

Inference on the mean of multivariate normal

Now we start on methodology of statistical inference, that is, using information in multivariate sample data to draw conclusions on population properties of interest.

Unless otherwise stated, we consider the dimension of the random vectors p is fixed, of moderate size, while the sample size $n > p$.

In this section we discuss statistics inference methods used for one sample of multivariate variables. The sample data are treated as independent observations from a multivariate normal distribution.

The next section will consider inference methods on multiple samples, where sample data are treated as independent observations from several multivariate normal distributions.

1 One-sample test on the mean

We are interested in estimating the components μ_1, \dots, μ_p of the mean vector $\boldsymbol{\mu}$ jointly, taking the covariance relationship among the variables into consideration.

One common type of statistical inference on the mean is formulated as hypothesis testing. Although hypothesis testing appears to be a small part of statistical data analysis, the test statistic discussed here is closely related to more general inference such as providing confidence regions for $\boldsymbol{\mu}$.

The hypothesis test to be considered is

$$\begin{cases} H_o: & \boldsymbol{\mu} = \boldsymbol{\mu}_o \\ H_a: & \boldsymbol{\mu} \neq \boldsymbol{\mu}_o \end{cases}$$

2 Review: univariate hypothesis test and t -statistic

First recall the one sample t -test and its t -statistic on testing the mean in the univariate case $p = 1$, which you should have come across in previous courses.

For a random sample of size n , assume

$$X_1, \dots, X_n \text{ i.i.d. } \sim N(\mu, \sigma^2)$$

The hypotheses of interest is

$$H_o: \mu = \mu_o \text{ vs. } H_a: \mu \neq \mu_o$$

Then under H_o , the t -test statistic T_1 stated below is of the Student's t -distribution with degrees of freedom $n - 1$,

$$T_1 = \frac{\bar{X} - \mu_o}{s/\sqrt{n}} \sim t_{n-1} \quad \text{under } H_o$$

where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is the sample mean — an unbiased estimator of the population mean μ with $\mathbb{E}(\bar{X}) = \mu$, $s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is the sample variance — an unbiased estimator of the population variance σ^2 with $\mathbb{E}(s^2) = \sigma^2$, t_k represents the Student t distribution with k degrees of freedom, and \sim denotes the probability distribution.

The rationale of hypothesis testing inference is that, if the X_i 's are indeed independent draws from a $N(\mu_o, \sigma^2)$ random variable, then T_1 would be in the interval where the t_{n-1} distribution has high probability, that is, T_1 should be not too far away from 0. If the observed value of T_1 lands away from 0, then we conclude that the sample is not likely from a $N(\mu_o, \sigma^2)$ distribution. The cutoff of the decision boundary depends on our tolerance of the error that the data are from a $N(\mu_o, \sigma^2)$ distribution but fall to the far tail of the distribution just by chance.

Formally, at test level $\alpha = \mathbb{P}(\text{reject } H_o | H_o \text{ is true})$, which is the Type-I error that we decide to tolerate, H_o is rejected when the value of the sample test statistic $|T_1| \geq t_{n-1, \alpha/2}$, where $t_{n-1, \alpha/2} > 0$ is the upper $\frac{\alpha}{2}$ quantile of t_{n-1} distribution, defined by $\mathbb{P}(t_{n-1} > t_{n-1, \alpha/2}) = \frac{\alpha}{2}$.

Without the normality assumption of the X_i 's, the t test statistic T_1 is approximately of t_{n-1} distribution when n is moderately large, in the spirit of the Central Limit Theorem.

Remarks on t -distributions and t -statistic

- Recall that, a random variable Y is of t -distribution of degrees of freedom d if

$$Y = \frac{Z}{\sqrt{W/d}}$$

with

$$Z \sim N(0, 1), \quad W \sim \chi_d^2, \quad Z \perp\!\!\!\perp W,$$

where $\perp\!\!\!\perp$ denotes independent relationship, \sim the probability distribution of the random variable, $N(0, 1)$ the standard normal distribution, and χ_d^2 the chi-square distribution with degrees of freedom d .

- It is clear that t distribution is closely related to normal distribution and chi-square distribution.

Recall that, a random variable X is of chi-square distribution with degrees of freedom k , if

$$X = Z_1^2 + \dots + Z_k^2, \quad \text{with } Z_i \text{ i.i.d. } \sim N(0, 1)$$

That is, a chi-square random variable with degrees of freedom k can be expressed as the sum of squares of k independent standard normal random variables.

Consequently the density function of χ_k^2 -distribution can be derived from the density of standard normal (likely done in STAT245x0) by applying variable substitution method to the density functions.

χ_k^2 distribution is a special case of gamma distribution with the shape parameter $\frac{k}{2}$ and rate parameter $\frac{1}{2}$.

- From the density functions of normal distribution and chi-square distribution, the density of t distribution can be derived. The density function of t -distribution of degrees of freedom d is

$$f(t) = \frac{1}{\sqrt{d} B(\frac{1}{2}, \frac{d}{2})} \left(1 + \frac{t^2}{d} \right)^{-\frac{d+1}{2}}, \quad t \in \mathbb{R}.$$

In the above, $B(\cdot, \cdot)$ denotes the Beta function

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = \int_0^1 x^{a-1}(1-x)^{b-1} dx, \quad a, b > 0,$$

where $\Gamma(\cdot)$ is the Gamma function

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$

Some most useful properties of the Gamma function are: $\Gamma(1) = 1$, $\Gamma(z+1) = z\Gamma(z)$, thus $\Gamma(k) = (k-1)!$ for integer k , and $\Gamma(\frac{1}{2}) = \sqrt{\pi}$.

- The distribution of the test statistic t asymptotically approaches the standard normal distribution as $n \rightarrow \infty$. The proof of the asymptotic distribution uses the exponential limit formula

$$e^x = \lim_{x \rightarrow \infty} \left(1 + \frac{x}{n}\right)^x$$

and the Gamma function property

$$\lim_{n \rightarrow \infty} \frac{\Gamma(n+c)}{\Gamma(n)n^c} = 1$$

for scalar c .

- The above leads to the useful result that if $\{X_1, \dots, X_n\}$ is an i.i.d. sample from a general probability distribution with mild regularity conditions such as having second moment, then \bar{X} is asymptotically of normal distribution for large n , by the Central Limit Theorem, and the t statistic is asymptotically approaching $N(0, 1)$.

If we square the t -statistic then rewrite the test statistic as

$$T_1^2 = (\bar{x} - \mu_o) \left(\frac{s^2}{n}\right)^{-1} (\bar{x} - \mu_o) = (\bar{x} - \mu_o) [Var(\bar{x} - \mu_o)]^{-1} (\bar{x} - \mu_o)$$

We then obtain a convenient form suitable for generalization of the t -statistic to multidimensional cases.

3 Hotelling's T^2 statistic

The multivariate normal generalization of the t -test is to test $H_0: \mu = \mu_o$, when $\mu, \mu_o \in \mathbb{R}^p$ with $p \geq 2$.

Recall that observed multivariate data can be displayed as

$$\begin{array}{l} \text{observation 1} \\ \text{observation 2} \\ \vdots \\ \text{observation } j \\ \vdots \\ \text{observation } n \end{array} \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_j \\ \vdots \\ \mathbf{x}'_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2k} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{j1} & x_{j2} & \cdots & x_{jk} & \cdots & x_{jp} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} & \cdots & x_{np} \end{bmatrix}$$

The data can be viewed as realizations of n i.i.d. multivariate vector,

$$\begin{array}{l} \text{random vector 1} \\ \text{random vector 2} \\ \vdots \\ \text{random vector } j \\ \vdots \\ \text{random vector } n \end{array} \begin{bmatrix} \mathbf{X}'_1 \\ \mathbf{X}'_2 \\ \vdots \\ \mathbf{X}'_j \\ \vdots \\ \mathbf{X}'_n \end{bmatrix} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1k} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2k} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ X_{j1} & X_{j2} & \cdots & X_{jk} & \cdots & X_{jp} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{nk} & \cdots & X_{np} \end{bmatrix}$$

Assuming the random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. from an $N_p(\mu, \Sigma)$, a p -variate multivariate normal distribution. The sample mean vector is

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

and sample variance-covariance matrix is

$$\mathbf{S} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})'$$

Denote the proposed null mean vector

$$\mu_0 = \begin{bmatrix} \mu_{10} \\ \vdots \\ \mu_{p0} \end{bmatrix}$$

The test statistic is **Hotelling's T^2** ,

$$T^2 = (\bar{\mathbf{X}} - \mu_0)' [\widehat{Cov}(\bar{\mathbf{X}} - \mu_0)]^{-1} (\bar{\mathbf{X}} - \mu_0) = (\bar{\mathbf{X}} - \mu_0)' \left(\frac{\mathbf{S}}{n}\right)^{-1} (\bar{\mathbf{X}} - \mu_0)$$

To test the hypotheses

$$H_o: \mu = \mu_o \quad \text{vs.} \quad H_a: \mu \neq \mu_o$$

Under the multivariate normality assumption, the test statistic T^2 is related to the F distribution under the null,

$$T^2 \sim \frac{(n-1)p}{n-p} F_{p, n-p} \quad \text{under } H_o$$

With the observed data $x_1, \dots, x_n \in \mathbb{R}^p$ with sample mean \bar{x} , given a test level α , we would reject H_o if

$$\frac{n-p}{(n-1)p} T^2 > F_{p, n-p}(\alpha)$$

where $F_{p, n-p}(\alpha)$ is the upper 100α percentile of $F_{p, n-p}$, that is, $\mathbb{P}\{F_{n, n-p} > F_{p, n-p}(\alpha)\} = \alpha$.

Remarks on T^2

- Analogous to the t -statistic in the univariate test, T^2 is a sample statistic, its probability distribution does not involve μ or Σ , which consist of unknown parameters.
- T^2 is invariant under non-singular linear transformation of the variables $\mathbf{X}_{p \times 1} \rightarrow C_{p \times p} \mathbf{X}_{p \times 1} + d_{p \times 1}$ (exercise).
- The \mathbf{S} in Hotelling's T^2 is the $p \times p$ sample covariance matrix.

Elementwise,

$$\mathbf{S} = \left[\frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) \right]_{i,k=1, \dots, p} = [s_{ik}]_{i,k=1, \dots, p}$$

where s_{ik} is the sample covariance between the i th and k th variables, as defined at the beginning of the course.

In vector-matrix form of the observed sample data,

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

The p -variate vector \mathbf{x}_j is the j th observation. In other words, $(n-1)\mathbf{S}$ can be written as the sum of n matrices, each of the n matrices is of dimensions $p \times p$, the j th matrix has the form

$$(\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' = \begin{bmatrix} x_{j1} - \bar{x}_1 \\ x_{j2} - \bar{x}_2 \\ \vdots \\ x_{jk} - \bar{x}_k \\ \vdots \\ x_{jp} - \bar{x}_p \end{bmatrix} \begin{bmatrix} x_{j1} - \bar{x}_1, & x_{j2} - \bar{x}_2, & \cdots, & x_{jk} - \bar{x}_k, & \cdots, & x_{jp} - \bar{x}_p \end{bmatrix} = \left[(x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) \right]_{i,k=1, \dots, p}$$

where $j = 1, \dots, n$ index the observations.

- Write the sample covariance matrix as $S = [s_{ik}]_{p \times p}$, and assume that S is of full rank, that is, the p columns are linearly independent. Denote the inverse of the sample covariance matrix as

$$S^{-1} = [s^{ik}]_{p \times p}$$

Then Hotelling's T^2 can be written as a sum of p^2 univariate terms:

$$T^2 = (\bar{x} - \mu_0)' \left(\frac{S}{n} \right)^{-1} (\bar{x} - \mu_0) = n \sum_{i=1}^p \sum_{k=1}^p s^{ik} (x_i - \mu_{i0})(x_k - \mu_{k0})$$

- Hotelling's T^2 is of $\frac{(n-1)p}{n-p} F_{p, n-p}$ distribution under the null hypothesis.

The F -distribution of T^2 is derived from the independence between \bar{X} & S , and the distribution of S^{-1} (an inverse Wishart distribution, not covered in this class).

The proof is complicated and beyond the coverage of this course, thus omitted.

A complete, detailed proof can be found in Chapter 3 of Muirhead (1982).

- When $n \rightarrow \infty$, T^2 is asymptotically of χ_p^2 distribution.

However the test using T^2 is more useful when n is not very large.

- Review of F distribution

- A random variable is of F -distribution with degrees of freedom d_1 and d_2 , denoted as F_{d_1, d_2} , if it is the ratio of two independent Chi-square random variables divided by their degrees of freedom respectively,

$$\frac{Y_1/d_1}{Y_2/d_2}, \quad \text{with } Y_1 \sim \chi_{d_1}^2, \quad Y_2 \sim \chi_{d_2}^2$$

and the numerator is independent of the denominator,

$$Y_1 \perp\!\!\!\perp Y_2.$$

In the above, χ_d^2 denotes the Chi-square distribution with degrees of freedom d .

- The density of F_{d_1, d_2} is

$$f(y)_{d_1, d_2} = \frac{1}{B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)} \left(\frac{d_1}{d_2}\right)^{\frac{d_1}{2}} y^{\frac{d_1}{2}-1} \left(1 + \frac{d_1}{d_2}y\right)^{-\frac{d_1+d_2}{2}}, \quad y > 0.$$

- In particular, when $d_1 = p = 1$, $d_2 = n - p = n - 1$, the $F_{p, n-p}$ density becomes

$$f(y)_{1, n-1} = \frac{1}{B\left(\frac{1}{2}, \frac{n-1}{2}\right)} \left(\frac{1}{n-1}\right)^{\frac{1}{2}} y^{-\frac{1}{2}} \left(1 + \frac{y}{n-1}\right)^{-\frac{n}{2}}, \quad y > 0,$$

which is the density of $(t_{n-1})^2$, the square of a t-distribution random variable with degrees of freedom $n - 1$.

- Note that $p < n$ is assumed in the above derivation and application.

When $p > n$, inverse of covariance matrix does not exist. Even if pseudo-inverse matrix is used, T^2 test loses power (Bai and Saranadasa 1996). Various remedies have been studied and proposed.

4 Hotelling's T^2 and likelihood ratio tests

Hotelling's T^2 for multivariate normal is equivalent to a likelihood ratio test.

Let θ be a parameter of interest, $\theta \in \Theta$, the whole parameter space. To test if the true parameter value of θ lies inside a specific subspace Θ_0 , the hypothesis test can be expressed as

$$\begin{cases} H_0: & \theta \in \Theta_0 \subset \Theta \\ H_a: & \theta \in \Theta \setminus \Theta_0 \end{cases}$$

For example, we may test on $\Theta_0 = \{(0, 0)\}$, $\Theta = \mathbb{R}^2$.

Likelihood Ratio

The Likelihood Ratio test statistic is the ratio of the maximum likelihood under H_0 and the unrestricted maximum likelihood under a more general $H_0 \cup H_a$.

$$\Lambda = \frac{\max_{\theta \in \Theta_0} L(\theta)}{\max_{\theta \in \Theta} L(\theta)}, \quad \text{where } \Theta_0 \subset \Theta.$$

- Possible values of Λ are in $[0, 1]$.
- When H_0 is true, Λ is close to 1.
- When H_a is true, Λ is closer to 0.
- Consequently, H_0 is rejected if $\Lambda < c$, where c is a critical value depending on the test level α (*type-I error*).
- When the sample size n is large, often $-2 \ln \Lambda$ is used as the likelihood ratio test statistic, with asymptotic χ^2 distribution under the null,

$$-2 \ln \Lambda \sim \chi_{v-v_0}^2 \quad \text{under } H_0 \quad (\text{asymptotically})$$

where

$$v_0 = \dim(\Theta_0), \quad v = \dim(\Theta).$$

Maximum likelihood and Likelihood Ratio for multivariate normal

To test the mean parameter vector of multivariate normal, the likelihood ratio test has a simpler form.

We have derived (in lecture notes "The multivariate normal distribution" in week 2) that the maximum of the likelihood function for a multivariate normal sample of size n is

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

where the ML estimators for the population mean vector and variance-covariance matrix are

$$\hat{\mu} = \bar{x}, \quad \hat{\Sigma} = S_n = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})'$$

Therefore, under the null hypothesis $H_0: \mu = \mu_0$, the maximum likelihood has the form

$$\max_{\Sigma} L(\mu_0, \Sigma) = L(\hat{\mu}_0, \hat{\Sigma}_0) = \frac{1}{(2\pi)^{np/2} |\hat{\Sigma}_0|^{n/2}} e^{-np/2}$$

where the ML estimator for the variance-covariance matrix under the null H_0 is

$$\hat{\Sigma}_0 = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_0)(\mathbf{x}_j - \boldsymbol{\mu}_0)'$$

For $\boldsymbol{\mu}_0 \in \mathbb{R}^p$, consider the hypotheses

$$\begin{cases} H_0 : & \boldsymbol{\mu} = \boldsymbol{\mu}_0 \\ H_a : & \boldsymbol{\mu} \neq \boldsymbol{\mu}_0 \end{cases}$$

The resulting likelihood ratio has a simple form.

$$\Lambda = \frac{\max_{\Sigma} L(\boldsymbol{\mu}_0, \Sigma)}{\max_{\boldsymbol{\mu}, \Sigma} L(\boldsymbol{\mu}, \Sigma)} = \left(\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_0|} \right)^{n/2} = \left(\frac{\left| \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right|}{\left| \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_0)(\mathbf{x}_j - \boldsymbol{\mu}_0)' \right|} \right)^{n/2}$$

The statistics $\Lambda^{2/n}$ is called the *Wilks' lambda*.

Wilks' Lambda and Hotelling's T^2

The following gives the relation between the likelihood ratio and Hotelling's T^2 .

$$\text{Wilks' Lambda} \quad \Lambda^{2/n} = \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_0|} = \left(1 + \frac{T^2}{n-1} \right)^{-1}$$

$$\text{Hotelling's } T^2 \quad T^2 = (n-1) \left(\frac{|\hat{\Sigma}_0|}{|\hat{\Sigma}|} - 1 \right)$$

Proof. In the following we derive the relation between Λ and T^2 .

We use the following matrix algebra result on block matrix determinants:

$$\begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = |A_{11}| |A_{22} - A_{21}A_{11}^{-1}A_{12}| = |A_{22}| |A_{11} - A_{12}A_{22}^{-1}A_{21}| \quad (1)$$

under the conditions that the diagonal block are invertible square matrices: $|A_{11}| \neq 0, |A_{22}| \neq 0$.

The results in (1) follow from the Schur complements of block matrices.

Assuming sample variance-covariance matrix \mathbf{S} is of full rank, (1) can be applied to the block matrix

$$\begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = \begin{vmatrix} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' & \sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0) \\ \sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' & -1 \end{vmatrix},$$

We obtain

$$\begin{aligned} |A_{22}| |A_{11} - A_{12}A_{22}^{-1}A_{21}| &= (-1) \left| \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' \right| \\ &= (-1) \left| \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_0)(\mathbf{x}_j - \boldsymbol{\mu}_0)' \right| \end{aligned} \quad (2)$$

and

$$\begin{aligned} |A_{11}| |A_{22} - A_{21}A_{11}^{-1}A_{12}| &= \left| \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right| \\ &\times \left| -1 - n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' \left(\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right)^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0) \right| \end{aligned} \quad (3)$$

Notice that

$$\begin{aligned} &\left| -1 - n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' \left(\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right)^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0) \right| \\ &= -1 - n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' [(n-1)\mathbf{S}]^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' \\ &= (-1) \left(1 + \frac{1}{n-1} n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0) \right) \\ &= (-1) \left(1 + \frac{1}{n-1} T^2 \right) \end{aligned}$$

where Hotelling's $T^2 = n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)\mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)$.

From (1), The right hand sides of (2) and (3) are equal.

$$(-1) \left| \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_0)(\mathbf{x}_j - \boldsymbol{\mu}_0)' \right| = (-1) \left| \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right| \left(1 + \frac{1}{n-1} T^2 \right)$$

The equality above leads to the desired result

$$\Lambda^{2/n} = \frac{\left| \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right|}{\left| \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_0)(\mathbf{x}_j - \boldsymbol{\mu}_0)' \right|} = \left(1 + \frac{1}{n-1} T^2 \right)^{-1}$$

which shows the desired relationship between Λ and T^2 . □

Usefulness of the alternative expression of T^2

There is a need to compute the inverse matrix \mathbf{S}^{-1} in evaluating Hotelling's T^2 using $T^2 = n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)\mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)'$. Calculating \mathbf{S}^{-1} is computationally demanding, unless the dimensions is quite low, thus often avoided. The new formula $T^2 = (n-1)(|\hat{\Sigma}_0|/|\hat{\Sigma}| - 1)$ is often easier to evaluate.

5 Confidence regions

Let $\boldsymbol{\theta}$ be the vector of parameters of interest, $\boldsymbol{\theta} \in \Theta$, the parameter space. In the univariate case, a parameter estimate $\hat{\theta} = \hat{\theta}(X)$ based on a random sample $X = \{X_1, \dots, X_n\}$ is often presented along with a confidence interval $\mathbb{P}(L(X) < \theta < U(X)) = 1 - \alpha$, where L and U can be calculated from the sample X , similar to $\hat{\theta}$ which can be evaluated from sample values of X . The region (L, U) is a $100(1 - \alpha)\%$ confidence interval for the unknown parameter θ .

Analogously, in multivariate inference, a $100(1 - \alpha)\%$ **confidence region** $R(\mathbf{X})$ for $\boldsymbol{\theta}$ based on a p -variate random sample $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ is defined as

$$\mathbb{P}(\boldsymbol{\theta} \in R(\mathbf{X})) = 1 - \alpha,$$

in the sense that among all possible sample data \mathbf{X} , the parameter vector $\boldsymbol{\theta}$ is in the the region $R(\mathbf{X})$ about $(1 - \alpha)100\%$ of the time. The probability in the above equality is with respect to the true (unknown) parameter $\boldsymbol{\theta}$.

Confidence region of the mean vector by T^2

Recall

$$n(\bar{\mathbf{X}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) = T^2 \sim \frac{(n-1)p}{n-p} F_{p, n-p} \quad \text{under } H_0.$$

where $F_{d_1, d_2}(\alpha)$ is the upper $\alpha \times 100\%$ percentile of the F distribution F_{d_1, d_2} , so $\mathbb{P}[F_{d_1, d_2} \geq F_{d_1, d_2}(\alpha)] = \alpha$.

The $(1 - \alpha)100\%$ T^2 confidence region of the mean μ based on the sample mean \bar{X} is derived from

$$\mathbb{P}\left(n(\bar{X} - \mu)'S^{-1}(\bar{X} - \mu) \leq \frac{(n-1)p}{n-p}F_{p, n-p}(\alpha)\right) = 1 - \alpha$$

Remarks on T^2 confidence region

- When \bar{X} is multivariate normal and μ is known (only in theory), the \bar{x} values satisfying the inequality

$$n(\bar{x} - \mu)'S^{-1}(\bar{x} - \mu) \leq \frac{(n-1)p}{n-p}F_{p, n-p}(\alpha) \quad (4)$$

gives an elliptically shaped region centered at μ , in which sample mean vector \bar{x} falls with probability $1 - \alpha$.

- When \bar{x} is from an i.i.d. sample of multivariate normal and μ is unknown (as in practice), plausible values of μ satisfying (4) form a random region of ellipsoid centered at sample mean \bar{x} , and the region includes the true population mean vector μ with confidence level $(1 - \alpha)100\%$. This is the T^2 **confidence region** of the mean vector.
- When $p = 1$, the inequality gives the usual $(1 - \alpha)100\%$ confidence interval constructed by t -test.
- The Euclidean distance between vectors $x = (x_1, \dots, x_p)'$ and $\mu = (\mu_1, \dots, \mu_p)'$ is

$$\sqrt{(x_1 - \mu_1)^2 + \dots + (x_p - \mu_p)^2} = \sqrt{(x - \mu)'(x - \mu)}$$

The Mahalanobis distance (a.k.a. **statistical distance**), a type of standardized distance, of variable vector x to constant vector μ is defined as

$$\sqrt{(x - \mu)' \widehat{Cov}(x)^{-1}(x - \mu)}$$

where $\widehat{Cov}(x)$ denotes the sample covariance matrix of x .

- The Mahalanobis distance takes into consideration the covariance structure of the data. Hotelling's T^2 uses the Mahalanobis distance between the mean \bar{x} and μ_0 to test $H_0: \mu = \mu_0$. The T^2 Confidence Region uses the statistical distance to construct a confidence region for the population mean μ .
- For multivariate normal $X \sim N_p(\mu, \Sigma)$, the contours of constant density function

$$(x - \mu)' \Sigma^{-1}(x - \mu) = c^2$$

are ellipsoids in \mathbb{R}^p , centered at μ with axes proportional to $\sqrt{\lambda_i}e_i$, where $\Sigma e_i = \lambda_i e_i, i = 1, \dots, p$.

Consequently, the constant contours of T^2 confidence region are also ellipsoids in \mathbb{R}^p , centered at \bar{x} with axes proportional to $\sqrt{\lambda_i} \sqrt{\frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha)} e_i$, where $S e_i = \lambda_i e_i$ for $i = 1, \dots, p$.

6 Simultaneous confidence intervals of component means

Other than the above ellipsoidal confidence region formed by \bar{x} , in many cases it is still desirable to have a statement

" $\mu_k \in (a_k, b_k)$ for all $k = 1, \dots, p$ "

with a reasonable level of confidence.

Such statement must involve all \bar{x}_k , forming a hyper-rectangular region $\in \mathbb{R}^p$. Compared with the advantages of ellipsoidal confidence region using statistical distance, these hyper-rectangles (orthotopes) are often easier to form and to compute.

Confidence interval for a linear combination of component means

To relate a component mean \bar{x}_k with the mean vector \bar{x} , notice that \bar{x}_k can be expressed as $a'\bar{x}$, where $a = [0 \dots 0 \ 1 \ 0 \dots 0]'$ has the k th component = 1 and 0 for other components.

First we will derive a confidence interval for $a'\bar{\mu}$, a linear combination of the component means.

If $X \sim N_p(\mu, \Sigma)$, then any linear combination of its components (written as $a'X$) is a univariate normal,

$$a'X \sim N(\mu, \sigma^2), \quad \text{with } \mu = a'\mu, \quad \sigma^2 = a'\Sigma a$$

Especially, the mean vector \bar{X} of a random sample has the same property.

$$\bar{X} \sim N_p\left(\mu, \frac{1}{n}\Sigma\right) \Rightarrow a'\bar{X} \sim N\left(a'\mu, \frac{1}{n}a'\Sigma a\right)$$

If x_1, \dots, x_n is a p -variate random sample from a common distribution with sample mean \bar{x} and sample covariance matrix S , then for any $a \in \mathbb{R}^p \setminus \{0\}$,

$$\begin{cases} y_1 = a'x_1 \\ y_2 = a'x_2 \\ \vdots \\ y_n = a'x_n \end{cases}$$

is a random sample from the univariate normal variable $Y = a'X$, with univariate sample mean and sample variance

$$\bar{y} = \frac{1}{n} \sum_{j=1}^n y_j = a'\bar{x}, \quad s_y^2 = a'Sa.$$

For a fixed vector $a \in \mathbb{R}^p$, the t -statistic

$$t = \frac{\bar{y} - \mu_y}{\sqrt{s_y^2/n}} = \frac{a'\bar{x} - a'\mu}{\sqrt{a'Sa/n}} \sim t_{n-1}$$

since $\mu_y = a'\mu$ is the true mean of Y . Therefore a $(1 - \alpha)100\%$ confidence interval of parameter $\mu_y = a'\mu$ is

$$a'\bar{x} - t_{n-1, \alpha/2} \frac{\sqrt{a'Sa}}{\sqrt{n}} \leq a'\mu \leq a'\bar{x} + t_{n-1, \alpha/2} \frac{\sqrt{a'Sa}}{\sqrt{n}} \quad (5)$$

In other words, for an observed \bar{x} and fixed $a \in \mathbb{R}^p$, the $(1 - \alpha)100\%$ confidence interval consists of plausible μ values for which $\left| \frac{a'\bar{x} - a'\mu}{\sqrt{a'Sa/n}} \right| \leq t_{n-1, \alpha/2}$, or equivalently,

$$\frac{n[a'(\bar{x} - \mu)]^2}{a'Sa} \leq t_{n-1, \alpha/2}^2 \quad (6)$$

Thus, if we only want to construct only one confidence interval for a specific μ_k , we choose $a = [0 \dots 0 \ 1 \ 0 \dots 0]'$ with the k th component = 1 and 0 otherwise. Then the above inequality (6) becomes

$$\frac{n[a'(\bar{x} - \mu)]^2}{a'Sa} = \frac{(\bar{x}_k - \mu_k)^2}{s_{kk}/n} \leq t_{n-1, \alpha/2}^2$$

which yields a confidence interval for a single μ_k .

6.1 Hyper-rectangle by marginal confidence intervals of component means (naive)

If we ignore the multivariate nature and treat the component variables as independent, a marginal confidence interval of each component mean μ_k can be obtained from the t-distribution,

$$\bar{x}_k - t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}} \quad (7)$$

This can be derived from univariate one-sample t -statistic. A single such interval contains its mean parameter μ_k with $(1-\alpha)100\%$ confidence, ignoring or without considering probability statements on other components. However in general we do not know the confidence level that the p intervals contain all μ_k simultaneously.

The only exception is when the component variables are independent, then

$$\begin{aligned} & P \left\{ \bar{x}_k - t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}}, \forall k = 1, \dots, p \right\} \\ &= \prod_{k=1}^p P \left\{ \bar{x}_k - t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + t_{n-1, \alpha/2} \sqrt{\frac{s_{kk}}{n}} \right\} \\ &= (1-\alpha)^p < 1-\alpha \end{aligned}$$

That is, the p marginal intervals in the form of (7) contain all μ_k simultaneously would have confidence level $(1-\alpha)^p 100\%$ when independence holds. The coverage probability of the p simultaneous confidence intervals would be $(1-\alpha)^p < 1-\alpha$, and independence among component variables is not true or interesting in multivariate studies.

6.2 Bonferroni simultaneous confidence intervals of component means

In order to have a desirable confidence level for the p confidence intervals in (7) to contain all μ_k simultaneously, we need to pay the price of making the intervals wider, according to the commonly used method of Bonferroni.

Let C_k = the event that the confidence interval of μ_k contains μ_k with confidence level $(1-\alpha_k)100\%$. So

$$P(C_k = \text{true}) = 1 - \alpha_k$$

for a given $k = 1, \dots, p$. To achieve the desired confidence level for simultaneous confidence intervals, we need

$$P(C_k = \text{true}, k = 1, \dots, p) \geq 1 - \alpha$$

The Bonferroni method is based on the following probability derivation.

$$\begin{aligned} P(C_k = \text{true}, k = 1, \dots, p) &= 1 - P(C_k = \text{false}, \text{for some } k) \\ &\geq 1 - \sum_{k=1}^p P(C_k = \text{false}) = 1 - \sum_{k=1}^p [1 - P(C_k = \text{true})] \\ &= 1 - \sum_{k=1}^p \alpha_k = 1 - (\alpha_1 + \dots + \alpha_p) \end{aligned}$$

The desired $(1-\alpha)100\%$ confidence level for p simultaneous confidence intervals can be guaranteed if $\sum_{i=1}^p \alpha_i \leq \alpha$.

$$\alpha_1 + \dots + \alpha_p \leq \alpha \implies P(C_k = \text{true}, k = 1, \dots, p) \geq 1 - \alpha$$

A common, convenient choice is to pick each $\alpha_k = \alpha/p$, then we obtain the Bonferroni confidence intervals

$$\bar{x}_k - t_{n-1, \alpha/2p} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + t_{n-1, \alpha/2p} \sqrt{\frac{s_{kk}}{n}}, \quad k = 1, \dots, p.$$

The statement is true simultaneously for all $k = 1, \dots, p$, with confidence level at least $(1-\alpha)100\%$.

$$P \left(\bar{x}_k - t_{n-1, \alpha/2p} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + t_{n-1, \alpha/2p} \sqrt{\frac{s_{kk}}{n}}, k = 1, \dots, p \right) \geq 1 - \alpha$$

6.3 Asymptotic approximation and simultaneous confidence intervals

When the sample size n is large relative to p (i.e. when $n-p$ is large), the Central Limit Theorem leads to

$$n(\bar{X} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{X} - \boldsymbol{\mu}) \sim \chi_p^2 \quad (\text{approximately asymptotically})$$

without assuming normality of the distribution.

The χ^2 approximation gives the asymptotic simultaneous confidence intervals

$$\bar{x}_k - \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{s_{kk}}{n}} \leq \mu_k \leq \bar{x}_k + \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{s_{kk}}{n}}, \quad k = 1, \dots, p.$$

The above asymptotic statement comes from the multivariate normal property (exercise)

$$\mathbf{X} \sim N_p(\boldsymbol{\mu}, \Sigma) \implies (\mathbf{X} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}) \sim \chi_p^2$$

and the more general confidence interval result for any linear combination $\mathbf{a}'\boldsymbol{\mu}$, that when $n-p$ is sufficiently large,

$$\mathbf{a}'\bar{\mathbf{X}} \pm \sqrt{\chi^2(\alpha)} \sqrt{\frac{\mathbf{a}'\mathbf{S}\mathbf{a}}{n}}$$

is a $(1-\alpha)100\%$ confidence interval of $\mathbf{a}'\boldsymbol{\mu}$ asymptotically approximately.

6.4 Simultaneous confidence intervals of several component means by T^2

The inequality (6) is developed for a fixed \mathbf{a} . In the following we will use Hotelling's T^2 test statistic to construct an analogous result for all \mathbf{a} , which leads to hyper-rectangular simultaneous confidence intervals of $\bar{\mu}_k, k = 1, \dots, p$.

The left hand side of (6) is in the form

$$\frac{(\mathbf{v}'\mathbf{w})^2}{\mathbf{v}'\mathbf{B}\mathbf{v}}, \quad \text{with } \mathbf{v} = \mathbf{a}, \mathbf{w} = \bar{\mathbf{x}} - \boldsymbol{\mu}, \mathbf{B} = \frac{1}{n}\mathbf{S}.$$

A useful extension of the Cauchy-Schwarz Inequality

$$(\mathbf{v}'\mathbf{w})^2 \leq \|\mathbf{v}\|^2 \|\mathbf{w}\|^2 = (\mathbf{v}'\mathbf{v})(\mathbf{w}'\mathbf{w})$$

is the Maximization Lemma,

$$\frac{(\mathbf{v}'\mathbf{w})^2}{\mathbf{v}'\mathbf{B}\mathbf{v}} \leq \mathbf{w}\mathbf{B}^{-1}\mathbf{w}$$

where the equality holds if and only if $\mathbf{v} \propto \mathbf{B}^{-1}\mathbf{w}$ (\propto means proportional to).

The Maximization Lemma implies that the left side of (6) achieves maximum when we choose $\mathbf{a} \propto n\mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu})$.

Therefore for all $\mathbf{a} \in \mathbb{R}^p$, we have

$$\frac{n[\mathbf{a}'(\bar{\mathbf{x}} - \boldsymbol{\mu})]^2}{\mathbf{a}'\mathbf{S}\mathbf{a}} \leq \max_{\mathbf{a}} \frac{n[\mathbf{a}'(\bar{\mathbf{x}} - \boldsymbol{\mu})]^2}{\mathbf{a}'\mathbf{S}\mathbf{a}} = n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) = T^2 \quad (8)$$

The last expression above is in the form of Hotelling's T^2 statistic. Since $\boldsymbol{\mu}$ is the true mean,

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) = T^2 \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

Thus the T^2 in (8) satisfies the probability property

$$P\left(T^2 \leq \frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)\right) = 1 - \alpha,$$

where $F_{d_1, d_2}(\alpha)$ denotes the upper α quantile of the F_{d_1, d_2} distribution, as defined earlier.

Since the left side inequality in (8) holds for all $\mathbf{a} \in \mathbb{R}^p$,

$$\left| \frac{\mathbf{a}'\bar{\mathbf{x}} - \mathbf{a}'\boldsymbol{\mu}}{\sqrt{\mathbf{a}'\mathbf{S}\mathbf{a}/n}} \right| \leq T^2, \quad \forall \mathbf{a} \in \mathbb{R}^p,$$

from the probability property of T^2 we obtain

$$P\left(\left| \frac{\mathbf{a}'\bar{\mathbf{x}} - \mathbf{a}'\boldsymbol{\mu}}{\sqrt{\mathbf{a}'\mathbf{S}\mathbf{a}/n}} \right| \leq \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)}\right) \geq 1 - \alpha, \quad \forall \mathbf{a} \in \mathbb{R}^p.$$

In other words,

$$P\left(\left| \mathbf{a}'\bar{\mathbf{x}} - \mathbf{a}'\boldsymbol{\mu} \right| \leq \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)} \frac{\sqrt{\mathbf{a}'\mathbf{S}\mathbf{a}}}{n}\right) \geq 1 - \alpha, \quad \forall \mathbf{a} \in \mathbb{R}^p.$$

The probability statement gives an at least $(1 - \alpha)100\%$ confidence interval of parameter $\mathbf{a}'\boldsymbol{\mu}$, for any $\mathbf{a} \in \mathbb{R}^p$,

$$P\left(\mathbf{a}'\bar{\mathbf{x}} - \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)} \frac{\sqrt{\mathbf{a}'\mathbf{S}\mathbf{a}}}{n} \leq \mathbf{a}'\boldsymbol{\mu} \leq \mathbf{a}'\bar{\mathbf{x}} + \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)} \frac{\sqrt{\mathbf{a}'\mathbf{S}\mathbf{a}}}{n}\right) \geq 1 - \alpha. \quad (9)$$

In particular, when we choose vector \mathbf{a} to be $\mathbf{a} = [0 \cdots 0 \ 1 \ 0 \cdots 0]'$ with the k th component = 1 and 0 otherwise, then

$$\mathbf{a}\bar{\mathbf{x}} = \bar{x}_k, \quad \mathbf{a}\boldsymbol{\mu} = \mu_k, \quad \mathbf{a}'\mathbf{S}\mathbf{a} = s_{kk} = s_k^2,$$

we obtain a $(1 - \alpha)100\%$ confidence interval for μ_k , the populations mean parameter for the k th variable. Because (9) holds for any and all $\mathbf{a} \in \mathbb{R}^p$, we can obtain such confidence interval for $\mu_k, k = 1, \dots, p$ simultaneously,

$$\bar{x}_k - \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)} \frac{s_{kk}}{n} \leq \mu_k \leq \bar{x}_k + \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)} \frac{s_{kk}}{n}, \quad \text{for all } k = 1, \dots, p. \quad (10)$$

The intervals in (10) are called the simultaneous T^2 confidence intervals for the means. The simultaneous confidence intervals form a hyper-rectangle in \mathbb{R}^p .

Remarks on the T^2 confidence intervals for the means

- In fact, any confidence interval on a linear combination of the μ_k 's can be obtained "for free" from the general statement in (9). For example, we can obtain confidence intervals for differences of the means of two component variables $\mu_k - \mu_i$.
- The simultaneous confidence intervals in (10) derived using T^2 is a hypercube or hyper-rectangle in \mathbb{R}^p , compared with the T^2 confidence region (4) which is an ellipsoid in \mathbb{R}^p . Both are useful.
- Because the statement in (9) is so general, the result for any specific $\mathbf{a} \in \mathbb{R}^p$ is on the conservative side, ending up with wider confidence intervals, sometime too wide to be practically informative.

6.5 Summary of simultaneous confidence intervals of component mean

We have discussed four types of hyper-rectangle confidence regions for component means. The hyper-rectangle is formed by p simultaneous confidence intervals for component mean parameter $\mu_k, k = 1, \dots, p$, in the form

$$\bar{x}_k \pm (\text{"Multiplier"}) \sqrt{\frac{s_{kk}}{n}}, \quad k = 1, \dots, p.$$

Different methods are based on different probability distributions, thus use different multipliers. The multipliers are multiple of critical values of various probability distributions.

- Marginal confidence intervals using t statistics, ignoring dependence among component variables.

$$\text{Multiplier} = t_{n-1, \frac{\alpha}{2}}$$

- Bonferroni simultaneous confidence intervals using t statistics

$$\text{Multiplier} = t_{n-1, \frac{\alpha}{2p}}$$

- Asymptotic simultaneous confidence intervals applies χ^2 statistic,

$$\text{Multiplier} = \sqrt{\chi_p^2(\alpha)}$$

- Hotelling's T^2 induced simultaneous confidence intervals use F statistic,

$$\text{Multiplier} = \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)}$$

These simultaneous confidence intervals are useful for small to moderate p . When p is in the hundreds or higher, these hyper-rectangles simultaneous confidence intervals are likely too large to be useful. Other criteria are sought and developed, under the subject of multiple testing (not covered in this course).

Note: Related chapter in Johnson and Wichern: Chapter 5.