

## Introduction

• Concept of comparison of quantities:

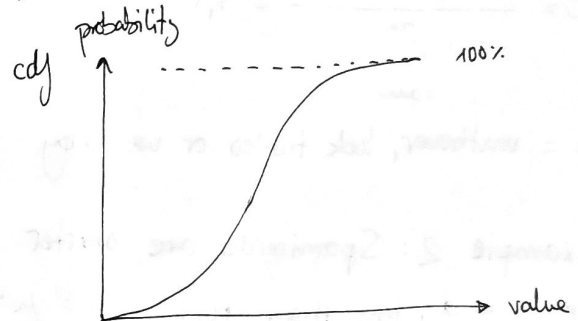
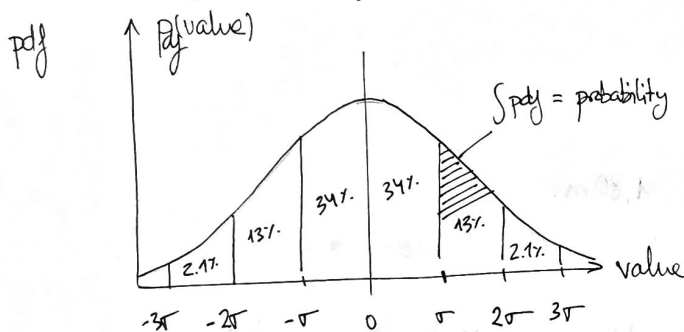
- + Do AVE trips from Madrid to Seville cost  $> 60 \text{ €}$ ?
- + Do people who sleep more have better health?
- + ~~Is~~ Is effective a certain treatment against some illness?

## Key concepts to understand previously

Population  $\rightarrow$  All existing elements in a dataset. All spaniards. It is not possible to know the truth of a population without examining all its elements.

Sample  $\rightarrow$  Subset of a population usually selected randomly by a statistician to infer some information about.

Distribution  $\rightarrow$  mathematical function that provides the probabilities of occurrence of different outcomes of an experiment: pdf (probability density function), and cdf (cumulative distribution function), for continuous distributions.



Normal Distribution, T Student Distribution  $\rightarrow$

• Standardization (normal distribution):  $z =$

$$z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\bar{x}$   $\rightarrow$  mean of sample  
 $\mu$   $\rightarrow$  mean of distribution  
 $\sigma$   $\rightarrow$  standard deviation of distribution  
 $\sqrt{n}$   $\rightarrow$  number of elements in sample

• Standardization (T student):

Degrees of freedom

Number of decisions that can be made while computing a statistic.

- computing a sample mean:  $n$  observations - 1
- choosing hat every weekday: all weekdays - 1

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$\bar{x}$   $\rightarrow$  mean of sample  
 $\mu$   $\rightarrow$  mean of distribution  
 $s$   $\rightarrow$  estimated standard deviation from sample  
 $\sqrt{n}$   $\rightarrow$  number of elements in sample

$\bar{z}$  &  $t$  have no dimensions

# Types of Hypothesis tests

+ One sample vs constant : now

+ Two or more samples to each other

Related samples

Independent samples

ANOVA (>2 samples)

} after

Example 1: Spaniards are taller than 1,70m.

One sided, greater than test

a) Get a sample of Spaniards, imagine

$\bar{x} = 1,72$  with a sample size of  $n = 5000$  and  $s = 0,2$

$H_0: \mu \leq \mu_0 \rightarrow$  "Spaniards are not taller than 1,70m"

$H_1: \mu > \mu_0$

$$b) t = \frac{1,72 - 1,70}{\frac{0,20}{\sqrt{5000}}} = 7,07$$

$p =$  whatever, look tables or use scipy

Example 2: Spaniards are shorter than 1,80m.

One sided, less than test.  $H_0: \mu \geq \mu_0$ ,  $H_1: \mu < \mu_0$ , idem but with  $t < 0$ .

Example 3: Spaniards' Height is significantly different than 1,80.

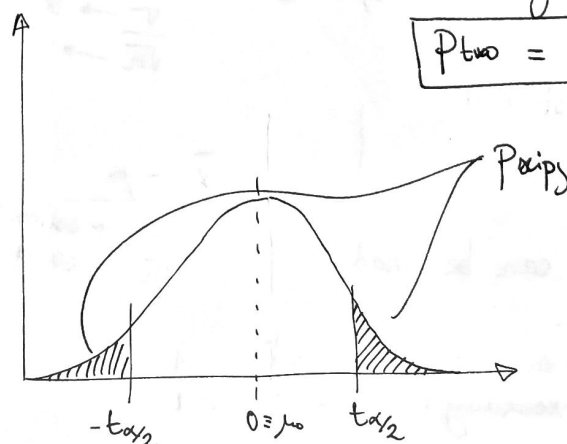
Two sided test

$H_0: \mu = \mu_0$

$H_1: \mu \neq \mu_0$

+ In scipy, all test are two sided by default:

$$P_{two} = 2 \cdot P_{one}$$



test significance, is a threshold we set,  $p$  obtained is the probability of a result obtained by chance.

If that  $p$  is less than 5%, test is successful and null hypothesis  $H_0$  is rejected:  $t > 0$  is mandatory as well as  $p < \alpha$

+ Confidence interval:

$$\bar{x} \pm t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

contains population mean with  $\alpha\%$  confidence (95% for example).

Test assumptions, must be true to be valid

- Observations must be independent of each other in the sample.
- Data distribution is normal.
- Sample size is at least 30.
- For z test,  $\sigma$  is known (never happens!), otherwise use t test.