

# Two Sample $t$ -Test

Grinnell College

November 14, 2025

Suppose that we have a diagnostic test for an infectious disease which has a Type I error rate of 5% and a Type II error rate of 1%. Answer the following:

1. The null hypothesis for any diagnostic test is that the individual does not have the disease. What constitutes a Type I and Type II error?
2. Suppose we use it to test for the disease on a population of 1,000 people where 40% have the disease. Construct a table showing the number of correct and incorrect conclusions based on the truth of  $H_0$
3. Of individuals with a positive test, what percentage actually had the disease?
4. For testing for an infectious disease, is Type I or Type II error more important to control?

# The process is the same

What we did before, we will do today:

1. Construct a null hypothesis,  $H_0$
2. Collect data and compute our statistic (i.e.,  $\bar{x}$ )
3. Evaluate that statistic in the context of a null distribution, i.e.,

$$t = \frac{\bar{X} - \mu_0}{\hat{\sigma}/\sqrt{n}}$$

4. Reject or fail to reject hypothesis
  - ▶ Type I errors
  - ▶ Type II errors

# Group Differences

Often in statistical inference, we are interested in investigating the *difference* between two or more groups

For example, we may have two groups,  $A$  and  $B$ , with a mean value for each group,  $\mu_A$  and  $\mu_B$

Expressed in our null hypothesis, this equates to

$$H_0 : \mu_A = \mu_B \quad \text{or} \quad H_0 : \mu_A - \mu_B = 0$$

# Two-sampled t-test

Just as in the univariate case for testing the mean, we can use a  $t$ -test to evaluate the difference in means between two groups

There are a number of various assumptions about our data, all resulting in slightly different tests (degrees of freedom and standard error):

1. Independent, groups same size and have same variance
2. Independent, groups have unequal sizes and similar variance
3. Independent, groups have different sizes and different variances
4. Paired testing

In general, we will concern ourselves with (3) and (4)

# Two-sampled t-test

Just as in the univariate case for testing the mean, we can use a  $t$ -test to evaluate the difference in means between two groups

There are a number of various assumptions about our data, all resulting in slightly different tests (degrees of freedom and standard error):

1. ~~Independent, groups same size and have same variance~~
2. ~~Independent, groups have unequal sizes and similar variance~~
3. Independent, groups have different sizes and different variances
4. Paired testing

In general, we will concern ourselves with (3) and (4)

## Example

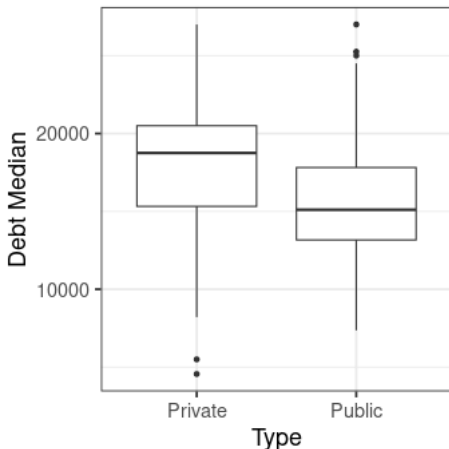
Consider our college data, where we might investigate the differences in median debt upon graduate for public and private schools

### ▶ Private Schools

- ▶  $\bar{x}_1 = 18,028$
- ▶  $\hat{\sigma}_1 = 3,995$
- ▶  $n_1 = 647$

### ▶ Public Schools

- ▶  $\bar{x}_2 = 15,627$
- ▶  $\hat{\sigma}_2 = 3,111$
- ▶  $n_2 = 559$



## t-test, Independent samples, heterogeneous groups

Our  $t$ -statistic takes the form

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}}$$

This  $t$ -statistic only approximately follows a  $t$ -distribution, making the calculation of its degrees of freedom non-trivial, usually approximated using  $n_1 + n_2 - 2$  (or with software)

Otherwise, the process for constructing confidence intervals or testing hypotheses is exactly the same



## Example

Again, we will use R to compute this, utilizing a special “formula” syntax when using data.frames (will cover in lab)

```
1 > t.test(Debt_median ~ Private, college)
2
3   Welch Two Sample t-test
4
5 data:  Debt_median by Private
6 t = 11.2, df = 1075, p-value <0.00000000000000002
7 alternative hypothesis: true difference in means between group
   Private and group Public is not equal to 0
8 95 percent confidence interval:
9   1981.0 2820.6
10 sample estimates:
11 mean in group Private
12                   18028
13 mean in group Public
14                   15627
```

## Paired t-test

The **paired t-test** or **paired difference test** is a test for assessing differences in group means where the groups consist of the same subjects with multiple observations

While it ostensibly shares many characteristics with a two-sample t-test, in practice it more closely resembles that of a one-sample test:

$$t_{\text{paired}} = \frac{\overline{X}_D - \mu_0}{\hat{\sigma}_D / \sqrt{n}}$$

where  $n$  represents the number of *unique* subjects and  $\overline{X}_D$  and  $\hat{\sigma}_D$  represent the mean and standard deviation of the *difference*, respectively

# Paired t-test

Just as with the unpaired case, our null hypothesis is typically that

$$H_0 : \mu_0 = 0$$

Paired testing between groups allows us to control for within-subject variation, effectively reducing variation and making it easier to detect a true difference (power)

This comes at a cost, however – for  $n$  subjects we are required to make  $2n$  unique observations

## Example – French Institute

Consider the results of a summer institute program sponsored by the National Endowment for the Humanities to improve language abilities in foreign language high school teachers

Twenty teachers were given a listening test of spoken French before and after the program, with a maximum score of 36. We are interested in determining the efficacy of the summer institute

## Example – French Institute

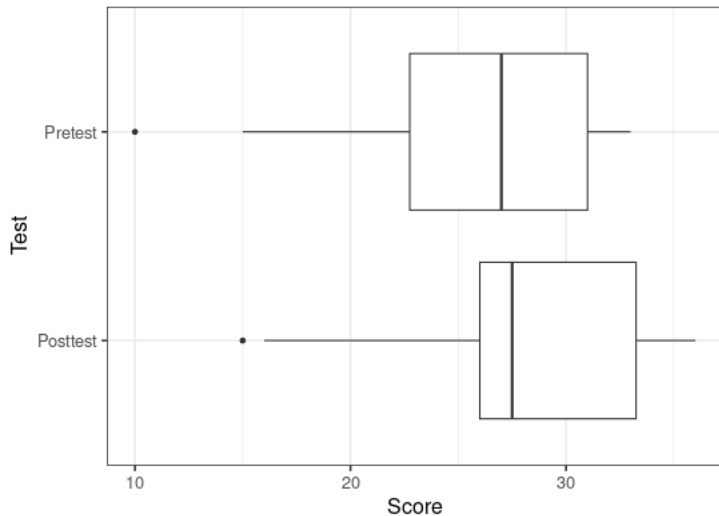
1. What is the null hypothesis for this study?
  - ▶ What would be a Type I error?
  - ▶ A Type II error?
2. How many total subjects do we have?
3. How many recorded observations do we have?

## Example – French Institute

The results of the tests are as follows:

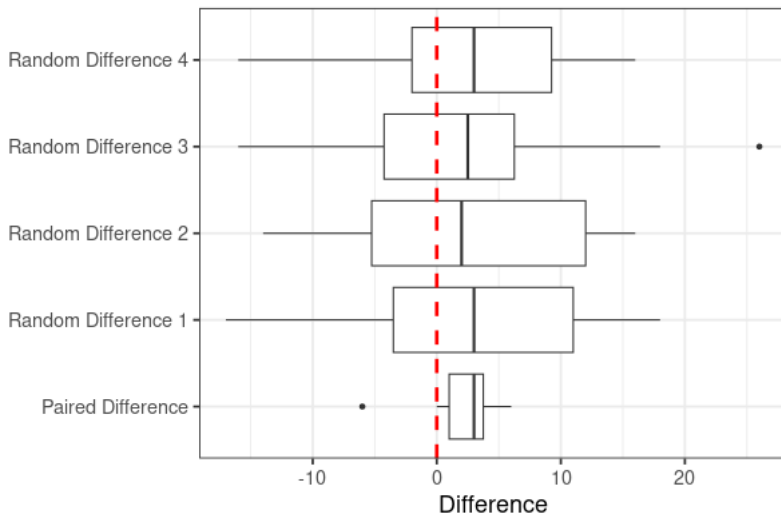
ID	Pretest	Posttest	Difference	ID	Pretest	Posttest	Difference
1	32	34	2	11	30	36	6
2	31	31	0	12	20	26	6
3	29	35	6	13	24	27	3
4	10	16	6	14	24	24	0
5	30	33	3	15	31	32	1
6	33	36	3	16	30	31	1
7	22	24	2	17	15	15	0
8	25	28	3	18	32	34	2
9	32	26	-6	19	23	26	3
10	20	26	6	20	23	26	3

## Example – French Institute



## Example – French Institute

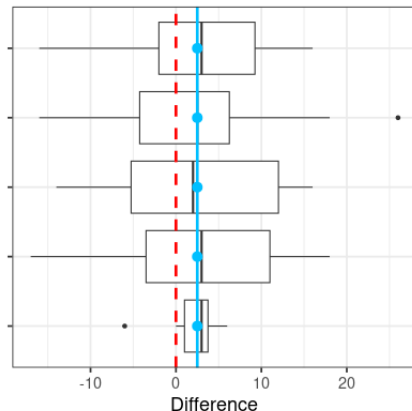
Plotted below is a boxplot of observed differences if people were randomly shuffled and repaired in each group





There are a few things to notice here:

- ▶ The mean value for each arrangement is *identical*
- ▶ The groups that were randomly assigned show far greater variability
- ▶ Less variability = more power



## Example – French Institute

Results of the *paired t-test*

```
1 > t.test(post, pre, paired = TRUE)
2
3   Paired t-test
4
5 data:  post and pre
6 t = 3.86, df = 19, p-value = 0.001
7 alternative hypothesis: true mean difference is
   not equal to 0
8 95 percent confidence interval:
9  1.1461 3.8539
10 sample estimates:
11 mean difference
12           2.5
```

## Example – French Institute

Results of the unpaired t-test, no power to find difference

```
1 > t.test(post, pre, paired = FALSE)
2
3   Welch Two Sample t-test
4
5 data:  post and pre
6 t = 1.29, df = 37.9, p-value = 0.2
7 alternative hypothesis: true difference in
   means is not equal to 0
8 95 percent confidence interval:
9  -1.424  6.424
10 sample estimates:
11 mean of x mean of y
12    28.3    25.8
```

- ▶ Hypothesis testing works nearly identically for two groups as it did with one group
- ▶ CLT applies for both difference in proportions as well as difference in group means
- ▶ Two-sample t-tests have a paired version
  1. Reduces variability
  2. Also reduces degrees of freedom
- ▶ We can use R to do most of these for us

# Multiple Comparisons

One prevalent issue in hypothesis testing is that of **multiple comparisons** whereby several hypothesis tests are conducted simultaneously

As the number of hypothesis tests conducted grows in number, so to does the probability of one of those tests being decided in error

# Multiple Comparisons

Consider conducting 2 hypothesis tests, each with a Type I error rate of 5%

For any given test, the probability of *not* making an error is

$$P(\text{No type I error}) = 0.95$$

1. What is the probability that neither test has a Type I error?
2. What is the probability that *at least* one test has a Type I error?

## Example

Suppose that I am interested in testing if there is a non-zero correlation between cost and average faculty salary in each of the 8 regions of our college dataset

Suppose further we are testing for significance at the level  $\alpha = 0.05$

	Region	$p$ -value
1	Far West	0.7667
2	Great Lakes	0.0085
3	Mid East	0.0001
4	New England	0.0061
5	Plains	0.9487
6	Rocky Mountains	0.7394
7	South East	0.0143
8	South West	0.0344

# Family-wise error rates (FWER)

For a collection of independent hypothesis tests, the **family-wise error rate (FWER)** describes the probability of making one or more Type I errors

For  $m$  independent tests with a Type I error rate of  $\alpha$ , the FWER is defined as

$$\text{FWER} = 1 - (1 - \alpha)^m$$



## Example

Suppose that I am interested in testing if there is a non-zero correlation between cost and average faculty salary in each of the 8 regions of our college dataset

If my Type I error rate for each test is 5%, what is the probability that I make at least one Type I error?

$$\begin{aligned}P(\text{At least one Type I error}) &= 1 - P(\text{Probability of no Type I errors}) \\&= 1 - (1 - 0.05)^8 \\&= 33.6\%\end{aligned}$$

That is, instead of making a Type I error 1 in 20 times, we are now making it 1 in 3 times

# FWER Correction

Just as we control the Type I error rate of a single hypothesis test with  $\alpha$ , we also have an interest in controlling the FWER

For  $m$  hypothesis tests controlled at level  $\alpha$ , the correction  $\alpha^* = \alpha/m$  is known as the **Bonferonni Adjustment**

If instead for a series of  $m$  tests we reject the null hypothesis when  $p < \alpha^*$ , we will control the FWER at level  $\alpha$

Assuming the 8 regions of our hypothesis test are independent, our Bonferonni adjustment for  $\alpha = 0.05$  should be

$$\alpha^* = 0.05/8 = 0.00625$$

where our new FWER is

$$\begin{aligned}\text{FWER} &= 1 - (1 - 0.00625)^8 \\ &= 0.04892\end{aligned}$$

Testing $p < \alpha$		
	Region	$p$ -value
1	Far West	0.7667
2	Great Lakes	0.0085
3	Mid East	0.0001
4	New England	0.0061
5	Plains	0.9487
6	Rocky Mountains	0.7394
7	South East	0.0143
8	South West	0.0344

Testing $p < \alpha^*$		
	Region	$p$ -value
1	Far West	0.7667
2	Great Lakes	0.0085
3	Mid East	0.0001
4	New England	0.0061
5	Plains	0.9487
6	Rocky Mountains	0.7394
7	South East	0.0143
8	South West	0.0344

Based on the evidence observed, we will ultimately make one of two decisions:

1. Reject  $H_0$
2. Fail to reject  $H_0$

Depending on the true state of  $H_0$ , we can be incorrect in two ways:

1. Type I Error ( $\alpha$ ):  $H_0$  is true, yet we reject anyway
2. Type II Error ( $\beta$ ):  $H_0$  is false, yet we fail to reject it

Finally, there is the issue of *multiple comparisons*

1. Family-wise error rate
2. Bonferonni correction