# STATISTICS FOR DATA SCIENCE
# HYPOTHESIS and INFERENCE

**Dr. Deepa Nair**

Department of Science and Humanities

# STATISTICS FOR DATA SCIENCE

**UNIT-4      HYPOTHESIS and INFERENCE**
**Session-7**
**Large - Sample tests for Difference between two means**

**Dr. Deepa Nair**
Department of Science and Humanities

**Example**

Tennis Elbow

**To operate or not to operate?**

Mean % of the people cured by Surgery

$$\mu_X = 72, \sigma_X = 8, n_X = 32$$

Mean % of the people cured by Physiotherapy
$$\mu_Y = 75, \sigma_Y = 6, n_Y = 36$$

**Example:**

Can we conduct a Hypothesis test ?

**Can we say that Surgery is better than Physiotherapy ?**

- To determine whether the means of two populations are equal.

- The data will consist of two samples, one from each population.

- Let $X_1, \ldots X_{n_X}$ $and$ $Y_1, \ldots . Y_{n_Y} (n_X > 30, n_Y > 30)$ be large sample from a population with means $\mu_X and \mu_Y$, and standard deviations $\sigma_X$ $and$ $\sigma_Y$.

- We will compute the difference of the sample means.

$$H_0: \mu_X - \mu_Y = \Delta_0,$$

$$H_0: \mu_X - \mu_Y > \Delta_0,$$

$$H_0: \mu_X - \mu_Y < \Delta_0$$

- **Compute the $z$-score,**

$$z = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{\sqrt{\sigma^2_X / n_X + \sigma^2_Y / n_Y}}$$

Type equation here.

- $\bar{X} - \bar{Y} \sim N(\Delta_0, \sigma^2_X / n_X + \sigma^2_Y / n_Y)$

- **If $\sigma_X$ and $\sigma_Y$ are unknown they may be approximated with $s_X$ and $s_Y$ respectively.**

- Compute the $P$-value

- The $P$-value is an area under the normal curve, which depends on the alternate hypothesis as shown in the table:

| Alternate Hypothesis | *P*-value |
|---|---|
| $H_1: \mu_X - \mu_Y > \Delta_0$ | *Area to the right of z* |
| $H_1: \mu_X - \mu_Y < \Delta_0$ | *Area to the left of z* |
| $H_1: \mu_X - \mu_Y \neq \Delta_0$ | *Sum of the areas in the tails cut off by z and* $-z$ |

**Example:**

The article "Effect of Welding Procedure on Flux Cored Steel Wire Deposits" (N. Ramini de Rissone, I. de S. Bott, et al., *Science and Technology of Welding and Joining*, 2003:113–122) compares properties of welds made using carbon dioxide as a shielding gas with those of welds made using a mixture of argon and carbon dioxide.

**Example:**

- One property studied was the diameter of inclusions, which are particles embedded in the weld.

- A sample of 544 inclusions in welds made using argon shielding averaged 0.37 $\mu$m in diameter, with a standard deviation of 0.25 $\mu$m.

- A sample of 581 inclusions in welds made using carbon dioxide shielding averaged 0.40 $\mu$m in diameter, with a standard deviation of 0.26 $\mu$m. (Standard deviations were estimated from a graph.)

- Can you conclude that the mean diameters of inclusions differ between the two shielding gases?

**Solution:**

$$H_0: \mu_X - \mu_Y = 0$$

$$H_1: \mu_X - \mu_Y \neq 0$$

$$\overline{X} = 0.37 \ , \overline{Y} = 0.40$$

$$s_X = 0.25, S_Y = 0.26, n_X = 544, n_Y = 581$$

$$\overline{X} - \overline{Y} \sim N(0. 0.01521^2)$$

Solution:

$$z = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{\sqrt{\sigma^2{}_X/n_X + \sigma^2{}_Y/n_Y}} = -1.97$$

Solution:

- This is a two-tailed test, and the $P$-value is $0.0488$.

- A follower of the $5\%$ rule would reject the null hypothesis.

- It is certainly reasonable to be skeptical about the truth of $H_0$

**Example:**

- In a random sample of 100 tube lights produced by company A, the mean lifetime (mlt) of tube light is 1190 hours with standard deviation of 90 hours.

- Also in a random sample of 75 tube lights from company B the mean lifetime is 1230 hours with standard deviation of 120 hours.

- Is there a difference between the mean lifetime of the two brands of tube lights at a significance level of 0.05 ?

**Solution:**

- Let $X_A$ , $X_B$ denote the lifetime(in hours) of tube lights produced by company $A$ and B respectively.

- It is given that the mean lifetime of tube lights of company $A$ is $\overline{X_A}$=1190, standard deviation for tube lights of $A$ is $s_A = 90$.

- Similarly $\overline{X_B}$ =1230, $s_B$= 120, $n_A =$ sample size of tube lights from $A$= 100, $n_B =$ sample size from B = 75

**Solution:**

$H_0: \mu_A - \mu_B = 0$ i.e., no difference.

Alternate hypothesis: $H_1: \mu_A - \mu_B \neq 0$ i.e., there is difference.

$$\mu_{\overline{X}_A - \overline{X}_B = \Delta_0 = 0}$$

$$\sigma_{\overline{X}_A - \overline{X}_B} = \sqrt{\sigma_{\overline{X}_A}^2 + \sigma_{\overline{X}_B}^2} = \sqrt{\frac{S_A^2}{n_A} + \frac{S_B^2}{n_B}}$$

$$= \sqrt{\frac{(90)^2}{100} + \frac{(120)^2}{75}} \quad = 16.5227$$

**Solution:**

**Test statistic:**

$$z = \frac{(\bar{X}_A - \bar{X}_B) - \Delta_0}{\sqrt{\frac{S_A^2}{n_A} + \frac{S_B^2}{n_B}}} = \frac{1190 - 1230}{16.5227} = -2.421$$

For $\alpha = 0.05$.

Reject N.H. since $P - Value \, 0.0156 < 0.05$

i.e., yes, there is difference between the mean lifetimes of the tube lights produced by

$A$ and $B$.

**Solution:**

For $\alpha = 0.01$

   Accept N.H. since Reject N.H. since $P - Value\ 0.0156 >$ $0.05$

So there is  no difference between $\overline{X}_A$ and $\overline{X}_B$.

**Example:**

- To test the effects a new pesticide on rice production, a farm land was divided into 60 units of equal areas, all portions having identical qualities as to soil, exposure to sunlight etc.
- The new pesticide is applied to 30 units while old pesticide to the remaining 30.
- Is there reason to believe that the new pesticide is better than the old pesticide if the mean number of kg. of rice harvested/unit using new pesticide (N.P)is 496.31with s.d of 17.18 kgs while for old pesticide (O.P) is 485.41kgs and 14.73kgs.Test at a level of significance (a) $\alpha = 0.005$

**Example:**

$$H_0 : \mu_X - \mu_Y \leq 0$$

$H_1 : \mu_X - \mu_Y > 0$ i.e., new pesticide is superior to (better than) old pesticide.

$$\bar{X} = 496.31, \bar{Y} = 485.41, s_X = 17.18, s_Y = 14.73, n_X = 30, n_Y = 30$$

Test statistic is

$$z = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{\sqrt{\sigma^2{}_X/n_X + \sigma^2{}_Y/n_{\bar{Y}}}} = \frac{(496.31 - 485.41) - 0}{\sqrt{\frac{(17.18)^2}{30} + \frac{(14.73)^2}{30}}} = 2.63814$$

**Example:**

**Decision**

$$\alpha = 0.05$$

Reject N.H. since $P - \text{Value} \, 0.041 < 0.05$ i.e., accept A.H.

or new pesticide is superior to old pesticide.

**Dr.  Deepa Nair**

Department of Science and Humanities

**deepanair@pes.edu**