

# Analysis of technical and economic<sup>al</sup> data

- What is the subject all about?

↳ what kind of consequences can be drawn from measured or observed data?

↳ how can we "compress" information?

( measurement with 1000 data  $\mapsto$  how can we compress it and get valuable information? )

- The way of teaching the subject

↳ usually: basic definitions, statements, methods, applications

↳ in this subject:

- simple methods of data processing

↳ these require no new mathematical knowledge

↳ • just a few new definitions: descriptive statistics

↳ we will pose some questions, they will be answered

later  
↓

the general form of these answers will lead to the axioms of probability theory and the basis of statistics

## Guideline of the subject:

- descriptive statistics (2h)
  - probability theory and basics of statistics (2h)
  - applications in measurement techniques (1h)  
(error definitions, direct and indirect measurement)
  - confidence interval (1h)
  - quality control and reliability (1h)
  - statistical tests (parametric; non-parametric) (2h)
  - analysis of variance (1÷2h)
  - relationship between variables
    - ↳ correlation coefficient (2÷3h)
    - ↳ regression analysis
    - ↳ Wald method
  - Conclusion (1h)
- 
- 14h

## Elementary methods of data processing

- let's do some data collection: Mathematics A2 ~~no~~ marks

mark	$V_i$	$\frac{V_i}{n}$
$\overset{!}{z}_2 = 2$	250	0,5
$\overset{!}{z}_3 = 3$	150	0,3
$\overset{!}{z}_4 = 4$	70	0,14
$\overset{!}{z}_5 = 5$	30	0,06
$\Sigma:$	500 "n	1

$V_i$ : frequency

$\frac{V_i}{n}$ : relative frequency

$$\sum_{i=1}^n V_i = n = 500$$

$$\sum_{i=1}^n \frac{V_i}{n} = 1 \quad \square$$

- the average can be calculated:

$$\bar{z} = \frac{1}{n} \cdot \sum_{j=1}^{500} z_j \quad ; \quad \text{where:}$$

$$z_j: \quad 2, 2, 5, 3, 2, 4, 2, 5, 4, \dots$$

$$j: \quad 1, 2, 3, \dots$$

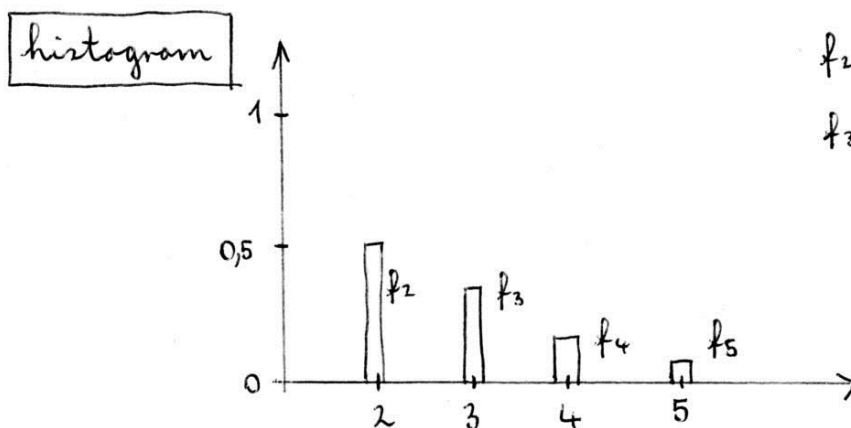
- let's create four sub-intervals, they should be disjoint  
(non-overlap)

↳ the sum of the frequencies should be equal to the  
sum of the individuals

$$\bar{z} = \frac{1}{n} \cdot \sum_{j=1}^{500} z_j = \frac{1}{n} \cdot \sum_{i=2}^5 v_i \cdot z_i = \sum_{i=2}^5 \frac{v_i}{n} \cdot z_i = 0,5 \cdot 2 + 0,3 \cdot 3 + 0,14 \cdot 4 + 0,06 \cdot 5 = 2,76$$

↑  
changing the index!

- the average carries not enough information, let's show the distribution  
of the marks:



$$f_2 = \frac{v_2}{n} \quad f_4 = \frac{v_4}{n}$$

$$f_3 = \frac{v_3}{n} \quad f_5 = \frac{v_5}{n}$$

- the mark is a discrete variable, here we get a natural order

- if the variable is non-numerical but a quality or categorical variable

↳ how can we construct a  $\mathcal{M}$  histogram? (car types, countries, etc.)

↓  
**PARETO DIAGRAM** (Vilfredo Pareto) In this case the variables don't define an exact order

- in case of the Pareto diagram we begin with the highest value!

## Continuous variables

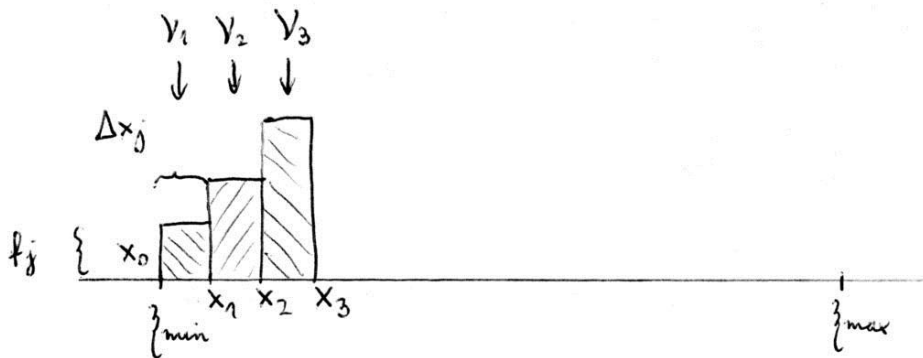
- for creating a histogram the data has to be assigned to groups

data:

$$\{x_1\} < \{x_2\} < \dots < \{x_n\} \quad (\text{ordered data})$$

$$\{x_{\min} = x_0 = \min \{x_1, \dots, x_n\}$$

$$\{x_{\max} = x_{15} = \max \{x_1, \dots, x_n\}$$



$$I = \{x_{\max} - x_{\min} \rightarrow \text{data range}$$

- creating discrete intervals:  $[x_0, x_1); [x_1, x_2); [x_2, x_3); \dots; [x_{15-1}, x_{15}]$

$I_5$ : the number of intervals

- let's calculate the frequencies:  $v_1, v_2, \dots, v_{I_5}$

$$\text{relative frequencies: } \frac{v_1}{n}, \frac{v_2}{n}, \dots, \frac{v_{I_5}}{n}$$

$$f_j \cdot \Delta x_j = \frac{v_j}{n} \rightarrow \text{area of the rectangle}$$

- how can we select  $IS$ ?

↳ if  $IS$  is too high, the histogram gets disturbed

↳ if  $IS$  is too low, the nature of the data is spoiled

therefore:

$$n \leq 100$$

$$IS \cong \sqrt{n}$$

$$n \geq 100$$

$$IS \cong \log_2 n + 1$$

↳ "logarithm to base two"  
of  $n$

$n$	$\sqrt{n}$	$\log_2 n + 1$
100	10	8
200	14	9
500	22	10
1000	32	11

How should we determine  $\Delta x_j$ ?

① uniform subdivision

↳ basic level (secretary, Excel)

② non-uniform subdivision

↳ the goal is, that every  $\Delta x_j$  should contain the same amount of data ( $V_j$ )  
nearby const.

Example:

body heights  $\mapsto$  the students should guess the height of the teacher!

Queries:  $\overbrace{170; 170; 173; 173; 175; 175; 175; 176; 177; 178; 178; 181; 183; 185}$

$$n = 14$$

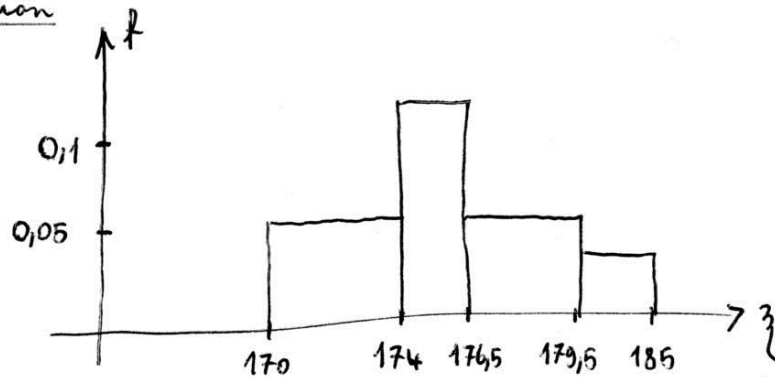
$$x_{\min} = 170$$

$$IS = \sqrt{n} \cong 4$$

$$x_{\max} = 185$$

$x_{Li}$	$x_{Ui}$	$v_i$	$v_i/n$	$\Delta x_i$	$f_i$
[ 170	174 )	4	0,286	4	0,071
[ 174	176,5 )	4	0,286	2,5	0,114
[ 176,5	179,5 )	3	0,214	3	0,071
[ 179,5	185 ]	3	0,214	5,5	0,039
	$\Sigma:$	14	1		

### Empirical density function



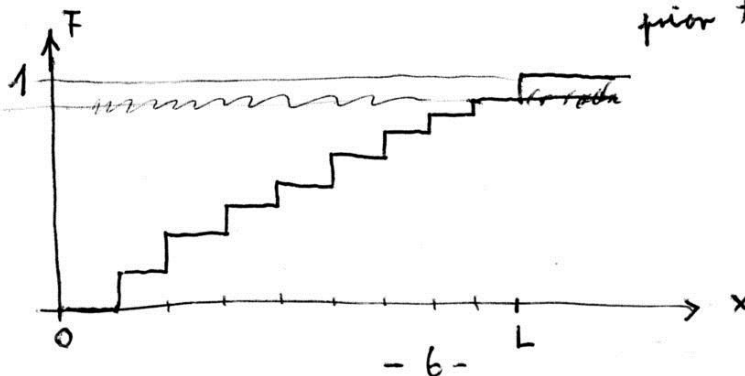
### Questions:

- how does the relative frequency behave if  $n \rightarrow \infty$
- how does the histogram behave if  $n \rightarrow \infty$  (answers will follow later)

2012.08.06

### Cumulative distribution function (CDF) $\rightarrow$ empirical!

- what percent of a data-series is below a certain value?
- pulling a homogenous beam:
  - several measurements are carried out  $\rightarrow$  in what percent of the measurements does the beam break for example prior to  $\frac{4}{5} \cdot L$ ?

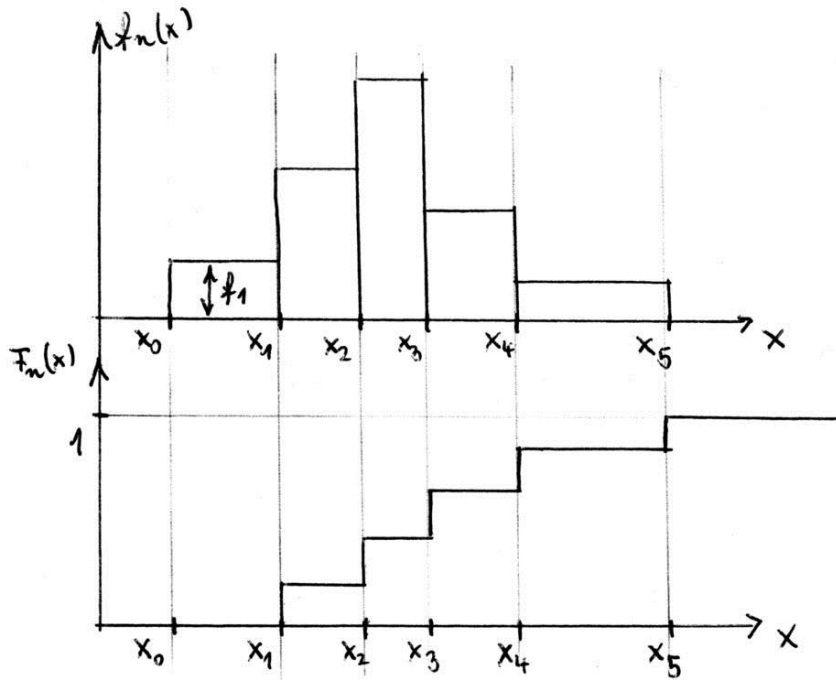


$$F_n(x) = \frac{h}{n} \quad ; \quad \text{with} \quad \{x_2 \leq x \leq x_{2+1}\}$$

$\{1\}, \{2\}, \dots, \{n\}$  (sample units)

- the CDF shows the relative frequency of the sample units left from  $x$ .

- What is the relationship between the empirical density function and the empirical cumulative distribution function?



$x_0 - x_1$	$v_1$
$x_1 - x_2$	$v_2$
$x_2 - x_3$	$v_3$
$x_3 - x_4$	$v_4$
$x_4 - x_5$	$v_5$

$$f_i = \frac{v_i}{n \cdot \Delta x_i}$$

where:

$$F_n(x_1) = \frac{v_1}{n}$$

$$F_n(x_2) = \frac{v_1}{n} + \frac{v_2}{n}$$

$\vdots$

$$F_n(x_5) = \frac{1}{n} \cdot (v_1 + v_2 + \dots + v_5)$$

general formula:

$$F_n(x_h) = \sum_{x_i \leq x_h} f_n(x_i) \cdot \Delta x_i$$

(!)

↳ the sum of the areas left from the actual point

### Question:

- what happens with  $\left. \begin{matrix} f_n(x) \\ F_n(x) \end{matrix} \right\}$  if  $n \rightarrow \infty$ ?
- what kind of relationship is between  $F_n(x)$  and  $f_n(x)$  if  $n \rightarrow \infty$ ?

↓

the answer comes later on!

### Median

- the median is the "middle of the sample". Half of the sample is above, half below it.

• let's have a measurement:  $\{z_1, z_2, \dots, z_n\}$

• arrange the data:  $z_1^* \leq z_2^* \leq \dots \leq z_n^*$  "ordered sample"

① if  $n$  is odd

$$m = \frac{n+1}{2}; \quad m \text{ must be even now}$$

$$\hookrightarrow \tilde{z} = z_m^*, \text{ the median}$$

② if  $n$  is even

$$m_1 = \frac{n}{2}$$

$$\hookrightarrow \tilde{z} = \frac{1}{2} (z_{m_1}^* + z_{m_1+1}^*)$$

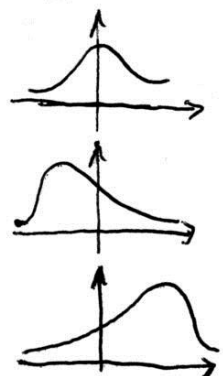
$\sim$  : "tilde"

- relationship between median and average:

• if the histogram is about symmetrical  $\Rightarrow \tilde{z} \approx \bar{z}$

• if  $-||-$  leans to the left  $\rightarrow \tilde{z} < \bar{z}$

$-||-$  leans to the right  $\rightarrow \tilde{z} > \bar{z}$

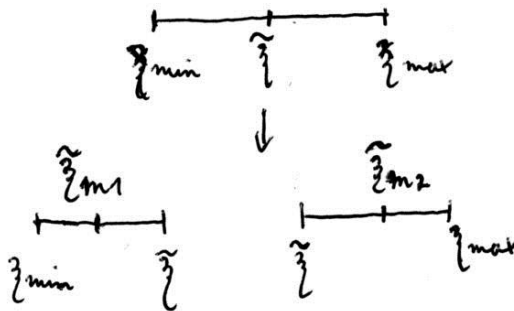




## Lower and upper quartiles

- the median splits the sample into two parts (upper and lower)

↳ let's determine the medians of the upper and lower parts!



example:  $n=19$

$z_1^* \leq z_2^* \leq \dots \leq z_{19}^*$      $\mapsto$  the median is  $\tilde{z} = z_{10}^*$

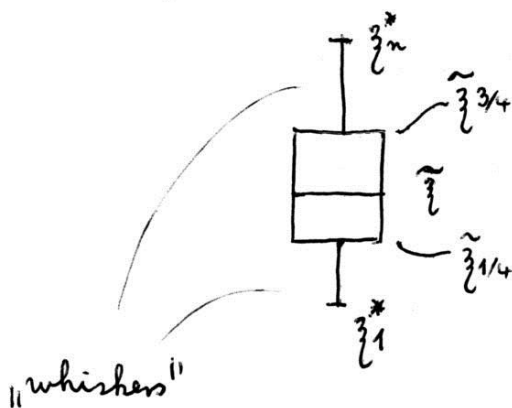
lower quartile:  $\tilde{z}_{1/4} = z_5^*$

upper quartile:  $\tilde{z}_{3/4} = z_{15}^*$

how can we display the median and quartiles?

**Boxplot**

- global representation of the measurement sample



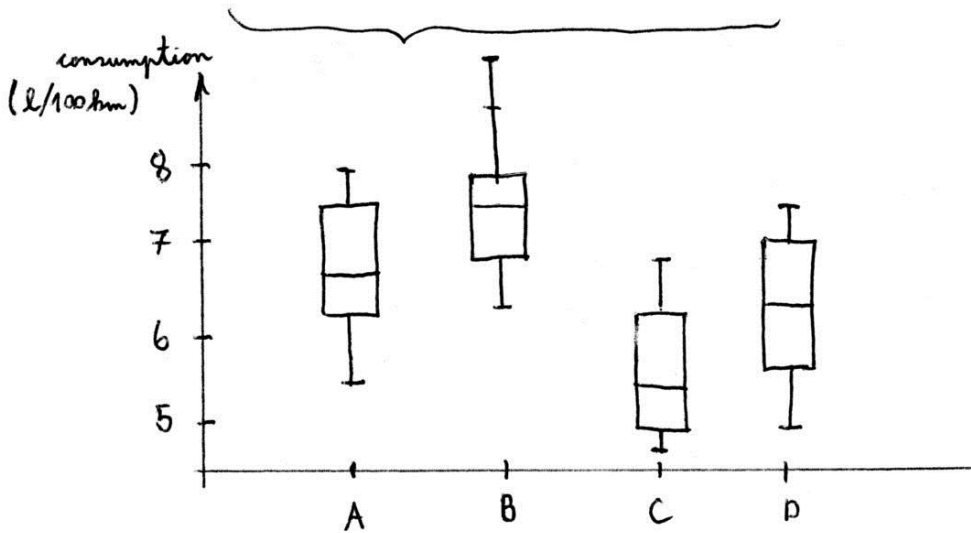
advantage of using box plots:

- different samples can be compared easily

example:

- comparison of fuel consumption of different car types

	$\tilde{x}$	$\tilde{x}_{1/4}$	$\tilde{x}_{3/4}$	$\tilde{x}_{max}$	$\tilde{x}_{min}$
A					
B					
C					
D					



observations from the ~~the~~ box-plot:

- fuel consumption of C is significantly lower, than the others
- type B can consume very much fuel under certain circumstances



how to deal with the outliers?

## Outliers

- definition: "Inter-quartile range" = IQR

$$IQR = \tilde{z}_{3/4} - \tilde{z}_{1/4}$$

- the data is an outlier if it is further away from the box than  
~~1,5~~  $1,5 \cdot IQR$

- extreme outlier, if further than  $3 \cdot IQR$

calculation:

• upper outliers:  $z > \tilde{z}_{3/4} + 1,5 \cdot IQR$

• upper extreme outliers:  $z > \tilde{z}_{3/4} + 3 \cdot IQR$

• lower outliers:  $z < \tilde{z}_{1/4} - 1,5 \cdot IQR$

• lower extreme outliers:  $z < \tilde{z}_{1/4} - 3 \cdot IQR$

example  $\rightarrow$  we will see it in the practice !!

---

## Standard deviation of the sample

- how can we describe the fluctuation of the data?

variation:  $z_i - \bar{z}$ , for every "i" point

let's calculate the average variation:

$$\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z}) = \frac{1}{n} \sum_{i=1}^n z_i - \frac{1}{n} \sum_{i=1}^n \bar{z} = \bar{z} - \frac{1}{n} \cdot n \bar{z} = 0 \quad \square$$

- something is wrong

↳ the variation  $(z_i - \bar{z})$  might be +/-, therefore  
the average is zero

- what can we use instead?

$$\sum |z_i - \bar{z}| \quad \text{or} \quad \boxed{\sum (z_i - \bar{z})^2}$$

↳ the absolute value is not ideal, because it is  
difficult to ~~the~~ differentiate!

therefore:

$$s^2 = \frac{1}{n} \cdot \sum_{i=1}^n (z_i - \bar{z})^2 \quad ; \quad \text{"standard deviation of the sample"}$$

but the formula contains  $(n-1)$  independent terms instead of  $n$

/ because  $\sum (z_i - \bar{z}) = 0$  /

⇓

sample standard deviation:

$$s^{*2} = \frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})^2 \quad \text{"bracket"}$$

↳ this is corrected, adjusted

$$s^* = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})^2}$$

• the unit is the same as the original variable!

- a different formula might be used also:

$$\begin{aligned}
 s^{*2} &= \frac{1}{n-1} \cdot \sum_{i=1}^n (z_i - \bar{z})^2 = \frac{1}{n-1} \sum (z_i^2 + \bar{z}^2 - 2 \cdot z_i \bar{z}) \stackrel{\nabla}{=} \\
 &= \frac{1}{n-1} \cdot \left[ \sum z_i^2 + n \cdot \bar{z}^2 - 2 \cdot \underbrace{\sum z_i \bar{z}}_{2 \cdot n \cdot \bar{z}^2} \right] = \frac{1}{n-1} \left( \sum z_i^2 - n \cdot \bar{z}^2 \right) = \\
 &= \frac{n}{n-1} \left[ \frac{1}{n} \cdot \sum z_i^2 - \bar{z}^2 \right] \geq 0
 \end{aligned}$$

"the avg. of the squares is always higher than the square of the avgs!"

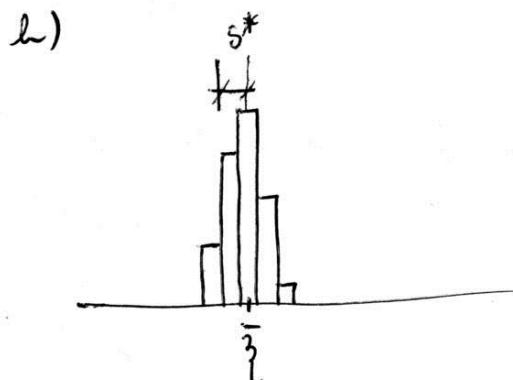
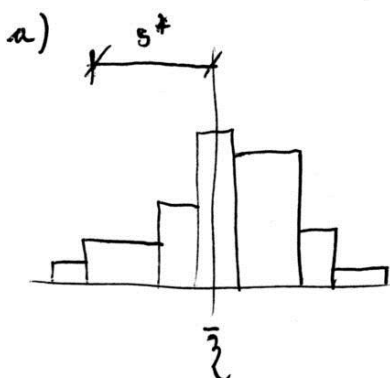
(Schwartz - Cauchy - Bunyakovski)

- comparing  $s^*$  and  $\tilde{z}_{1/4}$  ;  $\tilde{z}_{3/4}$ :

- few outliers increase  $s^*$  significantly, but  $\tilde{z}_{1/4}$  and  $\tilde{z}_{3/4}$  are not influenced
- if outliers are present, the box plot gives better information

practical use:

$\bar{z}$  and  $s^*$  describe the measurement data



- let  $x$  be the expected value of a production process

( $D=20$  mm diameter,  $V=0,5$  l beer, 1 kg bread, 500 g detergent)

- let's make measurements; draw a histogram, calculate  $\bar{x}$  and  $s^*$

if  $|\bar{x} - x| = H > 0 \rightarrow$  adjustment problems

if  $s^*$  is high  $\rightarrow$  inaccurate production

### Data acquisition

- measurement, observation
- data sampling

#### ① measurement, observation

- globally at similar conditions performed or occurred processes, their parameters are recorded.

#### ② data sampling

- the goal is to form a statement on the multitude / population

$\hookrightarrow$  but we do not want to investigate every single item

- it would be expensive, time consuming
- there are cases, where the item is destroyed during investigation (stress test, tensile test)

$\uparrow$  "nahitósóba"

- a sample should be taken from the population

$\hookrightarrow$  the statements for the sample should be true for the whole population

- in order to do this, the following points have to be ensured:

- the population should be homogeneous regarding the investigated property  
(every item should have the same properties)
- every item should get into the sample with the same probability

---

- later on we will state that the sample is a collection of independent items with the same distribution

↳ some distribution is ensured by the correct definition of the population  
(same production line, same rice)

↳ the identical probability of selection is ensured with random sampling.  
(either with or without putting the item back)

---

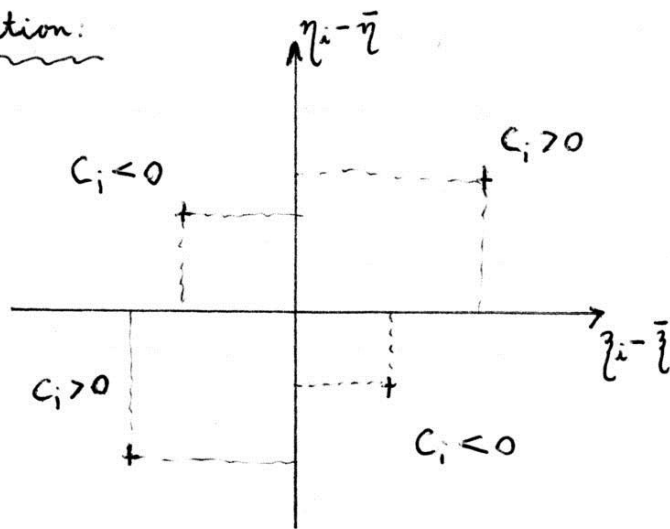
### Empirical correlation coefficient

- to this point we only dealt with one variable

↳ let's assume two variables. Is there a relationship between them?

$\left. \begin{array}{l} z_1; z_2; \dots; z_n \\ \eta_1; \eta_2; \dots; \eta_n \end{array} \right\} \text{pairing the } n \text{ points:}$   
 $(z_1, \eta_1); (z_2, \eta_2); \dots; (z_n, \eta_n)$

graphical representation:



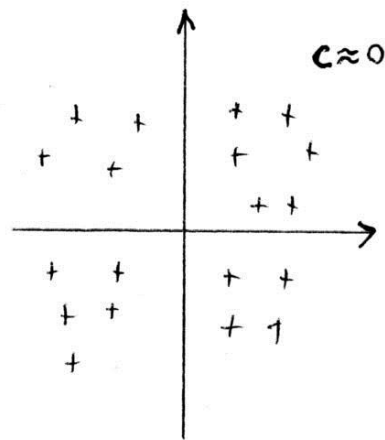
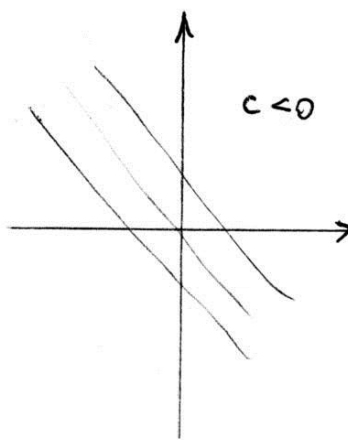
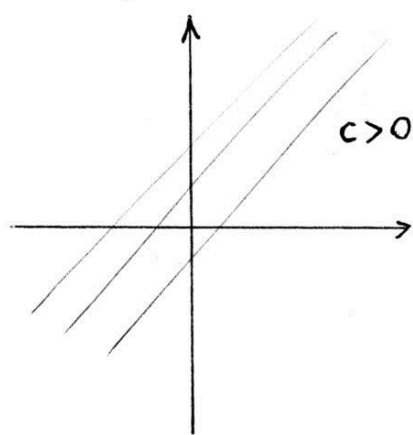
- let's define  $c_i$ :

$$c_i = (z_i - \bar{z})(\eta_i - \bar{\eta})$$

- create the sum:

$$c = \sum_{i=1}^N c_i = \sum_{i=1}^N (z_i - \bar{z})(\eta_i - \bar{\eta})$$

- what will we get for  $c$ ?



$c$  depends on:

• amount of points: let's take the average,  $\frac{1}{N} \cdot c$

• standard deviation of the variables: let's normalize,  $\frac{1}{N} \cdot c \cdot \frac{1}{s_z} \cdot \frac{1}{s_\eta}$

↓

$$s(z, \eta) = \frac{\frac{1}{N} \cdot \sum (z_i - \bar{z})(\eta_i - \bar{\eta})}{s_z \cdot s_\eta}$$



$S(\xi, \eta)$ : empirical correlation coefficient

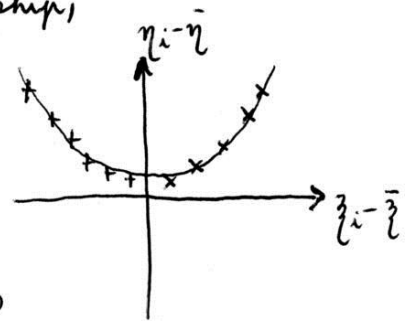
- we can determine that:

- if  $S(\xi, \eta) > 0$ ; the relationship is progressive
- if  $S(\xi, \eta) < 0$ ; —||— is degressive
- if  $S(\xi, \eta) \approx 0$ ; no relationship

questions:

- what happens to  $S(\xi, \eta)$  if  $n \rightarrow \infty$ ?
- what values can  $S(\xi, \eta)$  have? ( $|S| \leq 1$ )
- does it really represent a general relationship, or less than that?

↳ counterexample:  $\eta = \xi^2$



$S \approx 0$ , but the relationship is obvious!

### Rank correlation

- can we simplify the formula for correlation, if the observed variable is discrete?

Example:

- a jury of two judge something (wines, ice skating, sausage competition)

the contenders: A, B, C, D

	A	B	C	D
$\xi$ : J1	1	4	2	3
$\eta$ : J2	3	4	1	2

→ ranking of the jury

- question: are the juries in concordance?

jury J1:  $\zeta_1; \zeta_2; \zeta_3; \zeta_4$   
 J2:  $\eta_1; \eta_2; \eta_3; \eta_4$  } calculate  $\bar{\zeta}; \bar{\eta}; s_{\zeta}; s_{\eta}$

$$r(\zeta, \eta) = \frac{\frac{1}{N} \cdot \sum (\zeta_i - \bar{\zeta})(\eta_i - \bar{\eta})}{s_{\zeta} \cdot s_{\eta}}$$

↑  
rank correlation coeff

if  $r \approx 1 \rightarrow$  concordance

if  $r \approx -1 \rightarrow$  opposite opinion

- correlation between the ranks = rank correlation  
 (Spearman)

- the formula can be simplified:

$$d_1 = \zeta_1 - \eta_1$$

$$d_2 = \zeta_2 - \eta_2$$

$$r = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

in the example:

$$d_1 = 1 - 3 = -2$$

$$d_2 = 4 - 4 = 0$$

$$d_3 = 2 - 1 = 1$$

$$d_4 = 3 - 2 = 1$$

$$\sum d_i^2 = 4 + 0 + 1 + 1 = 6$$

$$r = 1 - \frac{6 \cdot 6}{4(16 - 1)} = 0,4$$

- if the jury has more than two members:

2012.08.20.

Kendall's W or Kendall's concordance

J1, J2, ..., J<sub>NJ</sub>

NJ: number of juries

C1, C2, ..., C<sub>NC</sub>

NC: number of contenders

	J1	J2	J3	...	J <sub>NJ</sub>	Σ
C1	r <sub>11</sub>	r <sub>21</sub>	r <sub>31</sub>		r <sub>NJ1</sub>	R <sub>1</sub>
C2	r <sub>12</sub>	r <sub>22</sub>				R <sub>2</sub>
⋮	r <sub>13</sub>					⋮
⋮						
C <sub>NC</sub>	r <sub>1NC</sub>				r <sub>NJNC</sub>	R <sub>NC</sub>

$R_i$ : the sum of  $r_{ij}$  for the " $i^{\text{th}}$ " contender; this is the resultant rank number

$$R_1 = \sum_{j=1}^{NJ} r_{j1} \quad ; \quad R_2 = \sum_{j=1}^{NJ} r_{j2}$$

- let's arrange the resultant rank numbers:

$$R_1^* \leq R_2^* \leq \dots \leq R_{NC}^* \quad \mapsto \text{final order}$$

$$\bar{R} = \frac{1}{NC} \sum_{i=1}^{NC} R_i \quad ; \quad \text{let's calculate the squared deviation}$$

$$S^* = \sum_{i=1}^{NC} (R_i - \bar{R})^2$$

- what will be the value of  $S^*$ ?

↳ if the members of the jury are in concordance:

• the order is the same:

$$\left. \begin{array}{l} R_1^* = 1 \cdot NJ \\ R_2^* = 2 \cdot NJ \\ \vdots \end{array} \right\} \Delta R = NJ!$$

↳ if there is no concordance:

$\Delta R$  is getting smaller! (maybe 0)

the maximum of  $S$ :

(if  $\Delta R = NJ$ )

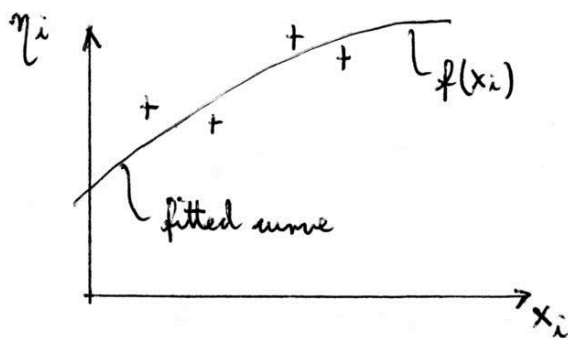
$$S_{\max} = \frac{NJ^2 \cdot NC(NC^2 - 1)}{12}$$

- the actual  $S$  should be normalized with  $S_{\max}$ :

$$W = \frac{12 \cdot \sum_{i=1}^{NC} (R_i - \bar{R})^2}{NJ^2 \cdot (NC^3 - NC)}$$

$0 \leq W \leq 1$  ; if there is concordance:  $W = 1$   
(absolute)

Approximation of the relationship



$[x_i, \eta_i]_{i=1}^n$  measurement

- the idea comes from Gauss (1777 - 1855) ; around 1795 LNM

(Daniel Kehlmann: Measuring the World)

$$\min. \left\{ \sum_{i=1}^n [\eta_i - f(x_i)]^2 \right\} \rightarrow \text{that is what we are looking for!}$$

-  $f(x)$  should possess some free parameters

for example:

$$f(x) = d_0 + d_1 x + d_2 x^2 \rightarrow \text{non-linear function}$$

let's determine:

$$\min \left\{ \sum_{i=1}^n (\eta_i - d_0 - d_1 x_i - d_2 x_i^2)^2 \right\} = D(d_0, d_1, d_2)$$

$$\left\{ \eta_i, x_i \right\}_{i=1}^n \text{ is given}$$

- let's assume that there exists a local minimum, the derivative there is zero.

$$\frac{\partial D}{\partial d_0} = 0 \quad ; \quad \frac{\partial D}{\partial d_1} = 0 \quad ; \quad \frac{\partial D}{\partial d_2} = 0$$

↳ 3 unknowns, 3 equations

$$\frac{\partial D}{\partial d_0} = \sum_{i=1}^n 2 \cdot (\eta_i - d_0 - d_1 x_i - d_2 x_i^2) (-1) = 0$$

$$\frac{\partial D}{\partial d_1} = \sum_{i=1}^n 2 \cdot (\eta_i - d_0 - d_1 x_i - d_2 x_i^2) (-x_i) = 0$$

$$\frac{\partial D}{\partial d_2} = \sum_{i=1}^n 2 \cdot (\eta_i - d_0 - d_1 x_i - d_2 x_i^2) (-x_i^2) = 0$$

Exposition:

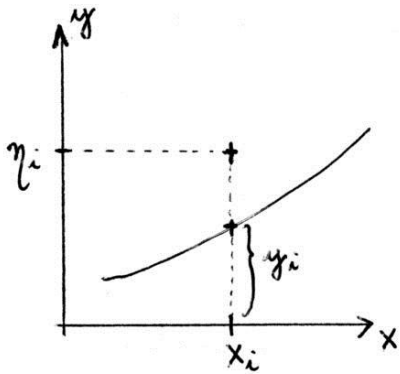
$$\sum \eta_i - n \cdot d_0 - d_1 \sum x_i - d_2 \sum x_i^2 = 0$$

$$\sum x_i \eta_i - d_0 \sum x_i - d_1 \sum x_i^2 - d_2 \sum x_i^3 = 0$$

$$\sum x_i^2 \eta_i - d_0 \sum x_i^2 - d_1 \sum x_i^3 - d_2 \sum x_i^4 = 0$$

}  $d_0, d_1, d_2$  can be determined  
(Excel uses the same method)

## Coefficient of determination



- total standard deviation of  $\eta$

$$\frac{1}{n} \sum (\eta_i - \bar{\eta})^2 = S_{\text{tot}}^2$$

- residual standard deviation

$$\frac{1}{n} \sum (\eta_i - y_i)^2 = S_{\text{err}}^2$$

coeff. of determination:

$$R^2 = \frac{S_{\text{tot}}^2 - S_{\text{err}}^2}{S_{\text{tot}}^2} \begin{cases} 1, & \text{if } S_{\text{err}} = 0 \\ 0, & \text{if } S_{\text{tot}} = S_{\text{err}} \end{cases}$$

- if  $\eta_i = y_i$ , then the curve passes through the measurement points

$$\hookrightarrow \text{in this case: } S_{\text{err}} = 0 \Rightarrow R^2 = 1 \quad \square$$

- if  $y_i = \bar{\eta}$ , then the curve is selected to be the average of the points

$$\hookrightarrow \text{in this case } S_{\text{err}}^2 = S_{\text{tot}}^2 \Rightarrow R^2 = 0 \quad \square$$

$\Downarrow$

the coeff. of determination shows the "goodness" of the fitting

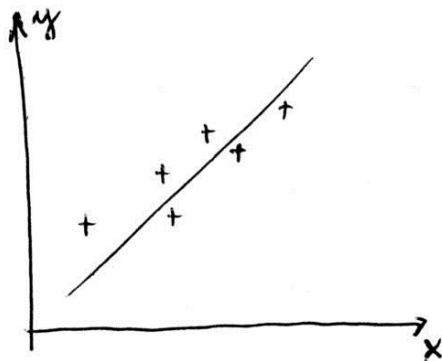
- the nearer is  $R^2$  to 1 the nearer goes the curve to the points

## Caution in the interpretation

- if  $R^2 = 1$ , thus  $\eta_i = y_i$ , the fitted curve passes through the points, this is not regression but interpolation

$\hookrightarrow$  this will not describe the process properly!

it can be shown:



$$r(y, x) = \frac{\frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})}{s_x \cdot s_y}$$

$$R^2 = r^2(y, x)$$

Definitions to this point: (with the posed questions)

- relative frequency  $\frac{V}{n}$       what if  $n \rightarrow \infty$ ?
- average  $\frac{1}{n} \sum z_i$       —||—
- empirical density function      —||—
- empirical cumulative distribution function (CDF)      —||—
- relationship between  $f_n(x)$  and  $F_n(x)$ , if  $n \rightarrow \infty$
- sample standard deviation  $s^*$       —||—
- empirical correlation coefficient  $\rho$       —||—

/.

# Probability theory

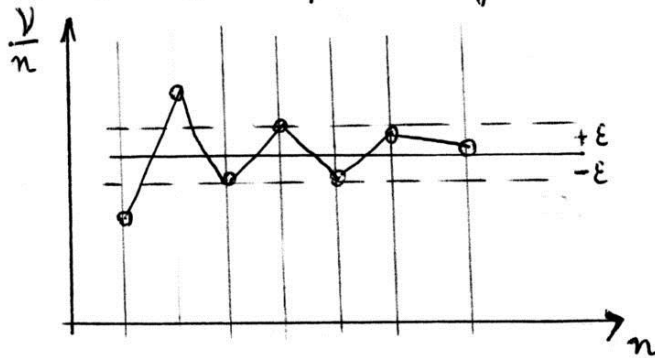
- history: • the Greek did not deal with the probability theory
- 17th century: Pascal and Fermat (basics)
- beginning of the 20th century: Kolmogorov

## random event

- the result of an event (experiment, etc.) is not unequivocally determined by the circumstances taken into account. (clearly)

↳ this is incidence

- relative frequency  $\rightarrow$  probability



$v$ : frequency (number of <sup>„successful“</sup> events)  
 $n$ : number of all events

↳ relative frequency varies randomly

↳ even as  $n \rightarrow \infty$ , the rel. frequ. „stabilizes“. The probability that it jumps out from a  $+\epsilon/-\epsilon$  interval, decreases

## probability of an event

A: event (e.g. to throw a dice and get a five)

$$P(A) = p \quad 0 \leq p \leq 1$$

$P(A) = 1$ : certain event

$P(A) = 0$ : impossible event



## deterministic variable

example:  $V = \frac{s}{t}$ , it does not depend on coincidence

## random variable

- the circumstances do not determine the variable unequivocally.  
taken into account

- it can be discrete or continuous

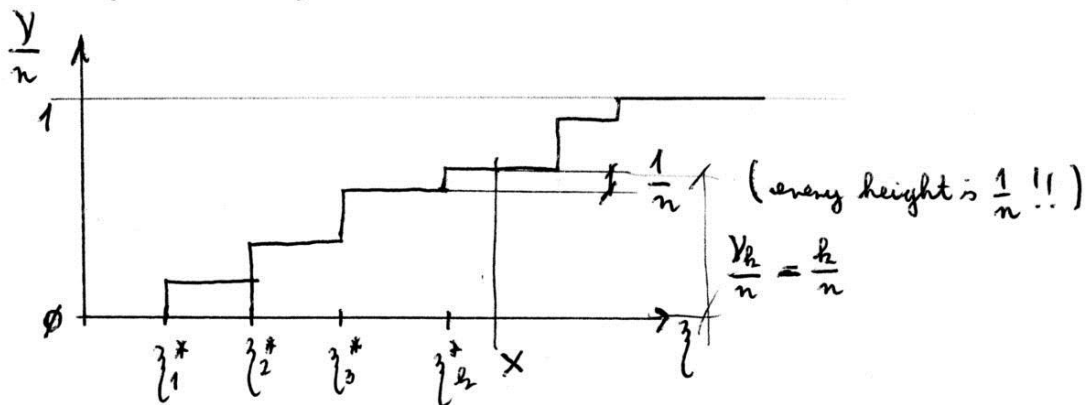
examples: dice, roulette, result of a measurement, body height

the empirical cumulative distribution function  $\rightarrow$  cumulative distribution function

ordered observations:

$\{z_1^*, z_2^*, \dots, z_n^*\}$

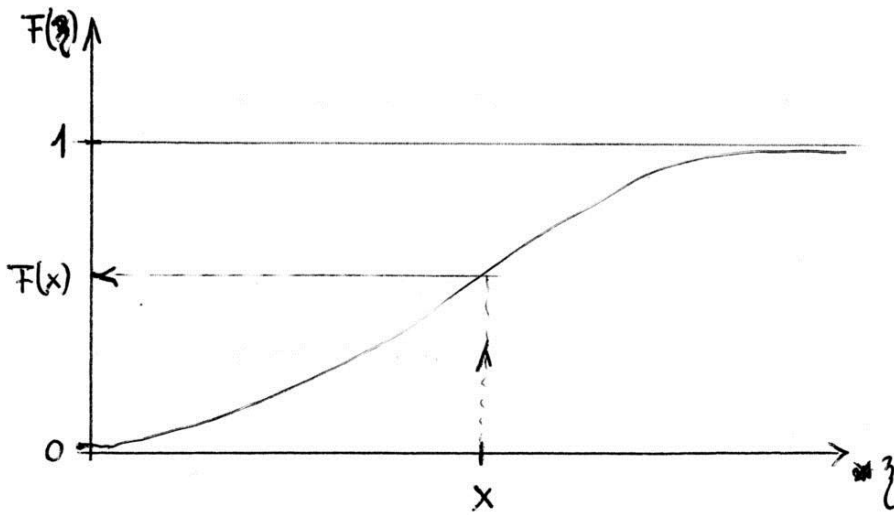
(discrete variable!)



- for an arbitrary  $x$  there belongs " $h$ " observations, which are less than " $x$ ":

$$F_n(x) = \frac{h}{n}$$

- if  $n \rightarrow \infty$ , then  $F_n(x) \rightarrow F(x)$



event "A":

$$A: \{z < x\}$$

$$P(A) = P\{z < x\} = F(x)$$

Properties:

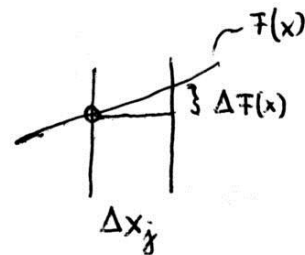
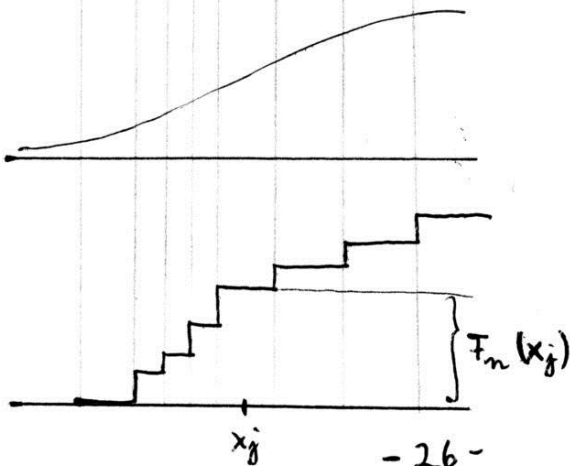
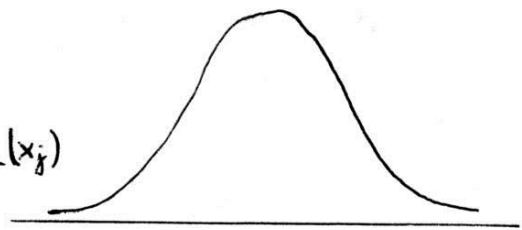
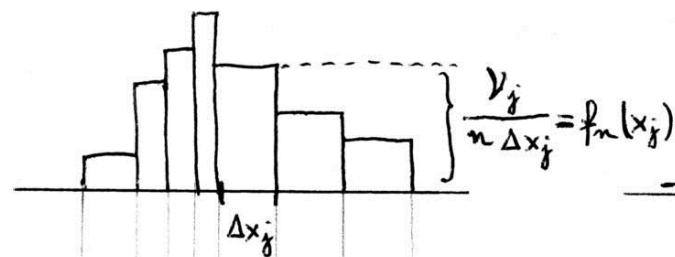
$$F(x) \rightarrow 0 ; \text{ if } x \rightarrow -\infty$$

$$F(x) \rightarrow 1 ; \text{ if } x \rightarrow +\infty$$

probability of belonging to an interval:

$$\left. \begin{aligned} P\{z < a\} &= F(a) \\ P\{z < b\} &= F(b) \end{aligned} \right\} P(a < z < b) = F(b) - F(a)$$

empirical density function  $\rightarrow$  density function



$$\frac{\Delta F(x)}{\Delta x} \cong \frac{F(x+\Delta x) - F(x)}{\Delta x} = \frac{\frac{v_j}{n}}{\Delta x}$$

$$F_n(x_j) = \sum_{x < x_j} \frac{V_j}{n} = \sum_{x < x_j} \underbrace{\frac{V_j}{n \Delta x_j}}_{f_n(x_j)} \Delta x_j = \sum f_n(x_j) \Delta x_j$$

if  $n \rightarrow \infty$

$$F(x) = \int_{-\infty}^x f(t) dt$$

properties of the density function

•  $f(x) \rightarrow 0$ ; if  $x \rightarrow \pm \infty$

•  $f(x) \geq 0$

•  $P(a \leq z \leq b) = \int_a^b f(x) dx$

•  $f(x) = \frac{dF}{dx}$

examples:  
- normal  
- uniform

average  $\rightarrow$  mean value

$$\{z_1, z_2, \dots, z_n\}$$

number of intervals!!

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i = \frac{1}{n} \sum_{j=1}^K x_j V_j = \frac{1}{n} \sum_{j=1}^K x_j \frac{V_j}{\Delta x_j} \Delta x_j = \sum x_j \underbrace{\frac{V_j}{n \Delta x_j}}_{f_n(x_j)} \Delta x_j =$$

$$= \sum x_j f_n(x_j) \Delta x_j \xrightarrow{n \rightarrow \infty} \int x f(x) dx$$

("first-order moment")

definition of the mean value:

$$M(z) = \int x f(x) dx$$

properties

$$M(z+\eta) = M(z) + M(\eta)$$

$$\eta = az + b$$

$$M(az+b) = a \cdot M(z) + b$$

$$M(a) = a$$

sample standard deviation  $\rightarrow$  standard deviation (and variance)

$$s^{*2} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$\downarrow$

$$D^2(z) = \sigma^2 = M[(z - M(z))^2], \text{ variance (mean squared difference)}$$

$$\sigma = \sqrt{D^2(z)}, \text{ standard deviation}$$

properties

•  $z$  ;  $D^2(z) = \sigma_z^2$

$$\eta = az + b \rightarrow D^2(\eta) = D^2(az + b) = a^2 D^2(z) + D^2(b) = a^2 \cdot \sigma_z^2$$

•  $z, \eta$

$$D^2(z + \eta) = D^2(z) + D^2(\eta) \quad \text{if } z \text{ and } \eta \text{ are independent}$$

but:

$$D(z + \eta) \neq D(z) + D(\eta)$$

Important:

in case of a random variable  $z$ ,  $D(z)$  and  $M(z)$  are constant values (not random variables anymore)!

(deterministic)

Variance of the average

$$\bar{z} = \frac{1}{n} \sum z_i$$

$$D^2(\bar{z}) = D^2\left(\frac{1}{n} \sum z_i\right) = \frac{1}{n^2} D^2\left(\sum z_i\right) = \frac{1}{n^2} \sum D^2(z_i) = \frac{n \cdot D^2(z)}{n^2} = \frac{\sigma^2}{n}$$

2012.10.04.

# Standardization

$$Z; m = M(Z); D^2(Z) = \sigma_Z^2$$

- standardized variable:

$$\eta = \frac{Z - m}{\sigma} \quad (\text{standard score})$$

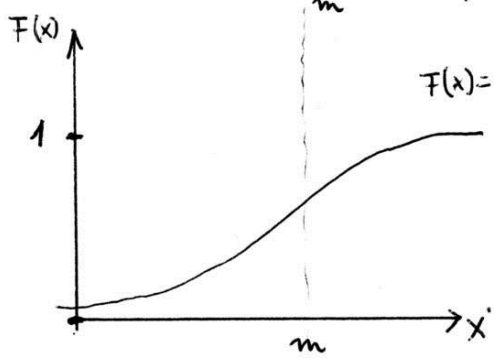
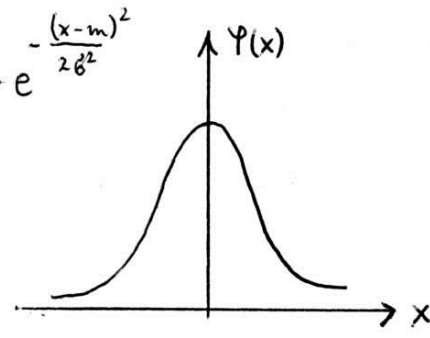
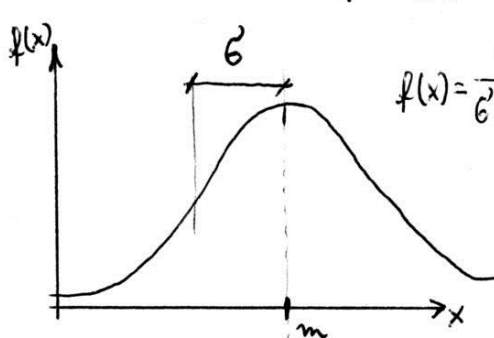
$$M(\eta) = M\left(\frac{Z - m}{\sigma}\right) = \frac{1}{\sigma} \cdot M(Z - m) = \frac{1}{\sigma} \cdot (M(Z) - M(Z)) = 0$$

$$D^2(\eta) = D^2\left(\frac{Z - m}{\sigma}\right) = \frac{1}{\sigma^2} \cdot D^2(Z - m) = \frac{D^2(Z)}{\sigma^2} = 1$$

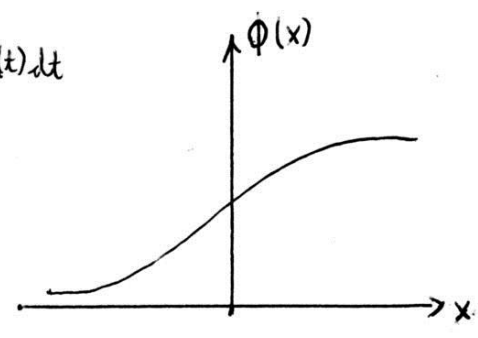
↓

the standardized variable has a mean value of zero and a variance of one.

$$\begin{pmatrix} M(\eta) = 0 \\ D(\eta) = 1 \end{pmatrix}$$



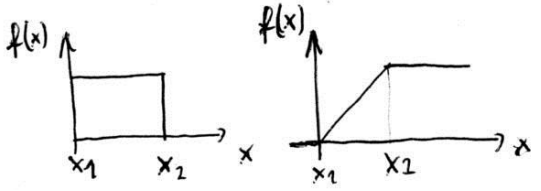
$$F(x) = \int f(t) dt$$



normal distribution  
 $N(m, \sigma)$

standard normal distribution  
 $N(0, 1)$

uniform distribution?



## Law of large numbers

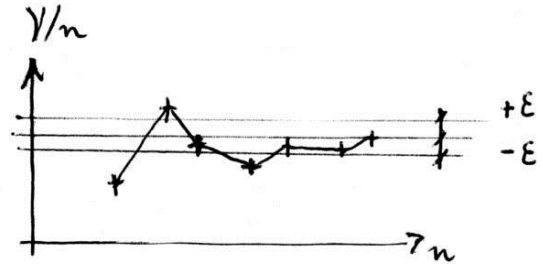
- in case of an experiment which was carried out several times:

relative frequency  $\rightarrow$  probability

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{Y}{n} - p\right| \geq \varepsilon\right) = 0$$

- the convergence is stochastic

$$\bar{z} \xrightarrow{\text{st.}} M(\bar{z})$$

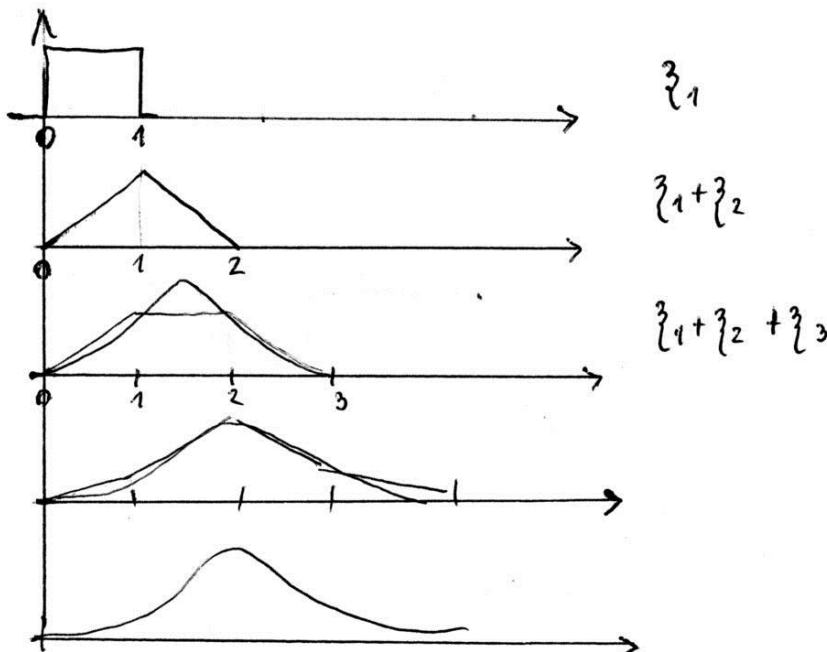


## Central limit theorem (CLT)

- let's take "n" independent random variables, the distribution of their sum will be a normal distribution

$$\lim_{n \rightarrow \infty} P\left(\frac{z_1 + z_2 + \dots + z_n - n \cdot m}{\sigma \cdot \sqrt{n}}\right) = \phi(x)$$

example:  $z_i, i=1, 2, \dots, n$  ; the distribution of each variable is uniform



## Reminder:

- data sampling

↳ the data sample should be a collection of independent items with the same distribution

↳ the data sample should be "sufficiently large"

↳ a data sample can be described:

- average
- sample standard distribution
- histogram

## Estimation

- the parameters (mean, variance, etc.) are calculated from the sample

## estimator:

- the value which estimates the particular parameter

## Properties

2012.10.11

### ① unbiased estimation

actual parameter:  $a$

estimator:  $d$

- if  $M(d) = a$ , then the estimation is unbiased
- the average  $\bar{x}$  is the unbiased estimator of the mean:

$$M(\bar{z}) = M\left(\frac{1}{n} \sum z_i\right) = \frac{1}{n} \cdot \sum M(z_i) = \frac{1}{n} \cdot n \cdot M(z) = M(z)$$

- standard deviation, sample standard deviation

$$M(s^2) \neq D^2(z) \quad \text{but} \quad M(s^{*2}) = D^2(z)$$

## ② consistent estimation

$d_i$ : estimator calculated from a sample of "i" elements (units)

$d_1, d_2, \dots, d_n$

- the estimation is consistent if " $d_i$ " tends to " $a$ " ~~if  $n \rightarrow \infty$~~  with increasing  $n$

## ③ efficient estimation

$d_1$  estimation

$d_2$  estimation

}  $d_1$  is the more efficient estimator if:

$$D^2(d_1) < D^2(d_2)$$

## Measurement errors

- measurement errors and error propagation

- the results of a measurement:  $z_i$  random variable

- every measurement has an actual value:  $X$

$$z = X + \epsilon$$

$\epsilon$ : the error of the measurement

↳ this might be the error in the measuring system, or the noise of the observed value

(example: 1 liter bottles)

$$M(z) = M(X + \epsilon) = M(X) + M(\epsilon) = X + M(\epsilon)$$

$$M(\epsilon) = \Delta x \begin{cases} = \emptyset & \rightarrow \epsilon \text{ is a random error} \\ \neq \emptyset & \rightarrow \Delta x \text{ is a systematic error} \end{cases}$$

- 32 -  $\sigma = \epsilon - \Delta x$  is the random error



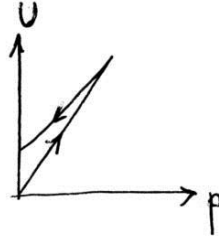
## Systematic error

- the systematic error of the measuring device or the measuring method

- examples:

• zero error (incorrect zeroing, displaced scale, etc.)

• hysteresis of the device



• dynamic behaviour of the device (for example an unsteady <sup>process is</sup> ~~process is~~ ~~measured~~ measured with a static device)

## Random error

- several effects which cannot be taken into account

influence the measurement (temperature of the air, voltage of power supply, material properties, etc.)

- since  $M(E) = 0 \Rightarrow M(\xi) = X$

- the evaluation of the measurement is performed with statistical methods (confidence intervals, hypothesis testing, regression)

## Error propagation

- classification of measurements:

① direct measurements

• the measured quantity is directly compared to the etalon (standard)

↳ a device representing the unit of the quantity

## ② indirect measurements

- the required quantity is calculated from measured quantities

$$\hookrightarrow \text{torque: } M = l \cdot F \quad \text{"lever arm"}$$

$$\text{power: } P = U \cdot I$$

- every measured quantity has an error (systematic, random)



these errors "propagate" through the calc. process

2012.10.18.

### 1.1 Propagation of systematic errors

- different quantities:  $x_1 ; x_2 ; \dots ; x_n$   
sys errors:  $\Delta x_1 ; \Delta x_2 ; \dots ; \Delta x_n$

- the value of  $y$  has to be calculated:

$$y = f(x_1, x_2, \dots, x_n)$$

- how can we determine the systematic error of  $y$ ? \*

$$\Delta y = ?$$

$$\Delta y = f(x_1 + \Delta x_1, x_2 + \Delta x_2, \dots, x_n + \Delta x_n) - f(x_1, x_2, \dots, x_n)$$

↳ if every  $\Delta x_i$  is known, this can be calculated!

↳ but: how do each  $\Delta x_i$  influence the sys. error of  $y$ ?



if we want to decrease  $\Delta y$ , which  $\Delta x_i$  should be reduced?

(i.e.: what parameter should be measured more accurately?)

How does the value of  $\Delta x_i$  influence  $\Delta y$ ?

(first order Taylor series)

- Let us create the Taylor series for  $y + \Delta y$ :

$$y + \Delta y \cong f(x_1, x_2, \dots, x_n) + \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Delta x_i + \dots$$

only the linear part is kept!!

$$H_y = \Delta y \cong \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Delta x_i$$

↑

the resultant systematic error  
(absolute error)

sensitivity parameter

↳ this shows in what extent each  $\Delta x_i$  influences the resultant systematic error

the maximum of the res. sys. error:

(absolute error)

$$\Delta y_{\max} = \sum_{i=1}^n \left| \frac{\partial f}{\partial x_i} \right| \cdot |\Delta x_i|$$

relative systematic error:

$$H_{\text{rel}} = \frac{\Delta y}{y} \cong \sum_{i=1}^n \frac{\partial f}{\partial x_i} \cdot \frac{\Delta x_i}{y}$$

specific case:

$$y = C \cdot x_1^{h_1} \cdot x_2^{h_2} \cdot \dots \cdot x_n^{h_n}$$

$$\frac{\partial f}{\partial x_i} = C \cdot x_1^{h_1} \cdot x_2^{h_2} \cdot \dots \cdot h_i \cdot x_i^{h_i-1} \cdot \dots \cdot x_n^{h_n} = h_i \cdot \frac{y}{x_i}$$

therefore:

$$\frac{\Delta y}{y} = \sum_{i=1}^n h_i \cdot \frac{\Delta x_i}{x_i}$$

short example:

the volume of a cylinder is calculated as follows:

$$V = \frac{D^2 \pi h}{4} = \frac{\pi}{4} \cdot D^2 \cdot h$$

$$\textcircled{1} \quad H_{\text{rel}} = \frac{\Delta V}{V} = 2 \cdot \frac{\Delta D}{D} + 1 \cdot \frac{\Delta h}{h} //$$

$$\textcircled{2} \quad \frac{\partial V}{\partial D} = \frac{\pi}{2} \cdot D \cdot h \quad ; \quad \frac{\partial V}{\partial h} = \frac{\pi}{4} \cdot D^2$$

$$H_{\text{abs}} = \Delta V = \frac{\pi}{2} \cdot D \cdot h \cdot \Delta D + \frac{\pi}{4} \cdot D^2 \cdot \Delta h \quad /: V = \frac{D^2 \pi}{4} \cdot h$$

$$\frac{\Delta V}{V} = 2 \cdot \frac{\Delta D}{D} + \frac{\Delta h}{h} //$$

## 2. Propagation of random errors

$x_1, x_2, \dots, x_n$

random variables

$\Rightarrow y$  is also a random variable

$$\varepsilon_i = X - x_i$$

- let us use the previous formula, instead of  $\Delta x_i$  we have  $\varepsilon_i$

$$\boxed{\varepsilon_y \cong \sum_{i=1}^n \frac{\partial f}{\partial x_i} \cdot \varepsilon_i}$$

calculating the mean:

$$M(\varepsilon_y) = M\left(\sum_{i=1}^n \frac{\partial f}{\partial x_i} \cdot \varepsilon_i\right) = \sum_{i=1}^n M\left(\frac{\partial f}{\partial x_i} \cdot \varepsilon_i\right) = \sum_{i=1}^n \frac{\partial f}{\partial x_i} M(\varepsilon_i) = 0$$

$\Downarrow$   
therefore  $\varepsilon_y$  is also random type error

- the variance of  $E_y$ :

$$D^2(\overset{\text{actual value}}{y} + E_y) = D^2(E_y) = D^2\left(\sum_{i=1}^n \frac{\partial f}{\partial x_i} \varepsilon_i\right) = \sum_{i=1}^n \left(\frac{\partial f}{\partial x_i}\right)^2 D^2(\varepsilon_i)$$

$$\sigma_y^2 = \sum_{i=1}^n \left(\frac{\partial f}{\partial x_i}\right)^2 \sigma_i^2$$

$$\hookrightarrow \sigma_y = \sqrt{\sigma_y^2} ; \text{ standard deviation}$$

- relative variance:

$$\left(\frac{\sigma_y}{y}\right)^2 = \sum_{i=1}^n \left(\frac{\partial f}{\partial x_i}\right)^2 \left(\frac{\sigma_i}{y}\right)^2$$

\* special case:

$$y = C \cdot x_1^{h_1} \cdot x_2^{h_2} \cdot \dots \cdot x_n^{h_n}$$

$$\left(\frac{\sigma_y}{y}\right)^2 = \sum_{i=1}^n \left(h_i \frac{\sigma_i}{x_i}\right)^2$$

- in case of a measurement the variance / st. dev. is unknown

$\hookrightarrow$  the sample standard deviation ( $s^*$ ) is calculated

$\hookrightarrow$  the formulas above are used with " $s^*$ "

example:

- volume of a cylinder;

$$V = \frac{D^2 \pi}{4} \cdot h$$

- measuring diameter ( $D_1, D_2, \dots, D_N$ ) and height ( $h_1, h_2, \dots, h_N$ )

$$\bar{D} = \frac{1}{N} \sum_{j=1}^N D_j \quad \mapsto \quad s_D^* \quad s_{\bar{D}}^* = \frac{s_D^*}{\sqrt{N}}$$

$$\bar{h} = \frac{1}{N} \sum_{j=1}^N h_j \quad \mapsto \quad s_h^* \quad s_{\bar{h}}^* = \frac{s_h^*}{\sqrt{N}}$$

- but it is better to calculate each  $V_i = \frac{D_i^2 \pi}{4} h_i$

(if we don't want to investigate the effect of the components)

### Confidence interval

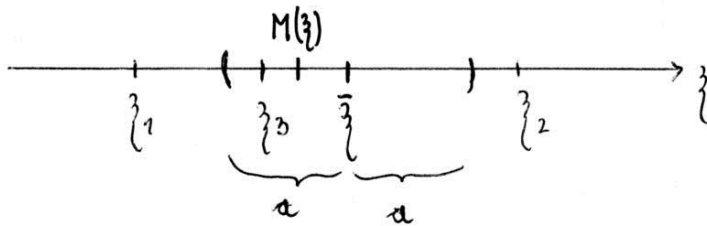
- let us assume that we deal with random errors

- a measurement is carried out:  $\{z_1, z_2, \dots, z_n\}$

• the average  $\bar{z}$  can be calculated

question:

how precisely do we estimate the mean  $M(z)$   
using the average?



remember:  $\bar{z} \rightarrow M(z)$   
 $n \rightarrow \infty$

- let us create an interval around  $\bar{z}$ ! The average is in the middle, the radius is "a"

- the interval should contain  $M(z)$ :

$$M(z) \in [\bar{z} - a; \bar{z} + a] \text{ with a probability of "p"}$$

- how is "a" determined?

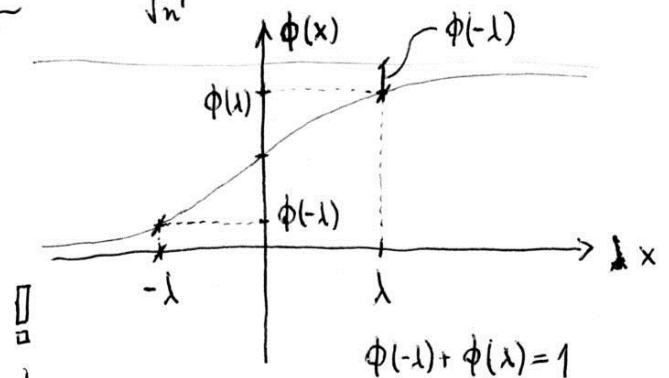
$$P(\bar{z} - a \leq M(\bar{z}) \leq \bar{z} + a) = p$$

$$P(-a \leq \bar{z} - M(\bar{z}) \leq +a) = P\left(-a \cdot \frac{\sqrt{n}}{\sigma} \leq \frac{\bar{z} - M(\bar{z})}{\frac{\sigma}{\sqrt{n}}} \leq +a \cdot \frac{\sqrt{n}}{\sigma}\right)$$

let us create the standard score of the avg!

$$M(\bar{z}) = M(\bar{z})$$

$$D(\bar{z}) = \frac{D(\bar{z})}{\sqrt{n}}$$



$$P(-\lambda \leq \eta \leq \lambda) = p = \Phi(\lambda) - \Phi(-\lambda) = 2\Phi(\lambda) - 1$$

$$\Downarrow$$

$$\Phi(\lambda) = \frac{p+1}{2} ; \lambda = \Phi^{-1}\left(\frac{p+1}{2}\right)$$

example:

if the significance level is  $p = 95\%$ , then  $\lambda = 1,96$

$$\lambda = a \cdot \frac{\sqrt{n}}{\sigma} \Rightarrow \boxed{a = \lambda \cdot \frac{\sigma}{\sqrt{n}}}$$

, the radius of the confidence interval

$$\bar{z} \pm a ; \bar{z} - a : \text{lower confidence limit}$$

$$\bar{z} + a : \text{upper confidence limit}$$

- what if the st. deviation is unknown? (This is usually the case)

• let us calculate the sample st. dev.:

$$s^{*2} = \frac{1}{n-1} \sum_{i=1}^n (\bar{z} - z_i)^2$$

st. in this case we have:

$$\eta^* = \frac{\bar{z} - M(\bar{z})}{\frac{s^*}{\sqrt{n}}} \notin N(0,1)$$

instead:

Student distribution

(William Gosset, worker of the Guinness  
brewery)

- Student's "t" distribution

⇓

$\lambda_{st}$  instead of  $\lambda$ , therefore:

$$a = \lambda_{st} \cdot \frac{s^*}{\sqrt{n}}$$

---

STATISTICAL HYPOTHESIS TESTING