

Lezione 14 - Hypothesis Testing on two populations.

In the last two class we built a procedure to test a hypothesis on a value of a population.

We want to see if two samples share similar characteristic.

We will only see how to compare the mean of two populations.

First Case:

with known σ_i^2

Let's assume the two samples are Gaussian, and they are independent

Let $X_{11}, \dots, X_{n_1} \sim N(\mu_1, \sigma_1^2)$ of size n_1

let $X_{12}, \dots, X_{n_2} \sim N(\mu_2, \sigma_2^2)$ of size n_2 (usually $n_1 \neq n_2$)

$H_0: \mu_1 - \mu_2 = \Delta_0$ vs. $H_1: \mu_1 - \mu_2 \neq \Delta_0$. which Δ_0 is fixed real.

How do we build a rejection region?

$\underbrace{\bar{X}_1 - \bar{X}_2}$ is an unbiased estimator of $\mu_1 - \mu_2$.

sample means
of i-th

$$\text{Var}(\bar{X}_1 - \bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

(Unbiased since

$$E(\bar{X}_1 - \bar{X}_2) = \bar{X}_1 - \bar{X}_2$$

already unbiased

$$\rightarrow \text{Test Statistic} \rightarrow \frac{\bar{X}_1 - \bar{X}_2 - \Delta_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \stackrel{H_0}{\sim} N(0, 1)$$

Reject H_0 at level α if $\frac{|\bar{x}_1 - \bar{x}_2 - \Delta_0|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} > z_{1-\frac{\alpha}{2}}$

$$p\text{-value} = 2(1 - \phi(|z_0|))$$

What if the hypothesis is different and one-sided?

The rejection region will change (on slides)

CI for $\mu_1 - \mu_2$ of level $1 - \alpha$ if σ_1^2 and σ_2^2 known:

$$\left(\bar{x}_1 - \bar{x}_2 - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

Exercises:

Compare the mean tensile strength of two steel alloys

alloy A:

engineers want to determine if there is a statistically significant difference between the two means at 5% significance level.

alloy B:

$$n_A = 10$$

x_{1A}, \dots, x_{10A} = tensile strength [MPa] in specimen i for alloy A.

$$\bar{x}_A = 620 \text{ MPa}$$

$$n_B = 12, x_{1B}, \dots, x_{12B} \rightarrow \bar{x}_B = 605 \text{ MPa}$$

$$H_0: \mu_A = \mu_B \quad H_1: \mu_A \neq \mu_B \iff H_0: \mu_A - \mu_B = 0 \quad H_1: \mu_A - \mu_B \neq 0$$

$$X_{1n}, \dots, X_{10A} \stackrel{iid}{\sim} N(\mu_A, \sigma_A^2) \quad \sigma_A = 15$$

$$X_{1B}, \dots, X_{12B} \stackrel{iid}{\sim} N(\mu_B, \sigma_B^2) \quad \sigma_B = 20$$

Reject H_0

At significance level $\alpha = 5\%$.

$$\iff |z_0| = \frac{|\bar{x}_A - \bar{x}_B - 0|}{\sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}} \\ 2.01$$

$$\geq z_{1-\frac{\alpha}{2}} = z_{.975} = 1.96$$

Yes, the equality is satisfied, so we reject H_0 .

This is a strong conclusion, since if we were wrong, the probability of doing so is no greater than 5%.

p-value = 4,4% \rightarrow moderate evidence of H_1 ,

\hookrightarrow we can write the table for this on the form below.

If $\alpha = 1\% \Rightarrow$ we cannot reject H_0 .

If $\alpha = 10\% \Rightarrow$ we reject H_0 since $\alpha \geq p\text{-value}$.

Second Case:

Two Gaussian, with unknown but equal variance

$$\sigma_1^2 = \sigma_2^2 = \sigma^2 \text{ and } \sigma^2 \text{ is unknown}$$

$H_0: \mu_1 - \mu_2 = \Delta_0$ vs. $H_1: \mu_1 - \mu_2 \neq \Delta_0$, where Δ_0 is a fixed real value.

$\bar{X}_1 - \bar{X}_2$ is an unbiased estimator for $\mu_1 - \mu_2$

$$\text{Var}(\bar{X}_1 - \bar{X}_2) = \sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

\hookrightarrow but unknown

So to estimate σ^2 , we use an estimator called the *pooled estimator*:

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

Reject H_0 at level α if $|t_0| := \frac{|x_1 - x_2 - \Delta_0|}{\sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} > t_{1-\frac{\alpha}{2}} (n_1 + n_2 - 2)$

$$p\text{-value} = 2 \left(1 - F_{t, n_1 + n_2 - 2}(|t_0|) \right)$$

Third Case:

What if the hypothesis of equality variances is NOT reasonable?

We are able to give an answer for LARGE samples only.

Suppose $\sigma_1^2 \neq \sigma_2^2$ are unknown, if n_1 and n_2 are large ($n_1, n_2 \geq 50$)

(Slides)

Example:

Plants grown in 2 adjacent labs.

The growth time of two plant varieties are known,

days

$X \rightarrow$ growth time of plant \leq

$Y \rightarrow$ " " " "

$$n_1 = 18$$

$$\sum_{i=1}^{18} x_i = 137$$

2

$$\sum_{i=1}^{18} x_i^2 = 1120$$

$$n_2 = 14$$

$$\sum_{j=1}^{14} y_j = 82$$

$$\sum_{j=1}^{14} y_j^2 = 610$$

$$\text{Independent} \begin{cases} X_1, \dots, X_{18} \stackrel{iid}{\sim} N(\mu_X, \sigma_X^2) \\ Y_1, \dots, Y_{14} \stackrel{iid}{\sim} N(\mu_Y, \sigma_Y^2) \end{cases} \quad \sigma_X^2 = \sigma_Y^2 = \sigma^2$$

a) $\hat{\mu}_x = \bar{x} = 7.611$ $\hat{\mu}_y = \bar{y} = 5.8571$

$$S_x^2 = 4.5458$$

$S_y^2 = 9.9780 \rightarrow$ These are different because

$$\text{b) } S^2 = S_p^2 = \frac{(18-1)S_x^2 + (14-1)S_y^2}{18+14-2}$$

$$\simeq 6.8997$$

it's an observed sample that has been subjected to random error and it's possible that this is different from the real value.

c) Verify the hypothesis that the two mean times differ by exactly one day vs. the hypothesis that plant y grows slower than plant x by at least one day

$$H_0: \mu_x = \mu_y + 1 \quad H_1: \mu_x > \mu_y + 1$$

$$H_0: \mu_x - \mu_y = 1 = \Delta_0 \quad H_1: \mu_x - \mu_y > 1$$

Reject H_0

$$\text{at } \alpha \iff t_0 = \frac{\bar{x} - \bar{y} - 1}{\sqrt{s_p^2 \left(\frac{1}{18} + \frac{1}{14} \right)}} > t_{1-\alpha, 18+14-2}$$

Verify that $\Rightarrow H_0$

From that, the data shows, the violence is $\Rightarrow H_1$

} Read slowly.

$$p\text{-value} = 1 - F_{t(30)}(t_0) = 21.34\%$$

\hookrightarrow computed with R,
not the tables.

There is no evidence against H_0 .

Example:

$$X_{11}, \dots, X_{1n_1} \stackrel{iid}{\sim} F_1 \quad \mu_1 \rightarrow \bar{X}_1 \quad n_1, n_2 = \text{large}$$

$$X_{21}, \dots, X_{2n_2} \stackrel{iid}{\sim} F_2 \quad \mu_2 \rightarrow \bar{X}_2$$

$$H_0: \mu_1 = \mu_2 \text{ vs. } H_1: \mu_1 \neq \mu_2$$

$$\mu_1 - \mu_2 = \Delta_0$$

$$\bar{X}_1 \stackrel{\text{approx}}{\sim} N(\mu_1, \frac{\sigma_1^2}{n_1})$$

$$\bar{X}_2 \stackrel{\text{approx}}{\sim} N(\mu_2, \frac{\sigma_2^2}{n_2})$$

Two populations proportionate (Slides, since she went quick)

$p_1 = \text{number of successes in sample 1 / total attempts}$

$p_2 = \text{u 2 / total}$

Two Gaussian paired populations

we have not seen

2 populations collected in PAIRS and:

$(X_{1,1}, X_{2,1}), (X_{1,2}, X_{2,2}), \dots, (X_{1,n}, X_{2,n})$ iid bivariate Gaussian distribution

E.g. energy consumption before and after lubricant.

$$\mu_x > \mu_y$$

We want to see if our average the consumption is different.
We can't assume they are independent since they measure the same thing, just after one change.

$$X_{1,1}, X_{1,2}, \dots, X_{1,n} \stackrel{\text{iid}}{\sim} N(\mu_1, \cdot)$$

$$X_{2,1}, X_{2,2}, \dots, X_{2,n} \stackrel{\text{iid}}{\sim} N(\mu_2, \cdot)$$

Paired t-test

Define: $D_j = X_{1,j} - X_{2,j}$ $\xrightarrow{\text{if } \mu_1 - \mu_2}$
 $\Rightarrow D_1, \dots, D_n \stackrel{iid}{\sim} N(\mu_D, \sigma_D^2)$
↳ can't show now \hookrightarrow unknown.

We can use a t-test.

Paired vs. Unpaired Comparison

(Slide)

Example in slides