

STT 843: Multivariate Analysis

4. Inferences about a Mean Vector (Chapter 5.5-5.8)

Guanqun Cao

Department of Statistics and Probability
Michigan State University

Spring 2026

Outline

- 1 Large Sample Inference about a Population Mean Vector
- 2 Paired Comparisons
- 3 Comparing Mean Vectors from Two Populations
- 4 The Two-Sample Situation When $\Sigma_1 \neq \Sigma_2$
- 5 Control regions for future individual observations

Large Sample Inference about a Population Mean Vector

Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be a random sample from a population with mean $\boldsymbol{\mu}$ and positive definite covariance matrix $\boldsymbol{\Sigma}$. When $n - p$ is large, the hypothesis $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ is rejected in favor of $H_1 : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$, at a level of significance approximately α , if the observed

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0) > \chi_p^2(\alpha)$$

Here $\chi_p^2(\alpha)$ is the upper (100α) th percentile of a chi-square distribution with p.d.f.

Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be a random sample from a population with mean $\boldsymbol{\mu}$ and positive definite covariance matrix $\boldsymbol{\Sigma}$. If $n - p$ is large,

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

will contain $\mathbf{a}^T \boldsymbol{\mu}$, for every \mathbf{a} , with probability approximately $1 - \alpha$.

Example 4.6 (Constructing large sample simultaneous confidence intervals)

A music educator tested thousands of Finnish students on their native musical ability in order to set national norms in Finland. Summary statistics for part of the data set are given in Table 5.5. These statistics are based on a sample of $n = 96$ Finnish 12 th graders. Construct 90% simultaneous confidence intervals for individual mean components $\mu_i, i = 1, 2, \dots, 7$.

Paired Comparisons

In the single response (univariate) case, let X_{j1} denote the response to treatment 1, and let X_{j2} denote the response to treatment 2 for the j th trial. That is, (X_{j1}, X_{j2}) are measurements recorded on the j th unit or j th pair of like units. By design, the n differences

$$D_j = X_{j1} - X_{j2}, \quad j = 1, 2, \dots, n$$

Should reflect only the differences D_j represent independent observations from an $N(\delta, \sigma_d^2)$ distribution. the variable

$$t = \frac{\bar{D} - \delta}{s_d / \sqrt{n}}$$

where $\bar{D} = \frac{1}{n} \sum_{j=1}^n D_j$ and $s_d^2 = \frac{1}{n-1} \sum_{j=1}^n (D_j - \bar{D})^2$ has a t-distribution with $n - 1$ d.f.

- An α -level test of

$$H_0 : \delta = 0 \quad \text{vs} \quad H_1 : \delta \neq 0$$

may be conducted by comparing $|t|$ with $t_{n-1}(\alpha/2)$ -the upper $100(\alpha/2)$ th percentile of a t-distribution with $n - 1$ d.f.

- A $100(1 - \alpha)\%$ confidence interval for the mean difference $\delta = E(X_{j1} - X_{j2})$ is provided the statement

$$\bar{D} - t_{n-1}(\alpha/2) \frac{s_d}{\sqrt{n}} \leq \delta \leq \bar{D} + t_{n-1}(\alpha/2) \frac{s_d}{\sqrt{n}}$$

- Multivariate extension of the paired-comparison procedure to distinguish between p response, two treatments, and n experimental units. The p paired difference random variables become

$$D_{j1} = X_{1j1} - X_{2j1}$$

$$D_{j2} = X_{1j2} - X_{2j2}$$

$$\vdots$$

$$D_{jp} = X_{1jp} - X_{2jp}$$

Let $\mathbf{D}_j^\top = [D_{j1}, D_{j2}, \dots, D_{jp}]$, and assume for $j = 1, 2, \dots, n$ that

$$E(\mathbf{D}_j) = \boldsymbol{\delta} = [\delta_1, \delta_2, \dots, \delta_p]^\top \quad \text{and} \quad \text{Cov}(\mathbf{D}_j) = \Sigma_d$$

If, in addition, $\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_n$ are independent $N_p(\boldsymbol{\delta}, \Sigma_d)$ random vectors, inference about the vector of mean differences $\boldsymbol{\delta}$ can be based upon a T^2 statistics

$$T^2 = n(\bar{\mathbf{D}} - \boldsymbol{\delta})^\top \mathbf{S}_d^{-1} (\bar{\mathbf{D}} - \boldsymbol{\delta})$$

where $\bar{\mathbf{D}} = \frac{1}{n} \sum_{j=1}^n \mathbf{D}_j$ and $\mathbf{S}_d = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{D}_j - \bar{\mathbf{D}}) (\mathbf{D}_j - \bar{\mathbf{D}})^\top$.

Let the differences D_1, D_2, \dots, D_n be a random sample from $N_p(\delta, \Sigma_d)$ population. Then

$$T^2 = n(\bar{D} - \delta)^T \mathbf{S}_d^{-1} (\bar{D} - \delta)$$

is distributed as an $[(n-p)p/(n-p)F_{p,n-p}]$ random variable, whatever the true δ and Σ_d .

Given the observed differences

$$\mathbf{D}_j^\top = [D_{j1}, D_{j2}, \dots, D_{jp}], j = 1, 2, \dots, n,$$

- an α -level test of

$$H_0 : \delta = 0 \quad \text{versus} \quad H_1 : \delta = 0$$

for an $N_p(\delta, \Sigma_d)$ population rejects H_0 if the observed

$$T^2 = n \overline{\mathbf{D}}^\top \mathbf{S}_d^{-1} \overline{\mathbf{D}} > \frac{(n-p)p}{n-p} F_{p, n-p}(\alpha)$$

where $F_{p, n-p}(\alpha)$ is the upper $(100\alpha)\%$ th percentile of an F-distribution with p and $n-p$ d.f.

- A $100(1 - \alpha)\%$ confidence region for δ consists of all δ such that

$$(\bar{\mathbf{D}} - \delta)^T \mathbf{S}_d^{-1} (\bar{\mathbf{D}} - \delta) > \frac{(n - p)p}{n(n - p)} F_{p, n-p}$$

- Also $100(1 - \alpha)\%$ simultaneous confidence intervals for the individual mean differences δ_i are given by

$$\delta_i : \bar{D}_i \pm \sqrt{\frac{(n - 1)p}{(n - p)} F_{p, n-p}(\alpha)} \sqrt{\frac{s_{D_i}^2}{n}}$$

where \bar{D}_i is the i th element of $\bar{\mathbf{D}}$ and $s_{D_i}^2$ is the i th diagonal element of \mathbf{S}_d^{22} .

Example 4.7 (Checking for a mean difference with paired observations) Municipal wastewater treatment plants are required by law to monitor their discharges into rivers and streams on a regular basis. Concern about the reliability of data from one of these self-monitoring programs led to a study in which samples of effluent were divided and sent to two laboratories for testing. One-half was sent to a private commercial laboratory routinely used in the monitoring program. Measurements of biochemical oxygen demand (BOD) and suspended solid (SS) were obtained, for $n = 11$ sample splits, from the two laboratories.

Do the two laboratories's chemical analyses agree? If differences exist, what is their nature?

Alternatively, think of each observation

$$\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_{1i} \\ \mathbf{x}_{2i} \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{pre-tests} \\ \leftarrow \text{post-tests} \end{array}$$

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix}$$

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix}$$

Interest is in $\mathbf{C}\mathbf{x}_i$, where

$$\mathbf{C}_{p \times 2p} = \begin{bmatrix} 1 & & 0 & -1 & & 0 \\ & 1 & & & -1 & \\ & & \ddots & & & \ddots \\ 0 & & & 1 & 0 & -1 \end{bmatrix}$$

Note

$$\mathbf{d}_i = \mathbf{C}\mathbf{x}_i$$

$$\bar{\mathbf{d}} = \mathbf{C}\bar{\mathbf{x}}$$

$$\mathbf{S}_d = \mathbf{C}\mathbf{S}\mathbf{C}'$$

$$T^2 = n\bar{\mathbf{x}}'\mathbf{C}'(\mathbf{C}\mathbf{S}\mathbf{C}')^{-1}\mathbf{C}\bar{\mathbf{x}} \sim T_{p,n-1}^2$$

and

$$\begin{aligned} &= \frac{(n-1)p}{(n-1-p+1)} F_{p,n-1-p+1} \\ &\stackrel{q}{=} \frac{(n-1)p}{n-p} F_{p,n-p} \end{aligned}$$

An extension to a comparison of p treatments given to each subject over time

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{evaluation after day 2 dosage} \\ \leftarrow \text{evaluation after day } j \text{ dosage} \end{array}$$

Interest may lie in comparisons of treatment means

$$\mathbf{C}_{(p-1) \times p} \boldsymbol{\mu} = \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & -1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & & -1 & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} = \begin{bmatrix} \mu_2 - \mu_1 \\ \mu_3 - \mu_2 \\ \vdots \\ \mu_p - \mu_{p-1} \end{bmatrix}$$

$$\begin{aligned}
 T^2 &= n(\mathbf{C}\bar{\mathbf{x}})'(\mathbf{CSC}')^{-1}\mathbf{C}\bar{\mathbf{x}} \sim T_{p-1, n-1}^2 \\
 &\stackrel{q}{=} \frac{(n-1)(p-1)}{(n-1-(p-1)+1)} F_{(p-1), n-1-(p-1)+1} \\
 &\stackrel{q}{=} \frac{(n-1)(p-1)}{n-p+1} F_{p-1, n-p+1}
 \end{aligned}$$

e.g., if comparing 3 days, we might use

$$\underset{2 \times 3}{\mathbf{C}} = \begin{bmatrix} -1 & 0 & 1 \\ 1 & -2 & 1 \end{bmatrix} \begin{array}{l} \leftarrow \text{linear} \\ \leftarrow \text{quadratic} \end{array}$$

e.g., if comparing 4 days, we might use

$$\mathbf{C}_{3 \times 4} = \begin{bmatrix} -3 & -1 & 1 & 3 \\ 1 & -1 & -1 & 1 \\ 1 & -3 & 3 & -1 \end{bmatrix} \begin{array}{l} \leftarrow \text{linear} \\ \leftarrow \text{quadratic} \\ \leftarrow \text{cubic} \end{array}$$

Sample	Summary statistics	
(Population 1)	$\mathbf{x}_{11}, \mathbf{x}_{12}, \dots, \mathbf{x}_{1n_1}$	
(Population 2)	$\mathbf{x}_{21}, \mathbf{x}_{22}, \dots, \mathbf{x}_{2n_2}$	
	$\bar{\mathbf{x}}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbf{x}_{1j}$	$\mathbf{S}_1 = \frac{1}{n_1-1} \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1) (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T$
	$\bar{\mathbf{x}}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} \mathbf{x}_{2j}$	$\mathbf{S}_2 = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2) (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T$

Two-sample comparison

- 1 The sample $\mathbf{X}_{11}, \mathbf{X}_{12}, \dots, \mathbf{X}_{1n_1}$ is a random sample of size n_1 from a p -variate population with mean $\boldsymbol{\mu}_1$ and covariance matrix $\boldsymbol{\Sigma}_1$.
- 2 The sample $\mathbf{X}_{21}, \mathbf{X}_{22}, \dots, \mathbf{X}_{2n_2}$ is a random sample of size n_2 from a p -variate population with mean $\boldsymbol{\mu}_2$ and covariance matrix $\boldsymbol{\Sigma}_2$.
- 3 Also, $\mathbf{X}_{11}, \mathbf{X}_{12}, \dots, \mathbf{X}_{1n_1}$ are independent of $\mathbf{X}_{21}, \mathbf{X}_{22}, \dots, \mathbf{X}_{2n_2}$.

$$H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{0} \quad \text{vs} \quad H_1 : \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 \neq \mathbf{0}$$

Further Assumption When n_1 and n_2 Are Small

- ① Both populations are multivariate normal.
- ② Also, $\Sigma_1 = \Sigma_2$ (same covariance matrix)

Set the estimate of Σ as

$$\begin{aligned} \mathbf{S}_{\text{pooled}} &= \frac{\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^T + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^T}{n_1 + n_2 - 2} \\ &= \frac{n_1 - 1}{n_1 + n_2 - 2} \mathbf{S}_1 + \frac{n_2 - 1}{n_1 + n_2 - 2} \mathbf{S}_2 \end{aligned}$$

$$\text{Cov}(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2) = \text{Cov}(\mathbf{X}_1) + \text{Cov}(\mathbf{X}_2) = \frac{1}{n_1}\Sigma + \frac{1}{n_2}\Sigma = \left(\frac{1}{n_1} + \frac{1}{n_2}\right)\Sigma$$

Hence

$$\left(\frac{1}{n_1} + \frac{1}{n_2}\right)\Sigma$$

is an estimator of $\text{Cov}(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)$.

If $\mathbf{X}_{11}, \mathbf{X}_{12}, \dots, \mathbf{X}_{1n_1}$ is a random sample of size n_1 from $N_p(\boldsymbol{\mu}_1, \Sigma)$ and $\mathbf{X}_{21}, \mathbf{X}_{22}, \dots, \mathbf{X}_{2n_2}$ is an independent random sample size n_2 from $N_p(\boldsymbol{\mu}_2, \Sigma)$, then

$$T^2 =$$

$$[\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)]^T \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{\text{pooled}} \right]^{-1} [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)]$$

is distributed as

$$\frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}$$

Consequently,

$$P(T^2 \leq c^2) = 1 - \alpha$$

where

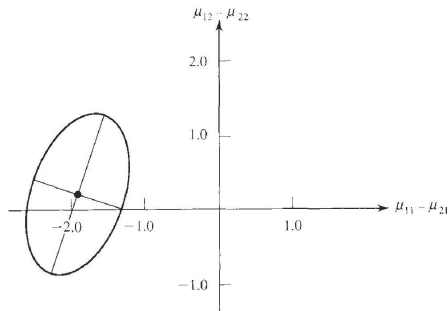
$$c^2 = \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}(\alpha)$$

Example 4.8 (Constructing a confidence region for the difference of two mean vectors)

Fifty bars of soap are manufactured in each of two ways. Two characteristics, $X_1 = \text{lather}$ and $X_2 = \text{mildness}$, are measured. The summary statistics for bars produced by method 1 and 2 are

$$\mathbf{x}_1 = \begin{bmatrix} 8.3 & 4.1 \end{bmatrix}^T, \mathbf{x}_2 = \begin{bmatrix} 10.2 & 3.9 \end{bmatrix}^T$$
$$\mathbf{S}_1 = \begin{bmatrix} 2 & 1 \\ 1 & 6 \end{bmatrix}, \quad \mathbf{S}_2 = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$$

Obtain a 95% confidence region for $\mu_1 - \mu_2$.



Simultaneous Confidence Intervals

Let $c^2 = [(n_1 + n_2 - 2)p / (n_1 + n_2 - p - 1)] F_{p, n_1 + n_2 - p - 1}(\alpha)$. With probability $1 - \alpha$.

$$\mathbf{a}^T (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \pm c \sqrt{\mathbf{a}^T \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{\text{pooled}} \mathbf{a}}$$

will cover $\mathbf{a}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ for all \mathbf{a} . In particular $\mu_{1i} - \mu_{2i}$ will be covered by

$$(\bar{X}_{1i} - \bar{X}_{2i}) \pm c \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2} \right) s_{ii, \text{pooled}}} \quad \text{for } i = 1, 2, \dots, p$$

The Two-Sample Situation When $\Sigma_1 \neq \Sigma_2$

Let the sample sizes be such that $n_1 - p$ and $n_2 - p$ are large. Then, an approximate $100(1 - \alpha)\%$ confidence ellipsoid for $\mu_1 - \mu_2$ is given by all $\mu_1 - \mu_2$ satisfying

$$T^2 = [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2)]^T \left(\frac{1}{n_1} \mathbf{S}_1 + \frac{1}{n_2} \mathbf{S}_2 \right)^{-1} [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2)] \leq \chi_p^2(\alpha)$$

where $\chi_p^2(\alpha)$ is the upper (100α) th percentile of a chi-square distribution with p d.f. Also $100(1 - \alpha)\%$ simultaneous confidence interval for all linear combinations $\mathbf{a}^T (\mu_1 - \mu_2)$ are provided by

$$\mathbf{a}^T (\mu_1 - \mu_2)$$

belongs to $\mathbf{a}^T (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\mathbf{a}^T \left(\frac{1}{n_1} \mathbf{S}_1 + \frac{1}{n_2} \mathbf{S}_2 \right) \mathbf{a}}$

Testing for Equality of Covariance Matrices

With g populations, the null hypothesis is

$$H_0 : \Sigma_1 = \Sigma_2 = \cdots = \Sigma_g = \Sigma$$

where Σ_l is the covariance matrix for the l th population, $l = 1, 2, \dots, g$, and Σ is the presumed common covariance matrix. The alternative hypothesis is that at least two of the covariance matrices are not equal.

- Assuming multivariate normal populations, a likelihood ratio statistic for testing above is given by

$$\Gamma = \prod_l \left(\frac{|\mathbf{S}_l|}{|\mathbf{S}_{\text{pooled}}|} \right)^{(n_l-1)/2}$$

Here n_l is the sample size for the l th group, \mathbf{S}_l is the l th group sample covariance matrix and \mathbf{S}_{pool} is the pooled sample covariance matrix given by

$$\mathbf{S}_{pool} = \frac{1}{\sum_l (n_l - 1)} \{ (n_1 - 1) \mathbf{S}_1 + (n_2 - 1) \mathbf{S}_2 + \cdots + (n_g - 1) \mathbf{S}_g \}$$

- Box's test is based on this χ^2 approximation to the sampling distribution of $-2 \ln \Gamma$. Setting $-2 \ln \Gamma = M$ (Box's M statistics) gives

$$M = \left[\sum_l (n_l - 1) \right] \ln |\mathbf{S}_{pooled}| - \sum_l [(n_l - 1) \ln |\mathbf{S}_l|]$$

Set

$$u = \left[\sum_I \frac{1}{(n_I - 1)} - \frac{1}{\sum_I (n_I - 1)} \right] \left[\frac{2p^2 + 3p - 1}{6(p + 1)(g - 1)} \right]$$

where p is the number of variables and g is the number of groups.

Then $C = (1 - u)M = (1 - u) \{ [\sum_I (n_I - 1)] \ln |\mathbf{S}_{\text{pooled}}| - \sum_I [(n_I - 1) \ln |\mathbf{S}_I|] \}$ has an approximate χ^2 distribution with

$$\nu = g \frac{1}{2} p(p + 1) - \frac{1}{2} p(p + 1) = \frac{1}{2} p(p + 1)(g - 1)$$

degrees of freedom. At significance level α , reject H_0 if

$$C > \chi_{p(p+1)(g-1)/2}^2(\alpha).$$

Example 4.9 (Testing equality of covariance matrices-nursing homes)

The Wisconsin Department of Health and Social Services reimburse nursing homes in the state for the services provided. The department develops a set of formulas for the rates for each facility, based on factors such as level of care, mean wage rate, and average wage rate in the state.

Nursing homes can be classified on the basis of ownership (private party, nonprofit organization, and government) and certification (skilled nursing facility, intermediate care facility, or combination of the two).

One purpose of a recent study was to investigate the effects of ownership or certification (or both) on cost.s. Four costs, computed on a per-patient-day basis and measured in hours per patient day, were selected for analysis X_1 = cost of nursing labor, X_2 = cost dietary labor, X_3 = cost of plant operation and maintenance labor, and X_4 = cost of housekeeping and laundry labor. A total of $n = 516$ observations on each of the $p = 4$ cost variables were initially separated according to ownership.

Summary statistics for each of the $g = 3$ groups are given in the following table.

Sample covariance matrices

$$\mathbf{S}_1 = \begin{bmatrix} .291 & & & & \\ -0.001 & .011 & & & \\ .002 & .000 & .001 & & \\ .010 & .003 & .000 & .010 & \\ .561 & & & & \\ -0.011 & .025 & & & \\ .001 & .004 & .005 & & \\ .037 & .007 & .002 & .019 & \end{bmatrix}; \mathbf{S}_3 = \begin{bmatrix} .261 & & & & \\ -0.030 & .017 & & & \\ .003 & -.000 & .004 & & \\ .018 & .006 & .001 & .013 & \end{bmatrix}$$

Assuming multivariate normal data, test hypothesis

$$H_0 : \Sigma_1 = \Sigma_2 = \Sigma_3 = \Sigma.$$

Variable	Raw score	
	Mean (\bar{x}_i)	Standard deviation ($\sqrt{s_{ii}}$)
$X_1 = \text{melody}$	28.1	5.76
$X_2 = \text{harmony}$	26.6	5.85
$X_3 = \text{tempo}$	35.4	3.82
$X_4 = \text{meter}$	34.2	5.12
$X_5 = \text{phrasing}$	23.6	3.76
$X_6 = \text{balance}$	22.0	3.93
$X_7 = \text{style}$	22.7	4.03

Table 6.1 Effluent Data

Sample j	Commercial lab		State lab of hygiene	
	x_{1j1} (BOD)	x_{1j2} (SS)	x_{2j1} (BOD)	x_{2j2} (SS)
	6	27	25	15
2	6	23	28	13
3	18	64	36	22
4	8	44	35	29
5	11	30	15	31
6	34	75	44	64
7	28	26	42	30
8	71	124	54	64
9	43	54	34	56
10	33	30	29	20
11	20	14	39	21

Control regions for future individual observations

Let $\mathbf{X}_i \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, i.i.d., $i = 1, \dots, n$. Let \mathbf{X} be a future observation from the same distribution. Then

$$T^2 = \frac{n}{n+1} (\mathbf{X} - \bar{\mathbf{X}})^\top \mathbf{S}^{-1} (\mathbf{X} - \bar{\mathbf{X}}) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

and a $100(1 - \alpha)\%$ p -dim prediction ellipsoid is given by all \mathbf{x} satisfying

$$\frac{n}{n+1} (\mathbf{x} - \bar{\mathbf{X}})^\top \mathbf{S}^{-1} (\mathbf{x} - \bar{\mathbf{X}}) \leq \frac{(n^2 - 1)}{n(n-p)} F_{p, n-p}(\alpha)$$