# 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

# 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 on Mean Vectors

**This Week's Topics:**

- **Single Mean Vector:** inference on $\boldsymbol{\mu}$ (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 $\boldsymbol{\mu}$
- Paired $t$-test $\rightarrow$ paired mean vectors
- Two-sample $t$-test $\rightarrow$ two mean vectors
- ANOVA $\rightarrow$ MANOVA

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

Inference and Comparison of Mean Vectors

CLEMS❄N
U N I V E R S I T Y

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

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

- $\bar{X}_n \overset{\cdot}{\sim} \mathrm{N}(\mu, \frac{\sigma^2}{n})$ 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., $X \sim \mathrm{N}(\mu, \sigma^2)$

- The standard error $\sqrt{\mathrm{Var}(\bar{X}_n)} = \frac{\sigma}{\sqrt{n}}$ provides a measure estimation precision. In practice, we use $\frac{s}{\sqrt{n}}$ instead where $s$ is the sample standard deviation

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

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

- $\bar{X}_n \overset{\cdot}{\sim} N(\mu, \frac{1}{n}\Sigma)$ 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 $\Sigma$ by its estimate $S$, the sample covariacne matrix

# Review: Interval Estimation of Univariate Population Mean $\mu$

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 $\mu$ 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, $\alpha$, the sample size $n$

# Constructing Confidence Intervals for Mean Vector

We will still use the general recipe

Sample mean $\pm$ multiplier $\times$ 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 $\mu_j$ is computed as

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

- Bonferroni Method: a CI for $\mu_j$ is computed as

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

- Simultaneous CIs: a CI for $\mu_j$ is computed as

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

**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

CLEMS🐯N
U N I V E R S I T Y

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

# Mineral Content Measurements Example Cont'd

Inference and
Comparison of Mean
Vectors

CLEMSON
U N I V E R S I T Y

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

- One at a Time CIs: $\bar{x}_j \pm t_{n-1,\alpha/2} \frac{s_j}{\sqrt{n}}, \quad j = 1, \cdots, p$. Therefore 95% CIs for $\mu_1$ and $\mu_2$ are:

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

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

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

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

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

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

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

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

# 95 % CIs Based on Three Methods

Inference and
Comparison of Mean
Vectors

CLEMS🐾N
U N I V E R S I T Y

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

# Confidence Ellipsoid

A confidence ellipsoid for $\boldsymbol{\mu}$ is the set of $\boldsymbol{\mu}$ satisfying

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

# Hypothesis Testing for Mean

Inference and
Comparison of Mean
Vectors

CLEMSON
U N I V E R S I T Y

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

- Recall: for univariate data, $t$ statistic

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

Under $H_0 : \mu = \mu_0$

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

- Extending to multivariate by analogy:

$$T^2 = n\left(\bar{\boldsymbol{X}}_n - \boldsymbol{\mu}_0\right)^T \boldsymbol{S}^{-1}\left(\bar{\boldsymbol{X}}_n - \boldsymbol{\mu}_0\right)$$

Under $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$

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

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

# Hypothesis Testing for Mean Vector $\mu$

1. State the null

$$H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$$

and the alternative

$$H_a : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$$

2. Compute the test statistic

$$F = \frac{n-p}{(n-1)p} n \left( \bar{\boldsymbol{X}}_n - \boldsymbol{\mu}_0 \right)^T \boldsymbol{S}^{-1} \left( \bar{\boldsymbol{X}}_n - \boldsymbol{\mu}_0 \right)$$

3. **Compute the P-value**. Under $H_0$ : $\quad F \sim F_{p,n-p}$

4. **Draw a conclusion**: We do (or do not) have enough statistical evidence to conclude $\boldsymbol{\mu} \neq \boldsymbol{\mu}_0$ at $\alpha$ significant level

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

Inference and Comparison of Mean Vectors

CLEMS🐾N
U N I V E R S I T Y

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

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

| Variable | Recommended Intake ($\boldsymbol{\mu}_0$) | Sample Mean ($\bar{\boldsymbol{x}}_n$) |
|---|---|---|
| Calcium | 1000 $mg$ | 624.0 $mg$ |
| Iron | 15 $mg$ | 11.1 $mg$ |
| Protein | 60 $g$ | 65.8 $g$ |
| Vitamin A | 800 $\mu g$ | 839.6 $\mu g$ |
| Vitamin C | 75 $mg$ | 78.9 $mg$ |

Here we would like to test, at $\alpha = 0.01$ level, if the $\boldsymbol{\mu} = \boldsymbol{\mu}_0$

# Women's Dietary Intake Example Analysis

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

1. State the null

$$H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$$

and the alternative

$$H_a : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$$

2. Compute the test statistic
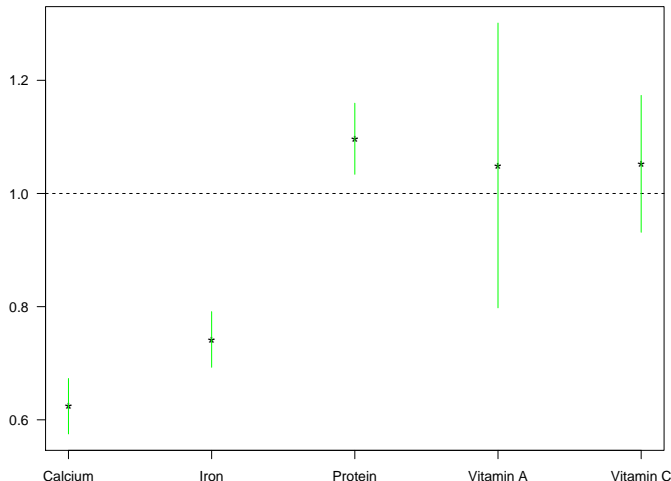
$$F = \frac{n-p}{(n-1)p} n \left(\bar{\boldsymbol{x}}_n - \boldsymbol{\mu}_0\right)^T \boldsymbol{S}^{-1} \left(\bar{\boldsymbol{x}}_n - \boldsymbol{\mu}_0\right) = 349.80$$

3. **Compute the P-value**. Under $H_0 : \quad F \sim F_{p,n-p} \Rightarrow$ p-value $= \mathbb{Pr}(F_{p,n-p} > 349.80) = 3 \times 10^{-191} < \alpha = 0.01$

4. **Draw a conclusion**: We do have enough statistical evidence to conclude $\boldsymbol{\mu} \neq \boldsymbol{\mu}_0$ at $\alpha = 0.01$ significant level

**Inference and Comparison of Mean Vectors**

CLEMS☙N
U N I V E R S I T Y

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

## 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

## 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
U N I V E R S I T Y

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

# Multivariate Paired Hotelling's T-Square

Inference and Comparison of Mean Vectors

CLEMSON
U N I V E R S I T Y

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

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 $H_a : \mu_D \neq \mathbf{0}$

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 $H_0 :$ $\quad F \sim F_{p,n-p}$

4. **Draw a conclusion**: We do (or do not) have enough statistical evidence to conclude $\mu_D \neq \mathbf{0}$ at $\alpha$ significant level

# Spouse Survey Data Example Analysis

1. State the null

$$H_0 : \boldsymbol{\mu}_D = \mathbf{0}$$

and the alternative

$$H_a : \boldsymbol{\mu}_D \neq \mathbf{0}$$

2. Compute the test statistic

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

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

4. **Draw a conclusion**: We do have enough statistical evidence to conclude $\boldsymbol{\mu}_D \neq \mathbf{0}$ at $0.05$ significant level

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

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

- The first population is the population of Genuine Bank Notes

- 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

# Review: Two Sample t-Test

Suppose we have data from a single variable from population 1: $X_{11}, X_{12}, \cdots, X_{1n_1}$ and population 2: $X_{21}, X_{22}, \cdots, X_{2n_2}$. Here we would like to draw inference about their population means $\mu_1$ and $\mu_2$.

**Assumptions**:

- Homoscedasticity: The data from both populations have common variance $\sigma^2$

- 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

- 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 $H_a : \mu_1 \neq \mu_2$

## 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

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

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

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

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

## The Two Sample Problem: The Multivariate Case

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

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

to infer the relationship between $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$, where

$$\boldsymbol{\mu}_i = \begin{bmatrix} \mu_{i1} \\ \mu_{i2} \\ \vdots \\ \mu_{ip} \end{bmatrix}$$

### Assumptions

- Both populations have common covariance matrix, i.e., $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$

- Independence: The subjects from both populations are independently sampled

- Normality: Both populations are normally distributed

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

## The Multivariate Two-Sample Problem

Here we are testing

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

Under the common covariance assumption we have

$$\boldsymbol{S}_p = \frac{(n_1 - 1)\boldsymbol{S}_1 + (n_2 - 1)\boldsymbol{S}_2}{n_1 + n_2 - 2},$$

where

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

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

Inference and Comparison of Mean Vectors

CLEMSON
U N I V E R S I T Y

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

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:

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

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

# Two-Sample Test for Swiss Bank Notes

**Inference and Comparison of Mean Vectors**

CLEMS🐯N
U N I V E R S I T Y

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

```
> (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 $\Rightarrow$ which ones?

# Simultaneous Confidence Intervals

Inference and
Comparison of Mean
Vectors

CLEMS☴N
U N I V E R S I T Y

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

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

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 | $(1.81, 2.33)$ |

# Checking Model Assumptions

**Assumptions**:

- Homoscedasticity: The data from both populations have common covariance matrix $\Sigma$

  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

# 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 $k \in \{1, 2, \cdots, p\}$, 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:

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

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

Inference and Comparison of Mean Vectors

CLEMS**N
U N I V E R S I T Y

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

- Pottery shards are collected from four sites in the British Isles:
  - Llanedyrn (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

# Review: (Univariate) Analysis of Variance (ANOVA)

**Inference and Comparison of Mean Vectors**

CLEMS☙N
U N I V E R S I T Y

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

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

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

- Test Statistic: $F^* = \frac{\text{MSTr}}{\text{MSE}}$. 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. $\sigma_1^2 = \sigma_2^2 = \cdots = \sigma_g^2$)

  - Responses for a given group are independent to each other

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

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

- **Notation**: $\boldsymbol{Y}_{ij}$ is the vector of variables for subject $j$ in group $i$; $n_i$ is the sample size in group $i$; $N = n_1 + n_2 + \cdots + n_g$ the total sample size

- **Assumptions**: 1) common covariance matrix $\boldsymbol{\Sigma}$; 2) Independence; 3) Normality

# Test Statistics for MANOVA

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

$$H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \cdots = \boldsymbol{\mu}_g.$$

The alternative hypothesis:

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

- **Mean vectors**:
  - Sample Mean Vector: $\bar{\boldsymbol{y}}_{i\cdot} = \frac{1}{n_i} \boldsymbol{Y}_{ij}, \quad i = 1, \cdots, g$
  - Grand Mean Vector: $\bar{\boldsymbol{y}}_{\cdot\cdot} = \frac{1}{N} \sum_{i=1}^{g} \sum_{j=1}^{n_i} \boldsymbol{Y}_{ij}$

- **Total Sum of Squares**:

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

# MANOVA Decomposition and MANOVA Table

Inference and Comparison of Mean Vectors

CLEMSON
U N I V E R S I T Y

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

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

$$= \sum_{i=1}^{g} \sum_{j=1}^{n_i} \left[ (\boldsymbol{Y}_{ij} - \bar{\boldsymbol{y}}_{i.}) + (\bar{\boldsymbol{y}}_{i.} - \bar{\boldsymbol{y}}_{..}) \right] \left[ (\boldsymbol{Y}_{ij} - \bar{\boldsymbol{y}}_{i.}) + (\bar{\boldsymbol{y}}_{i.} - \bar{\boldsymbol{y}}_{..}) \right]^T$$

$$= \underbrace{\sum_{i=1}^{g} \sum_{j=1}^{n_i} (\boldsymbol{Y}_{ij} - \bar{\boldsymbol{y}}_{i.})(\boldsymbol{Y}_{ij} - \bar{\boldsymbol{y}}_{i.})^T}_{\boldsymbol{E}} + \underbrace{\sum_{i=1}^{g} n_i (\bar{\boldsymbol{y}}_{i.} - \bar{\boldsymbol{y}}_{..})(\bar{\boldsymbol{y}}_{i.} - \bar{\boldsymbol{y}}_{..})^T}_{\boldsymbol{H}}$$

MANOVA Table

| **Source** | df | SS |
|---|---|---|
| Treatment | $g-1$ | $\boldsymbol{H}$ |
| Error | $N-g$ | $\boldsymbol{E}$ |
| Total | $N-1$ | $\boldsymbol{T}$ |

Reject $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \cdots = \boldsymbol{\mu}_g$ if the matrix $\boldsymbol{H}$ is "large" relative to the matrix $\boldsymbol{E}$

# Test Statistics for MANOVA

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

- Wilks Lambda

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

  Reject $H_0$ if $\Lambda^*$ is "small"

- Hotelling-Lawley Trace

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

  Reject $H_0$ if $T_0^2$ is "large"

- Pillai Trace

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

  Reject $H_0$ if $V$ is "large"

# 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
```

$\Rightarrow$ at least one of the chemicals differs among the sites

## 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