# 1 Key Concepts and Basic Statistics

1. Key concepts
2. Descriptive statistics
3. Frequency and probability-distributions
4. Hypothesis testing
5. P-values
6. Interval estimation
7. Suggested exercises

# ① Key concepts

Pop.

Sample

**Def. Population: All** the units to be studied

**Def. Sample: A subset** of the population

Some common sampling methods:

* Simple random sampling: All
- - - - - - - - - - - - - - - - -
subsets of the same size have the same probability of being drawn

↑ Usually means: All units have an equal probability of being selected

* Systematic sampling: Selecting the units at regular intervals from a "List" of all the population members
  ↳ E.g. an alphabetically ordered list


* Stratified sampling: Classify the population into strata (i.e. categories) and then sample from each stratum (i.e. category)

* Convenience sampling: Select those that are close at hand/ easy to investigate

# ② Descriptive statistics

## Measures of central tendency

→ Intended to indicate where the majority of the values are

→ The sample mean:

$$\bar{x} = \frac{\Sigma x}{n}$$

Example: 8, 10, 9, 7, 8

$$\bar{x} = \frac{8 + 10 + 9 + 7 + 8}{5} = 8.4$$

→ The median: The middle value that separates the highest values from the lowest

$\rightarrow$ The mode : The most frequent value

$\rightarrow$ The weighted mean

## Measures of variability:

$\rightarrow$ Intended to indicate the degree of variability or dispersion

Example: Sample 1 = 8, 10, 9, 7, 8   $\bar{x} = 8.4$

— " — 2 = 6, 7, 11, 8, 10   $\bar{x} = 8.4$

Sample 1:   x   xx   x   x

— " - 2:

   x   x   x         x         x

+——+———+———+——+———+———+——

6   7   8   9   10   11

$\rightarrow$ The sample variance:

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

Example:

sample 1

sample 2

| X | $(x-\bar{x})$ | $(x-\bar{x})^2$ | X | $(x-\bar{x})$ | $(x-\bar{x})^2$ |
|---|---|---|---|---|---|
| 8 | -0.4 | 0.16 | 6 | -2.4 | 5.76 |
| 10 | 1.6 | 2.56 | 7 | -1.4 | 1.96 |
| 9 | 0.6 | 0.36 | 11 | 2.6 | 6.76 |
| 7 | -1.4 | 1.96 | 8 | -0.4 | 0.16 |
| 8 | -0.4 | 0.16 | 10 | 1.6 | 2.56 |

SUM      5.2          17.2

$$s^2 = \frac{\sum (x-\bar{x})^2}{n-1} = \frac{5.2}{4} = 1.3$$

$$s^2 = \frac{17.2}{4} = 4.3$$

→ The sample standard deviation:

$$s = \sqrt{s^2} = \sqrt{\frac{\sum(x-\bar{x})^2}{n-1}}$$

Examples: Sample 1 = $s = \sqrt{1.3} = 1.14$

2 = $s = \sqrt{4.3} = 2.07$

→ The sample range: Highest value minus the lowest

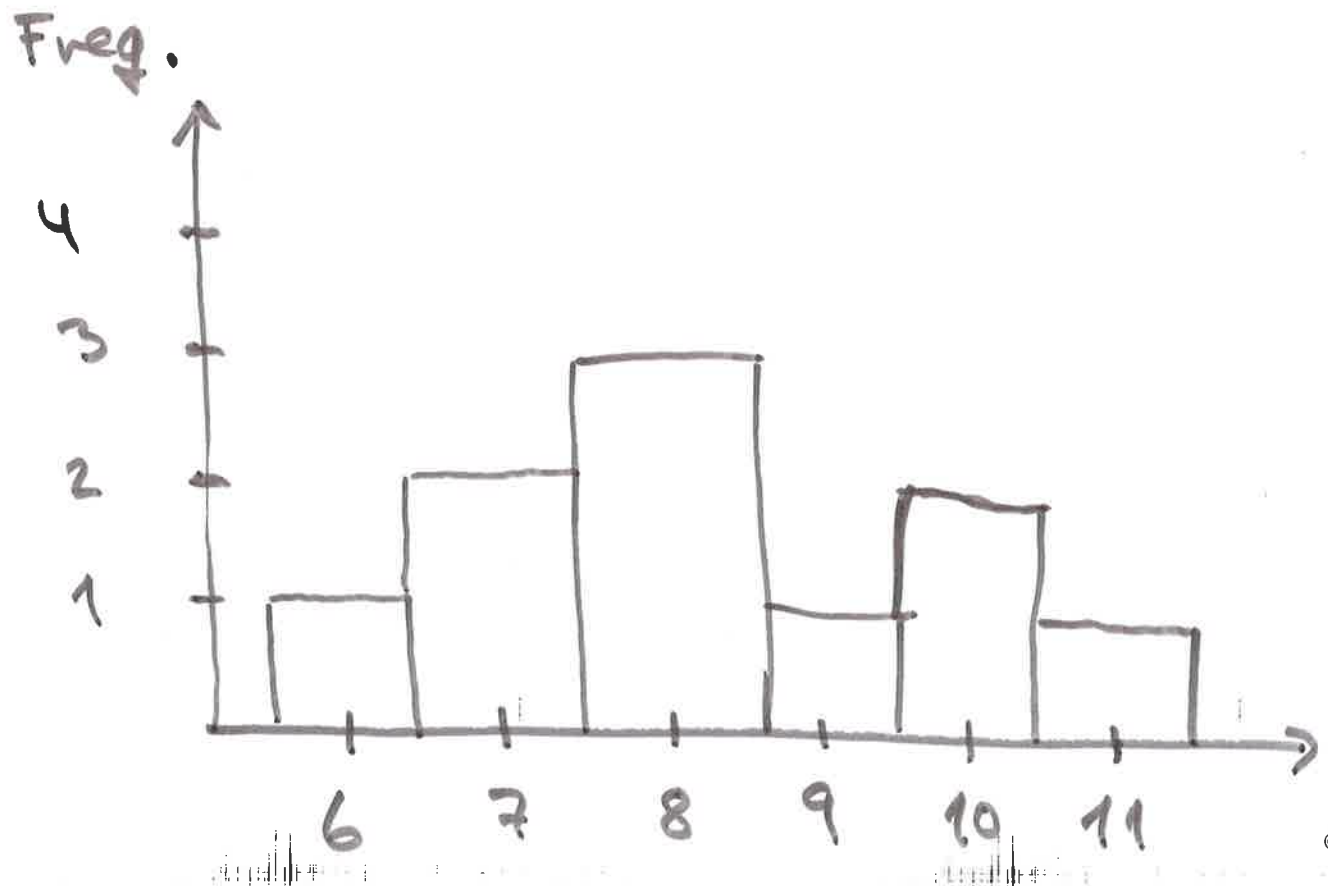Examples: Sample 1: $10 - 7 = \underline{\underline{3}}$

— 1 — 2: $11 - 6 = \underline{\underline{5}}$

# ③ Frequency and probability distributions

Def. <u>Frequency distribution</u>: A description of the number of times each value of a variable appears in the sample

Example: A histogram, i.e. a bar graph with no space between the bars
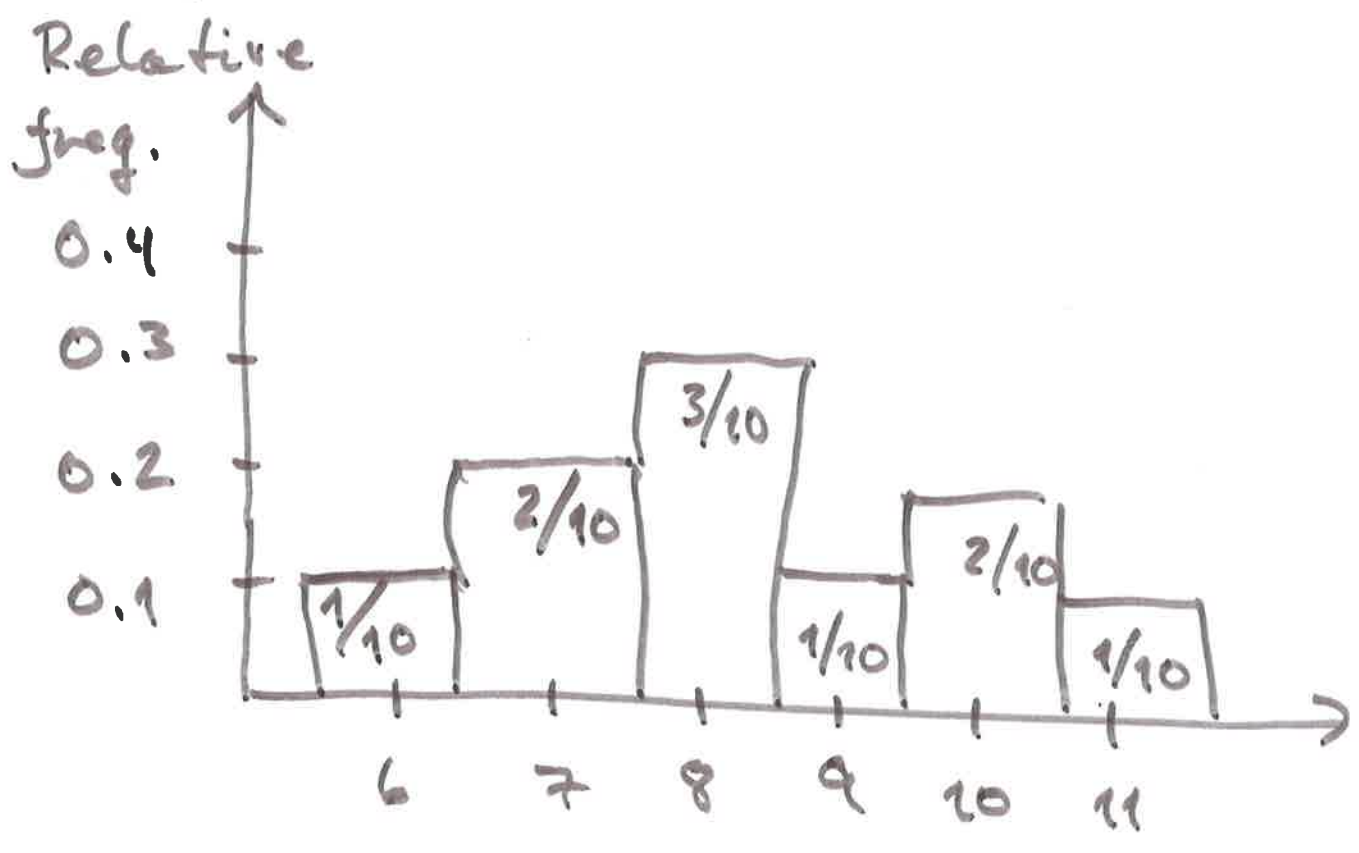
Sample = 8, 10, 9, 7, 8, 6, 7, 11, 8, 10

Freq.

# Def. Relative frequency distribution:

A description of the proportion of times each value appears in the sample

Example: A histogram

Sample = Same as previous ($n = 10$)



# Def. Probability distribution:
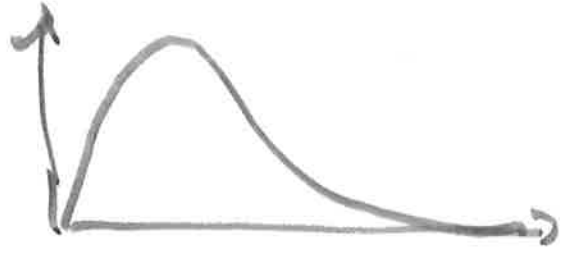
A histogram of relative frequencies
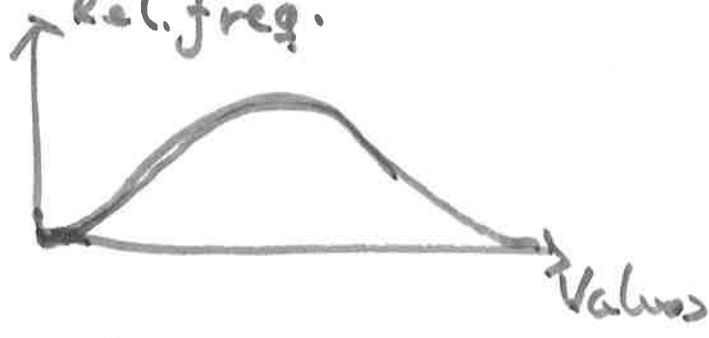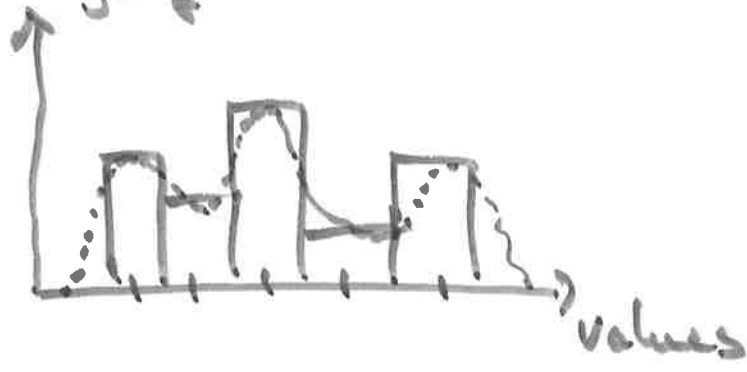
Example: Previous!

Discrete (categorical) probability distributions

vs.

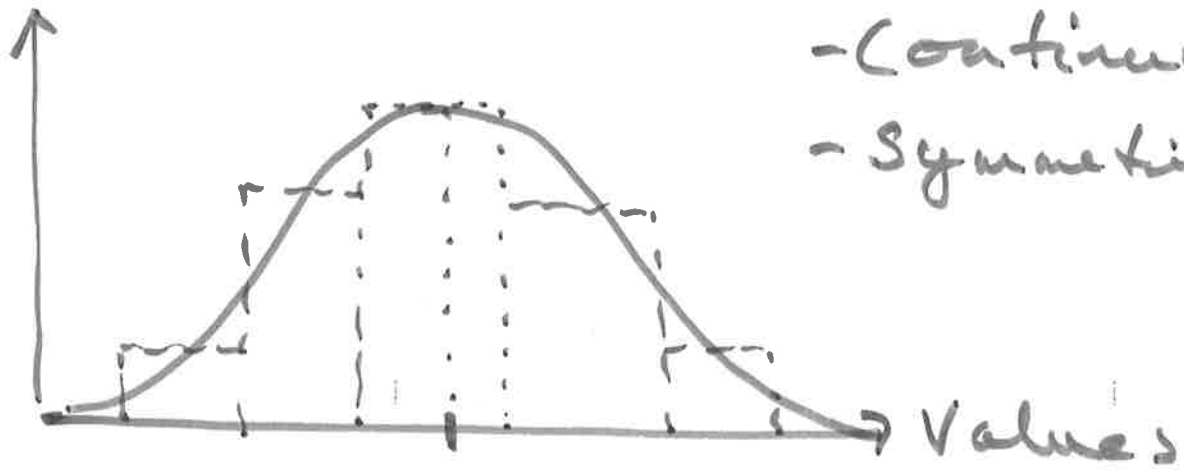Continuous Probability distributions

Rel. freq.



Values

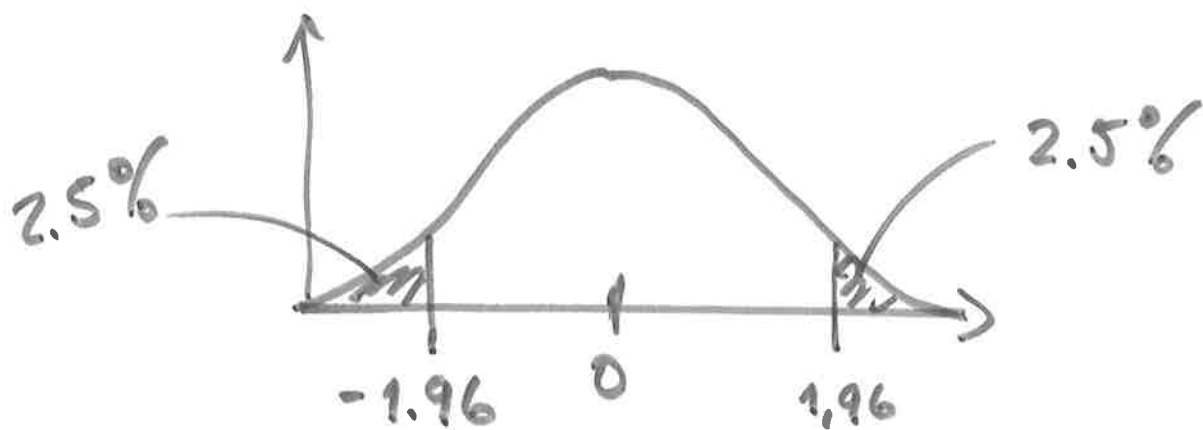Rel. freq.



Values



The normal distribution:

Rel. freq.



Mean

Values

- Continuous
- Symmetric

Some properties of the standard normal distribution:


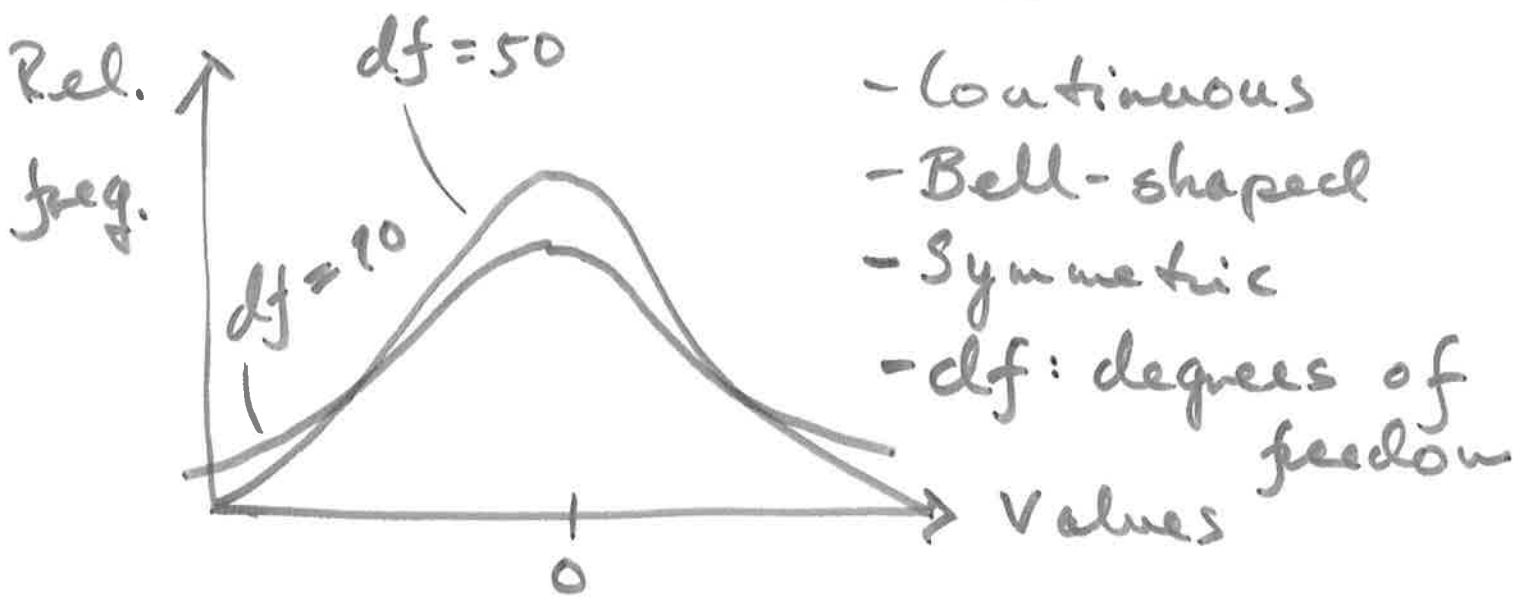
→ mean is 0

→ Area under the curved line is 1 or 100%

→ Symmetry:

- Area to the left of 0 is 0.5 or 50%

- Area to the right of 0 is 0.5

- Area to the left of -1.96 is 2.5%

# The t-distribution :



- Continuous
- Bell-shaped
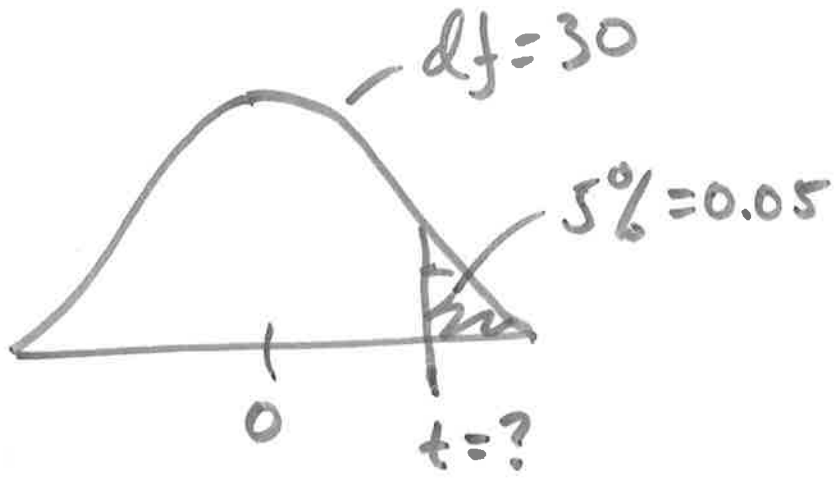- Symmetric
- df: degrees of freedom

Some properties :

→ There is in fact an infinite number of t-distributions, one for each df

→ df is usually approximately equal to the number of observations

→ The "flatness" of the t-distribution depends on df

$\rightarrow$ When $df \rightarrow \infty$, then the t-distribution becomes a standard normal distribution
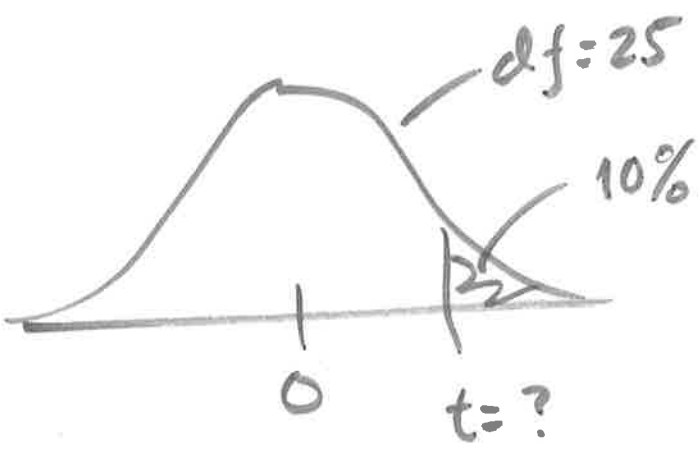
Examples:

What is t?

$$t_{0.05}(30) = 1.697$$



df = 30

5% = 0.05

0    t = ?

What is t?

$$t_{0.10}(25) = 1.316$$

df = 25

10%

0    t = ?

What is t?

$$t_{0.05}(20) = 1.725$$

df = 20

5%

0    t = ?

# Other important probability distributions:

- Chi-squared dist.
- F-distribution

}

Rel. freq.



0 → Values

## ④ Hypothesis testing



Pop.

... — ☁
Sample

**Def. Hypothesis:** A claim about a population

**Def. Null hypothesis $(H_0)$:** A claim about the population in terms of an equality

Examples:

$\rightarrow \mu = 20 : H_0$ (population mean equal to 20)

$\rightarrow H_0 : \mu = 9$

$\rightarrow H_0 : \mu = 2$

Def. <u>Alternative hypothesis $(H_A)$:</u>

A claim about the population of the "not $H_0$" type

Examples:

$\rightarrow H_A : \mu \neq 20 \qquad H_A : \mu > 20 \qquad H_A : \mu < 20$

$\rightarrow H_A : \mu \neq 9 \qquad H_A : \mu > 9 \qquad H_A : \mu < 9$

$\rightarrow$

NOTE: The investigator chooses
$H_0$ and $H_A$

Some conventions:

→ $H_A$ is usually the claim you would like to test or investigate

→ $H_0$ is what previous knowledge or findings suggest

→ Moral and/or health-reasons that come into play

<u>Def.</u> <u>Significance level ($\alpha$)</u>: The probability of rejecting $H_0$ when it is true

NOTE: The investigator chooses $\alpha$

Examples: $\alpha = 0.01 = 1\%$ } most common ones
$\alpha = 0.05 = 5\%$
$\alpha = 0.10 = 10\%$

**Def. Test expression (statistic):**

A formula whose value indicates whether to keep or reject $H_0$

Test expressions that are t-distributed typically have the following form:
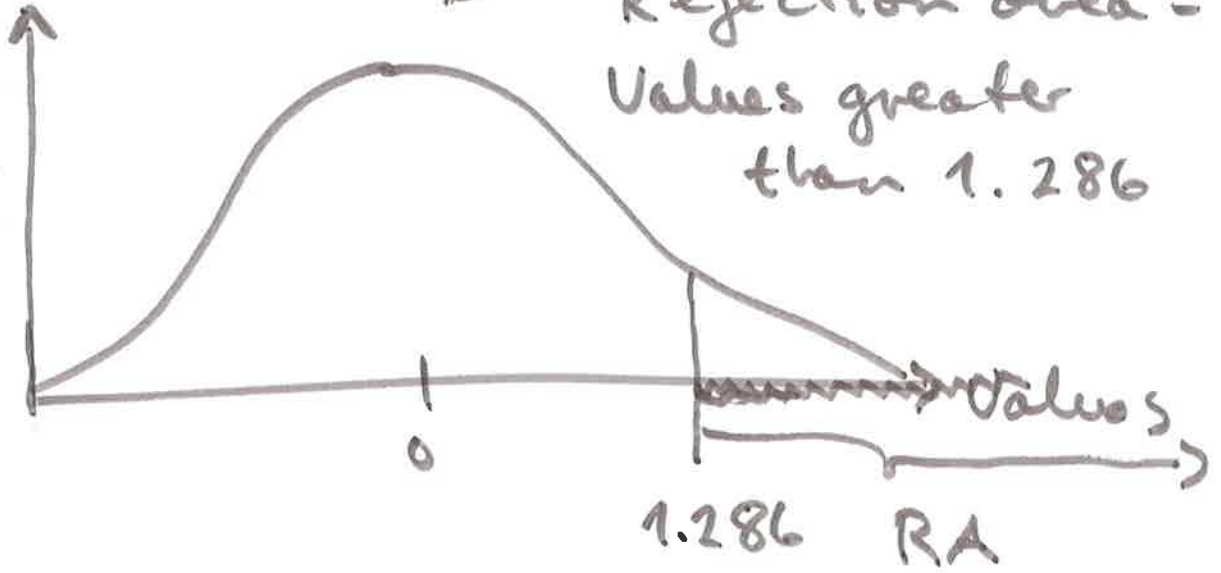
$$\overset{\text{sample estimate}}{\nearrow}\quad \frac{\bar{X} - H_0 \text{ value}}{\underset{\nwarrow \text{ sample standard dev.}}{S}}$$

**Def. Rejection area:** The values of the test-expression that makes us reject $H_0$

Example:
Relative
freq.

t-distribution
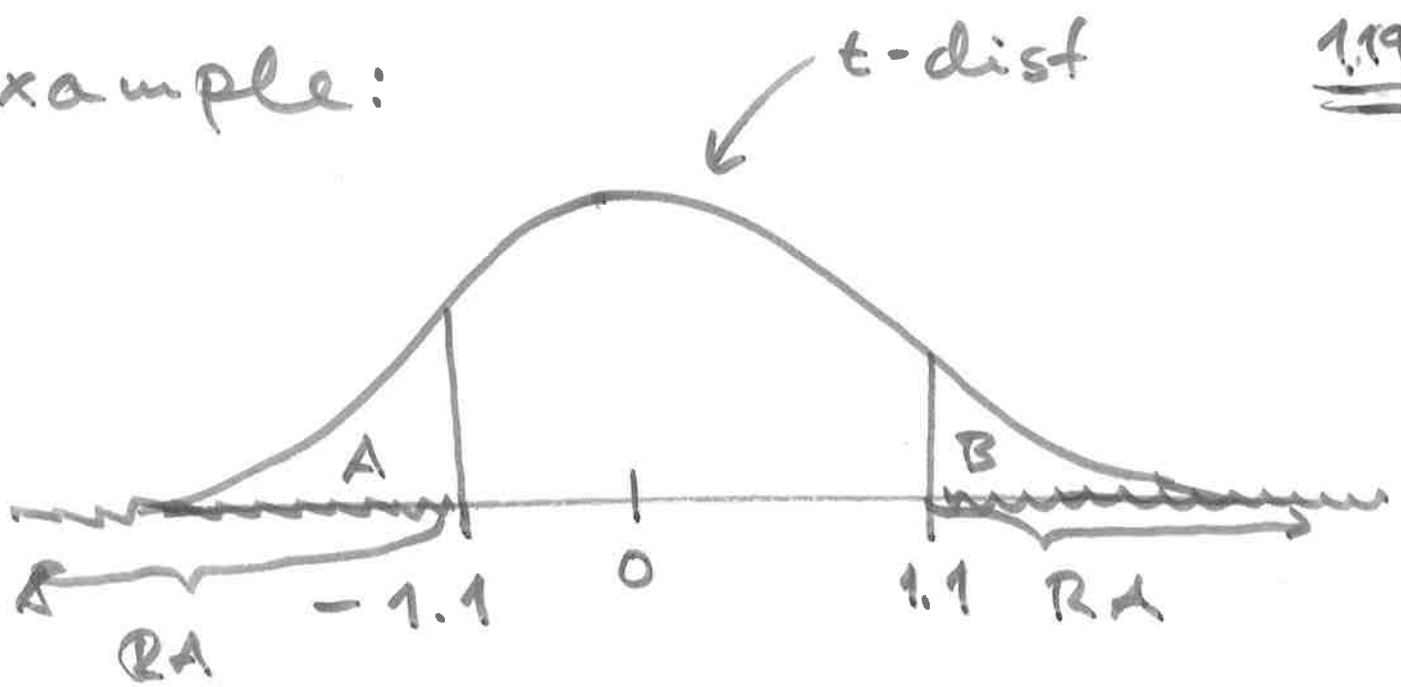Rejection area =
values greater
than 1.286



0

1.286    RA

values

Def. Critical value(s): The boundary (or boundaries) of the rejection area

Example: Above (i.e. previous example) we have 1 critical value, and it is equal to 1.286

Example:

t-dist



A   B

R   -1.1   0   1.1   RA
RA

Rejection area = Values greater than 1.1 and values smaller than -1.1. Critical values = 1.1 and -1.1.

Testing in 4 steps:

Step 1: Choose $\alpha$, $H_0$ and $H_A$

2: Find the critical value(s) and identify the rejection area

3: Compute the value of the test-expression

4: Conclude: Reject $H_0$ if the test-value lies in the rejection area, otherwise keep $H_0$

Type 1 and 2 errors:

|  | $H_0$ True | $H_A$ true |
|---|---|---|
| Reject $H_0$ | Type 1 error | Correct! |
| Keep $H_0$ | Correct! | Type 2 error |

## Testing the value of a population mean

Step 1: Choose $\alpha$, $H_0$ and $H_A$

2: Identify the rejection area using a t-distribution with $n-1$ df

3: Compute the test-expression

$$\frac{\bar{x} - H_0 \text{ value}}{s/\sqrt{n}}$$

4: Conclude: Reject $H_0$ if the test-value lies in the rejection area, otherwise keep $H_0$

Example: Travel time (in minutes)

Claim: Average travel time is greater than 20 minutes

Sample: 20, 15, 17, 55, 20, 10, 20, 90, 15 18, 40

$n = 11$         $\bar{x} = 29.09$         $s = 24.0394$
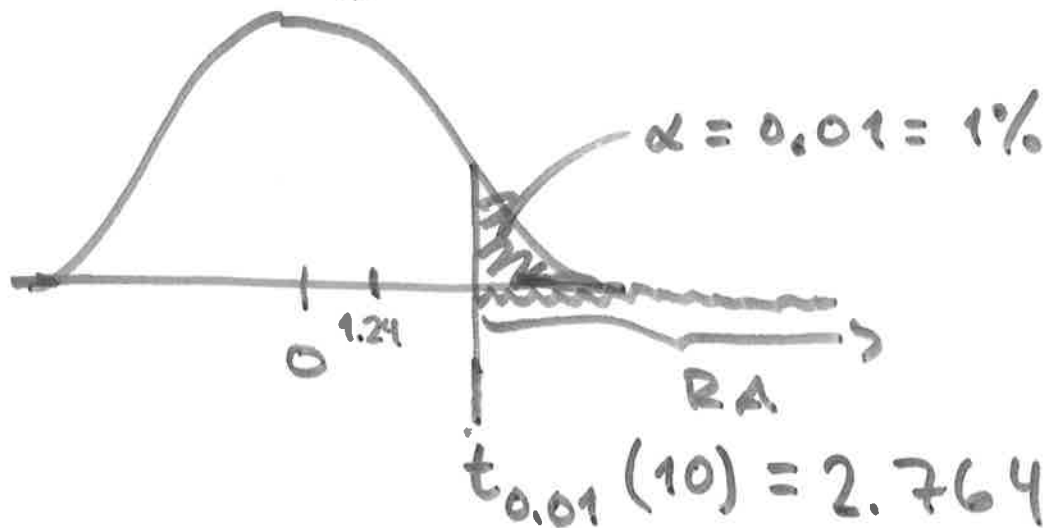
Step 1: $\alpha = 1\%$     $H_0: M = 20$

$H_A: M > 20$

2: Identify the rejection area:

$df = n - 1 = 11 - 1 = 10$



$\alpha = 0.01 = 1\%$

$t_{0.01}(10) = 2.764$

3: Compute the test-value:

$$29.09 \rightarrow \underbrace{\frac{\bar{x} - H_0 \text{ value}}{s / \sqrt{n}}}_{} = 1.2417$$

$20$

$24.0394 \rightarrow$ $s / \sqrt{n}$ $\leftarrow 11$

4: Conclusion: We keep $H_0$. That is,

we have found support in
favour of the claim that the
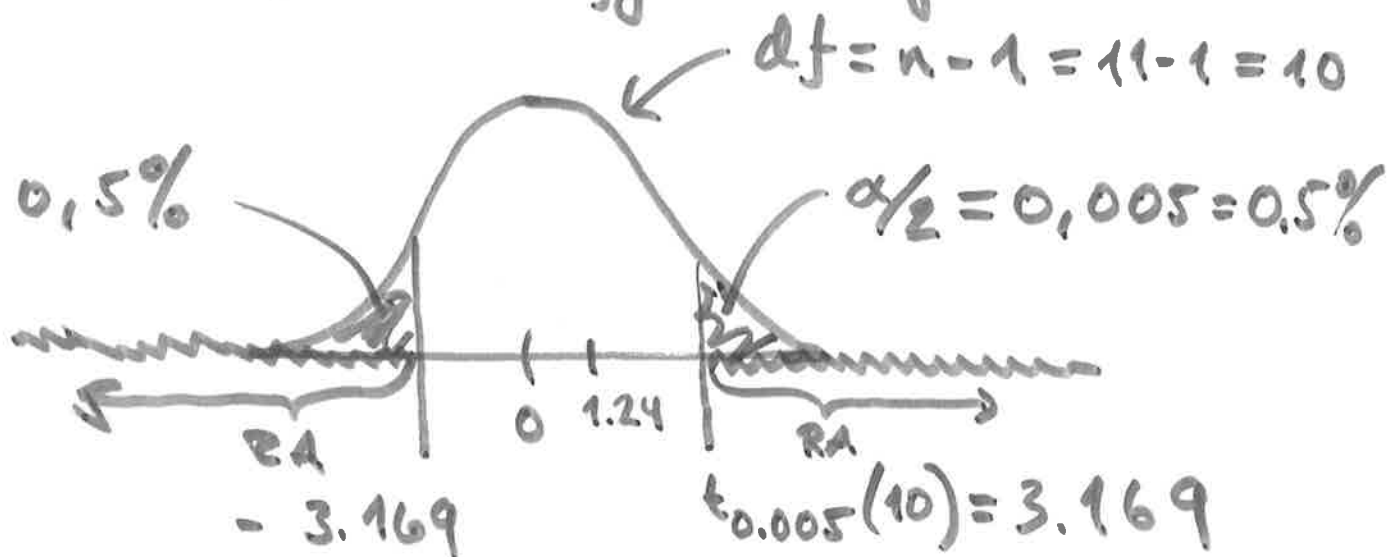average travel time is
20 minutes

Example:

- Claim: Average travel time is
  different from 20 minutes

- Sample: Same as previous

Step 1: $\alpha = 0.01 = 1\%$    $H_0: \mu = 20$

$H_A: \mu \neq 20$

2: Identify the rejection area:

$df = n - 1 = 11 - 1 = 10$

0,5%                    $\alpha/2 = 0,005 = 0.5\%$

RA                    RA
0  1.24
-3.169          $t_{0.005}(10) = 3.169$

RA: Values higher than 3.169 and values
lower than -3.169

3: Value of test-expression:

$$29.09 \nearrow \quad \frac{\bar{X} - H_0 \text{ value} \quad \leftarrow 20}{s/\sqrt{n} \quad \leftarrow 11} \quad = 1.2417$$

$$24.0394 \rightarrow$$

4: Conclusion: We keep $H_0$

# (5) P-values

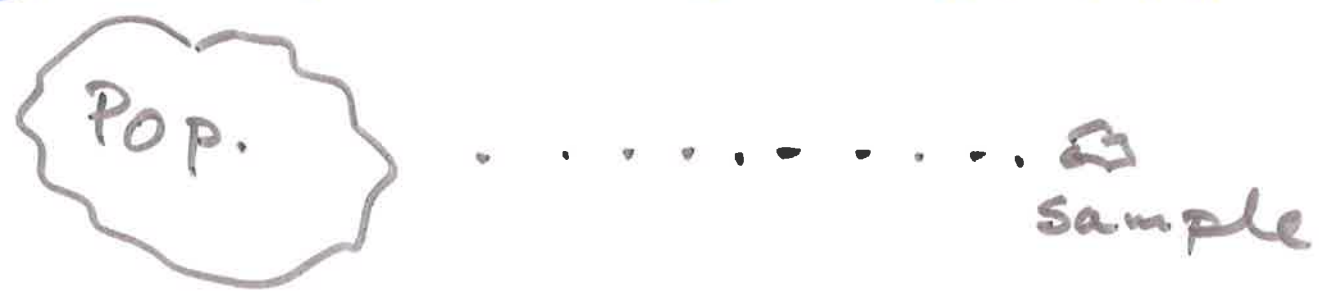Intuition: A number between 0 and 1 that summarises the information from a hypothesis test

Def. P-value: The smallest significance level at which the null hypothesis can be rejected

In practice: Reject $H_0$ if the p-value is smaller than $\alpha$

↑ The p-value game
    & dance !

NOTE: In order to compute a p-value "by hand", $H_0$ and $H_A$ must already have been defined, and the test-value must already have been computed

(b) Interval estimation



POP. . . . . . . . . . . . sample

$$\left.\begin{matrix} \mu \\ \sigma^2 \\ \sigma \end{matrix}\right\} \text{sample counterparts:} \begin{matrix} \bar{X} \\ s^2 \\ s \end{matrix}$$

Def. $(1-\alpha)\cdot 100\%$ confidence interval: An interval that contains the population value of interest (e.g. $\mu, \sigma^2, \sigma$, etc.) with $(1-\alpha)\cdot 100\%$ degree of certainty

Example: 90% confidence interval for the mean time it takes from you wake up until you leave the house (Q2)

Sample = 110, 45, 75, 90, 70, 90, 40, 30, 45, 60

$n = 10$     $\bar{x} = 65.5$     $s_x^2 = 674.72$     $s_x = 25.98$

$$L = \bar{x} - t_{\alpha/2}(df) \cdot s_x/\sqrt{n} = 50.43$$

65.5          0.05      10-1=9          10

1.833

$$U = \bar{x} + t_{\alpha/2}(df) \cdot s_x/\sqrt{n} = 80.57$$

Interpretation: The interval [50.43, 80.57] contains the population mean with 90% degree of certainty
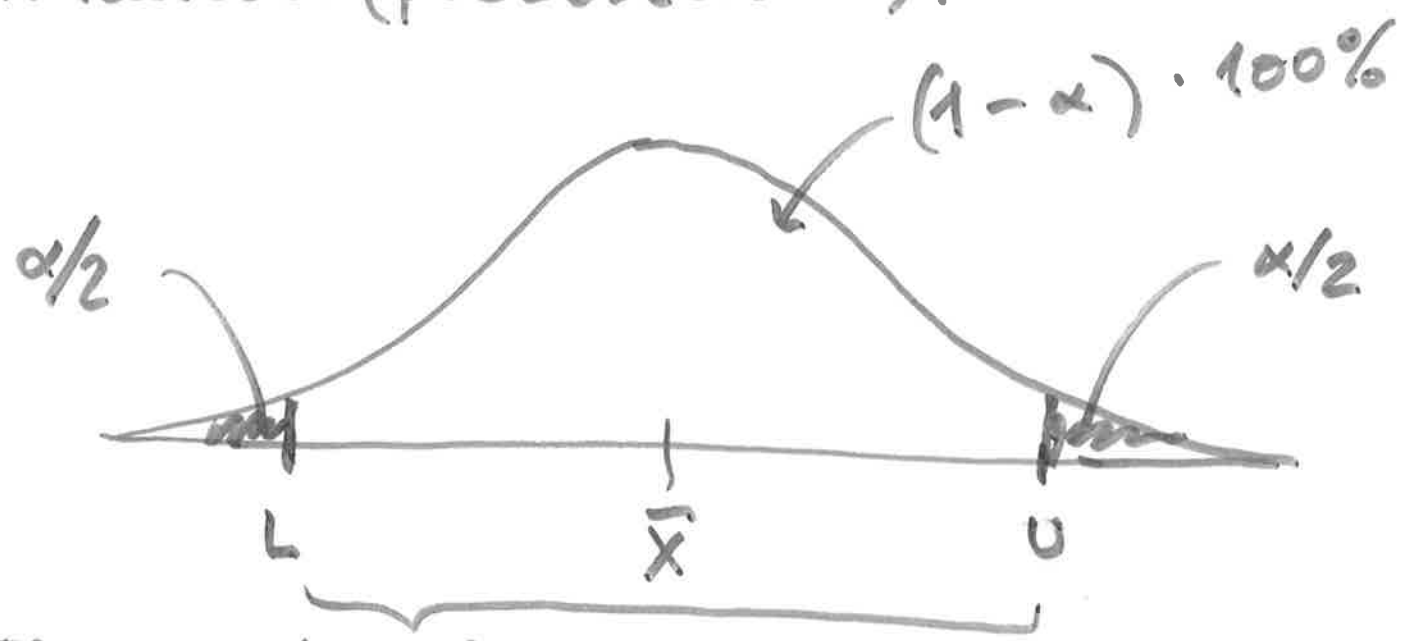
Example: Suppose [6, 11] is a 95% confidence interval for $\mu$, then $\mu$ lies between 6 and 11 with 95% degree of certainty

## Computation of confidence interval for $\mu$

Lower bound $= \bar{x} - t_{\alpha/2}(df) \cdot S_x / \sqrt{n}$

$$\uparrow_{n-1}$$

Upper bound $= \bar{x} + t_{\alpha/2}(df) \cdot S_x / \sqrt{n}$

Intuition (probabilistic):



$(1-\alpha) \cdot 100\%$

$\alpha/2$

$\alpha/2$

L

$\bar{x}$

U

This is the $(1-\alpha) \cdot 100\%$ confidence interval

# (7) Suggested exercises

Exercise set 1: 1a)–g), 4a)iv),
4b)iv), 4c)i), 5a)i), 5a)iii),
5b)i), 5b)iii), 5c)i), 7