

Experiment Design for Computer Sciences (0AL0400)

Topic 04 - Paired Comparison

Claus Aranha

caranha@cs.tsukuba.ac.jp

University of Tsukuba, Department of Computer Sciences

Version 2021.1 (Updated May 18, 2022)

Outline

In this lecture, we study two more cases of **Null Hypothesis Statistical Testing** that are quite common in experiments on algorithms:

- Hypothesis testing on the difference between two treatments. (**Comparison Testing**)
- Hypothesis testing when there is a strong correlation factor between the observations of the two treatments. (**Paired Testing**)

Treatments?

The word "treatment" is from the medicine literature, but here we use to indicate two different things that we want to compare. It could be two algorithms, or two parameter settings, or two experimental conditions, etc.

Part I – Two Sample Testing

Last Lecture Recap

In the last lecture, we studied the **Null Hypothesis Method of Statistical Inference**.

In this method, we determine a "Null Hypothesis" and an "Alternate Hypothesis" about the mean of a population of interest. Then, using data from an experiment, we determine the likelihood that this data corresponds to the Null hypothesis or to the Alternate hypothesis.

The statistic test we studied in the last week can be used for hypothesis involving **a single sample and an estimated mean**. How should we proceed if we want to compare two means, possibly from different populations?

Comparison of two processes

The comparison between two different approaches is a very common situation in scientific research:

- The efficacy of a new drug is compared against a control group;
- The precision of a new algorithm is compared against an old one;
- Two different website design proposals are compared regarding user preference;
- etc;

How can we adapt the Hypothesis testing procedure studied in the last lecture to these situations?

The analysis of these situations involves the calculation of statistics based on data from two different samples, so we will call it **two sample testing**.

Example: Comparing Methods for cutting steel rods

We will use the following situation to illustrate the hypothesis testing method:



A critical aspects of manufacturing steel rods is cutting the bars with a precise length.

Errors when cutting the bars will cause costs for reprocessing the rods.

An engineer is interested in comparing the current cutting process with a new method that could potentially improve the performance of the system by reducing the cutting error.

Comparing cutting methods

Quiz



We have two methods for cutting steel rods (old and new), and we want to find out which one has the smallest cutting error. Consider the following questions:

- How do we calculate / measure the cutting error of one of the methods?
- What is the observation / sample necessary to estimate this value?
- What is the variable that measures the cutting error difference between the two methods?
- What is a *statistical hypothesis* that represents the question of interest for this experiment?

Pause the video! Take some time to seriously answer these questions before you continue the material!

Modeling the cutting process

What is the cutting error?



Let's look at the first two questions:

- How do we calculate / measure the cutting error of one of the methods?
- What is the observation / sample necessary to estimate this value?

Let's consider a cutting error to be the difference between the length of a rod i and the target length l : $(|x_i - l|)$.

Assuming that the cutting error is a property of the method, we can estimate the **mean cutting error** using a sample X of n rods:

$$\hat{\mu}_e \text{ estimated by } e_X = \frac{\sum |x_i - l|}{n}$$

Modeling the cutting process

Cutting Error and Hypothesis



$$\hat{\mu}_e \text{ estimated by } e_X = \frac{\sum |x_i - l|}{n}$$

Using e_X as an estimate of the cutting error, it is possible to perform statistical inference about **one of the methods**. For example:

- Is the error of method Y equal or under a required value r ?
- $H_0 : e_Y \leq r$
- $H_1 : e_Y \geq r$

We can use the technique from the last lecture to solve this problem. But if we want to compare two methods: Y_1 and Y_2 , what do we do?

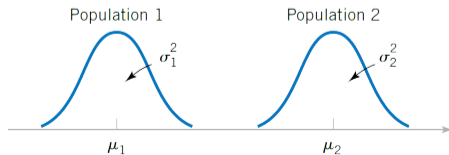
Modeling the cutting process

Comparing Errors



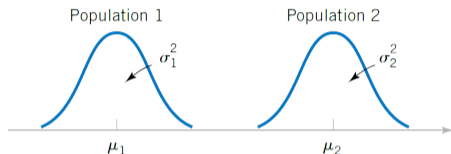
- What is the variable that measures the cutting error difference between the two methods?
- What is a *statistical hypothesis* that represents the question of interest for this experiment?

Let's consider these two questions. Remember that the error from each method is modeled as a random variable following a normal distribution.



Modeling the cutting process

Comparing Errors



The sum of two normal variables also follow a normal distribution. So, we can describe the difference between the cutting errors as the random variable $e_{\text{diff}} = (e_{\text{old}} - e_{\text{new}})$.

Because e_{diff} also follows a normal distribution, we can use the null hypothesis method to test the difference of the two methods:

- $H_0 : e_{\text{diff}} = 0$
- $H_1 : e_{\text{diff}} \neq 0$

A General Model for Comparing two Samples

Using the ideas from the previous example, let's describe a general statistical model to use when we want to test if two methods are quantitatively different.

Consider that we measure some observed value (y) taken from one of several methods ($i = 1, 2, \dots$), we understand that the value comes from some distribution with mean μ_i , at it will also have an error (ϵ) away from that mean, which is different for each observation. So we describe the j -th observation taken from the i -th method as

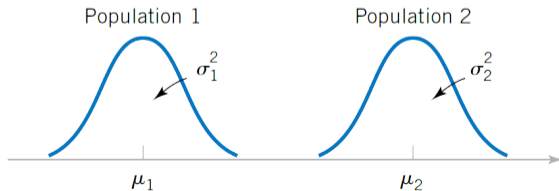
$$y_{ij} = \mu_i + \epsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, n_i \end{cases}$$

Statistical Models

Two population Model

$$y_{ij} = \mu_i + \epsilon_{ij} \begin{cases} i = 1, 2 \\ j = 1, \dots, n_i \end{cases}$$

Under this model for the observed variable (y_{ij}), we assume that the residuals ϵ_{ij} are independent and follow $\mathcal{N}(0, \sigma_i^2)$. Under these assumptions, the populations of the two samples look like this:



Comparison of two means

Null and Alternate Hypotheses

What should be the observed variable y ? The goal of this experiment is to measure if the new method produces steel rods closer to the nominal value. In this case, a possible response variable would be the **absolute error**, e.g., $y = |\ell - \ell_{nominal}|$.

Keeping in mind our statistical model, we can build the hypothesis around the **mean** of the absolute error (μ_j). In that case, we can state the null and alternate hypotheses as:

$$\begin{cases} H_0 : \mu_1 - \mu_2 = 0 \\ H_1 : \mu_1 - \mu_2 < 0 \end{cases} \quad \text{or, equivalently,} \quad \begin{cases} H_0 : \mu_1 = \mu_2 \\ H_1 : \mu_1 < \mu_2 \end{cases}$$

Comparison of two means

Calculating the statistic

Lets assume (for the moment) that the variance of the process is unknown but similar for both systems. Since it is unknown, we have to estimate the variance from the sample data. As assume $\sigma_1^2 \approx \sigma_2^2$, we can use the pooled variance estimator:

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

Based on this estimator and the stated assumptions, we calculate the T statistic:

$$T = \frac{(\bar{y}_1 - \bar{y}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t^{(n_1+n_2-2)}$$

Back to the Steel Rods Example

Calculation of the Rejection threshold

If we recall our working hypotheses for the steel rod example:

$$\begin{cases} H_0 : \mu_1 - \mu_2 = 0 \\ H_1 : \mu_1 - \mu_2 < 0 \end{cases}$$

we have that, under H_0 :

$$t_0 = \frac{(\bar{y}_1 - \bar{y}_2) - \cancel{(\mu_1 - \mu_2)}^0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{(\bar{y}_1 - \bar{y}_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t^{(n_1+n_2-2)}$$

We'll reject H_0 at the $(1 - \alpha)$ confidence level if $t_0 \leq t_{\alpha/2}^{(n_1+n_2-2)}$

Back to the Steel Rods Example

Statistic Test Parameters

Remember that we need to decide **three parameters** that will specify the statistical test:

- **Significance level:** The probability of a **Type I error**. Let's assume that the desired significance level is $\alpha = 0.05$.
- **Power:** The probability of a **Type II Error**. Let's assume that the desired sensitivity is $1 - \beta = 0.8$.
- **Meaningful difference:** What is the minimum difference between the two methods that we are interested in detecting? Let's assume 15cm .

The values for these variables depend on the needs of the specific experiment and/or application.

Calculating the Statistic

Computationally, we can perform the t-test for comparing the means of two independent populations by:

```
> y <- read.table("steelrods.csv", header = TRUE)
> t.test(y$Length.error ~ y$Process, alternative = "less",
+        mu = 0, var.equal = TRUE, conf.level = 0.95)
```

```
data:  y$Length.error by y$Process
t = -14.312, df = 32, p-value = 9.244e-16
alternative hypothesis: true difference in means is less than 0
95 percent confidence interval:
 -Inf -7.156884
sample estimates:
mean in group new mean in group old
      7.782353      15.900000
```

Comparison of two means

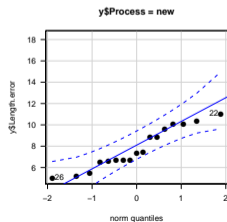
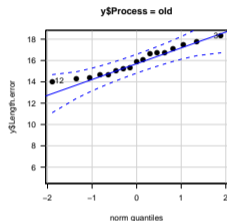
Testing the assumptions

The assumptions of the test must be verified. In this particular case:

- Normality;
- Equality of variances;
- Independence.

```
> qqPlot(y$Length.error, groups = y$Process,
         cex = 1.5, pch = 16, las = 1,
         layout = c(2, 1))
> shapiro.test(y$Length.error[y$Process == "new"])
# W = 0.92269, p-value = 0.164
> shapiro.test(y$Length.error[y$Process == "old"])
# W = 0.94971, p-value = 0.4519
```

Reminder: the t-test is quite robust to mild to moderate violations of the normality of the residuals / groups.



Comparison of two means

Testing the assumptions

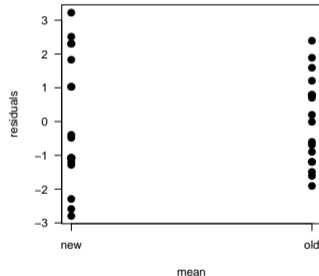
The assumptions of the test must be verified. In this particular case:

- Normality;
- Equality of variances;
- Independence;

```
> fligner.test(Length.error ~ Process, data = y)
# Fligner-Killeen:med chi-squared = 1.6837,
# df = 1, p-value = 0.1944

> residuals <- tapply(X = y$Length.error,
  INDEX = y$Process,
  FUN = function(x) {x - mean(x)})

> stripchart(x = residuals, vertical = TRUE,
  pch = 16, cex = 1.5, las = 1,
  xlab = "mean",
  ylab = "residuals")
```



Comparison of two means

Testing the assumptions

The assumptions of the test must be verified. In this particular case:

- Normality;
- Equality of variances;
- **Independence**;

As mentioned in the last class, there is no general test for the independence assumption, and it has to be guaranteed in the design phase.

One can at most test for serial autocorrelation in the residuals using Durbin-Watson's test, but this test is absolutely dependent on the ordering of the observations - very useful to detect ordering-related trends in the residuals, but not much more than that.

Comparison of two means

Unequal variances

Suppose now a more general case, in which the variances of the two populations are unknown and cannot be assumed equal.

For this cases, a modification on the t-test called *Welch's t test* is usually employed. The Welch statistic can be calculated as:

$$t_0^* = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Under the null hypothesis t_0^* is distributed approximately as a $t^{(\nu)}$ distribution, with:

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

Comparison of two means

Unequal variances

Let's illustrate the calculation of a comparison test with unequal variances in R. We will use the same data as before.¹

```
> t.test(y$Length.error ~ y$Process, alternative = "two.sided",
+       mu = 0,
+       var.equal = FALSE,           %% <- We only change this.
+       conf.level = 0.95))
Welch Two Sample t-test
data:  Length.error by Process
t = -14.312, df = 28.386, p-value = 1.645e-14
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.09278780 -0.06956515
sample estimates:
mean in group new mean in group old
 0.07782353      0.15900000
```

¹Note that this would not have been necessary, since we already checked the "equal variances" assumption.

Comparison of two means

Summary

To compare an estimator from samples of two populations that follow a normal distribution, we set our statistic and the corresponding hypotheses to be the difference of the target variables.

This technique for comparison testing is simple and extremely versatile.

Of course, there are cases where this approach does not apply. Next we will see a relatively common case where using the difference of the target variables would lead to a wrong inferential result.

Part II – Paired Testing

Paired Comparison of Two Samples

Outline

In the last part, we studied how to apply the statistical inference method using hypothesis testing to the situation where we want to compare two samples.

In this part, we will study a common special situation, where there is a strong dependency between observations in the samples.

The change in the calculation is very minor, but the results can be very different!

Examples of Paired Design

Paired Design happens when the observations in both samples have strong dependencies.

Example 1: Football shoes

Out of two brand of football shoes, you want to measure which one wears out faster.

- You make a team play two games, one with shoe A, one with shoe B. You measure the amount of wear for each player's shoes.
- You know that the Foward's shoes will wear much more than the Goalkeeper's shoes.

Example 2: Fuel Efficiency

You want to measure if a new kind of fuel is more efficient than an old one.

- You choose 10 cars, fill then with each type of fuel, make then run until they are out of fuel, and measure the distance.
- You know that different car types consume fuel at very different rates.

Computer Science Example

Comparison of Two Optimization Methods

A researcher develops a new optimization algorithm (A), and wants to compare its convergence speed against a method that represents the state-of-art (B).

The researcher believes that the proposed algorithm has a theoretical advantage on a **specific family of optimization problems**, so she selects a set of benchmark problems from that family.

Both methods are executed on the benchmark set, and the time-to-convergence is measured for each problem. The measurements are made under homogeneous conditions (same computer, same operating conditions, etc).

Computer Science Example

Consideration for Experiment Design

In this example, we are taking several problem instances, and running each of the two algorithms in all instances. Because of the expected variation in running time, we might want to run one "algorithm-instance" multiple times.

This problem has some important questions worth considering:

- What is the **estimator** that should be measured in this experiment?
- What is one **independent observation** for this experiment?
- What is the relevant sample size for the experiment?

Think about the difference between considering **individual runs** as observations and **individual problems** as observations.

Paired Experimental Design

Why is Pairing Necessary?

When we consider observations with strong dependencies (players, cars types, benchmark problems), the difference between the observations is a strong source of variation (noise) that is not related to the objective of the experiment.

This variation can, and must, be controlled in the experiment design.

An elegant solution to eliminate the influence of this nuisance parameter is the *pairing* of the measurements by problem:

- Observations are considered in pairs (A, B) for each benchmark problem;
- Hypothesis testing is done on the sample of "*differences for a benchmark*";

Paired Experimental Design

Statistical Model

Let y_{Aj} and y_{Bj} be the paired observations of the average time for methods A and B, for a problem instance j . The *paired difference* of an observation is simply $d_j = y_{Aj} - y_{Bj}$.

If we model our observations as an additive process:

$$y_{ij} = \underbrace{\mu + \tau_i}_{\mu_i} + \beta_j + \varepsilon_{ij}$$

where μ is the grand mean, τ_i is the effect of the i -th method on the mean (A or B), β_j is the effect of the j -th problem, and ε_{ij} is the model residual, then:

$$\begin{aligned} d_j &= y_{Aj} - y_{Bj} \\ &= \mu + \tau_A + \beta_j + \varepsilon_{Aj} - (\mu + \tau_B + \beta_j + \varepsilon_{Bj}) \\ &= (\cancel{\mu + \beta_j} - \cancel{\mu - \beta_j}) + \tau_A - \tau_B + \varepsilon_{Aj} - \varepsilon_{Bj} \\ &= \mu_D + \varepsilon_j \end{aligned}$$

Paired Experimental Design

Hypotheses

The hypotheses of interest can now be defined in terms of μ_D , e.g.:

$$\begin{cases} H_0 : \mu_D = 0 \\ H_1 : \mu_D \neq 0 \end{cases}$$

And now, we are back to our traditional "single sample hypothesis test". The population of interest is the difference in convergence time for the family of problems under investigation. The test statistic is given by:

$$T_0 = \frac{\bar{D}}{S_D/\sqrt{N}} \sim t^{(N-1)}$$

Where \bar{D} is the average of the paired differences, and N is the number of benchmark problems in the experiment.

Paired Experimental Design

Other Considerations

Some other important questions worth considering:

- In this example the minimally interesting effect size δ^* must be expressed in terms of *average time gains across problems* (not within individual instances).
- The most important sample size to consider in this situation refers to the *number of problem instances*, and not necessarily to the number of within-problems repeated measures;
- The number of repetitions within each problem will have an impact on the uncertainty associated to each observation (that is, to each value of mean time to convergence for each algorithm on each problem), which will propagate down to the residual variance.

Paired Experimental Design

Summary

- The Paired Design removes the effects of **controllable** nuisance factors from the analysis. (Problem type, personal characteristics, etc)
- It is strongly indicated in cases with **strong correlations between samples** (e.g., heterogeneous experimental conditions).

Paired Comparison Example

Going back to our example, assume the following facts about the desired comparison:

- The benchmark set is composed of seven problem instances ($N = 7$);
- The researcher is interested in finding differences in mean time to convergence greater than ten seconds ($\delta^* = 10$) with a power of at least $(1 - \beta) = 0.8$, using a significance level $\alpha = 0.05$;
- The researcher performs $n = 30$ repeated runs¹ of each algorithm in each problem, from random initial conditions.

¹ Not that this number is necessarily good, but it is generally an easy alternative if you don't want to keep justifying your choices to less statistically-savvy reviewers.

Executing the Paired Analysis

Step 1: load and precondition the data

```
> # Read data from CSV file
> data <- read.table("benchmark.csv", header=T)

# Change the type of the "Problem" variable
#           from "number" to "Factor"
> data$Problem <- as.factor(data$Problem)

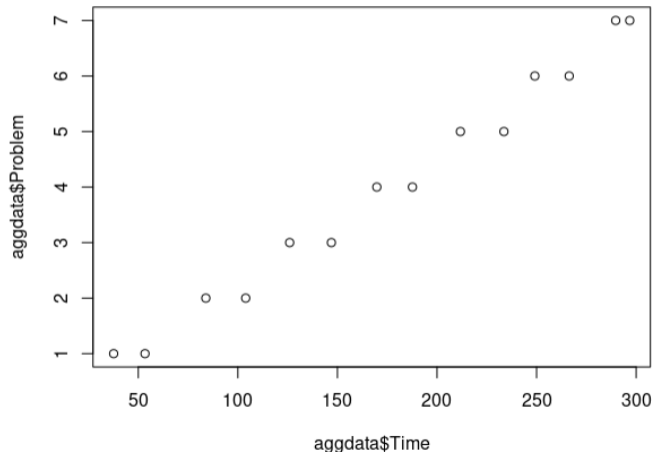
# Summarize within-problem observations by mean
> aggdata <- aggregate(Time ~ Problem:Algorithm,
+                       data = data, FUN = mean)
> summary(aggdata)
Problem Algorithm      Time
1:2      A:7      Min.   : 37.63
2:2      B:7      1st Qu.:109.45
3:2                      Median :178.73
4:2                      Mean   :175.48
5:2                      3rd Qu.:245.25
6:2                      Max.   :296.79
7:2
```

Problem	Algorithm	Time
1	A	37.62860
1	B	53.38058
2	A	83.93047
2	B	103.93440
3	A	126.01392
3	B	146.95828
4	A	169.77983
4	B	187.67396
5	A	211.75310
5	B	233.59827
6	A	249.12982
6	B	266.38702
7	A	289.73297
7	B	296.78549

Executing the Paired Analysis

Influence of the problem type in the results

Note that the difference between problems is bigger than the difference between methods.



Executing the Paired Analysis

Step 2: analysis

```
> # Perform paired t-test
> t.test(Time ~ Algorithm, data = aggdata,
+        paired = TRUE)          # <-- To do a paired test, just change here.
```

Paired t-test

data: Time by Algorithm

t = -9.1585, df = 6, p-value = 9.54e-05

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

-21.85862 -12.64118

sample estimates:

mean of the differences

-17.2499

Executing the Paired Analysis

Step 2: Alternate calculation

```
# Create an array with the difference per problem, and perform one-sample test.  
> difTimes <- aggdata$Time[1:7] - aggdata$Time[8:14])  
> t.test(difTimes)
```

One Sample t-test

data: difTimes

t = -9.1585, df = 6, p-value = 9.54e-05

alternative hypothesis: true mean is not equal to 0

95 percent confidence interval:

-21.85862 -12.64118 # Same result!

sample estimates:

mean of x

-17.2499

Check your understanding: Why is the paired test on two samples equivalent to the one sample test on the difference vector of the samples?

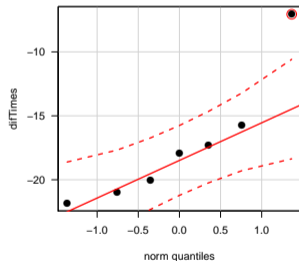
Paired Analysis

Step 3: Testing the assumptions

As usual, you should test the normality and variance of your data, and guarantee the independence of the observations. Let's check the normality test.

```
> shapiro.test(difTimes)
Shapiro-Wilk normality test
data:  difTimes
W = 0.8387, p-value = 0.09655

# Redo test without outlier
> indx <- which(difTimes == max(difTimes))
> t.test(difTimes[-indx])$p.value
[1] 6.179743e-06
> t.test(difTimes[-indx])$conf.int
[1] -21.41856 -16.48037
```



The normality test showed one big outlier. It does not invalidate the test, but it should be examined. You might learn something important!

Why is Pairing Important?

What happens if we ignore the dependency between observations?

```
> t.test(Time ~ Algorithm, data = aggdata)
```

```
Welch Two Sample t-test
```

```
data: Time by Algorithm
```

```
t = -0.3609, df = 11.993, p-value = 0.7245
```

```
alternative hypothesis: true difference in means is not equal to 0
```

```
95 percent confidence interval:
```

```
 -121.40320    86.90341
```

```
sample estimates:
```

```
mean in group A mean in group B
```

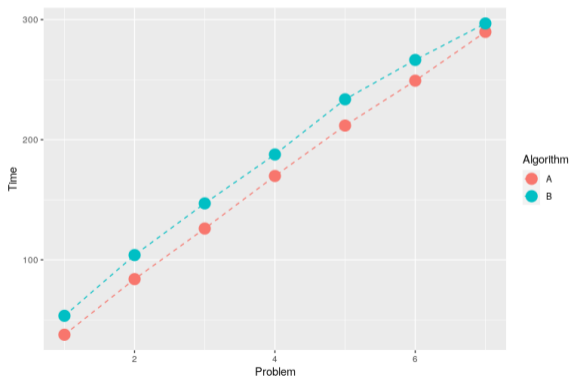
```
    166.8527      184.1026
```

If we don't take into account the large variation among problems, **it will hide variation between the two methods.**

Why is Pairing Important?

A visual Comparison

Paired Samples



Unpaired Samples



About these Slides

These slides were made by Claus Aranha, 2022. You are welcome to copy, re-use and modify this material.

These slides are a modification of "Design and Analysis of Experiments (2018)" by Felipe Campelo, used with permission.

Individual images in some slides might have been made by other authors. Please see the following references for those cases.

Image Credits I

[Page 6] Steel rod image: <http://www.shutterstock.com/pic-73207399/>

[Page 10] Two models image from D.C. Montgomery "Applied Statistics and Probability for Engineers", Wiley 2003

[Page 11] Two models image from D.C. Montgomery "Applied Statistics and Probability for Engineers", Wiley 2003