

1. Introduction

Welcome to the series of E-learning modules on Confidence Intervals for the difference between the means.

By the end of this session, you will be able to:

- Explain the confidence interval for the difference between the averages of two populations when the population standard deviations are known
- Explain the confidence interval for the difference between the averages of two populations when the population standard deviations are common but unknown
- Explain the confidence interval for the difference of means of two dependent samples

The problem for a quantitative population is the comparison of two population means, which is as important as the estimation of a single population mean μ .

We may be interested to make comparisons like the following:

- The average scores in the medical college admission test for students whose major was Biochemistry and those whose major was Biology
- The average yield in a chemical plant using new raw materials furnished by two different suppliers
- The average stem diameter of plants grown on two different types of nutrients

For each of these examples there are two populations:

The first with mean and variance μ_1 and σ_1^2 and the second with mean and variance μ_2 and σ_2^2 .

A random sample of size m is drawn from the population I and a sample of size n is drawn independently from the population II.

Intuitively, the difference between the two sample means would provide the maximum information about the actual difference between two population means. The best point estimator of the difference ($\mu_1 - \mu_2$) between the population means is ($\bar{x} - \bar{y}$)

The confidence interval formula yields a range (interval) within which the difference between two populations mean is located.

This topic describes how to construct a confidence interval for the difference between two means.

Estimation Requirements

The approach described in this topic is valid whenever the following conditions are met:

- Both samples are [simple random samples](#)
- The samples are [independent](#)
- Each [population](#) is at least 10 times larger than its respective [sample](#)
- The [sampling distribution](#) of the difference between means is approximately normally distributed

Generally, the sampling distribution will be approximately normally distributed if each sample is described by at least one of the following statements:

- The population distribution is normal
- The sampling distribution is [symmetric](#), [unimodal](#), without [outliers](#), and the sample size is fifteen or less
- The sampling distribution is moderately [skewed](#), unimodal, without outliers, and the sample size is between sixteen and forty
- The sample size is greater than forty, without outliers

The variability of the difference between sample means

To construct a confidence interval, we need to know the variability of the difference between sample means. This means we need to know how to compute the standard deviation of the sampling distribution of the difference.

- If the population standard deviations are known, the standard deviation of the sampling distribution is:

σ is equal to standard deviation of $(\bar{x} - \bar{y})$ which is equal to square root of $\sigma_1^2/m + \sigma_2^2/n$

Where, σ_1 is the standard deviation of the population 1, σ_2 is the standard deviation of the population 2, and m is the size of sample 1, and n is the size of sample 2

When the standard deviation of either population is unknown and the sample sizes (m and n) are large, the standard deviation of the sampling distribution can be estimated by the standard error, using the below equation:

Standard error of $(\bar{x} - \bar{y})$ is equal to square root of $s_1^2/m + s_2^2/n$

Where, s_1 is the standard deviation of the sample 1, s_2 is the standard deviation of the sample 2, m is the size of sample 1, and n is the size of sample 2.

2. CI's for the difference of Means of Two Populations with Known SD's

1. Confidence Intervals for the difference of means of two populations with known standard deviations

Assumptions:

- Populations standard deviations are known
- Population is normally distributed
- If population is not normal, use large sample

Suppose x_1, x_2, \dots, x_m and y_1, y_2, \dots, y_n are the random samples of size m and n drawn from the populations of size N_1 and N_2 with means μ_1 and μ_2 and variances σ_1^2 and σ_2^2 , we know that

\bar{x} is equal to $\sum x_i / m$ and \bar{y} is equal to $\sum y_i / n$

Z is equal to $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n}$ which follows normal with mean zero and variance 1

We can always find two quantities $-Z_{\alpha/2}$ and $+Z_{\alpha/2}$ from standard normal variate tables such as

Probability of $-Z_{\alpha/2} \leq Z \leq +Z_{\alpha/2}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n} \leq +Z_{\alpha/2}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n} \geq -Z_{\alpha/2}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n} \geq -Z_{\alpha/2}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n} \leq +Z_{\alpha/2}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma_1^2 / m + \sigma_2^2 / n} \geq -Z_{\alpha/2}$ is equal to $1 - \alpha$

Therefore, $1 - \alpha$ confidence interval for the difference of means of

two populations with known standard deviations is given by

$$\left[(\bar{x} - \bar{y}) \pm Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}} \right]$$

3. CI's for the difference of Means of Two Populations with Unknown but Common SD's

2. Confidence Interval for the difference of means of two populations with unknown but common standard deviations

Assumptions:

- Populations standard deviations are unknown but common
- Population is normally distributed
- If population is not normal, use large sample

Suppose x_1, x_2, \dots, x_m and y_1, y_2, \dots, y_n are the random samples of size m and n drawn from the normal populations of size N_1 and N_2 with unknown means μ_1 and μ_2 and common variance σ^2 ,

We know that

X_i follows Normal with mean μ_1 and variance σ^2 and

Y_i follows Normal with mean μ_2 and variance σ^2

Where, σ^2 is unknown

We know that \bar{x} is equal to $\sum x_i / m$ and \bar{y} is equal to $\sum y_i / n$

\bar{X} follows Normal with mean μ_1 and variance σ^2 / m and \bar{y} follows Normal with mean μ_2 and variance σ^2 / n

Therefore, $\bar{x} - \bar{y}$ follows Normal with mean $\mu_1 - \mu_2$ and variance $\sigma^2 / m + \sigma^2 / n$

Z is equal to $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma^2 / m + \sigma^2 / n}$ which follows normal with mean zero and variance 1

$(m - 1) \text{ into } s_1^2 / \sigma^2$ follows chi-square with $m - 1$ degrees of freedom and $(n - 1) \text{ into } s_2^2 / \sigma^2$ follows chi-square with $n - 1$ degrees of freedom

Which implies $(m - 1) \text{ into } s_1^2 / \sigma^2 + (n - 1) \text{ into } s_2^2 / \sigma^2$ follows chi-square with $m + n - 2$ degrees of freedom

t is equal to $Z / \sqrt{\chi^2 / \text{degrees of freedom}}$ which is equal to $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / \sqrt{\sigma^2 / m + \sigma^2 / n}$ divided by

$\sqrt{(m - 1) \text{ into } s_1^2 / \sigma^2 + (n - 1) \text{ into } s_2^2 / \sigma^2} / \sqrt{m + n - 2}$ follows t alpha distribution with $(m + n - 2)$ degrees of freedom

t is equal to $(\bar{x} - \bar{y}) - (\mu_1 - \mu_2) / S_p \sqrt{1 / m + 1 / n}$ follows t alpha $(m + n - 2)$

Where, S_p is equal to $\sqrt{(m - 1) \text{ into } s_1^2 + (n - 1) \text{ into } s_2^2} / \sqrt{m + n - 2}$

We can always find two quantities minus $t_{\alpha/2}(m+n-2)$ and $t_{\alpha/2}(m+n-2)$ from student's t variate table such that

Probability of minus $t_{\alpha/2}(m+n-2)$ less than or equal to t less than or equal to plus $t_{\alpha/2}(m+n-2)$ is equal to $1 - \alpha$

Probability of minus $t_{\alpha/2}(m+n-2)$ less than or equal to $(\bar{x} - \bar{y}) - \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ less than or equal to plus $t_{\alpha/2}(m+n-2)$ is equal to $1 - \alpha$

Probability of minus $t_{\alpha/2}(m+n-2)$ into $\frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ less than or equal to $(\bar{x} - \bar{y}) - \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ less than or equal to plus $t_{\alpha/2}(m+n-2)$ into $\frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ is equal to $1 - \alpha$

Probability of $t_{\alpha/2}(m+n-2)$ into $\frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ greater than minus $(\bar{x} - \bar{y}) + \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$

$(\mu_1 - \mu_2)$ greater than minus $t_{\alpha/2}(m+n-2)$ into $\frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) + \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ greater than

$(\mu_1 - \mu_2)$ greater than $(\bar{x} - \bar{y}) - \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ is equal to $1 - \alpha$

Probability of $(\bar{x} - \bar{y}) - \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ less than or equal to

$(\mu_1 - \mu_2)$ less than or equal to $(\bar{x} - \bar{y}) + \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}$ is equal to $1 - \alpha$

Therefore, hundred into $(1 - \alpha)$ percent Confidence Interval for the difference of means of two populations with unknown but common standard deviations is given by

$[(\bar{x} - \bar{y}) - t_{\alpha/2}(m+n-2) \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}, (\bar{x} - \bar{y}) + t_{\alpha/2}(m+n-2) \frac{s_p}{\sqrt{\frac{1}{m} + \frac{1}{n}}}]$

4. CI's for the Mean in case of Correlated Variables

3. Confidence Interval for the mean in case of correlated (dependent) variables

Let (x_i, y_i) , is equal to 1, 2, up to n be n pairs of observations which are correlated. Let the deviation d_i is equal to x_i minus y_i ,

Assume that d_i follows normal with mean θ and variance σ_d^2 .

Then, \bar{d} follows normal with mean θ and variance σ_d^2 by n

Z is equal to $\bar{d} - \theta$ by σ_d by \sqrt{n} which is equal to $\sqrt{n} \text{ into } (\bar{d} - \theta) \text{ by } \sigma_d$

We know that $(n - 1) \text{ into } s_d^2 \text{ by } \sigma_d^2$ follows chi square with $(n - 1)$ degrees of freedom, where s_d^2 is equal to summation $(d_i - \bar{d})^2$ whole square by $(n - 1)$

We have Z by root of chi square into $(n - 1)$ by $n - 1$ which is equal to $\sqrt{n} \text{ into } (\bar{d} - \theta) \text{ by } \sigma_d$ divided by square root of $(n - 1) \text{ into } s_d^2 \text{ by } \sigma_d^2$ into $(n - 1)$ which is equal to $\sqrt{n} (\bar{d} - \theta) \text{ by } s_d$ which follows t alpha with $(n - 1)$ degrees of freedom

We can always find two quantities minus $t_{\alpha}(n - 1)$ and $t_{\alpha}(n - 1)$ from students t -distribution table such that

Probability of minus $t_{\alpha}(n - 1)$ less than or equal to t less than or equal to plus $t_{\alpha}(n - 1)$ is equal to $1 - \alpha$

Probability of minus $t_{\alpha}(n - 1)$ less than or equal to $\sqrt{n} (\bar{d} - \theta) \text{ by } s_d$ less than or equal to plus $t_{\alpha}(n - 1)$ is equal to $1 - \alpha$

Probability of minus $t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}$ less than or equal to $\bar{d} - \theta$ less than or equal to plus $t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}$ is equal to $1 - \alpha$

Probability of $\bar{d} - t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}$ less than or equal to θ less than or equal to $\bar{d} + t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}$ is equal to $1 - \alpha$

Therefore, hundred into $(1 - \alpha)$ percent confidence interval for the population mean θ in case of correlated variables is given by

$[\bar{d} - t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}, \bar{d} + t_{\alpha}(n - 1) \text{ into } s_d \text{ by } \sqrt{n}]$, where s_d^2 is equal to summation $(d_i - \bar{d})^2$ whole square by $n - 1$ and \bar{d} is equal to summation d_i by n

5. Steps to Construct CI's for the Difference between Means

Steps to construct confidence interval for the difference between means

- Identify a sample statistic. Use the difference between sample means to estimate the difference between population means
- Select a confidence level. Often, researchers choose ninety percent, ninety-five percent or ninety nine percent confidence levels, but any percentage can be used
- Find the margin of error which is equal to the table or critical value into standard error of $(\bar{x} - \bar{y})$ and then use the derived rule to obtain the confidence interval

When the sample size is large, you can use a t score or a z score for the critical value. Since it does not require computing degrees of freedom, the z score is a little easier. When the sample sizes are small (less than thirty), use a t score for the critical value.

Here's a summary of our learning in this session, where we understood:

- The derivation of confidence limits for the difference of population means when the variances are known and unknown
- The derivation of confidence limits for the population mean in case of correlated variables