# MC3020 - Comparing Two Population Parameters

T. Mayooran,
Department of Interdisciplinary Studies,
Faculty of Engineering,
University of Jaffna.
Email: mayooran@eng.jfn.ac.lk

# Comparing two population parameters

➢ Comparing Two Population Proportions.

➢ Comparing Two Independent Population Means.

➢ Comparing Two Dependent or Matched Population Means.

➢ Comparing Two Independent Population Variances.

MC3020

# Comparing Two Population Proportions

# Confidence interval for $p_1 - p_2$

Conditions:

1. Let be $X_1$ the number of successes in $n_1$ Bernoulli trials having proportion of success $p_1$.

2. Let $X_2$ be the number of successes in $n_2$ Bernoulli trials having proportion of success $p_2$.

3. The two experiments are independent.

www.fppt.info

The point estimate for the difference between the proportions is $p_1 - p_2$,

$$\hat{p}_1 - \hat{p}_2 = \frac{X_1}{n_1} - \frac{X_2}{n_2}$$

The population mean of $\hat{p}_1 - \hat{p}_2$ is,

$$E(\hat{p}_1 - \hat{p}_2) = p_1 - p_2$$

The population variance of $\hat{p}_1 - \hat{p}_2$ is,

$$\frac{p_1(1 - p_1)}{n_1} - \frac{p_2(1 - p_2)}{n_2}$$

www.fppt.info

By the central limit theorem, for large $n_1$ and $n_2$, $\hat{p}_1 - \hat{p}_2$ has approximate normal distribution with mean $p_1 - p_2$ and variance $\frac{p_1(1-p_1)}{n_1} - \frac{p_2(1-p_2)}{n_2}$

Then $(1-\alpha) * 100\%$ confidence interval estimate for $p_1 - p_2$ is computed as

$$(\hat{p}_1 - \hat{p}_2 - E, \hat{p}_1 - \hat{p}_2 + E)$$

where, $E = Z_{\alpha/2} * \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$

www.fppt.info

# Example 1:

It is claimed that in the 2008 Democratic Presidential Nomination Primaries in USA, Senator Barack Obama was preferred by the black voters. To test the claim, a research firm sampled 600 black democrats and found that 384 support the senator and in another sample of 720 non-black democrats 417 support the senator. Construct a 97% confidence interval for the difference between the two populations proportions.

# Testing for the difference between two independent population proportions

Case 1: $H_0 : p_1 \geq p_2 \quad H_1 : p_1 < p_2$

Case 2: $H_0 : p_1 \leq p_2 \quad H_1 : p_1 > p_2$

Case 3: $H_0 : p_1 = p_2 \quad H_1 : p_1 \neq p_2$

The corresponding test statistic is

$$Z = \frac{\hat{p}_1 - \hat{p}_2 - (p_1 - p_2)}{\sqrt{\bar{p}(1 - \bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

where $\bar{p} = \frac{X_1 + X_2}{n_1 + n_2}$

Note that $p_1 - p_2 = 0$ in computations for all three cases above. But in general it is not necessarily zero as if we want to test that one proportion is at least an amount higher than the other then $p_1 - p_2$ is that least amount in proportion, and so on.

www.fppt.info

# Example 2:

It is claimed that in the 2008 Democratic Presidential Nomination Primaries in USA, Senator Barack Obama was preferred by the black voters. To test the claim, a research firm sampled 600 black democrats and found that 384 support the senator and in another sample of 720 non-black democrats 417 support the senator. Test the claim using 5% level of significance.

# Exercise:

In your Tutorial 5,

- 1
- 3

www.fppt.info

# Comparing Two Independent Population Means

www.fppt.info

Let $X_1, X_2, \ldots X_{n_1}$ be a random sample from a population with mean $\mu_1$ and variance $\sigma_1^2$.

Let $Y_1, Y_2, \ldots Y_{n_2}$ be a random sample from a population with mean $\mu_2$ and variance $\sigma_2^2$. The populations are independent.

The point estimate for the population mean difference $\mu_1 - \mu_2$ is the sample mean difference $\bar{X} - \bar{Y}$.

www.fppt.info

When the two population distributions are normal, the sampling distribution of $\bar{X} - \bar{Y}$ is normal with mean

$$\text{E}(\bar{X} - \bar{Y}) = \mu_1 - \mu_2$$

and variance

$$\text{Var}(\bar{X} - \bar{Y}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

# The confidence interval for the difference between the two-population means is computed as follows:

## Case 1:

When the two independent population distributions are normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are known, the $(1 - \alpha)*100$ % confidence interval for $\mu_1 - \mu_2$ is computed as,

$$\left( \bar{X} - \bar{Y} - Z_{\alpha/2} * \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \ , \bar{X} - \bar{Y} + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

## Case 2:

When the two independent population distributions are <span style="color:green">normal</span> and the <span style="color:green">population variances</span> $\sigma_1^2$ <span style="color:green">and</span> $\sigma_2^2$ <span style="color:green">are unknown and unequal</span>, the $(1-\alpha)*100$ % confidence interval for $\mu_1 - \mu_2$ is computed as,

$$\left( \bar{X} - \bar{Y} - t_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \ , \bar{X} - \bar{Y} + t_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right)$$

Where $s_1^2$ and $s_2^2$ are the corresponding sample variances, and the degrees of freedom for t is,

MC3020

$$df = \frac{(A + B)^2}{\dfrac{A^2}{n_1 - 1} + \dfrac{B^2}{n_2 - 1}}$$

Where $A = \dfrac{s_1^2}{n_1}$ and $B = \dfrac{s_2^2}{n_2}$

Case 3:

When the two independent population distributions are normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are unknown but equal, the $(1 - \alpha) * 100$ % confidence interval for $\mu_1 - \mu_2$ is computed as,

www.fppt.info

$$\left( \bar{X} - \bar{Y} - t\alpha_{/2} * S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \; , \bar{X} - \bar{Y} + t\alpha_{/2} * S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$$

Where $S_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$ is the pooled standard deviation from the two sample standard deviations, and the degrees of freedom for t is $n_1 + n_2 - 2$.

www.fppt.info

## Case 4:

When the two independent population distributions are not normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are known, and the sample size $n_1$ and $n_2$ are large, the $(1 - \alpha)*100$ % confidence interval for $\mu_1 - \mu_2$ is computed as,

$$\left( \bar{X} - \bar{Y} - Z_{\alpha/2} * \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \ , \bar{X} - \bar{Y} + Z_{\alpha/2} * \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

## Case 5:

When the two independent population distributions are not normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are unknown, and the sample size $n_1$ and $n_2$ are large, the $(1 - \alpha)*100$ % confidence interval for $\mu_1 - \mu_2$ is computed as,

$$\left( \bar{X} - \bar{Y} - Z_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \ , \bar{X} - \bar{Y} + Z_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right)$$

www.fppt.info

# Example 3:

Two different types of drugs 'A' and 'B' are tried on certain patients for increasing weight. Five randomly selected patients were given drug 'A', and 7 randomly selected patients were given drug 'B'. The increases in weight (in pounds) are given below:

Drug 'A':  8 12  13   9  3

Drug 'B': 10  8  12  15  6  8  11

Assume that the population distributions of the measurements are normal with equal variances. Construct a 95% confidence interval for the difference between the two means.

www.fppt.info

# Example 4:

To test effect of a fertilizer on rice production, 64 plots of land having equal areas were chosen. Half of these plots were treated with fertilizer and the other half were untreated. Other conditions were the same. The mean yield of rice on the untreated plots was 4.8 quintals with a standard deviation of 0.4 quintal, while the mean yield on the treated plots was 5.1 quintals with a standard deviation of 0.36 quintal. Construct a 94% confidence interval estimate for the mean difference between the untreated plots and treated plots.

www.fppt.info

# Exercise:

❖In your tutorial 5, 7.4

www.fppt.info

# Testing for the difference between two
# independent population means

Case 1: $H_0 : \mu_1 \geq \mu_2$    $H_1 : \mu_1 < \mu_2$

Case 2: $H_0 : \mu_1 \leq \mu_2$    $H_1 : \mu_1 > \mu_2$

Case 3: $H_0 : \mu_1 = \mu_2$    $H_1 : \mu_1 \neq \mu_2$

## Case 1:

When the two independent population distributions are normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are known, the test statistic is

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}}$$

## Case 2:

When the two independent population distributions are normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are unknown and unequal, the test statistic is

$$T = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\dfrac{S_1^2}{n_1} + \dfrac{S_2^2}{n_2}}}$$

where T has a t-distribution with degrees of freedom

$$df = \frac{(A+B)^2}{\dfrac{A^2}{n_1-1} + \dfrac{B^2}{n_2-1}}$$

Where $A = \dfrac{S_1^2}{n_1}$ and $B = \dfrac{S_2^2}{n_2}$

www.fppt.info

## Case 3:

When the two independent population distributions are normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are unknown but equal, the test statistic is

$$T = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{S_p \sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}}$$

where T has a t distribution with degrees of freedom $n_1 + n_2 - 2$ and $S_p = \sqrt{\dfrac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$ is the pooled standard deviation from the two sample standard deviations.

www.fppt.info

## Case 4:

When the two independent population distributions are not normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are known, and the sample size $n_1$ and $n_2$ are large, the test statistic is

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}}$$

www.fppt.info

**Case 5:**

When the two independent population distributions are not normal and the population variances $\sigma_1^2$ and $\sigma_2^2$ are unknown, and the sample size $n_1$ and $n_2$ are large, the test statistic is

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}}$$

www.fppt.info

Note that $\mu_1 - \mu_2$ for all three cases above. But in general it is not necessarily zero as if we want to test that one mean is at least an amount higher than the other then $\mu_1 - \mu_2$ is that least amount, and so on.

www.fppt.info

# Example 5:

Two different types of drugs 'A' and 'B' are tried on certain patients for increasing weight. Five randomly selected patients were given drug 'A' and 7 randomly selected patients were given drug 'B'. The increases in weight (in pounds) are given below:

Drug 'A': 8 12 13 9 3

Drug 'B': 10 8 12 15 6 8 11

Assume that the population distributions of the measurements are normal. Do the two drugs differ significantly with regard to their effect in increasing weight? Use 0.05 significance level.

www.fppt.info

# Example 6:

To test effect of a fertilizer on rice production, 64 plots of land having equal areas were chosen. Half of these plots were treated with fertilizer and the other half were untreated. Other conditions were the same. The mean yield of rice on the untreated plots was 4.8 quintals with a standard deviation of 0.4 quintal, while the mean yield on the treated plots was 5.1 quintals with a standard deviation of 0.36 quintal. Can we conclude that there is a significant improvement in rice production because of the fertilizer at 4% level of significance?

# Example :

An urban economist wanted to determine whether the mean price of a home in Lemont is less than the mean price of a home in Naperville. A random sample of homes sold in each neighborhood results in the following statistics, where the means and standard deviations are in thousands of dollars:
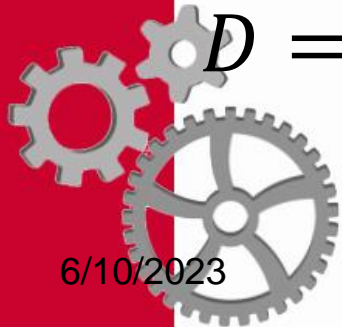
| Lemont | Naperville |
|---|---|
| $n_1 = 50$ | $n_2 = 50$ |
| $\overline{x}_1 = 200$ | $\overline{x}_2 = 300$ |
| $s_1 = 45$ | $s_2 = 75$ |

Test the claim that housing is less expensive in Lemont than in Naperville at the 5% level of significance.

www.fppt.info

# Comparing Two Dependent or Matched Population Means

www.fppt.info

Let $(X_1, Y_1), (X_2, Y_2), \ldots (X_n, Y_n)$ be pairs of random measurements from two dependent or matched populations. Such situations usually but not exclusively occur when measurements are taken on the same subjects before or after the experimentation. To consider all other factors same except the factor that measured by and , we combine the populations by finding the difference $D = X - Y$.

Let $D_1, D_2, \ldots . D_n$ be a random sample from a population with mean $\mu_D$ and variance $\sigma_D^2$.

The point estimate for the population mean difference is $\mu_1 - \mu_2$ the sample mean of the differences

$$\overline{D} = \frac{\sum D}{n} .$$

The mean of $\overline{D}$ is

$$\mathrm{E}(\overline{D}) = \mu_D$$

www.fppt.info

The variance of $\overline{D}$ is $V(\overline{D}) = \frac{\sigma_D^2}{n}$

When the population distribution of the differences $D$ is normal, the sampling distribution of $\overline{D}$ is normal with mean $\mathrm{E}(\overline{D}) = \mu_D$ and variance $V(\overline{D}) = \frac{\sigma_D^2}{n}$.

The $(1 - \alpha) * 100\%$ confidence interval for the difference between the two dependent population means $\mu_1 - \mu_2 = \mu_D$ is computed as follows,

Case 1

When the population distribution of the differences $D$ is normal and the population variance $\sigma_D^2$ is known, the $(1 - \alpha) * 100\%$ confidence interval for $\mu_1 - \mu_2 = \mu_D$ is computed as

$$\left( \bar{D} - Z_{\alpha/2} * \frac{\sigma_D}{\sqrt{n}} , \bar{D} + Z_{\alpha/2} * \frac{\sigma_D}{\sqrt{n}} \right)$$

<span style="color:red">Case 2</span>   When the population distribution of the differences $D$ <span style="color:green">is normal</span> and the <span style="color:green">population variance $\sigma_D^2$ is unknown</span>, the $(1 - \alpha) * 100\%$ confidence interval for $\mu_1 - \mu_2 = \mu_D$ is computed as

$$\left( \overline{D} - t\alpha_{/2} * \frac{S_D}{\sqrt{n}} \ , \overline{D} + t\alpha_{/2} * \frac{S_D}{\sqrt{n}} \right)$$

where the degrees of freedom for t is $n - 1$ and

$$S_D = \sqrt{\frac{\sum(D - \overline{D})^2}{n-1}} = \sqrt{\frac{n \sum D^2 - (\sum D)^2}{n(n-1)}} \quad \text{the sample}$$

standard deviation for the differences.

www.fppt.info

## Case 3

When the population distribution of the differences $D$ is not normal, $n$ is large, and the population variance $\sigma_D^2$ is known, the $(1 - \alpha) * 100\%$ confidence interval for $\mu_1 - \mu_2 = \mu_D$ is computed as,

$$\left( \bar{D} - Z_{\alpha/2} * \frac{\sigma_D}{\sqrt{n}} \;, \bar{D} + Z_{\alpha/2} * \frac{\sigma_D}{\sqrt{n}} \right)$$

## Case 4

When the population distribution of the differences $D$ is not normal, $n$ is large, and the population variance $\sigma_D^2$ is unknown, the $(1 - \alpha) * 100\%$ confidence interval for $\mu_1 - \mu_2 = \mu_D$ is computed as,

$$\left( \overline{D} - Z_{\alpha/2} * \frac{S_D}{\sqrt{n}} \ , \overline{D} + Z_{\alpha/2} * \frac{S_D}{\sqrt{n}} \right)$$

www.fppt.info

Note that since the measurements are on the same subjects, if there is no real difference in terms of the factor of interest, the distributions of the differences are often normal. Hence the situation 2 is a very common phenomenon.

# Testing for the difference between two dependent population means

Case 1: $H_0: \mu_1 \geq \mu_2 \qquad \equiv \quad H_0 : \mu_D \geq 0$

$\qquad\qquad H_1: \mu_1 < \mu_2 \qquad\qquad H_1: \mu_D < 0$

Case 2: $H_0: \mu_1 \leq \mu_2 \qquad \equiv \quad H_0 : \mu_D \leq 0$

$\qquad\qquad H_1 : \mu_1 > \mu_2 \qquad\qquad H_1 : \mu_D > 0$

Case 3: $H_0 : \mu_1 = \mu_2 \qquad \equiv \quad H_0 : \mu_D = 0$

$\qquad\qquad H_1 : \mu_1 \neq \mu_2 \qquad\qquad H_1 : \mu_D \neq 0$

www.fppt.info

Case 1: When the population distribution of the differences $D$ is normal and the population variance $\sigma_D^2$ is known, the test statistic is

$$Z = \frac{\overline{D} - \mu_D}{\sigma_D / \sqrt{n}}$$

Case 2 When the population distribution of the differences $D$ is normal and the population variance $\sigma_D^2$ is unknown, the statistic is

$$T = \frac{\overline{D} - \mu_D}{S_D / \sqrt{n}}$$

where T has a t distribution with $(n - 1)$ degrees of freedom

Case 3 When the population distribution of the differences $D$ is not normal, $n$ is large, and the population variance $\sigma_D^2$ is known, the test statistic is

$$Z = \frac{\bar{D} - \mu_D}{\sigma_D / \sqrt{n}}$$

Case 4 When the population distribution of the differences $D$ is not normal, $n$ is large, and the population variance $\sigma_D^2$ is unknown, the test statistic is

$$Z = \frac{\overline{D} - \mu_D}{S_D / \sqrt{n}}$$

www.fppt.info

Note that $\mu_1 - \mu_2 = \mu_D = 0$ for all three cases above. But in general it is not necessarily zero as if we want to test that one mean is at least an amount higher than the other then $\mu_1 - \mu_2$ is that least amount, and so on.

www.fppt.info

# Example 7:

To compare the demand for two different entrees, the manager of a cafeteria recorded the number of purchases for each entrée on seven consecutive days. The data are shown in the next table: Since the cafeteria demands depend on days of the week, the data are considered to be related. Assume that the differences of the measurements follow a normal distribution.

| | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday | Sunday |
|---|---|---|---|---|---|---|---|
| A | 420 | 374 | 434 | 395 | 637 | 594 | 679 |
| B | 391 | 343 | 469 | 412 | 538 | 521 | 625 |

www.fppt.info

a) Construct a 95% confidence interval for the mean difference.

b) Test the hypothesis that the demand for item A is higher than the demand for item B. Use 5% level of significance.

# Example:

A test preparation company claims that its SAT preparation course improves SAT math scores. The company administers the SAT to 9 randomly selected students and determines their scores. The same students then participate in the course. Upon completion, they retake the SAT. The results are presented below:

Before:  436      431      270      463      528      377      397
         413      525

After:   443      429      287      501      522      380      402
         450      548

Test the claim that the preparatory course improves SAT math scores at the 10% level of significance. (Assume that the differences between the scores have an approximate normal distribution.)

www.fppt.info

# Comparing Two Independent Population Variances

www.fppt.info

Let $X_1, X_2, \ldots X_{n_1}$ be a random sample from a population with mean $\mu_1$ and variance $\sigma_1^2$.
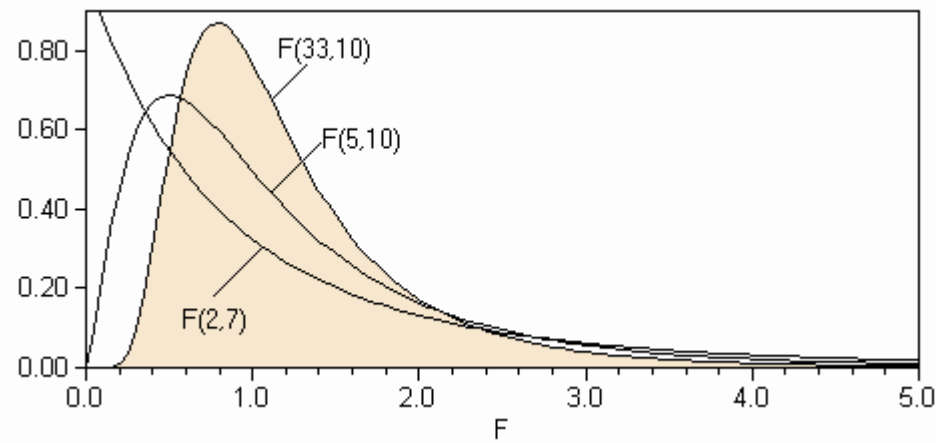
Let $Y_1, Y_2, \ldots Y_{n_2}$ be a random sample from a population with mean $\mu_2$ and variance $\sigma_2^2$. The populations are independent.

As the variances cannot be negative, we consider the ratio of the two variances instead of their difference in making comparison between them.

The point estimate for the ratio $\sigma_1^2 / \sigma_2^2$ is $S_1^2 / S_2^2$ where $S_1^2$ and $S_2^2$ are the respective sample variances.

www.fppt.info

It is known that when the populations are normal, $s_1^2 / s_2^2$ follows an F-distribution with degrees of freedom $n_1 - 1$ and $n_2 - 1$. The table for the F-distribution (Table 5) is given in the Appendix II.

www.fppt.info

Then $(1 - \alpha) * 100\%$ confidence interval for the $\sigma_1^2 / \sigma_2^2$ is computed as,

$$\frac{S_1^2}{S_2^2} * \frac{1}{F_{\alpha/2, n_1 - 1, n_2 - 1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{S_1^2}{S_2^2} * F_{\alpha/2, n_2 - 1, n_1 - 1}$$

Where $F_{\alpha/2, n_2 - 1, n_1 - 1}$ is the percentile value in the F-distribution such that the right side area is for degrees of freedom $n_2 - 1$ and $n_1 - 1$. And $F_{\alpha/2, n_1 - 1, n_2 - 1}$ is the percentile value in the F distribution such that the right side area is $\alpha/2$ for degrees of freedom $n_1 - 1$ and $n_2 - 1$.

www.fppt.info

# Testing for difference between two independent population variances,

### Case 1:

$$H_0 : \sigma_1^2 \geq \sigma_2^2 \quad H_1 : \sigma_1^2 < \sigma_2^2$$

### Case 2:

$$H_0 : \sigma_1^2 \leq \sigma_2^2 \quad H_1 : \sigma_1^2 > \sigma_2^2$$

### Case 3:

$$H_0 : \sigma_1^2 = \sigma_2^2 \quad H_1 : \sigma_1^2 \neq \sigma_2^2$$

www.fppt.info

Under the assumption that the null hypothesis $H_0$ is true, the test statistic

$$F = \frac{S_1^2}{S_2^2}$$

has F-distribution with degrees of freedom $n_1 - 1$ and $n_2 - 1$ .

www.fppt.info

# Example 7.8

Let us consider the final scores of the Author's two different sections of the elementary statistics courses:

Section 1: 38, 88, 91, 84, 97, 78, 51, 90, 72, 73, 73, 55, 83, 72, 97, 33, 78, 91, 93, 65, 86, 81, 87, 81, 28, 74

Section 2: 64, 36, 87, 73, 72, 43, 90, 81, 79, 43, 77, 89, 91, 72, 75, 68, 78, 72, 81, 72, 35, 72, 93, 74, 85.

Assume that the samples are from independent normal populations.

**Solution:** To construct a 90% confidence interval for $\frac{\sigma_1^2}{\sigma_2^2}$, we obtain $\frac{S_1^2}{S_2^2} *$

$$\frac{1}{F_{\alpha/2,n_1-1,n_2-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{S_1^2}{S_2^2} * F_{\alpha/2,n_2-1,n_1-1}$$

$$\frac{19.12^2}{16.48^2} * \frac{1}{F_{0.05,26-1,25-1}} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{19.12^2}{16.48^2} * F_{0.05,25-1,26-1}$$

$$\frac{19.12^2}{16.48^2} * \frac{1}{1.9750} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{19.12^2}{16.48^2} * 1.9643$$

$$0.6815 < \frac{\sigma_1^2}{\sigma_2^2} < 2.6440$$

www.fppt.info

To test the claim that the first section has higher variance compared to the second section using 5% level of significance, the hypotheses can be written as

$$H_0 : \sigma_1^2 \leq \sigma_2^2 \quad H_1 : \sigma_1^2 > \sigma_2^2$$

The test statistic

$$F = \frac{S_1^2}{S_2^2} = \frac{19.12^2}{16.48^2} = 1.3461$$

The 5% critical value for degrees of freedom 25 and 24 is 1.975.

So, we fail to reject $H_0$ at 0.05 level and conclude that the Section 1 variance is not significantly higher than the Section 2 variance.

Similarly, p-value= $P(F > 1.3461) > 0.05$

So we fail to reject at 0.05 level. Same conclusion!

www.fppt.info

# Exercises

In your Tutorial 5,

- 5
- 6
- 8
- 14

www.fppt.info

Don't hesitate to contact us if you have any questions about this course's teaching contents. Also, don't forget to check out the course page and Microsoft Team folder,

- course page Link: https://mayooran1987.github.io/MC3020/

- Course page's QR code

- Microsoft Team folder link