

P8130: Biostatistical Methods I

Lecture 7: Methods of Inference for Two-Means

Cody Chiuzan, PhD

Department of Biostatistics

Mailman School of Public Health (MSPH)

Lecture 6: Recap

- Sampling Distribution
- Central Limit Theorem (additional simulations in Recitation 2)
- Estimation/Confidence Interval for One-Sample Mean
- Hypothesis Testing for One-Sample Mean

Lecture 7: Outline

- Inference for two-samples means equal and unequal variances:
 - Confidence interval
 - Hypothesis testing
 - Test for equality of variances

Two-Sample Inference

- Thus far we were only interested in one parameter from a single population (one sample)
- What if we want to compare parameters in more than one population?
- In a two-sample setting, we want to make inferences about the parameters of two populations (both values unknown)
 - Paired vs Independent \leftrightarrow Longitudinal vs Cross-sectional study

The Paired t-test

Example: You are testing a new blood pressure medication on a sample of n randomly selected patients that came for a regular consult at CUMC in 2016. Systolic blood pressure (SBP) measurements were recorded at the first visit and 6 months later for all subjects.

Data structure: paired samples

Visit 1	Visit 2	Differences
X_{11}	X_{12}	$d_1 = X_{11} - X_{12}$
X_{21}	X_{22}	$d_2 = X_{21} - X_{22}$
\vdots	\vdots	\vdots
X_{n1}	X_{n2}	$d_n = X_{n1} - X_{n2}$

The Paired t-test

Assumptions:

SBP measurements are normally distributed with mean μ_1 and μ_2 , for visit 1 and 2, respectively.

It follows that the differences $d_i \sim N(\Delta, \sigma_d^2), i = 1, 2, \dots, n$.

The problem reduces to *one-sample t-test* based on the differences d_i (variance unknown), where:

Visit 1	Visit 2	Differences
X_{11}	$\rightarrow X_{12}$	$d_1 = X_{11} - X_{12}$
X_{21}	$\rightarrow X_{22}$	$d_2 = X_{21} - X_{22}$
	\vdots	\vdots
X_{n1}	$\rightarrow X_{n2}$	$d_n = X_{n1} - X_{n2}$

$$\bar{d} = \sum_{i=1}^n d_i / n \text{ and } s_d = \sqrt{\sum_{i=1}^n (d_i - \bar{d})^2 / (n - 1)}, \text{ n represents \# of pairs}$$

Two-Sided Paired t-test

Test for the mean of the differences with unknown variance

$$H_0: \mu_1 - \mu_2 = 0 \text{ or } \Delta = 0$$

vs.

$$H_1: \mu_1 - \mu_2 \neq 0 \text{ or } \Delta \neq 0$$

With significance level α pre-specified, compute the test statistic:

$$t = \frac{\bar{d} - 0}{s_d / \sqrt{n}}$$

Reject H_0 : if $|t| > t_{n-1, 1-\alpha/2}$

Fail to reject H_0 : if $|t| \leq t_{n-1, 1-\alpha/2}$

$t_{n-1, 1-\alpha/2}$ is called the critical value and it can be found in tables or calculated using software.

Confidence Interval for Δ

A $100(1 - \alpha)\%$ confidence interval for the true mean difference (Δ) for two paired samples is given by:

$$\bar{d} - t_{n-1, 1-\alpha/2} \frac{s_d}{\sqrt{n}} \leq \Delta \leq \bar{d} + t_{n-1, 1-\alpha/2} \frac{s_d}{\sqrt{n}}$$

Where:

\bar{d} is the point estimate of the mean difference

$\frac{s_d}{\sqrt{n}}$ is the estimated standard error of the differences

$t_{n-1, 1-\alpha/2}$ is the percentile of the t -distribution with $(n-1)$ degrees of freedom

The Paired t-test: Example

Test if a new diet has an effect on lowering the cholesterol levels. Use the data below recorded from 12 randomly selected subjects (time interval: 3mo).

Subject	Before (X_1)	After (X_2)	Difference (After-Before)
1	201	200	-1
2	231	236	5
3	221	216	-5
4	260	233	-27
5	228	224	-4
6	237	216	-21
7	326	296	-30
8	235	195	-40
9	240	207	-33
10	267	247	-20
11	284	210	-74
12	201	209	8

The Paired t-test: Example

Assume the observed differences constitute a random sample from a normally distributed population of differences.

In class practice.

Two-Sample t-test for Independent Samples

- Consider that our two samples are independent (there is no relation / correlation between the data points from the two groups) and that they are normally distributed:

$$X_1 \sim N(\mu_1, \sigma_1^2) \text{ and } X_2 \sim N(\mu_2, \sigma_2^2)$$

We can assume that the underlying variances of the two samples are equal:

$$\sigma_1^2 = \sigma_2^2 = \sigma^2$$

We can assume that the underlying variances of the two samples are unequal:

$$\sigma_1^2 \neq \sigma_2^2$$

Two-Sample Independent t-test: Equal Variances

- If we know the two population variances, then:

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

- Because $\sigma_1^2 = \sigma_2^2 = \sigma^2$, it follows that:

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}\right)$$

- Thus, testing $H_0: \mu_1 = \mu_2$ can use the test statistic:

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

Two-Sample Independent t-test: Equal Variances

- However, most of the times we do not know the common variance and we have to find one estimator using the sample variances s_1^2 and s_2^2 .
- The pooled estimate of the variance from two independent samples is given by:

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \quad (*)$$

- Explain the denominator.

Two-Sample Independent t-test: Equal Variances

$H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ (two-sided)

With significance level α pre-specified, compute the test statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}, \text{ where } s \text{ is given in (*)}.$$

Reject H_0 : if $|t| > t_{n_1+n_2-2, 1-\alpha/2}$

Fail to reject H_0 : if $|t| \leq t_{n_1+n_2-2, 1-\alpha/2}$

P-value = $2 \times P(t_{n_1+n_2-2} < t)$, if $t < 0$
= $2 \times P(t_{n_1+n_2-2} \geq t)$, if $t \geq 0$

Two-Sample Independent t-test: Unequal Variances

- If we know the two population variances, then:

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

- For testing $H_0: \mu_1 = \mu_2$ we can use the z-test statistic:

$$z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)}}$$

- Or (again) use s_1^2 and s_2^2 to estimate the two population variances (if these values are unknown).

Two-Sample Independent t-test: Unequal Variances

$H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ (two-sided)

With significance level α pre-specified, compute the test statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)}}$$

Compute the approximating degrees of freedom d' , and round down to the nearest integer d'' :

$$d' = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\left(\frac{s_1^2}{n_1}\right)^2 / (n_1 - 1) + \left(\frac{s_2^2}{n_2}\right)^2 / (n_2 - 1)} \quad (**)$$

Two-Sample Independent t-test: Unequal Variances

$H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ (two-sided)

With significance level α pre-specified, compute the test statistic:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)}}$$

Reject H_0 : if $|t| > t_{d'', 1-\alpha/2}$

Fail to reject H_0 : if $|t| \leq t_{d'', 1-\alpha/2}$

Where d'' represents the nearest integer of d' (degrees of freedom) (**).

Confidence Intervals for Two Independent Samples

Two-sided $100(1 - \alpha)\%$ confidence interval for:

- Two-independent samples (equal variance): $\mu_1 - \mu_2$ ($\sigma_1^2 = \sigma_2^2 = \sigma^2$):

$$(\bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2, 1-\alpha/2} s \sqrt{1/n_1 + 1/n_2}, \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2, 1-\alpha/2} s \sqrt{1/n_1 + 1/n_2})$$

- Two-independent samples (unequal variance): $\mu_1 - \mu_2$ ($\sigma_1^2 \neq \sigma_2^2$):

$$(\bar{X}_1 - \bar{X}_2 - t_{d'', 1-\alpha/2} \sqrt{s_1^2/n_1 + s_2^2/n_2}, \bar{X}_1 - \bar{X}_2 + t_{d'', 1-\alpha/2} \sqrt{s_1^2/n_1 + s_2^2/n_2})$$

Where d'' represents the nearest integer of d' (degrees of freedom) (**).

Test for Equality of Variances

Assume two independent random samples from $N(\mu_1, \sigma_1^2)$ and $N(\mu_2, \sigma_2^2)$.

Testing the equality of variances implies testing the hypotheses:

$$H_0: \sigma_1^2 = \sigma_2^2 \text{ vs } H_1: \sigma_1^2 \neq \sigma_2^2$$

The significance test is based on the relative magnitudes of the sample variances. (Why?)

With significance level α pre-specified, compute the test statistic:

$$F = \frac{s_1^2}{s_2^2} \sim F_{n_1-1, n_2-1}$$

The test statistic follows an F distribution with $n_1 - 1$ and $n_2 - 1$ degrees of freedom.

Test for Equality of Variances

Testing the hypotheses:

$$H_0: \sigma_1^2 = \sigma_2^2 \text{ vs } H_1: \sigma_1^2 \neq \sigma_2^2$$

With significance level α pre-specified, compute the test statistic:

$$F = \frac{s_1^2}{s_2^2} \sim F_{n_1-1, n_2-1}$$

Reject H_0 : if $F > F_{n_1-1, n_2-1, 1-\alpha/2}$ or $F < F_{n_1-1, n_2-1, \alpha/2}$

Fail to reject H_0 : if $F_{n_1-1, n_2-1, \alpha/2} \leq F \leq F_{n_1-1, n_2-1, 1-\alpha/2}$

$$\begin{aligned} \text{P-value} &= 2 \times P(F_{n_1-1, n_2-1} > F), \text{ if } F \geq 1 \\ &= 2 \times P(F_{n_1-1, n_2-1} < F), \text{ if } F < 1 \end{aligned}$$

Two-Samples Independent t-test

Example: To assess the impact of oral contraceptive use on bone mineral density (BMD), researchers in Canada carried out a study comparing BMD for women who had used oral contraceptives (OC) for at least 3 months to BMD for women who had never used oral contraceptives (OC). Summary values for BMD are given below:

	n	\bar{X}	s
Non OC users	10	1.08	0.16
OC users	10	1.00	0.14

- Assuming that the BMD is normally distributed in both OC and non-OC users, is there enough evidence to conclude that BMD levels differ between the two groups?
- Construct a 95% CI for the population mean difference.

In class practice.

Readings

Rosner, *Fundamentals of Biostatistics*, Chapter 8

- Sections: 8.1 – 8.7