p -variate mixture normal.

If \mathbf{X} is of p -variate spherical distribution and \mathbf{A} is a $p \times p$ matrix, \mathbf{v} is a p -vector, then the p -variate random vector $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{v}$ is said to have an elliptical distribution. If $E(\mathbf{X})$ exists, then $E(\mathbf{X}) = \mathbf{0}$ and $E(\mathbf{Y}) = \mathbf{v}$. If $Cov(\mathbf{X})$ also exists, then $Cov(\mathbf{X}) = c\mathbf{I}_p$ for some $c \geq 0$, and $Cov(\mathbf{Y}) = c\mathbf{A}\mathbf{A}'$. When \mathbf{X} has density function $f_{\mathbf{X}}(\mathbf{x}) = g(\mathbf{x}'\mathbf{x})$, \mathbf{Y} has density

$$f_{\mathbf{Y}}(\mathbf{y}) = |\det(\mathbf{A}^{-1})| f_{\mathbf{X}}[\mathbf{A}^{-1}(\mathbf{y} - \mathbf{v})] = |\det(\mathbf{B}^{-1})|^{1/2} g[(\mathbf{y} - \mathbf{v})' \mathbf{B}^{-1} (\mathbf{y} - \mathbf{v})]$$

where $\mathbf{B} = \mathbf{A}\mathbf{A}'$ is proportional to $Cov(\mathbf{Y})$ and is called the scale matrix, which completely determines the correlation matrix $corr(\mathbf{Y}) = D^{-1/2} \mathbf{b} D^{-1/2}$, where D is the diagonal matrix $D = diag(\mathbf{B})$.

The above definition and calculations of elliptic distributions are similar to that of multivariate normal. Elliptic distributions share many properties with the multivariate normal: marginal and conditional distributions are still elliptic, conditional means are linear in the means, etc.

The components of an elliptically symmetric, non-normal random vector \mathbf{X} are dependent. Ellipticity plus independence implies normality: If \mathbf{X} is elliptic, and if its covariance matrix is diagonal if and only if the components are independent, then \mathbf{X} is multivariate normal.

Note Relevant chapter in the text by Johnson and Wichern: Chapter 4.