



# Basics of Statistics

C - Probability and Statistics

2023

# Contents

1. Statistical inference. Introduction and basic concepts. Parameters
2. Point estimation. Estimators
  - a. Uni- and bivariate descriptive statistics
3. Estimation using confidence intervals
  - a. Statistics
  - b. Confidence and risk
  - c. Premises
4. Parameter estimation using confidence intervals
  - a. Confidence interval of  $\mu$  (known  $\sigma$  and unknown  $\sigma$ )
  - b. Confidence interval of  $\pi$
  - c. Confidence interval of  $\sigma$
5. Comparison of two parameters using CIs
  - a. Confidence interval of  $\mu_1 - \mu_2$  in paired samples
  - b. Confidence interval of  $\mu_1 - \mu_2$  in independent samples
  - c. Confidence interval of  $\pi_1 - \pi_2$  in independent samples
  - d. Confidence interval of  $\sigma_1^2 / \sigma_2^2$  in independent samples
6. Designs (how we obtain the data)
7. Functions in R for confidence intervals

# 1. Statistical inference

- We must provide evidence based on data.

For example, saying “my program works” requires evidence/data.

- It must be **reproducible**: only reproducible results might be of interest.

For example, a **miraculous** cure will not be useful for future patients.

- It must be **transparent**

to enable others to replicate the same results.

- **We infer** the characteristics of the **population** from a **random sample (RS)**.

For example, I can infer the population-wide connection speed from a random sample of speeds.



# 1. Statistical inference. Risks

- The scientific and technical (statistical) method:
  - by **deduction** → data collection design (population → RS)
  - by **induction** → inferring (estimating) results (RS → population)
- Statistical inference defines and **quantifies the risks** of this process. [E.g., the mean connection speed of the entire population cannot be known unless data are available for the entire population, but statistics allows us to estimate and **quantify the error** from a specific **random sample**.]
- The **evidence** provided by **data** ends with **the analysis**: e.g.,
  - “My program works well”
    - estimating a measure (e.g., **average** performance) and **its error**.
  - “My program improves the results of...”
    - estimating performance improvement (e.g., **mean difference**) and **its error**.

# 1. Statistical inference. Types of variables

To analyse the relationship between variables, we must establish the role of each one:

- **Response  $Y$ .** Measuring goal achievement – sometimes it can be an indirect measure.  
E.g., performance  $Y$  measured for a subject.
- **Decisions  $X$ .** We assign their values in experimental studies.  
They represent the potential to change the future: we want to measure the **effect** of  $X$  on  $Y$ .  
An experimental design allows the  $X$  to be independent of other variables.  
E.g., a teaching method based on **printed lists** of exercises ( $X=1$ ) compared with a method based on **e-status** ( $X=2$ ).
- **Co-variables  $Z$ .** These represent the conditions observed in *real* data.  
We can use  $Z$  to reduce the uncertainty of  $Y$  (we will have to quantify its success).  
We can obtain  $Z$  in both experimental and observational studies.  
 $Z$  are usually interrelated (*colinear* or *non-orthogonal*).  
E.g., the marks of two previous subjects ( $Z_1, Z_2$ ) usually have a certain relationship.

# 1. Statistical inference. Types of study

- **DO: Experimental studies**

We want **to change** the future **Y** through interventions in **X**.

In the analysis we estimate the **effects** of **X** on **Y**.

E.g., To try to improve the marks **Y**, we assign at random the students different work environments **X**.

- **SEE: Observational studies**

They allow us **to predict** **Y** from the observed values **Z**.

We will quantify the **capacity** of **Z** for **reducing** the **uncertainty** in the prediction of **Y**.

E.g., we compare the prediction of **Y** according to **Z<sub>1</sub>** or according to **Z<sub>2</sub>**, or depending on a certain **model m** of the two variables **m=f(Z<sub>1</sub>, Z<sub>2</sub>)**.

→ The group **Z<sub>1</sub>** reduces the uncertainty by 10%; **Z<sub>2</sub>** by 20%; and the model **m**, with both, by 25%.

**X** represents an **assignable** and well-defined cause.

The key to intervening is to be **owners** of **X**.

To guarantee independence from all **Z**, we assign **X** at random.

We assign respecting ethical and legal rights.

We are not **owners** of the **Z** variables (the units already come with the **Z** value).

We can establish **relationships** between **Z** and **Y**, which we can use to **predict** the values of **Y** from **Z**.

But the covariates **Z** may be related (**collinear**), so their *effects* on **Y** may be **confounded**.

Establishing **causality** requires many premises (which are beyond an introductory course).

# 1. Statistical inference. Basic concepts

- **Parameter:** an indicator of the population that we wish to know or estimate. E.g., the expectation ( $\mu$ ) of the heights of FIB students.
- **Statistic:** any indicator that is obtained as a function of the data of a sample. E.g., the sum of the heights of the students in a sample.
- **Estimator:** a statistic of a sample used to know the value of a parameter of the population. E.g., the average height in a random sample of FIB students is an estimator of the expectation ( $\mu$ ) of the heights of FIB students.

*Mean* may mean *expectation parameter* regarding the centre of gravity of the population distribution, or *statistical mean* regarding the average of a series of values obtained from a sample.

## 2. Point estimation

- An estimator  $\hat{\theta}$  of the unknown parameter  $\theta$  from the sample  $M(\omega_i)$   $(X_1, X_2, \dots, X_n)$  (a simple random sample defined in the appendix to Section B) is a function of the RVs:

$$\hat{\theta} = f(X_1, X_2, \dots, X_n)$$

- Point estimation:** the value that the estimator  $\hat{\theta}$  takes in a specific sample.

E.g.,  $\bar{x} = \frac{\sum x_i}{n}$  is the sample mean and is a point estimate of  $\mu$ .

Distinguish between the value  $\bar{x}$  (small letter) of a specific simple random sample and the sample mean random variable  $\bar{X}$  (capital letter).

- Standard error:** the variability of the estimator. In the above case of MEAN, the **standard error of the mean** (or *mean standard error*, or SE) is

$$se = \sqrt{V(\bar{X}_n)} = \sqrt{E[(\bar{X}_n - \mu)^2]} = \frac{\sigma}{\sqrt{n}}$$

Generally, the  $\sigma$  will be unknown and the standard error will have to be approximated using the corresponding estimator ( $\hat{\sigma}$ ) with the sample data:  $\widehat{se} = \frac{\hat{\sigma}}{\sqrt{n}} = \frac{s}{\sqrt{n}}$  (with  $s$  the point estimator of  $\sigma$ ).

Read the comments in point 5 of this section “Designs (how we obtain the data)” to use for Unit T



## 2. Point estimation. Cases

For the parameters we use letters of the Greek alphabet.

Parameter ( $\theta$ ) ( <b>POPULATION</b> )	Estimator ( $\hat{\theta}$ ) ( <b>SAMPLE</b> )
$\mu$ (expectation, population mean)	$\bar{x}$ (sample mean)
$\sigma^2$ (population variance) $\sigma$ (population standard deviation)	$s^2$ (sample variance) $s$ (sample standard deviation)
$\pi$ (probability)	$p$ (proportion)

MEAN:

$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$  The sample mean is a **point estimate of the parameter  $\mu$**  of central tendency.

STANDARD DEVIATION:

$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}} = \sqrt{\frac{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}}{n - 1}}$  The sample standard deviation is a **point estimate of the parameter  $\sigma$**  of dispersion.

PROPORTION:

$p = \sum_{i, x_i=1} 1/n$  The sample proportion is a **point estimate of the parameter  $\pi$** .

We must take into account the properties of the estimators (see the Appendix, along with other potential estimators).

## 2. Estimators and descriptive statistics

The above point estimators correspond to the functions of **descriptive statistics** for numerically summarising data (see more in the R section of the website).

The following table shows some (basic) **functions** in R for **descriptive statistics** in **univariate or bivariate numerical and categorical variables**:

	UNIVARIATE (numerical)	UNIVARIATE (categorical)	BIVARIATE
INDICATORS	<code>length( ) *</code> <code>mean( )</code> <code>var( )</code> <code>sd( )</code> <code>summary( )</code> <code>median( )</code>	<code>table( )</code>	<code>cov( , )</code> <code>cor( , )</code>
GRAPHICS	<code>hist( )</code> <code>boxplot( )</code>	<code>barplot(table( ))</code>	<code>plot( , )</code>

\* The sample size ( $n$ ) is not an estimator, but we include it in the list for practicality.

(More graph functions in R: <https://www.r-graph-gallery.com/>)

### 3. Estimation using confidence intervals

- We know how to calculate an “interval” that contains  $\bar{x}$  from  $\mu$ . But the real problem is **to approximate  $\mu$  from  $\bar{x}$**  (i.e., moving from an interval for the sample mean  $\bar{x}$  to one for the population mean  $\mu$ )
- From a probability  $1-\alpha$  between two (symmetric) values  $a$  and  $b$  (with known  $\sigma$ ):

$$P(a \leq \bar{X}_n \leq b) = 1 - \alpha \rightarrow P\left(\frac{a-\mu}{\sigma/\sqrt{n}} \leq \frac{\bar{X}_n-\mu}{\sigma/\sqrt{n}} \leq \frac{b-\mu}{\sigma/\sqrt{n}}\right) = 1 - \alpha \rightarrow P\left(z_{\frac{\alpha}{2}} \leq \frac{\bar{X}_n-\mu}{\sigma/\sqrt{n}} \leq z_{1-\frac{\alpha}{2}}\right) = 1 - \alpha$$

- we get the interval of the RV  $\bar{X}_n$  with **probability**  $1-\alpha$ :

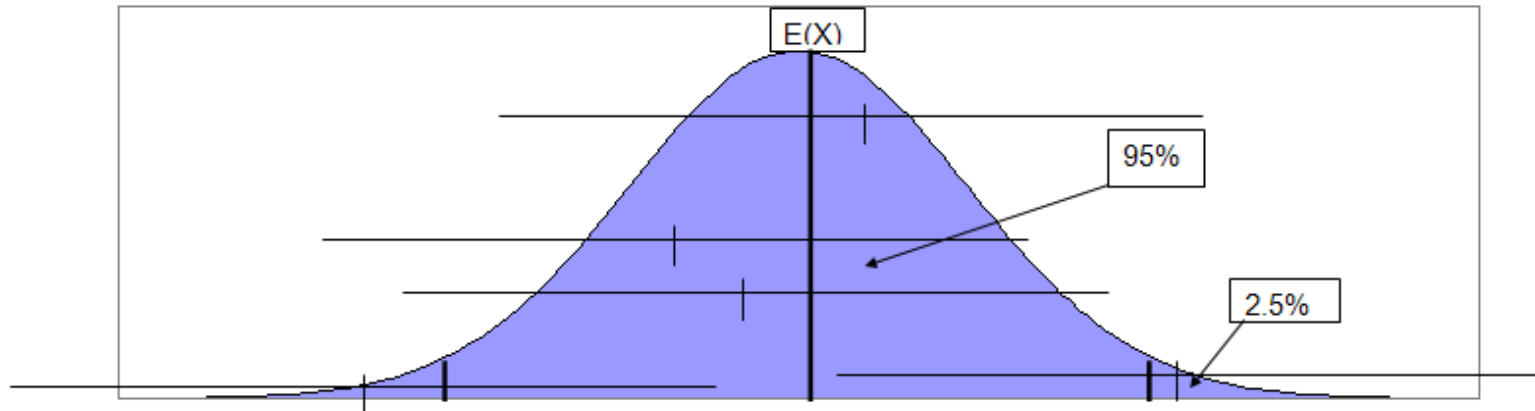
$$P\left(\mu + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{X}_n \leq \mu + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

- By rearranging, we get **the confidence interval (CI)  $1-\alpha$  of the parameter  $\mu$** :

$$P\left(\bar{X}_n + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

### 3. Estimation using confidence intervals

- $P\left(\bar{X}_n + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$  means that we can ensure that  $E(X) = \mu$  will be in the calculated range (with a confidence of  $1-\alpha$ )
- If  $1-\alpha$  is 95% ( $\alpha = 5\%$ ): **95% of the CIs will contain  $\mu$**  (see a simulation in the Appendix)



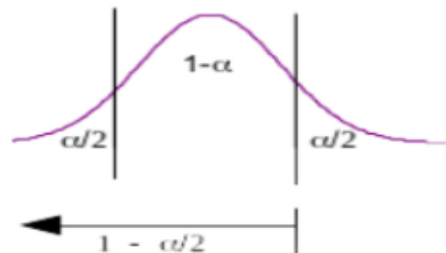
- This procedure is correct  $100 \cdot (1-\alpha)\%$  of the time!
- We call  **$CI(\mu, 1-\alpha)$**  the **CONFIDENCE INTERVAL**  $1-\alpha$  of  $\mu$

$$IC(\mu, 1 - \alpha) = \bar{x} \mp z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \quad (z_{\alpha/2} = -z_{1-\alpha/2} \text{ because } Z \text{ is symmetric})$$

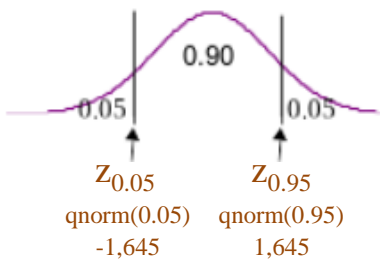
We will only observe one sample, and we will not know whether the found CI contains  $\mu$  or not, but we do know that in the long run this procedure gives true values  $100 \cdot (1-\alpha)\%$  of the time

# 3.a. Confidence and risk

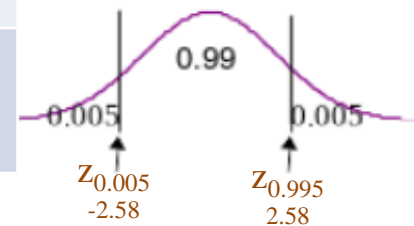
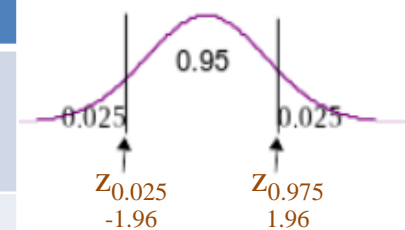
The calculation of a CI implies a confidence  $1-\alpha$  (and therefore a risk  $\alpha$ ), which we can represent as



And we can relate the confidence value to the quantile that we need to build the CI [E.g., the quantiles are indicated by a normal  $Z(0,1)$ , where we know that  $z_\alpha = -z_{1-\alpha}$  or  $z_{\alpha/2} = -z_{1-\alpha/2}$ ]



Confidence $1-\alpha$	Risk $\alpha$	$\alpha/2$	$1 - \alpha/2$
0.95	0.05	0,025	0,975
0.90	0.10	0.05	0.95
0.99	0.01	0,005	0,995



## 3.b. Statistics for inference

- We will see statistics of two types:
  - Ratio of **“signal” or “information”** (difference between a value  $\mu_0$  of the parameter and the sample value) to **“noise” or “error”** (standard error, SE).

These statistics are modelled following the Z or Student  $t^*$  model (in some cases we evaluate the “ $t$ -ratio” that quantifies by how many times the signal is greater than the noise).

$$\text{statistic } \hat{z} = \frac{(\bar{x} - \mu_0)}{\sigma/\sqrt{n}} = \frac{(\bar{x} - \mu_0)}{se} \quad \hat{z} \sim Z = N(0,1) \quad (\text{then the CI is } \mu \in \bar{x} \pm z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \text{ or } \bar{x} \pm z_{1-\frac{\alpha}{2}} \cdot se)$$

$$\text{statistic } \hat{t} = \frac{(\bar{x} - \mu_0)}{S/\sqrt{n}} = \frac{(\bar{x} - \mu_0)}{\widehat{se}} \quad \hat{t} \sim \text{T Student with degrees of freedom } \nu$$

- Ratio of variances.** These statistics are modelled following the F model.

$$\text{statistic } \hat{F} = \frac{S_A^2}{S_B^2} \quad \hat{F} \sim \text{F Fisher-Snedecor with degrees of freedom } \nu_1 \text{ and } \nu_2$$

*Student  $t_{\nu, 1-\alpha/2}$  ; Fisher  $F_{\nu_1, \nu_2, 1-\alpha/2}$  and chi squared  $\chi^2_{\nu, 1-\alpha/2}$  are defined in Section B (Appendix). Those models are derived from the normal distribution, and they are **parameterised with degrees of freedom ( $\nu$ ), depending on the sizes ( $n$ ) of the samples.***

## 3.c. Assumptions

- The fundamental assumption is to start from a random sample.

We say that values come from independent and identically distributed (IID) random variables.



- The premise of normality is necessary because CIs are based on the CLT theorem, which is based either on an original normal variable or a “large”  $n$ .

In small samples ( $n < 30?$ ), we will sustain the premise of normality with the prior knowledge of the response variable and with the graphic analysis with R.

See Part 7 on functions in R and Graphical Analysis of Normality in the Appendix to Section B.

## 4. Confidence interval for one parameter

Now we will see the **CI formulas** for 3 **single parameters**:

- The mean  $\mu$  (with or without known population variance)

*E.g., the mean mark of a subject*

- A proportion  $\pi$

*E.g., the proportion of passes of a subject*

- The variability  $\sigma^2$

*E.g., the deviation from the mean mark of a subject*



## 4.a. Confidence interval of $\mu$ (with known $\sigma$ )

- The confidence interval  $1-\alpha$  of  $\mu$  (with known  $\sigma$ ) is calculated as

$$CI(\mu, 1 - \alpha) = \bar{x} \pm z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

- Remember that we are using the CLT, which requires the random variable  $X$  to be normal or  $n$  to be “large”. Therefore, the requirement for performing this calculation is either  **$X \sim N$**  or  **$n$  “large”**
- This CI can be obtained by setting apart the parameter  $\mu$  from the statistic:  $\hat{Z} = \frac{(\bar{x}-\mu)}{\sigma/\sqrt{n}} =$   
 $\frac{(\bar{x}-\mu)}{se}$  whose distribution we know to be  $N(0,1)$
- Therefore the  $CI(\mu, 1-\alpha)$  can be seen as  $\bar{x} \pm z_{1-\frac{\alpha}{2}} \cdot se$ .

When  $n$  increases, the CI accuracy increases (narrower range). If the confidence increases (decreasing the risk  $\alpha$  of error), the accuracy of the CIs decreases (wider range).

To estimate  $\mu$ , we need to know  $\sigma$ , which is an unrealistic situation because  $\sigma$  is usually an unknown parameter (we can also assume a reasonable value from prior knowledge).

## 4.a. Confidence interval of $\mu$ with unknown $\sigma$

The previous confidence interval  $1-\alpha$  of  $\mu$  with unknown  $\sigma$  is calculated as

$$CI(\mu, 1 - \alpha) = \bar{x} \pm t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$$

This CI is obtained by isolating the parameter  $\mu$  from the statistic:  $\hat{t} = \frac{(\bar{x}-\mu)}{s/\sqrt{n}} = \frac{(\bar{x}-\mu)}{se}$

When  $\sigma$  is unknown, we replace  $\sigma$  by  $s$ ; and

the *Normal Z* by the **Student t** with  $n-1$  *df*

*being df=degrees of freedom*

In this case the initial variable X must be normal (the premise of normality) because the definition of the **Student t** model is based on normal variables.

The situation of not knowing  $\sigma$  is more realistic and frequent: no value is assumed but it is approximated by its point estimate  $s$ .

$t$  and  $N(0,1)$  are similar, increasingly so when  $n$  grows,  $t_{n \rightarrow \infty} \rightarrow N(0,1)$

For small values of  $n$ ,  $t$  has more variability, reflecting more uncertainty (as  $\sigma$  is approximated by  $s$ ).

Therefore, the CI with unknown  $\sigma$  will be wider than the equivalent assuming the true value of  $\sigma$ .

# 4.a Confidence interval of $\mu$ . Premises

To guarantee the confidence level of the CI, certain premises must be met.

The **fundamental** premise is that the origin of the sample must be **random**.

In addition:

- If sigma is known, one of the following two conditions is required:
  - **X~N**  $\rightarrow$  since the linear combination of normals is also normal ( $\bar{X} \sim N$ )
  - **The sample is “large”**  $\rightarrow$  by the CLT,  $\bar{X} \sim N$
- If sigma is unknown, one of the following conditions is required:
  - **X~N**  $\rightarrow (\bar{x} - \mu) / \sqrt{s^2/n} \sim t_{n-1}$
  - **The sample is large (“large” n)**  $\rightarrow$  by the CLT,  $\bar{X} \sim N$

In larger samples, the variation of  $s$  will be smaller ( $s$  estimates  $\sigma$  well), and we can consider that  $(\bar{x} - \mu) / \sqrt{s^2/n} \approx (\bar{x} - \mu) / \sqrt{\sigma^2/n} \sim N(0,1)$ .

In summary	... $\sigma$ Is known	... $\sigma$ is unknown
<i>If X is normal and...</i>	<i>We use the Normal</i>	<i>We use the Student t</i>
<i>If X is not normal but n is “large” and...</i>		

## 4.b. Confidence interval of $\pi$

$$\text{Let } X \sim B(n, \pi) \rightarrow E(X) = \pi \cdot n$$

$$V(X) = \pi \cdot (1 - \pi) \cdot n$$

$$\text{Then, } P = X/n \rightarrow E(P) = E(X/n) = E(X)/n = \pi \cdot n / n = \pi$$

$$V(P) = V(X/n) = V(X)/n^2 = \pi \cdot (1 - \pi) \cdot n / n^2 = \pi \cdot (1 - \pi) / n$$

By using the convergence from B to N,  $P \rightarrow N\left(\mu_P = \pi, \sigma_P = \sqrt{\frac{\pi(1-\pi)}{n}}\right)$

So, the statistic  $\hat{Z} = \frac{(P-\pi)}{\sigma_P} = \frac{(P-\pi)}{se}$  is distributed as N(0,1) provided

$n$  is “large” **and**  $\pi$  not extreme.

$$IC(\pi, 1 - \alpha) = P \pm z_{1-\frac{\alpha}{2}} se = P \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}$$

As a summary guide, check that

$$\pi \cdot n \geq 5 \text{ and } (1 - \pi) \cdot n \geq 5$$

The **paradox** that we need to know  $\pi$  to estimate the CI of  $\pi$  is usually solved in two ways:

a) by substituting  $\hat{\pi}$  with  $P$ : 
$$IC(\pi, 1 - \alpha) = P \pm z_{1-\alpha/2} \cdot \sqrt{(P(1 - P))/n}$$

b) by obtaining the maximum of  $\hat{\pi} \cdot (1 - \hat{\pi})$ , making  $\hat{\pi}$  equal to 0.5: 
$$IC(\pi, 1 - \alpha) = P \pm z_{1-\alpha/2} \cdot \sqrt{(0.5(1 - 0.5))/n}.$$

## 4.c. Confidence interval of $\sigma^2$

If  $X_i \rightarrow N$   $(n-1) \cdot \frac{s^2}{\sigma^2} = (n-1) \cdot \frac{(\sum_{i=1}^n (x_i - \bar{x})^2)/(n-1)}{\sigma^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma}\right)^2 \sim \chi_{n-1}^2$

We can relate the variance ratio statistic ( $S^2/\sigma^2$ ) to a  $\chi^2$   
as the sum of squared normal variables is  $\chi^2$

(see models derived from the normal in the Appendix to Section B).

Therefore, 
$$P\left(\chi_{n-1, \frac{\alpha}{2}}^2 \leq \frac{S^2 \cdot (n-1)}{\sigma^2} \leq \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) = 1 - \alpha$$

$$P\left(\frac{1}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} \leq \frac{\sigma^2}{S^2 \cdot (n-1)} \leq \frac{1}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = 1 - \alpha$$

$$P\left(\frac{S^2 \cdot (n-1)}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} \leq \sigma^2 \leq \frac{S^2 \cdot (n-1)}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = 1 - \alpha$$

$$IC(\sigma^2, 1 - \alpha) = \left[ \frac{s^2(n-1)}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \frac{s^2(n-1)}{\chi_{n-1, \frac{\alpha}{2}}^2} \right]$$

This is a CI for  $\sigma^2$ , not for  $\sigma$ !!

Since  $\chi^2$  is not symmetrical, it requires obtaining the upper and lower quantiles instead of doing  $\pm$ .

# 5. Confidence interval to compare two parameters

The **CI formulas** to...:

- Compare  $\mu_1$  and  $\mu_2$

*E.g., the CI of the differential effect ( $\mu_1 - \mu_2$ ) comparing averages between two subjects\**

We must differentiate between

- **paired samples\*\*** (each case results in two measures, pairs of measures)

*(the same students in both subjects,  $\mu_1 - \mu_2 = \mu_{\text{difference}} = \mu_d$ )*

- **independent samples** (each case is an independent measure)

*(different students in the two subjects)*

- Compare  $\pi_1$  and  $\pi_2$

*E.g., the CI of the differential effect ( $\pi_1 - \pi_2$ ) comparing averages between two subjects\**

- Compare  $\sigma^2_1$  and  $\sigma^2_2$

*E.g., the CI comparing deviations between two subjects\**

\* **The origin of the sample must be random.**

\*\* If possible, a design with paired data will be more efficient (as we will see below).

## 5.a. CI of $\mu_1 - \mu_2$ (or of $\mu_D$ ) in paired samples

If obtain a simple random **paired sample** of size  $n$ , and

we define  $\mathbf{d} = Y_1 - Y_2$  then  $E(d) = \mu_d$  and  $V(d) = \sigma_d^2$

and the  $n$  observed differences values have a mean  $\bar{d}$  and deviation  $s_d$ .

- The statistic  $\hat{t} = \frac{(\bar{d} - \mu_D)}{s_d / \sqrt{n}} = \frac{(\bar{d} - \mu_D)}{se}$  follows the  $t_{n-1}$  model

where  $(\bar{d} - \mu_D)$  is the **signal** and  $se = s_d / \sqrt{n}$  is the standard **error**.

- The CI of the population mean difference with confidence  $1 - \alpha$  is

$$IC(\mu_1 - \mu_2, 1 - \alpha) = IC(\mu_d, 1 - \alpha) = \bar{d} \pm t_{n-1, 1 - \frac{\alpha}{2}} \cdot \frac{s_d}{\sqrt{n}} = \bar{d} \pm t_{n-1, 1 - \frac{\alpha}{2}} \cdot se$$

It may be of practical interest to evaluate the  $t$ -ratio:  $t = \bar{d} / \frac{s_d}{\sqrt{n}} = \bar{d} / se$ ,

which says how many times the **signal** is greater than the **noise**.

## 5.b. CI of $(\mu_1 - \mu_2)$ independent samples

Be  $Y_1$  with  $E(Y_1) = \mu_1$ ,  $V(Y_1) = \sigma_1^2$ ; and  $Y_2$  with  $E(Y_2) = \mu_2$ ,  $V(Y_2) = \sigma_2^2$  with normal distributions ( **$\sigma_1$  and  $\sigma_2$  will be unknown values but must be assumed to be the same\***), from which we obtain two **independent** simple random samples of size  $n_1$  and  $n_2$  with means  $\bar{y}_1$ ,  $\bar{y}_2$  and deviations  $s_1$  and  $s_2$  as estimators of the common parameter  $\sigma$ .

- The statistic  $\hat{t} = \frac{(\bar{y}_1 - \bar{y}_2) - (\mu_1 - \mu_2)}{se}$  follows the distribution  $t_{n_1+n_2-2}$ , with standard error

$$se = s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \text{ where } s \text{ is the root of the "pooled" variance } s^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1-1) + (n_2-1)}$$

- The CI of the difference with confidence  $1-\alpha$  is

$$IC(\mu_1 - \mu_2, 1 - \alpha) = (\bar{y}_1 - \bar{y}_2) \pm t_{n_1+n_2-2, 1-\frac{\alpha}{2}} \cdot se =$$

$$(\bar{y}_1 - \bar{y}_2) \pm t_{n_1+n_2-2, 1-\frac{\alpha}{2}} \cdot s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

The **fundamental** condition is random samples. We **know** it.

The other **two assumptions** (normality of  $Y_1$ ,  $Y_2$  and  $\sigma_1^2 = \sigma_2^2$ ) will be analysed graphically.



## 5.c. CI of $\pi_1 - \pi_2$

Let  $P_1$  and  $P_2$  be the sample proportions of two binomial populations with  $\pi_1, \pi_2$ , from which we obtain two independent simple random samples of size  $n_1$  and  $n_2$ .

- The statistic  $\hat{z} = \frac{(P_1 - P_2) - (\pi_1 - \pi_2)}{se}$  follows the distribution  $N(0,1)$  with standard error  $se = \sqrt{P_1(1 - P_1)/n_1 + P_2(1 - P_2)/n_2}$ .
- The CI of the difference with confidence  $1 - \alpha$  is

$$IC(\pi_1 - \pi_2, 1 - \alpha) = (P_1 - P_2) \pm z_{1 - \frac{\alpha}{2}} \cdot se =$$

$$(P_1 - P_2) \pm z_{1 - \frac{\alpha}{2}} \cdot \sqrt{P_1(1 - P_1)/n_1 + P_2(1 - P_2)/n_2}$$

In this case, convergence requires “large” samples: usually  $P \cdot n > 5$  and  $(1 - P) \cdot n > 5$ .

## 5.d. CI of $\sigma_1^2/\sigma_2^2$

Let  $s_1$  and  $s_2$  be the sample deviations of two independent simple random samples of size  $n_1$  and  $n_2$  of two normal variables.

- The statistic  $\hat{F} = \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2}$  follows the distribution  $F_{(n_1-1, n_2-1)}$ .
- The CI of the ratio of variances with confidence  $1-\alpha$  is (following the same reasoning as the CI of  $\sigma^2$ )

$$IC(\sigma_1^2/\sigma_2^2, 1 - \alpha) = \left[ \frac{s_1^2/s_2^2}{F_{(n_1-1, n_2-1), 1-\frac{\alpha}{2}}}, \frac{s_1^2/s_2^2}{F_{(n_1-1, n_2-1), \frac{\alpha}{2}}} \right]$$

or (note the exchange of degrees of freedom of  $F$ )

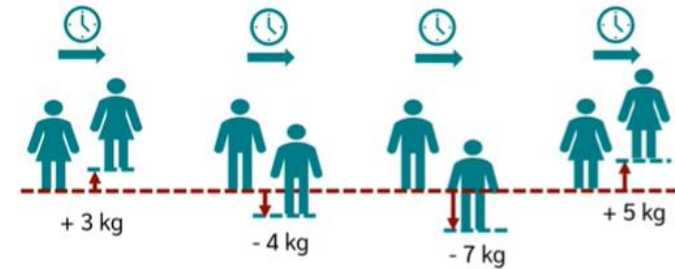
$$IC(\sigma_1^2/\sigma_2^2, 1 - \alpha) = \left[ s_1^2/s_2^2 F_{(n_2-1, n_1-1), \frac{\alpha}{2}}, s_1^2/s_2^2 F_{(n_2-1, n_1-1), 1-\frac{\alpha}{2}} \right]$$

## 6. Designs (how we obtain the data)

### Paired design:

One variable and two observations are taken from each unit (the two measures or responses),

**Requirement:** the first observation in a “pair” must not alter the state of the unit and therefore of the second observation.



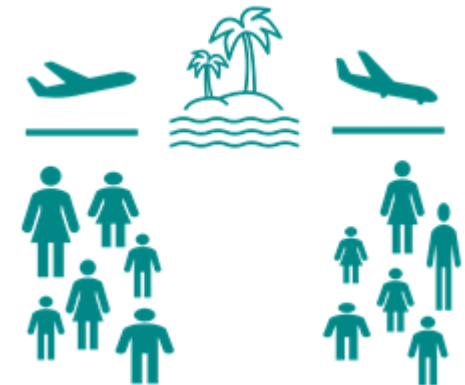
### Independent samples

For each unit, one observation and two variables are taken:

- (1) the outcome, the measure of the response and
- (2) the category to compare.

**Requirements:**

- It must be possible to assign the category to the unit (it cannot be a condition, such as sex).
- In observational studies, when the group is not assignable, the samples are selected separately.



**This is a simple approximation. The world of experiment design is much broader.**

The key is **random**:

Collecting data arbitrarily does not guarantee a random sample: **willy-nilly  $\neq$  at random.**

(1) **to plan** the random selection of the units to be measured; (2) to **carry out** the experiment correctly; (3) without **missing** values; and (4) to **document it** in a **reproducible** way.

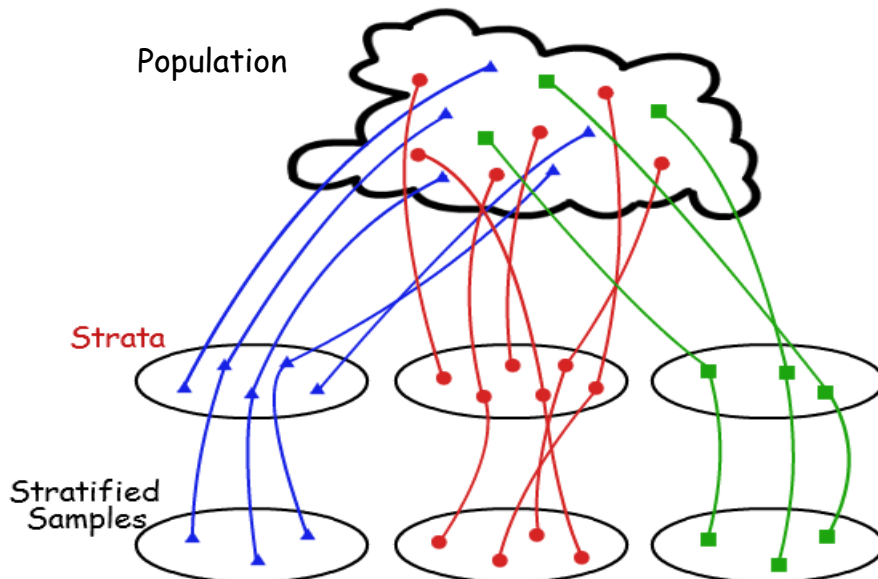
## 6. Designs (how we obtain the data)

The chosen design conditions the subsequent statistical analysis.

If the collection is complex, the statistical model used is also complex:

- Cases with **nested** data (*clusters*): groups from the top level (e.g., *school*) are first randomly selected, then groups from the lower level (*class*), until the individual (*student*) is reached.
- Cases with **stratified** data: all the strata can be seen, but within each stratum the individuals are randomly selected.

Then, the group of students chosen is not strictly a random sample; they must be analysed with appropriate techniques (not explained at a first level statistical course).



It is unusual to have the listing of the complete population; nor to be able to access any unit under the same conditions (both are requirements for a simple random sample).

Usually, some units will be more “visible” than others and will have a higher probability of being chosen (e.g., results that can only be obtained in order).

## 7. Functions in R for CIs

We will see the following:

- A list of functions in R.
- Functions in R for a sample. CI of  $\mu$ .
- Functions in R for two independent samples. CI of  $\mu_1 - \mu_2$ .
- Functions in R for paired samples. CI of  $\mu_D$ .
- Functions in R for paired samples. Graph of differences vs. means.
- Functions in R for comparing  $\sigma$ . CI of  $\sigma_1 - \sigma_2$ .
- Functions in R for  $\pi$ .

# A list of functions in R

Premise of normality (in Graphical Analysis of Normality in the Appendix to Section B):

`qqnorm(X)`

`qqline(X)`



CI of  $\mu$  with known  $\sigma$  (for this function you need the BSDA library):

`library(BSDA)`

`z.test(X, sigma.x= )`      # for a sample when sigma is known

CI of  $\mu$  (or  $\mu_s$ ) when  $\sigma$  (or  $\sigma_s$ ) is unknown:

`t.test(X)`      # for one sample

`t.test(XY)` or `t.test(X,Y,paired=T)`      # for paired samples

`t.test(X,Y,var.equal=T)`      # for two independent samples with equal variances

`t-test(X,Y,var.equal=F)`      # for two independent samples with different variances

CI of  $\sigma_s$  in two independent samples:

`var.test(X,Y)`

CI of  $\pi$ :

`prop.test` and `binom.test`

# In a sample

An example of nine values with positives and negatives (measurements above or below a threshold)

```
x <- c(-4,-2,-1,0,0,4,8,8,9) # mean = 2.4 SD = 4.9. We study normality with qqnorm(X), qqline(X)
library(BSDA)
z.test(X, sigma.x=4)      # CI assuming a population  $\sigma$  of 4.
```

```
z = 1.8333, p-value = 0.06675
Alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.1688409  5.0577298
Sample estimates: mean of x  2.444444
```

```
t.test(X)      # CI if we do not know the population  $\sigma$  but use the sample s.
```

```
t = 1.496, df = 8, p-value = 0.173
Alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -1.323423  6.212312
Sample estimates: mean of x  2.444444
```

You can check whether the CI limits coincide with those that would be calculated with the formulas. Apart from CIs, these functions in R provide the result of a *p-value*: P value assess the probability that the statistic is “extreme” in the distribution of the reference model (see more in the Appendix to section D with more functions that provide p-values)

# In two independent samples

An example of two samples to compare  $\mu_1$  and  $\mu_2$  with the CI of the differential effect ( $\mu_1 - \mu_2$ )

**X1** `<-c(1,2,3,5,6,6,7,7,8,8,9)` # mean = 5.6 SD = 2.62. We study normality with `qqnorm(X1)`, `qqline(X1)`

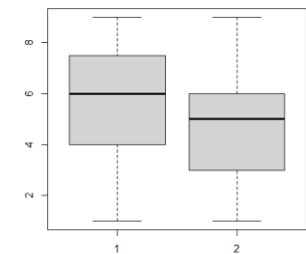
**X2** `<-c(1,1,3,4,5,5,6,7,9)` # mean = 4.5 SD = 2.65. We study normality with `qqnorm(X2)`, `qqline(X2)`

E.g., X1 and X2, two samples of marks with equal variability

**t.test(X1,X2,var.equal=T)**

```
t = 0.91335, df = 18, p-value = 0.3731
Alternative hypothesis: true difference is not equal to 0
95 percent confidence interval:
 -1.405312  3.566928
Sample estimates: mean of x mean of y  5.636364  4.555556
```

**boxplot(X1,X2)**



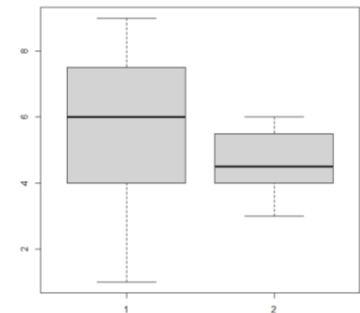
**X3** `<-c(3,4,4,4.5,4.5,5,6,6)` # mean = 4.6 SD = 1.0. We assume normality (or `qqnorm(X3)` and `qqline(X3)`).

E.g., X1 and X3, two samples of marks with non-equal variability

**t.test(X1,X3,var.equal=F)**

```
t = 1.1641, df = 13,793, p-value = 0.2641
Alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.8546781  2.8774054
Sample estimates: mean of x mean of y  5.636364  4.625000
```

**boxplot(X1,X3)**



**Interpretation:** the difference in means is up to 0.85 points in favour of group 3 or up to 2.9 points in favour of group 1, with a confidence of 95%



# In paired samples

An example of two samples with the CI of the differential effect ( $\mu_1 - \mu_2$ )

```
Y1 <- c(1,1,2,2.0,2,2.5,4,5,5.5,6,7.5,8,8,9.5,9,9.5)
```

```
Y2 <- c(1.5,1,2,1.0,3,3,3.5,5,6,6,8.5,8.5,9.5,8.5,9.1,9)
```

E.g., X1 and X2, two samples of marks, both came from the same student

```
t.test(Y1, Y2, paired=T)
```

```
t = -0.92936, df = 15, p-value = 0.3674
Alternative hypothesis: true difference is not equal to 0
95 percent confidence interval:
 -0.5351864  0.2101864
Sample estimates: mean of the differences -0.1625
```

In paired samples, it is possible to work with the difference of the values ( $D=Y1-Y2$ ), so it is like the case of one sample (instead of  $\mu_1 - \mu_2$ , we are interested in  $\mu_D$ ).

```
D <- Y1-Y2
```

```
-0.5  0.0  0.0  1.0 -1.0 -0.5  0.5  0.0 -0.5  0.0 -1.0 -0.5 -1.5  1.0 -0.1  0.5
```

```
mean(D)
```

```
-0.1625
```

```
sd(D)
```

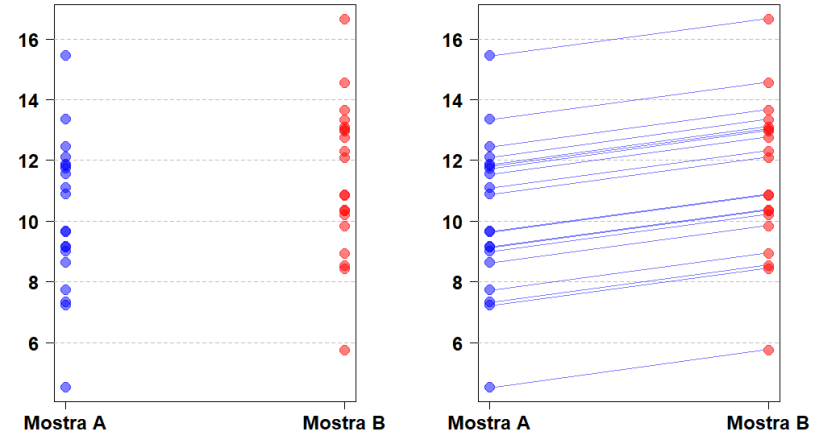
```
0.6994045
```

# In paired samples

It is very important not to do an analysis of paired data as independent data.

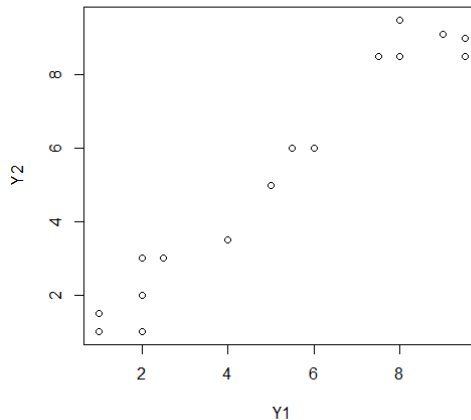
In the graph on the left (independent samples) we do not see that there could be differences in mean between the two populations.

In the graph on the right (paired samples) it is clearly seen that the mean is higher in sample B.



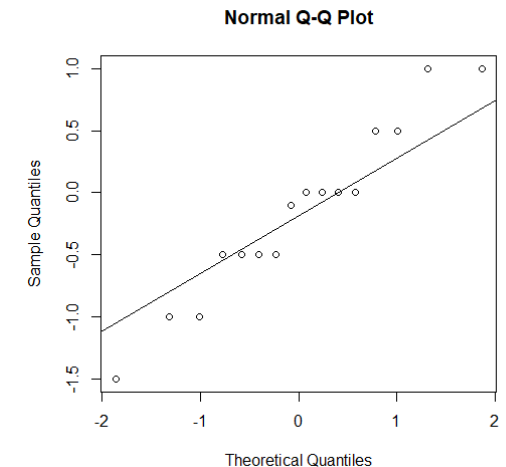
We have graphical descriptive functions in R to see the relationship between the two variables and the normality of the difference:

`plot(Y1, Y2)`



`qqnorm(Y1 - Y2)`

`qqline(Y1 - Y2)`

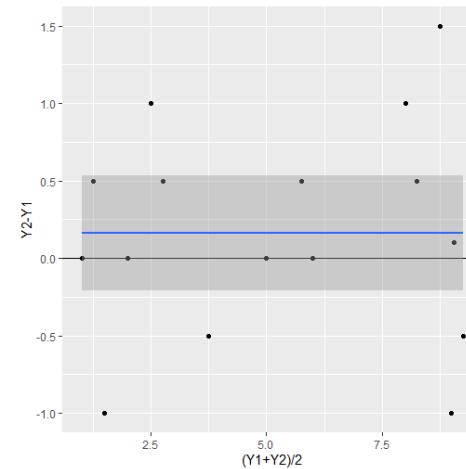


# In paired samples

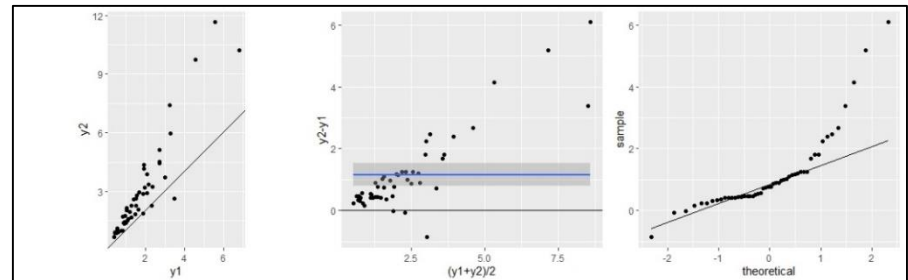
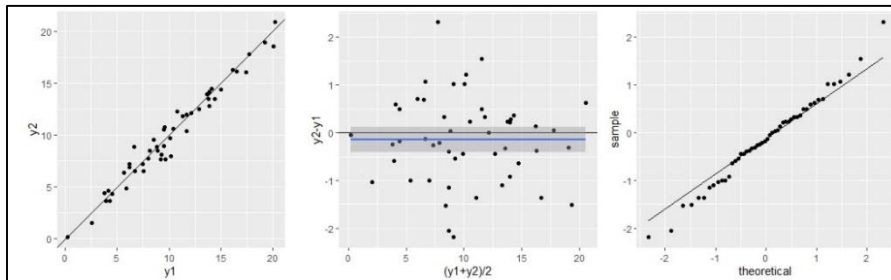
There are specific R functions for paired samples: the **Bland-Altman (BA) plot**, which represents the differences in the responses for each individual according to their means.

```
install.packages("PairedData")
library(PairedData)
p <- paired(Y1,Y2)
plot(p, type='BA')

( or plot((Y1+Y2)/2, Y2-Y1) )
```



The mean-difference plot (complementing the plot and the qqnorm) shows whether there is an additive (or multiplicative) effect and helps decide whether a transformation of the data would be appropriate (this will be seen in block D).



# Comparing variances

## An example with the CI of $\sigma_1^2/\sigma_2^2$

(as in the exercises comparing the variability in the duration of refills of ink cartridges of two brands).

```
A <- c(350, 361.9, 365, 365, 365, 370, 372, 377)
```

```
# mean(A)=365.7375 SD(A)=8.00231 var(A)=64.03696. We study normality with qqnorm(A), qqline(A)
```

```
B <- c(390, 391.7, 410, 412, 414, 418)
```

```
# mean(B)=405.95 SD(B)=12.00396 var(B)=144.095. We study normality with qqnorm(B), qqline(B)
```

```
var.test(B,A)
```

```
F test to compare two variances
data: B and A
F = 2.2502, num df = 5, denom df = 7, p-value = 0.3199
Alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 0.4257491 15.4206862
Sample estimates:
Ratio of variances
      2.250185
```

**Interpretation:**  $V(B)=2 \cdot V(A)$  although the 95% CI shows that the true population  $V(A)/V(B)$  ratio could be as small as 0.4, ( $V(A)=2.5 \cdot \text{Var}(B)$ ); and also as large as 15, ( $V(B)=15 \cdot V(A)$ ).

Final interpretation: high uncertainty, so more information should be considered.

There is no preference with the current data.

# For the CI of $\pi$

For example, tossing a coin 100 times and observing 56 heads.

```
prop.test(56,100) # requires convergence to Normal (large n)
```

```
1-sample proportions test with continuity correction
data: 56 out of 100, null probability 0.5
X-squared = 1.21, df = 1, p-value = 0.2713
Alternative hypothesis: true p is not equal to 0.5
95 percent confidence interval:
 0.4573588 0.6579781
Sample estimates: p      0.56
```

```
binom.test(56,100) # more appropriate if the sample is small
```

```
Exact binomial test
data: 56 and 100
Number of successes = 56, number of trials = 100, p-value = 0.2713
Alternative hypothesis: true probability of success is not equal to 0.5
95 percent confidence interval:
 0.4571875 0.6591640
Sample estimates:
Probability of success
                0.56
```

None of these CIs **exactly** match the one calculated with the formula approximating the normal distribution explained above. The agreement would increase with larger sample sizes and with proportions closer to  $\frac{1}{2}$ .

Although at a practical and interpretive level, all CIs agree with [46% to 66%]