

Two Sample t -Test

Grinnell College

April 17, 2026

Review

We are asked to evaluate a piece of diagnostic software that is used to sort through and categorize email as being spam or not spam. Email that is not spam is sent to your inbox, while email that is spam is sent to a junk folder.

- ▶ What is the null hypothesis for incoming email?
- ▶ What is a Type I and Type II error? Which is more important here?
- ▶ Suppose that the diagnostics of the spam software filter has a Type I error rate of 2% and a Type II error rate of 20%. In a typical month, the average user will receive 1,000 emails, with approximately 5% of them being spam. *Supposing that an email is marked as spam, what is the probability that a given email was not spam? (False positive rate)*

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

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

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, μ_A and μ_B

Expressed in our null hypothesis, this equates to

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

Test of Association

It is both helpful and correct to think of a two-sample t-test as a *test of association*

If two groups (categorical variable) have different group means (quantitative variable), then knowing any individual's group membership gives us better information than guessing an overall average

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

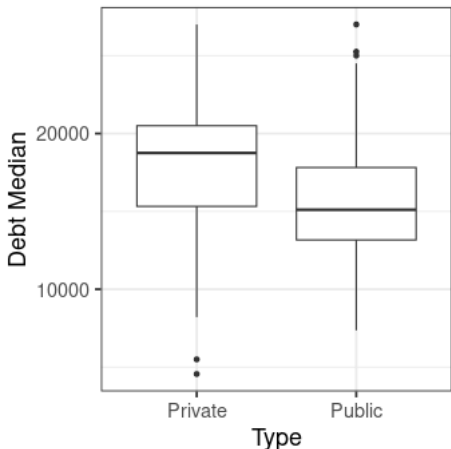
Consider our college data, where we might investigate the differences in median debt upon graduate for public and private schools

- ▶ 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:

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

where n represents the number of *unique* subjects and \bar{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

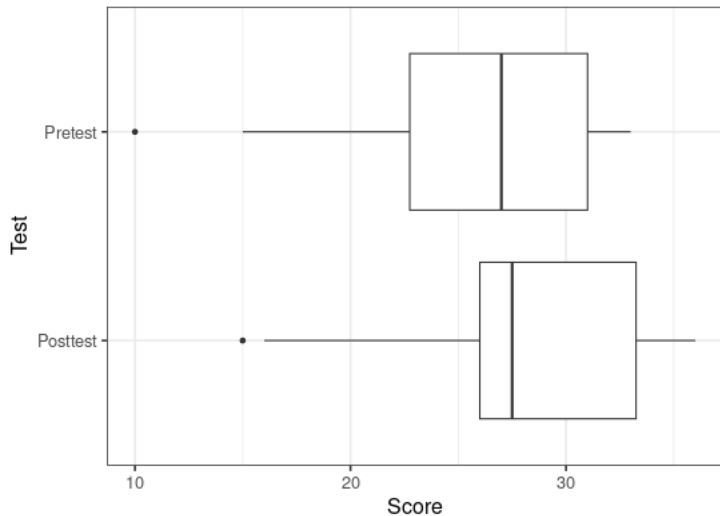
1. What is the null hypothesis for this study?
 - ▶ What would be a Type I error?
 - ▶ A Type II error?
2. How many total subjects do we have?
3. How many recorded observations do we have?

Example – French Institute

The results of the tests are as follows:

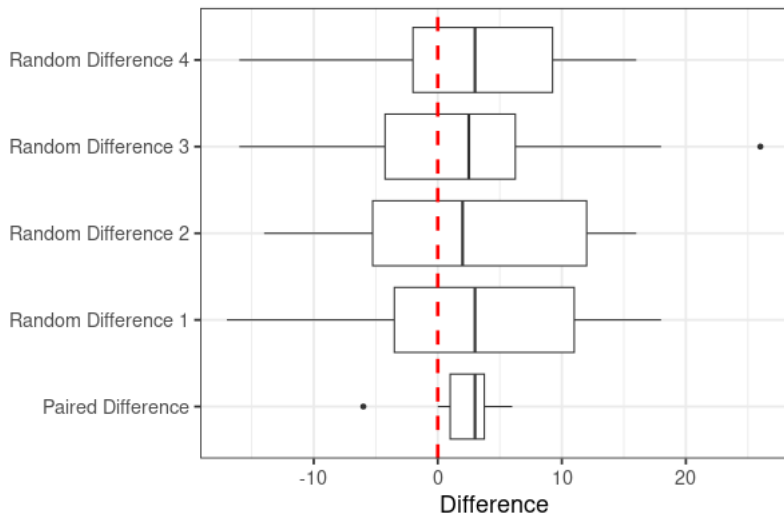
ID	Pretest	Posttest	Difference	ID	Pretest	Posttest	Difference
1	32	34	2	11	30	36	6
2	31	31	0	12	20	26	6
3	29	35	6	13	24	27	3
4	10	16	6	14	24	24	0
5	30	33	3	15	31	32	1
6	33	36	3	16	30	31	1
7	22	24	2	17	15	15	0
8	25	28	3	18	32	34	2
9	32	26	-6	19	23	26	3
10	20	26	6	20	23	26	3

Example – French Institute



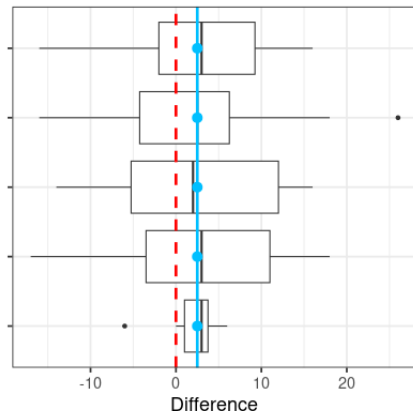
Example – French Institute

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