

Lecture 4

Inference and Comparison of Mean Vectors

Readings: Johnson & Wichern 2007, Chapter 5.1-5.4; 6.1-6.4; 6.8

DSA 8070 Multivariate Analysis

Whitney Huang
Clemson University

Inference and Comparison of Mean Vectors
CLEMSON UNIVERSITY
Confidence Intervals/Region for Population Means
Hypothesis Testing for Mean Vector
Multivariate Paired Hotelling's T-Square
Comparisons of Two Mean Vectors
Multivariate Analysis of Variance
4.1

Notes

Agenda

- 1 Confidence Intervals/Region for Population Means
- 2 Hypothesis Testing for Mean Vector
- 3 Multivariate Paired Hotelling's T-Square
- 4 Comparisons of Two Mean Vectors
- 5 Multivariate Analysis of Variance

Inference and Comparison of Mean Vectors
CLEMSON UNIVERSITY
Confidence Intervals/Region for Population Means
Hypothesis Testing for Mean Vector
Multivariate Paired Hotelling's T-Square
Comparisons of Two Mean Vectors
Multivariate Analysis of Variance
4.2

Notes

Inference on Mean Vectors

This Week's Topics:

- **Single Mean Vector:** inference on μ (multivariate one-sample t -test)
- **Paired Mean Vectors:** differences between paired observations \Rightarrow reduce to one-sample Hotelling's T^2 on differences
- **Two Independent Mean Vectors:** Hotelling's T^2 two-sample test
- **Several Mean Vectors:** MANOVA (multivariate extension of ANOVA)

Analogy with Univariate Methods:

- One-sample t -test \rightarrow single μ
- Paired t -test \rightarrow paired mean vectors
- Two-sample t -test \rightarrow two mean vectors
- ANOVA \rightarrow MANOVA

Inference and Comparison of Mean Vectors
CLEMSON UNIVERSITY
Confidence Intervals/Region for Population Means
Hypothesis Testing for Mean Vector
Multivariate Paired Hotelling's T-Square
Comparisons of Two Mean Vectors
Multivariate Analysis of Variance
4.3

Notes

Review: Sampling Distribution of Univariate Sample Mean \bar{X}_n

Suppose X_1, X_2, \dots, X_n is a random sample from a univariate population distribution with mean $\mathbb{E}(X) = \mu$ and variance $\text{Var}(X) = \sigma^2$. The sample mean \bar{X}_n is a function of random sample and therefore has a distribution

- $\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$ when the sample size n is “sufficiently” large \Rightarrow This is the central limit theorem (CLT)
- The result above is exact if the population follows a normal distribution, i.e., $X \sim N(\mu, \sigma^2)$
- The standard error $\sqrt{\text{Var}(\bar{X}_n)} = \frac{\sigma}{\sqrt{n}}$ provides a measure estimation precision. In practice, we use $\frac{s}{\sqrt{n}}$ instead where s is the sample standard deviation

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.4

Notes

Sampling Distribution of Multivariate Sample Mean Vector \bar{X}_n

Suppose X_1, X_2, \dots, X_n is a random sample from a multivariate population distribution with mean vector $\mathbb{E}(X) = \mu$ and covariance matrix Σ .

- $\bar{X}_n \sim N(\mu, \frac{1}{n}\Sigma)$ when the sample size n is “sufficiently” large \Rightarrow This is the multivariate version of CLT
- The result above is exact if the population follows a normal distribution, i.e., $X \sim N(\mu, \Sigma)$
- Again, the estimation precision improves with a larger sample size. Like the univariate case we would need to replace Σ by its estimate S , the sample covariacne matrix

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.5

Notes

Review: Interval Estimation of Univariate Population Mean μ

The general format of a confidence interval (CI) estimate of a population mean is

Sample mean \pm multiplier \times standard error of mean.

For variable X , a CI estimate of its population mean μ is

$$\bar{X}_n \pm t_{n-1, \frac{\alpha}{2}} \frac{s}{\sqrt{n}},$$

Here the multiplier value is a function of the confidence level, α , the sample size n

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.6

Notes

Constructing Confidence Intervals for Mean Vector

We will still use the general recipe

Sample mean ± multiplier × standard error of mean.

The multiplier value also depends the strategy used for dealing with the multiple inference issue

- One at a Time CIs: a CI for μ_j is computed as

$$\bar{x}_j \pm t_{n-1, \frac{\alpha}{2}} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p$$

- Bonferroni Method: a CI for μ_j is computed as

$$\bar{x}_j \pm t_{n-1, \frac{\alpha}{2p}} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p$$

- Simultaneous CIs: a CI for μ_j is computed as

$$\bar{x}_j \pm \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p, \alpha}} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p$$

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.7

Notes

Example: Mineral Content Measurements [source: Penn Stat Univ. STAT 505]

This example uses a dataset that includes mineral content measurements at two different arm bone locations for $n = 64$ women. We will determine confidence intervals for the two population means. The sample means and standard deviations for the two variables are:

Variable	Sample size	Mean	Std Dev
domradius (X_1)	$n = 64$	$\bar{x}_1 = 0.8438$	$s_1 = 0.1140$
domhumerus (X_2)	$n = 64$	$\bar{x}_2 = 1.7927$	$s_2 = 0.2835$

Let's apply the three methods we learned to construct 95% CIs

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.8

Notes

Mineral Content Measurements Example Cont'd

- One at a Time CIs: $\bar{x}_j \pm t_{n-1, \alpha/2} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p.$
Therefore 95% CIs for μ_1 and μ_2 are:

$$\mu_1 : 0.8438 \pm \underbrace{1.998}_{t_{63, 0.025}} \times \frac{0.1140}{\sqrt{64}} = [0.815, 0.872]$$

$$\mu_2 : 1.7927 \pm 1.998 \times \frac{0.2835}{\sqrt{64}} = [1.722, 1.864]$$

- Bonferroni Method: $\bar{x}_j \pm t_{n-1, \alpha/2p} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p.$

$$\mu_1 : 0.8438 \pm \underbrace{2.296}_{t_{63, 0.0125}} \times \frac{0.1140}{\sqrt{64}} = [0.811, 0.877]$$

$$\mu_2 : 1.7927 \pm 2.296 \times \frac{0.2835}{\sqrt{64}} = [1.711, 1.874]$$

- Simultaneous CIs:

$$\bar{x}_j \pm \sqrt{\frac{(n-1)p}{n-p} F_{p, n-p, \alpha}} \frac{s_j}{\sqrt{n}}, \quad j = 1, \dots, p$$

$$\mu_1 : 0.8438 \pm 2.528 \times \frac{0.1140}{\sqrt{64}} = [0.808, 0.880]$$

$$\mu_2 : 1.7927 \pm 2.528 \times \frac{0.2835}{\sqrt{64}} = [1.703, 1.882]$$

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

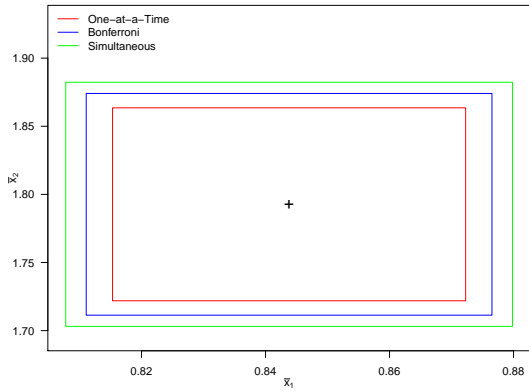
Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.9

Notes

95 % CIs Based on Three Methods

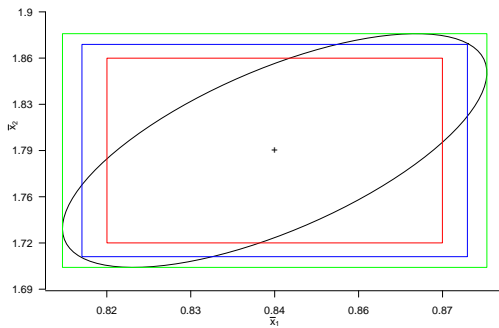


Notes

Confidence Ellipsoid

A confidence ellipsoid for μ is the set of μ satisfying

$$n(\bar{X}_n - \mu)^T S^{-1} (\bar{X}_n - \mu) \leq \frac{(n-1)p}{n-p} F_{p, n-p, \alpha}$$



Notes

Hypothesis Testing for Mean

- Recall: for univariate data, t statistic

$$t = \frac{\bar{X}_n - \mu_0}{s/\sqrt{n}} \Rightarrow t^2 = \frac{(\bar{X}_n - \mu_0)^2}{s^2/n} = n(\bar{X}_n - \mu_0)(s^2)^{-1}(\bar{X}_n - \mu_0)$$

Under $H_0 : \mu = \mu_0$

$$t \sim t_{n-1}, \quad t^2 \sim F_{1, n-1}$$

- Extending to multivariate by analogy:

$$T^2 = n(\bar{X}_n - \mu_0)^T S^{-1} (\bar{X}_n - \mu_0)$$

Under $H_0 : \mu = \mu_0$

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

Note: T^2 here is the so-called **Hotelling's T-Square**

Notes

Hypothesis Testing for Mean Vector μ

- 1 State the null
$$H_0 : \mu = \mu_0$$

and the alternative
$$H_a : \mu \neq \mu_0$$
- 2 Compute the test statistic
$$F = \frac{n-p}{(n-1)p} n (\bar{X}_n - \mu_0)^T S^{-1} (\bar{X}_n - \mu_0)$$
- 3 Compute the P-value. Under $H_0 : F \sim F_{p,n-p}$
- 4 Draw a conclusion: We do (or do not) have enough statistical evidence to conclude $\mu \neq \mu_0$ at α significant level

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.13

Notes

Example: Women's Dietary Intake [source: Penn Stat Univ. STAT 505]

The recommended intake and a sample mean for all women between 25 and 50 years old are given below:

Variable	Recommended Intake (μ_0)	Sample Mean (\bar{x}_n)
Calcium	1000 mg	624.0 mg
Iron	15 mg	11.1 mg
Protein	60 g	65.8 g
Vitamin A	800 μ g	839.6 μ g
Vitamin C	75 mg	78.9 mg

Here we would like to test, at $\alpha = 0.01$ level, if the $\mu = \mu_0$

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.14

Notes

Women's Dietary Intake Example Analysis

- 1 State the null
$$H_0 : \mu = \mu_0$$

and the alternative
$$H_a : \mu \neq \mu_0$$
- 2 Compute the test statistic
$$F = \frac{n-p}{(n-1)p} n (\bar{x}_n - \mu_0)^T S^{-1} (\bar{x}_n - \mu_0) = 349.80$$
- 3 Compute the P-value. Under $H_0 : F \sim F_{p,n-p} \Rightarrow$
p-value = $\mathbb{P}_F(F_{p,n-p} > 349.80) = 3 \times 10^{-191} < \alpha = 0.01$
- 4 Draw a conclusion: We do have enough statistical evidence to conclude $\mu \neq \mu_0$ at $\alpha = 0.01$ significant level

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

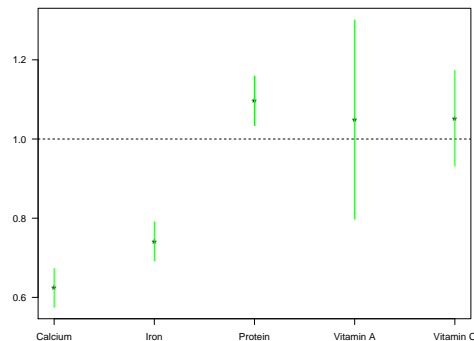
Multivariate Analysis of Variance

4.15

Notes

Profile Plots

- 1 Standardize each of the observations by dividing their hypothesized means
- 2 Plot either simultaneous or Bonferroni CIs for the population mean of these standardized variables



Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.16

Notes

Spouse Survey Data Example

A sample ($n = 30$) of husband and wife pairs are asked to respond to each of the following questions:

- 1 What is the level of passionate love you feel for your partner?
- 2 What is the level of passionate love your partner feels for you?
- 3 What is the level of companionate love you feel for your partner?
- 4 What is the level of companionate love your partner feels for you?

Responses were recorded on a typical five-point scale: 1) None at all 2) Very little 3) Some 4) A great deal 5) Tremendous amount.

We will try to address the following question: Do the husbands respond to the questions in the same way as their wives?

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.17

Notes

Multivariate Paired Hotelling's T-Square

Let X_F and X_M be the responses to these 4 questions for females and males, respectively. Here the quantities of interest are $\mathbb{E}(D) = \mu_D$, the average differences across all husband and wife pairs.

- 1 State the null $H_0 : \mu_D = 0$ and the alternative hypotheses $H_a : \mu_D \neq 0$
- 2 Compute the test statistic

$$F = \frac{n-p}{(n-1)p} n \bar{D}_n^T S_D^{-1} \bar{D}_n$$

- 3 Compute the P-value. Under $H_0 : F \sim F_{p,n-p}$
- 4 Draw a conclusion: We do (or do not) have enough statistical evidence to conclude $\mu_D \neq 0$ at α significant level

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.18

Notes

Spouse Survey Data Example Analysis

- 1 State the null

$$H_0 : \mu_D = 0$$

and the alternative

$$H_a : \mu_D \neq 0$$

- 2 Compute the test statistic

$$F = \frac{n-p}{(n-1)p} n \bar{D}_n^T S_D^{-1} \bar{D}_n = 2.942$$

- 3 **Compute the P-value.** Under $H_0 : F \sim F_{p,n-p} \Rightarrow$
p-value = $\Pr(F_{p,n-p} >) = 0.0394 < \alpha = 0.05$

- 4 **Draw a conclusion:** We do have enough statistical evidence to conclude $\mu_D \neq 0$ at 0.05 significant level



Notes

Motivating Example: Swiss Bank Notes (Source: PSU stat 505)

Suppose there are two distinct populations for 1000 franc Swiss Bank Notes:

- 1 The first population is the population of Genuine Bank Notes
- 2 The second population is the population of Counterfeit Bank Notes

For both populations the following measurements were taken:

- 1 Length of the note
- 2 Width of the Left-Hand side of the note
- 3 Width of the Right-Hand side of the note
- 4 Width of the Bottom Margin
- 5 Width of the Top Margin
- 6 Diagonal Length of Printed Area

We want to determine if counterfeit notes can be distinguished from the genuine Swiss bank notes



Notes

Review: Two Sample t-Test

Suppose we have data from a single variable from population 1: $X_{11}, X_{12}, \dots, X_{1n_1}$ and population 2: $X_{21}, X_{22}, \dots, X_{2n_2}$. Here we would like to draw inference about their population means μ_1 and μ_2 .

Assumptions:

- 1 **Homoscedasticity:** The data from both populations have common variance σ^2
- 2 **Independence:** The subjects from both populations are independently sampled $\Rightarrow \{X_{1i}\}_{i=1}^{n_1}$ and $\{X_{2j}\}_{j=1}^{n_2}$ are independent to each other
- 3 **Normality:** The data from both populations are normally distributed (not that crucial for "large" sample)

Here we are going to consider testing $H_0 : \mu_1 = \mu_2$ against $H_a : \mu_1 \neq \mu_2$



Notes

Review: Two Sample t-Test

We define the sample means for each population using the following expression:

$$\bar{x}_1 = \frac{\sum_{j=1}^{n_1} x_{1j}}{n_1}, \quad \bar{x}_2 = \frac{\sum_{j=1}^{n_2} x_{2j}}{n_2}.$$

We denote the sample variance

$$s_1^2 = \frac{\sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1)^2}{n_1 - 1}, \quad s_2^2 = \frac{\sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2)^2}{n_2 - 1}.$$

Under the **homoscedasticity** assumption, we can “pool” two samples to get the pooled sample variance

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

Test statistic

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \stackrel{H_0}{\sim} t_{n_1 + n_2 - 2}$$

We can use this result to construct confidence intervals and to perform hypothesis tests



Notes

The Two Sample Problem: The Multivariate Case

Now we would like to use two independent samples $\{X_{11}, \dots, X_{12}, \dots, X_{1n_1}\}$ and $\{X_{21}, \dots, X_{22}, \dots, X_{2n_2}\}$, where

$$X_{ij} = \begin{bmatrix} X_{ij1} \\ X_{ij2} \\ \vdots \\ X_{ijp} \end{bmatrix}$$

to infer the relationship between μ_1 and μ_2 , where

$$\mu_i = \begin{bmatrix} \mu_{i1} \\ \mu_{i2} \\ \vdots \\ \mu_{ip} \end{bmatrix}$$

Assumptions

- Both populations have common covariance matrix, i.e., $\Sigma_1 = \Sigma_2$
- Independence**: The subjects from both populations are independently sampled
- Normality**: Both populations are normally distributed



Notes

The Multivariate Two-Sample Problem

Here we are testing

$$H_0 : \begin{bmatrix} \mu_{11} \\ \mu_{12} \\ \vdots \\ \mu_{1p} \end{bmatrix} = \begin{bmatrix} \mu_{21} \\ \mu_{22} \\ \vdots \\ \mu_{2p} \end{bmatrix}, \quad H_a : \mu_{1k} \neq \mu_{2k} \text{ for at least one } k \in \{1, 2, \dots, p\}$$

Under the **common covariance** assumption we have

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

where

$$S_i = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)^T, \quad i = 1, 2$$



Notes

The Two-Sample Hotelling's T-Square Test Statistic

The two-sample t test is equivalent to

$$t^2 = (\bar{x}_1 - \bar{x}_2)^T \left[s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \right]^{-1} (\bar{x}_1 - \bar{x}_2).$$

Under H_0 , $t^2 \sim F_{1,n_1+n_2-2}$. We can use this result to perform a hypothesis test

We can extend this to the multivariate situation:

$$T^2 = (\bar{x}_1 - \bar{x}_2)^T \left[S_p \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \right]^{-1} (\bar{x}_1 - \bar{x}_2)$$

Under H_0 , we have

$$F = \frac{n_1 + n_2 - p - 1}{p(n_1 + n_2 - 2)} T^2 \sim F_{p,n_1+n_2-p-1}$$

We can use this result to perform inferences for multivariate cases

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.25

Notes

Two-Sample Test for Swiss Bank Notes

```
> (xbar1 <- colMeans(dat[real, -1]))
      V2      V3      V4      V5      V6      V7
214.969 129.943 129.720   8.305  10.168 141.517
> (xbar2 <- colMeans(dat[fake, -1]))
      V2      V3      V4      V5      V6      V7
214.823 130.300 130.193  10.530  11.133 139.450
> Sigma1 <- cov(dat[real, -1])
> Sigma2 <- cov(dat[fake, -1])
> n1 <- length(real); n2 <- length(fake); p <- dim(dat[, -1])[2]
> Sp <- ((n1 - 1) * Sigma1 + (n2 - 1) * Sigma2) / (n1 + n2 - 2)
> # Test statistic
> T.squared <- as.numeric(t(xbar1 - xbar2) %*% solve(Sp * (1 / n1 + 1 / n2)) %*% (xbar1 - xbar2))
> Fobs <- T.squared * ((n1 + n2 - p - 1) / ((n1 + n2 - 2) * p))
> # p-value
> pf(Fobs, p, n1 + n2 - p - 1, lower.tail = F)
[1] 3.378887e-105
```

Conclusion

The counterfeit notes can be distinguished from the genuine notes on at least one of the measurements => which ones?

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.26

Notes

Simultaneous Confidence Intervals

$$\bar{x}_{1k} - \bar{x}_{2k} \pm \sqrt{\frac{p(n_1 + n_2 - 2)}{n_1 + n_2 - p - 1} F_{p,n_1+n_2-p-1,\alpha}} \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2} \right) s_{k,p}^2},$$

where $s_{k,p}^2$ is the pooled variance for the variable k

Variable	95% CI
Length of the note	(−0.04, 0.34)
Width of the Left-Hand note	(−0.52, −0.20)
Width of the Right-Hand note	(−0.64, −0.30)
Width of the Bottom Margin	(−2.70, −1.75)
Width of the Top Margin	(−1.30, −0.63)
Diagonal Length of Printed Area	(1.81, 2.33)

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.27

Notes

Checking Model Assumptions

Assumptions:

- **Homoscedasticity:** The data from both populations have common covariance matrix Σ

Will return to this in next slide

- **Independence:**

This assumption may be violated if we have clustered, time-series, or spatial data

- **Normality:**

Multivariate QQplot, univariate histograms, bivariate scatter plots

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.28

Notes

Testing for Equality of Mean Vectors when $\Sigma_1 \neq \Sigma_2$

- **Bartlett's test** can be used to test if $\Sigma_1 = \Sigma_2$ **but** this test is **sensitive** to departures from normality
- As as crude rule of thumb: if $s_{1,k}^2 > 4s_{2,k}^2$ or $s_{2,k}^2 > 4s_{1,k}^2$ for some $k \in \{1, 2, \dots, p\}$, then it is likely that $\Sigma_1 \neq \Sigma_2$
- Life gets difficult if we cannot assume that $\Sigma_1 = \Sigma_2$ However, if both n_1 and n_2 are "large", we can use the following approximation to conduct inferences:

$$T^2 = (\bar{X}_1 - \bar{X}_2)^T \left[\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right]^{-1} (\bar{X}_1 - \bar{X}_2) \overset{H_0}{\sim} \chi_p^2$$

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.29

Notes

Comparing More Than Two Populations:
Romano-British Pottery Example (source: PSU stat 505)

- Pottery shards are collected from four sites in the British Isles:
 - Llanedyn (L)
 - Caldicot (C)
 - Isle Thorns (I)
 - Ashley Rails (A)
- The concentrations of five different chemicals were be used
 - Aluminum (Al)
 - Iron (Fe)
 - Magnesium (Mg)
 - Calcium (Ca)
 - Sodium (Na)
- **Objective:** to determine whether the chemical content of the pottery depends on the site where the pottery was obtained

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.30

Notes

Review: (Univariate) Analysis of Variance (ANOVA)

- $H_0 : \mu_1 = \mu_2 = \dots = \mu_g$
 H_a : At least one mean is different

Source	df	SS	MS	F statistic
Treatment	$g - 1$	$SSTr$	$MSTr = \frac{SSTr}{g-1}$	$F = \frac{MSTr}{MSE}$
Error	$N - g$	SSE	$MSE = \frac{SSE}{N-g}$	
Total	$N - 1$	$SSTo$		

- Test Statistic: $F^* = \frac{MSTr}{MSE}$. Under H_0 ,
 $F^* \sim F_{df_1=g-1, df_2=N-g}$
- Assumptions:**
 - The distribution of each group is normal with equal variance (i.e. $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_g^2$)
 - Responses for a given group are independent to each other

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.31

Notes

One-way Multivariate Analysis of Variance (One-way MANOVA)

Subject \ Group	1	2	...	g
1	$Y_{11} = \begin{bmatrix} Y_{111} \\ Y_{112} \\ \vdots \\ Y_{11p} \end{bmatrix}$	$Y_{21} = \begin{bmatrix} Y_{211} \\ Y_{212} \\ \vdots \\ Y_{21p} \end{bmatrix}$...	$Y_{g1} = \begin{bmatrix} Y_{g11} \\ Y_{g12} \\ \vdots \\ Y_{g1p} \end{bmatrix}$
2	$Y_{21} = \begin{bmatrix} Y_{121} \\ Y_{122} \\ \vdots \\ Y_{12p} \end{bmatrix}$	$Y_{22} = \begin{bmatrix} Y_{221} \\ Y_{222} \\ \vdots \\ Y_{22p} \end{bmatrix}$...	$Y_{g2} = \begin{bmatrix} Y_{g21} \\ Y_{g22} \\ \vdots \\ Y_{g2p} \end{bmatrix}$
\vdots	\vdots	\vdots	\vdots	\vdots
n_i	$Y_{1n_i} = \begin{bmatrix} Y_{1n_i1} \\ Y_{1n_i2} \\ \vdots \\ Y_{1n_ip} \end{bmatrix}$	$Y_{2n_i} = \begin{bmatrix} Y_{2n_i1} \\ Y_{2n_i2} \\ \vdots \\ Y_{2n_ip} \end{bmatrix}$...	$Y_{gn_i} = \begin{bmatrix} Y_{gn_i1} \\ Y_{gn_i2} \\ \vdots \\ Y_{gn_ip} \end{bmatrix}$

- Notation:** Y_{ij} is the vector of variables for subject j in group i ; n_i is the sample size in group i ;
 $N = n_1 + n_2 + \dots + n_g$ the total sample size
- Assumptions:** 1) common covariance matrix Σ ; 2) Independence; 3) Normality

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.32

Notes

Test Statistics for MANOVA

- We are interested in testing the null hypothesis that the group mean vectors are all equal

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_g.$$

The alternative hypothesis:

$$H_a : \mu_{ik} \neq \mu_{jk} \text{ for at least one } i \neq j \text{ and at least one variable } k$$

- Mean vectors:**
 - Sample Mean Vector:** $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$, $i = 1, \dots, g$
 - Grand Mean Vector:** $\bar{y}_{..} = \frac{1}{N} \sum_{i=1}^g \sum_{j=1}^{n_i} Y_{ij}$

- Total Sum of Squares:**

$$T = \sum_{i=1}^g \sum_{j=1}^{n_i} (Y_{ij} - \bar{y}_{..})(Y_{ij} - \bar{y}_{..})^T$$

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.33

Notes

MANOVA Decomposition and MANOVA Table

$$\begin{aligned} T &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \mathbf{y}_{..})(\mathbf{Y}_{ij} - \bar{\mathbf{y}})^T \\ &= \sum_{i=1}^g \sum_{j=1}^{n_i} [(\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.}) + (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})][(\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.}) + (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})]^T \\ &= \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.})(\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.})^T}_E + \underbrace{\sum_{i=1}^g n_i (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})(\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})^T}_H \end{aligned}$$

MANOVA Table

Source	df	SS
Treatment	$g - 1$	H
Error	$N - g$	E
Total	$N - 1$	T

Reject $H_0 : \mu_1 = \mu_2 = \dots = \mu_g$ if the matrix H is “large” relative to the matrix E

Inference and Comparison of Mean Vectors

CLEMSON
UNIVERSITY

Confidence
Intervals/Region
for Population
Means

Hypothesis Testing
for Mean Vector

Multivariate Paired
Hotelling's
T-Square

Comparisons of
Two Mean Vectors

Multivariate
Analysis of
Variance

4.34

Notes

Test Statistics for MANOVA

There are several different test statistics for conducting the hypothesis test:

- Wilks Lambda

$$\Lambda^* = \frac{|E|}{|H + E|}$$

Reject H_0 if Λ^* is “small”

- Hotelling-Lawley Trace

$$T_0^2 = \text{trace}(\mathbf{H}\mathbf{E}^{-1})$$

Reject H_0 if T_0^2 is “large”

- Pillai Trace

$$V = \text{trace}(\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1})$$

Reject H_0 if V is “large”

Inference and Comparison of Mean Vectors

CLEMSON
UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.35

Notes

[illegible]

Romano-British Pottery Example

```
> dat <- read.table("pottery.txt", header = F)
> out <- manova(cbind(V2, V3, V4, V5, V6) ~ V1, data = dat)
> summary(out, test = "Wilks")
```

	Df	Wilks	approx F	num Df	den Df	Pr(>F)
V1	3	0.012301	13.088	15	50.091	1.84e-12 ***
Residuals	22					

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
> summary(out)
```

	Df	Pillai	approx F	num Df	den Df	Pr(>F)
V1	3	1.5539	4.2984	15	60	2.413e-05 ***
Residuals	22					

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

⇒ at least one of the chemicals differs among the sites

Inference and Comparison of Mean Vectors

CLEMSON UNIVERSITY

Confidence Intervals/Region for Population Means

Hypothesis Testing for Mean Vector

Multivariate Paired Hotelling's T-Square

Comparisons of Two Mean Vectors

Multivariate Analysis of Variance

4.36

Notes

Summary

In this lecture, we learned about:

- Confidence Intervals/Regions for Mean Vector
- Hypothesis Testing for Mean Vector
- Multivariate Version of Paired Tests
- Hypothesis Testing for Two Mean Vectors
- MANOVA

In the next two lectures, we will learn about Multivariate Regression

Inference and Comparison of Mean Vectors
CLEMSON UNIVERSITY
Confidence Intervals/Region for Population Means
Hypothesis Testing for Mean Vector
Multivariate Paired Hotelling's T-Square
Comparisons of Two Mean Vectors
Multivariate Analysis of Variance
4.37

Notes

Notes

Notes
