

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

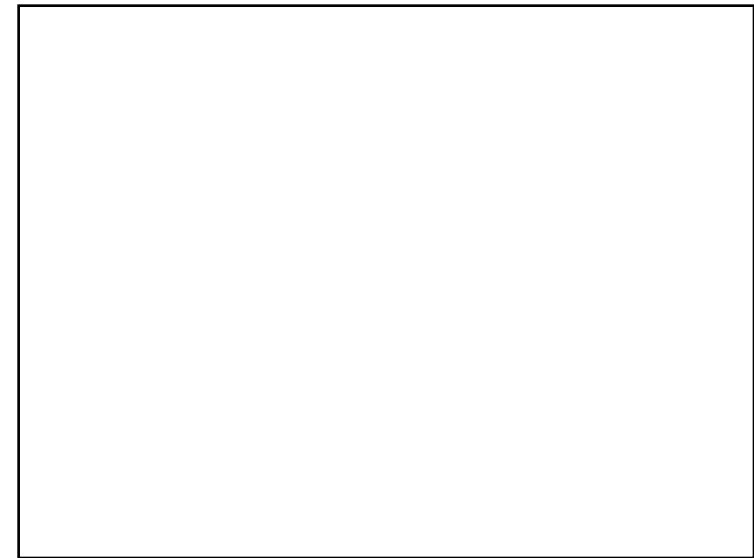
<http://www.bernhardkoester.de/vorlesungen/inhalt.html>



(Advanced) Statistics A

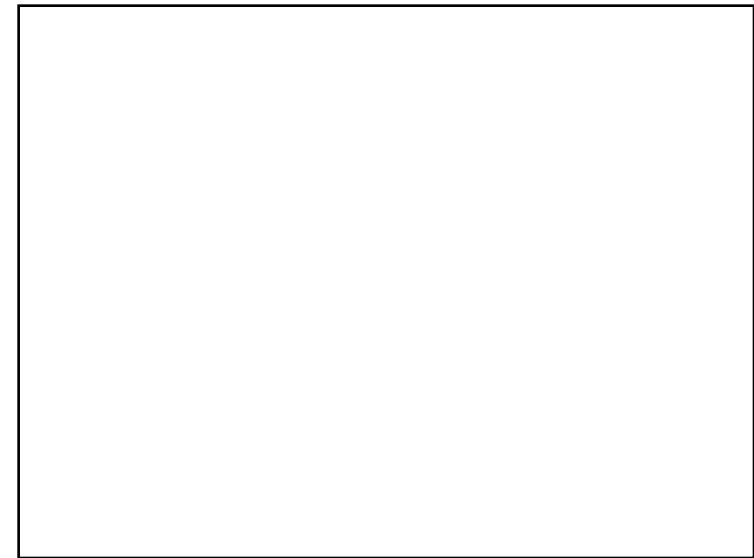
Sommersemester 2022

Prof. Dr. Bernhard Köster



- Due to the pandemic the organizational form of this lecture is not fixed right now!
- Start is Thursday, 2 March 2022 online via zoom (see moodle)
- Classroom lectures may be possible from 21 March 2022 on
- Further information will be given in the lecture. Times are volatile nowadays!
- During the pandemic, I always had ideas, who to deal with the uncertain situation, but always I hade to bow force majeure. Therefore I do not make any annoucements in advance.
- But I can assure you, that we will find a good teaching concept!
- Kind regards

Bernhard Köster



Prof. Dr. Bernhard Köster

Room: S 113

Street: Friedrich-Paffrath-Straße 101

location: 26389 Wilhelmshaven

Tel. +49 4421 985-2766

Email: bernhard.koester@jade-hs.de

Consultation hour: by arrangement
or just have a look into my office!
or [Webex/Zoom](#)...

Literature

- **Statistics for Business and Economics**, Anderson, Sweeney, Williams, Camm, Cochran
- **Mathematical Statistics for Economics and Business**, Mittelhammer
- **Statistics for Business and Economics, Global Edition**, Newbold, Thorne, Carlson
- **A Guide to Modern Econometrics**, Verbeek



Some open resources

[OpenIntro Statistics](#)

[Statistical Thinking for the 21st Century](#)

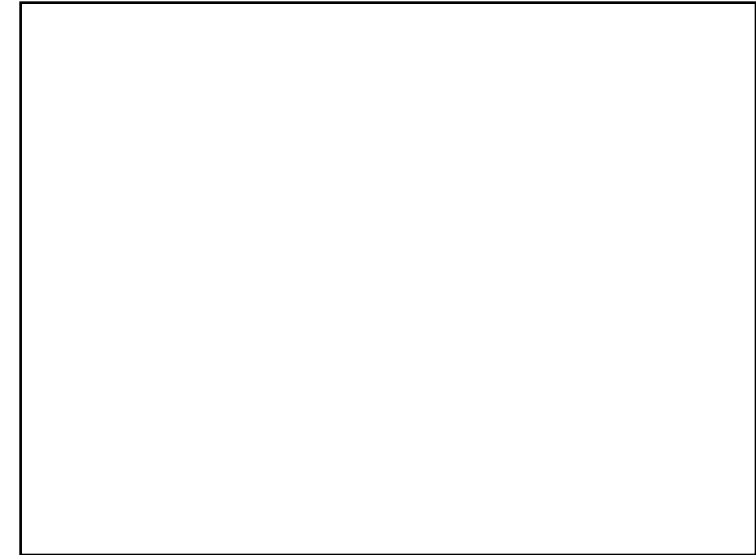
[MIT OPENCOURSEWARE](#)

-> there is much more!!!

Handwritten notes: \rightarrow Statistics

1	2
6	7
8	7
6	5

Handwritten notes: \rightarrow via calculator $\rightarrow Y = a + bX$



- Best unbiased estimator

- Efficiency

- Consistency

- Sufficiency

→ Confidence Interval

- One and two sided tests

- α (type I) / β (type II) error

→ Significance

- Hypothesis testing
t, F, χ^2

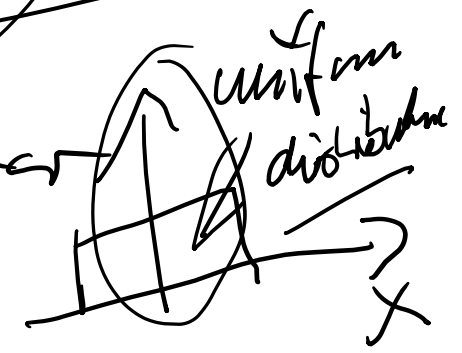
- ML - estimator

Decision Theory

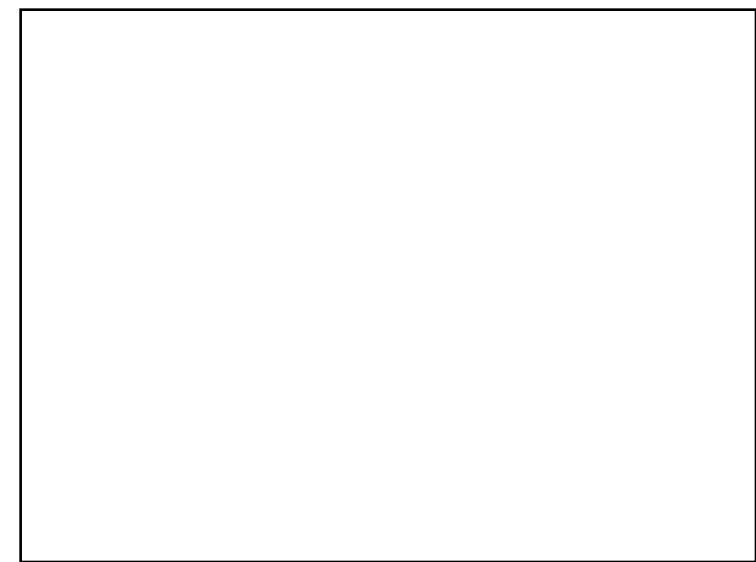


normal distribution

1. kind 2 kind error
Probability distribution



uniform distribution



Introduction and Revision

- **Means and Median**
- **Dispersion: MAD, Variance, Standard Deviation, Coefficient of Variation**
- **Contingency table and independence**
- **Bayes' Theorem**
- **Random Variable**
- **Probability Distributions, density function and cumulative distribution function**
- **Central limit theorem**

Example

$$\bar{x}_g = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$

$$\bar{x}_h = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

1) Generate in Excel a sample of the length of 15 male students around 175cm and a maximum deviation of 25cm.

2) Calculate the arithmetic, geometric, harmonic mean and the median of the distribution.

3) Calculate the MAD, variance, standard deviation, and the coefficient of variation of the distribution.

$$MAD = \frac{1}{n} \sum_{i=1}^n |x_a - x_i|$$

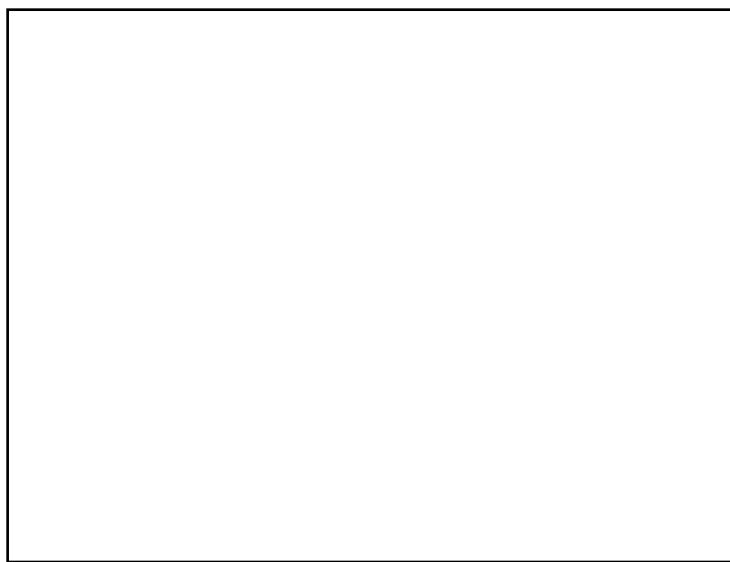
$$s^2 = \frac{1}{n} \left(\sum (x_a - x_i)^2 \right)$$

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$$

4) Generate in Excel a sample of the length of 11 female students around 165cm and a maximum deviation of 20cm.

5) Calculate the contingency table of the attributes gender and length ((0-155), [155,170), [170,185), [185,200] and the marginal distributions

1) Are the attributes gender and length statistically independent?



WHV

100 km/h
50 km/h

OL → What is the uncertainty?



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

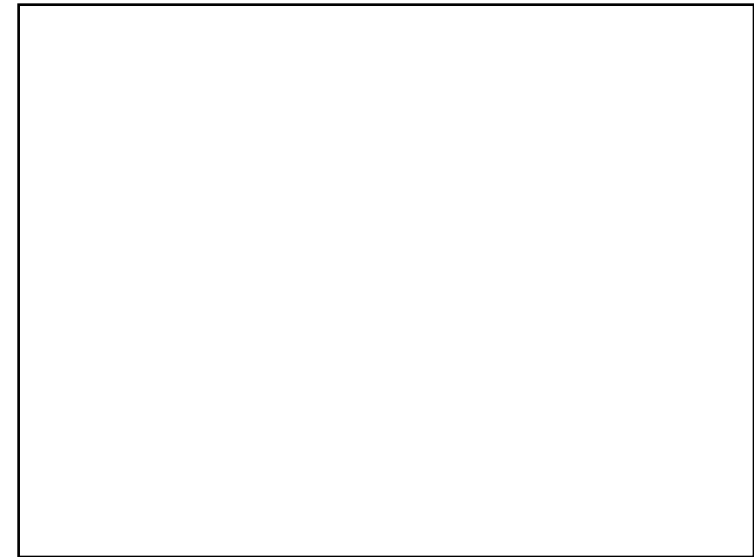
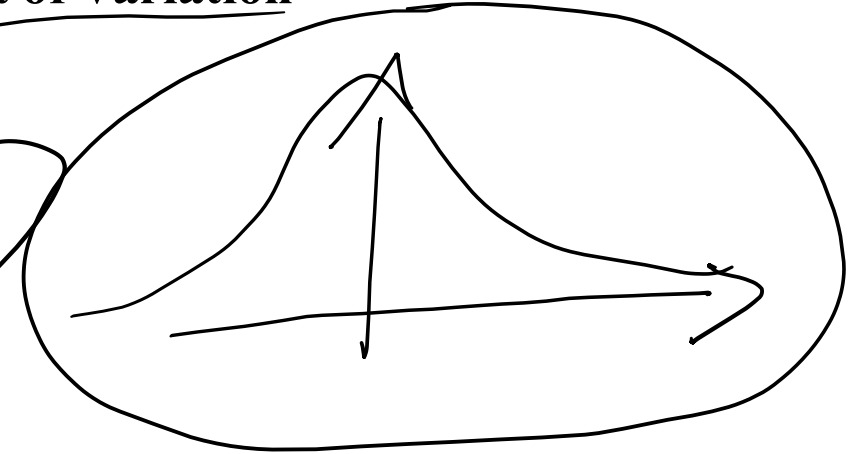
Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

applied, geometric

Introduction and Revision

- **Means and Median**
- **Dispersion: MAD, Variance, Standard Deviation, Coefficient of Variation**
- **Contingency table and independence**
- **Bayes' Theorem**
- **Random Variable**
- **Probability Distributions, density function and cumulative distribution function**
- **Central limit theorem**



Example

- 1) Generate in Excel a sample of the length of 15 male students around 175cm and a maximum deviation of 25cm.
 - 2) Calculate the arithmetic, geometric, harmonic mean and the median of the distribution.
 - 3) Calculate the MAD, variance, standard deviation, and the coefficient of variation of the distribution.
 - 4) Generate in Excel a sample of the length of 11 female students around 165cm and a maximum deviation of 20cm.
 - 5) Calculate the contingency table of the attributes gender and length ((0-155),[155,170),[170,185),[185,200] and the marginal distributions
- 1) Are the attributes gender and length statistically independent?

Conditional probability

The conditional probability of an event A is the probability of the occurrence of the event A given event B has happened (or happens simultaneously with A)

→ conditional probability of A given: $P(A | B)$.

Definition:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Bayes Theorem:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

Symmetrically

Goats and Cars

Suppose in a game show you are sitting in front of 3 doors. You know behind one door you'll win a car and behind the two other doors you get only a goat

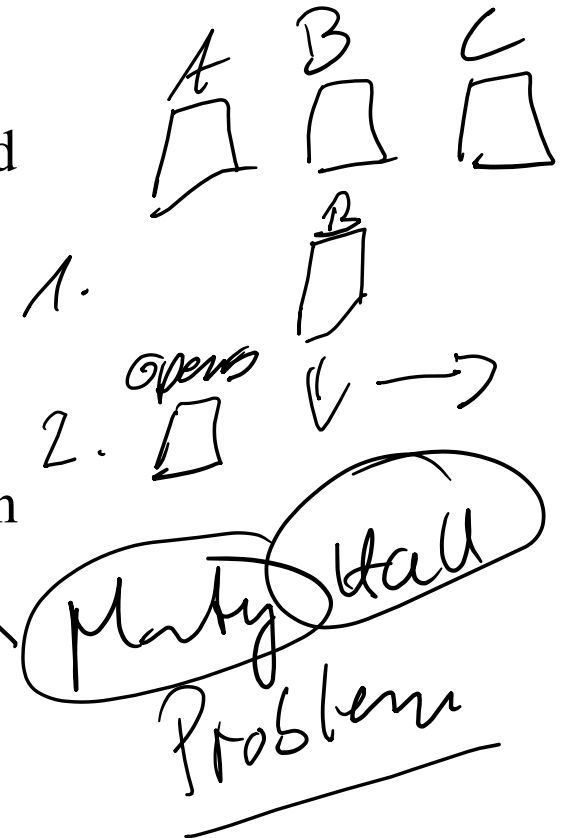
Round 1: You choose a door

After that, the showmaster, knowing the door with the car, opens a door with a goat

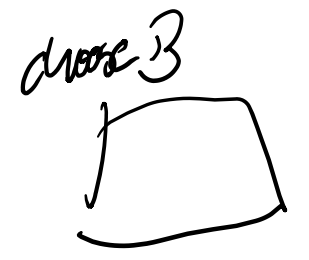
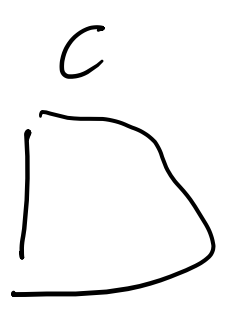
Round 2: You're asked if you want to change doors

What is your decision?

If you don't believe it, try to simulate this in Excel!



* \Rightarrow winning the by
switching the door
= $1 \cdot \frac{1}{3}$



1.
2.

$\frac{2}{3}$



Should I switch?
Shouldn't I switch?
It doesn't matter!

First Decision had a
winning probability of $\frac{1}{3}$

$$P(\text{Car}) = \frac{1}{3}$$

why $P(\text{behind not opened door}) = \frac{1}{2}$

$$P(\text{Car} | \text{not opened door in Round 2}) = \frac{P(B|K) \cdot P(A)}{P(B)}$$

what we need

$$P(\text{door is not opened} | \text{a car is behind the door}) = 1$$



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

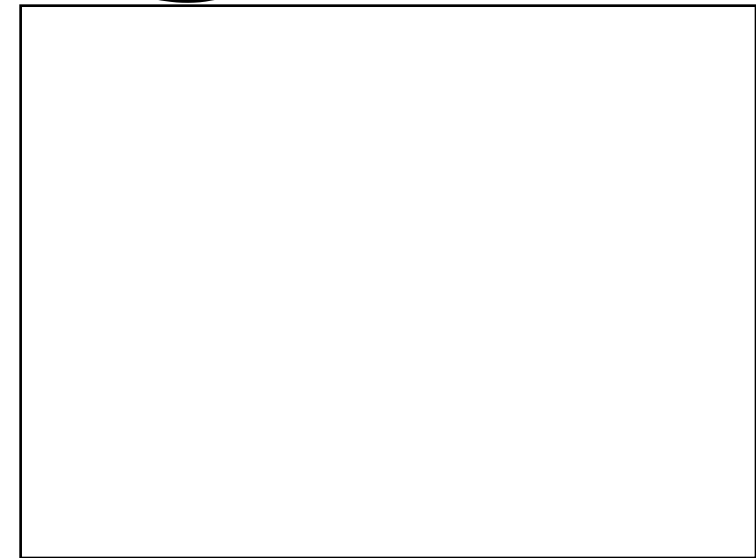
<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Random Variable

Definition:

A **Random Variable** X is the mapping of the sample space Ω into the real numbers \mathbb{R} . That means every elementary event A_i maps onto a real number x_i and the probability $p_i = P(A_i) = P(X = x_i) = p(x_i)$ is known.

→ p_i ist the probability that random variable X equals the outcome x_i .



Random Variable – Explanations

- Every possible event of a random experiment can be expressed via a random variable.

- Random variables are expressed in capital letters.

Example: dice with sample space X is element of $\{1,2,3,4,5,6\}$

- The realization of a random variable is expressed in lowercase letters:

Example rolling dice: X realizes the value $x = 5$

- Probability in the example of rolling dice:

$$P(x < 5) = ??$$

$$= \frac{4}{6} = \frac{2}{3} \approx 66\% = P(x=1) = \frac{1}{6}$$

Handwritten notes: $P(x=4) = P(x=3) = P(x=2)$

Discrete and continuous random variables (RV)

Discrete RV:

RV, which have only a finite or countably infinite realizations.

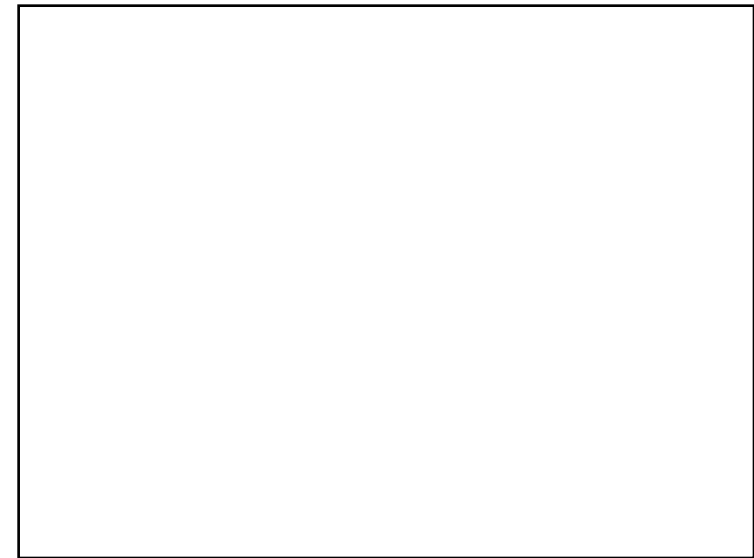
Continuous RV:

RV which can realize within an interval every possible real number.

Realizations in the Interval

12:15 and 12:45 → 30 numbers

~~The~~ A Bus is stopping within this time interval



Discret

Probability mass function:

$$P(X = x_i) = p_i$$

Cumulative distribution function:

$$P(X \leq x) = \sum_{x \leq x_i} p(x_i)$$

Expectedated value:

$$E(X) = \mu = \sum_{i=1}^n x_i p(x_i)$$

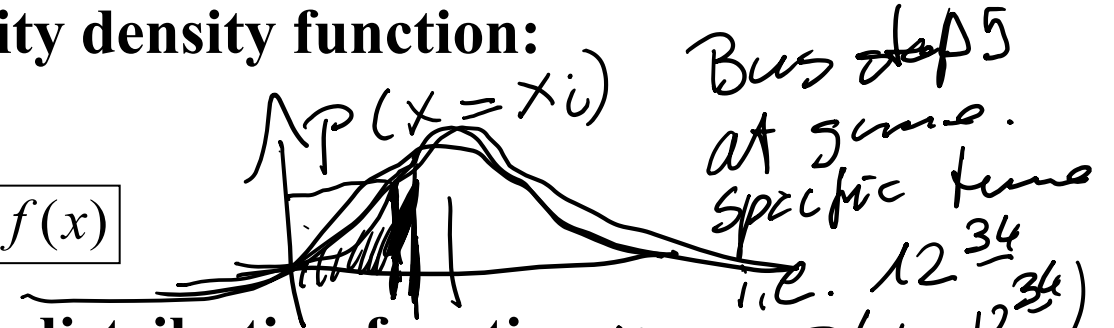
Variance:

$$Var(X) = \sigma^2 = \sum_{i=1}^n (x_i - \mu)^2 p(x_i)$$

Continous

Probability density function:

$$f(x)$$



Cumulative distribution function:

$$F(x) = \int_{x'=-\infty}^x f(x') dx'$$

Expectedated value :

$$E(X) = \mu = \int_{x=-\infty}^{x=\infty} x f(x) dx$$

Variance:

$$Var(X) = \sigma^2 = \int_{x=-\infty}^{x=\infty} (x - \mu)^2 f(x) dx$$

Probability can be interpreted as the Area below the density function = ?

Calculating with the Expected Value and the Variance

Constant:

$$\bar{E}(a) = a$$

$$\text{Var}(a) = 0$$

Constant factor:

$$E(a \cdot X) = a \cdot E(X)$$

$$\text{Var}(aX) = a^2 \cdot \text{Var}(X)$$

Linear Transformation:

$$E(a \cdot X + b) = a \cdot E(X) + b$$

$$\text{Var}(a \cdot X + b) = a^2 \cdot \text{Var}(X)$$

$$\begin{aligned} E(X) &= \sum x_i p_i \Rightarrow E(ax) \\ &= \sum ax_i p_i \\ &= a \sum x_i p_i = a\mu \end{aligned}$$

$$\begin{aligned} \text{Var}(aX) &= \sum (ax_i - a\mu)^2 p_i \\ &= \sum a^2 (x_i - \mu)^2 p_i \\ &= a^2 \sum (x_i - \mu)^2 p_i \end{aligned}$$

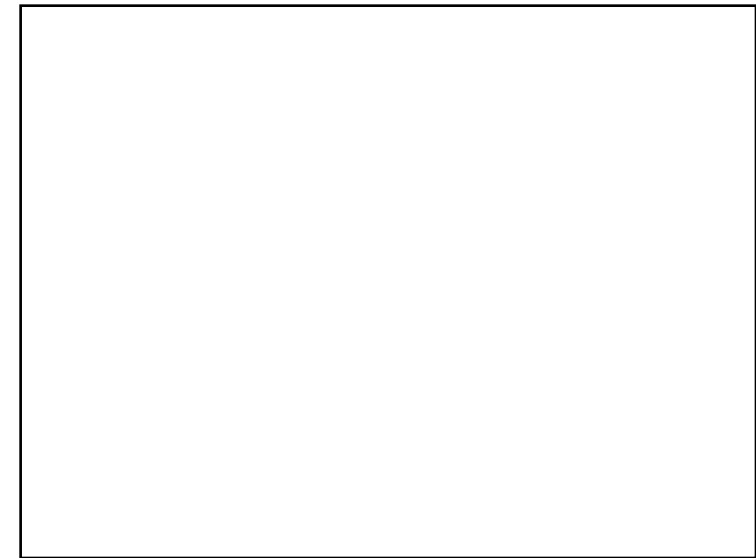
\Rightarrow Expected Value operator is a linear operator

Probability Axioms

Axiom 1: $P(A) \geq 0$ for all $A \in \Omega$
Every event has a non negative probability

Axiom 2: $P(\Omega) = 1$
The certain event has probability 1

Axiom 3: $P(A \cup B) = P(A) + P(B)$ if $A \cap \bar{B} = \emptyset$
Addition rule for disjoint events



Rules

because in general this could be different from the empty set

Rule 1: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

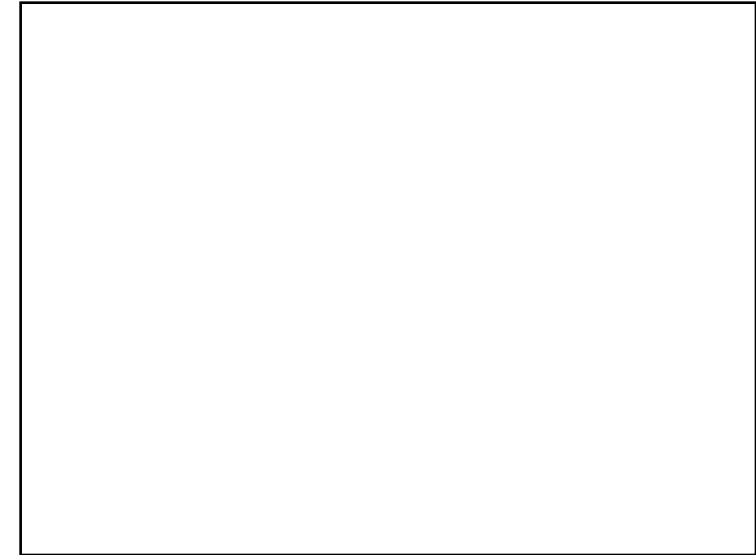
Addition rule for any events

Rule 2: $P(A \setminus B) = P(B) - P(A \cap B)$

A without B

Rule 3: $P(\bar{A}) = 1 - P(A)$

Probability of complementary events



Conditional probability

The conditional probability of an event A is the probability of the occurrence of the event A given event B has happened (or happens simultaneously with A)

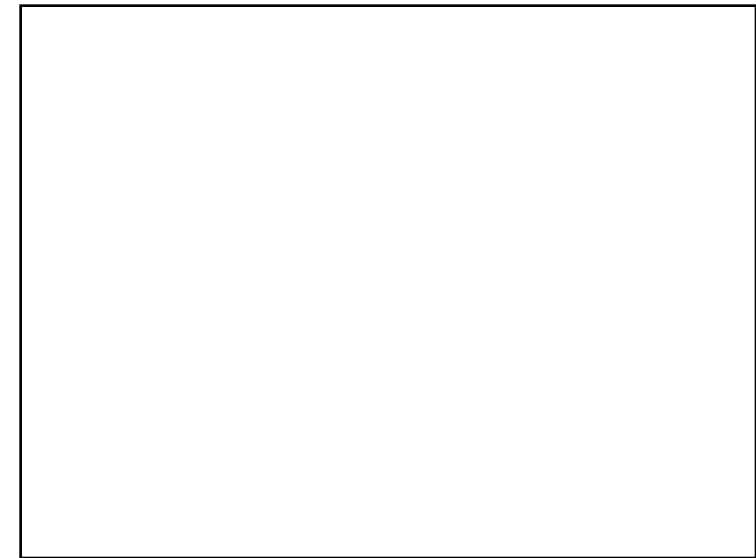
→ conditional probability of A given B: $P(A | B)$.

Definition:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Bayes Theorem:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$



Example

Suppose the following probabilities are known:

$$\underline{P(A) = 0,5;}$$

$$\underline{P(B) = 0,3;}$$

$$\underline{P(A \cap B) = 0,2}$$

Calculate

a) $P(A \cup B)$

b) $P(A|B)$

c) $P(\overline{A \cap B})$

d) $P(A \setminus B)$

$$a) P(A \cup B) = P(A) + P(B) - P(A \cap B) = 50\% + 30\% - 20\% = 60\%$$

$$b) P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0,2}{0,3} = \frac{2}{3} = 66\%$$

$$c) P(\overline{A \cap B}) = 1 - P(A \cap B) = 1 - 0,2 = 80\%$$

$$d) P(A \setminus B) = P(B) - P(A \cap B) = 0,3 - 0,2 = 10\%$$

A without B



Conditional Probability



Statistically independent events

$$\underline{P(A|B) = P(A)}$$

or

$$\underline{P(B|A) = P(B)}$$

→

$$P(A \cap B) = P(A) \cdot P(B)$$

Statistically independent events: Example

$n=2$ First student 365 second 365 $365 \cdot 365$



Sylvester 1988 in Casino in Konstanz at the roulette table, the event $A = \{0,3\}$ occurred 9 times in a row. Calculate the probability of this event

$p = \frac{2}{37}$ $P(X \text{ is } \{0,3\} \text{ 9 times in a row}) = \underbrace{\frac{2}{37} \cdot \frac{2}{37} \cdot \dots \cdot \frac{2}{37}}_{9 \text{ times}}$

$1 - 36 \times 0 \rightarrow 37$
 $\left(\frac{2}{37}\right)^9 = 4 \cdot 10^{-12}$ 9 numbers

u

Suppose in class two students have birthday at the same day. What is the the minimum number of students in class, that the probability for this event is 50%?

u: Number of students

365^u possible distributions of Birthday Just guess
 First was at some specific day this Birthday 182
 Then the second has only 364 possibilities left
 3rd student 363
 $365 \cdot 364 \cdot \dots \cdot 365 - (u - 1)$ possibilities that no students have at the same day their birthday
 $1 - \frac{365 \cdot 364 \cdot \dots \cdot 365 - (u - 1)}{365^u} = 50\%$ \Rightarrow we need 23 students

Total Probability

If the sample space Ω consists of k disjoint elementary events A_i , then the probability of event B is:

→

$$P(B) = \sum_{i=1}^k P(B | A_i) \cdot P(A_i)$$

Total Probability: Example

Every day a small village is visited by a postman. If he is in good mood he is in time with probability 90%. If he is in bad mood he is late with probability 40%.

What is the probability that he is late at any day, if he is on average in good at 7 out of 10 days?

$$\begin{aligned} P(\text{Postman is late}) &= P(\text{Postman is late} \mid \text{good}) \cdot P(\text{good}) \\ &\quad + P(\text{Postman is late} \mid \text{bad}) \cdot P(\text{bad}) \\ &= 0,07 + 0,12 = 19\% \end{aligned}$$

Total Probability and Bayes Theorem

Suppose you are testing for a rare illness and you have the following probabilities:

A = {patient is ill}

B = {test is positive}

Since we have a rare illness: $P(A) = 0,1\%$

Since a test does not have 100% accuracy, from surveys we have:

$$P(B|A) = 0,98$$

$$P(B|\bar{A}) = 0,03$$

Suppose your test result is positive. What is the probability, that you are really ill?

$$P(B) = P(B|A)P(A) + P(B|\bar{A})P(\bar{A}) \quad \text{Total probability}$$

What I want to know is

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{0,98 \cdot 0,001}{0,98 \cdot 0,001 + 0,03 \cdot 0,999}$$

↑
↑
Test positive

$$= 3,2\%$$

quite surprisingly
a very low probability

$$P(\bar{A}) = 1 - 0,001 = 0,999$$
$$P(A) = 0,001$$

just have a guess
around 90%

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Binomial distribution

- A random experiment has two possible outcomes (success and failure).
- The random experiment will be repeated n times.
- What is the probability to be k times successful within $n \geq k$ repetitions
- Every repetition is independent from each other. Within every repetition success has probability p and failure probability $1-p$.

Binomial Distribution

Suppose p of „success“ is known, define the random variable X = „number of success k “ within $n \geq k$ repetitions. $n! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot \dots \cdot n$

→ X is Binomial distributed $X \sim B(n; p)$ with

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Example

$$\binom{7}{3} = \frac{7!}{3!(7-3)!} = \frac{1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7}{1 \cdot 2 \cdot 3 \cdot (1 \cdot 2 \cdot 3 \cdot 4)} = \frac{5 \cdot 6 \cdot 7}{1 \cdot 2 \cdot 3} = 35$$

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$E(X) = np$$

$$Var(X) = np(1-p)$$

Binomial coefficients

1	2	3	4	5	6	7
1	3	3	1			
1	4	6	4	1		
1	10	10	5	1		

Examples Binomial Distribution

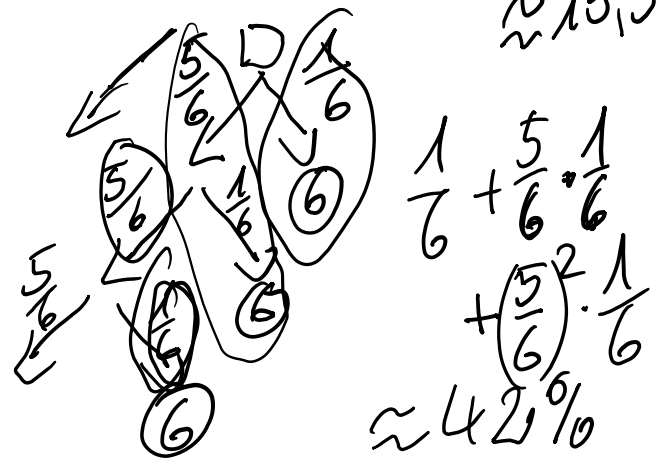
- Probability within 10 throws of a dice to obtain 3 times a 6?
- Probability to move the first time in the game „Mensch ärgere Dich nicht“ (Hombre, no te enfades)?
- You play four different statistically independent Lotteries. The probability of success is always 20%.
 - What is the expected value and the variance and standard deviation?
 - What is the expected payoff, if every success counts 10 Euros?
 - Calculate the probability to win exactly two times
 - Calculate the probability that you loose at least three times.

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$n = 10 \quad k = 3 \quad p = \frac{1}{6}$

$$\binom{10}{3} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^7 = \frac{1 \cdot 5^7}{6^{10}} \binom{10}{3}$$

$\approx 15,5\%$



a) $n = 4 \quad p = \frac{2}{10}$

$$E(X) = n \cdot p = \frac{8}{10} = 0,8$$

$$Var(X) = n p (1-p) = \frac{4 \cdot 2}{10} \cdot \frac{8}{10} = 0,64$$

$$std\ dev = \sqrt{0,64} = 0,8$$

b) $E(10X) = 10 E(X) = 8$

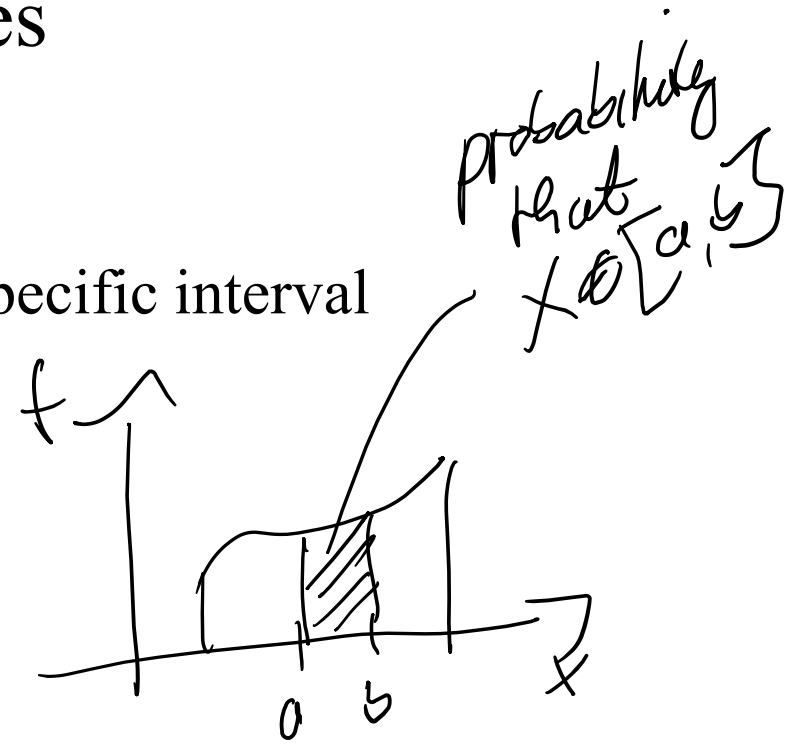
c) $k = 2 \quad \binom{4}{2} (0,2)^2 (0,8)^2 = \frac{1 \cdot 2 \cdot 3 \cdot 4}{1 \cdot 2 \cdot 1 \cdot 2} \frac{4}{100} \frac{64}{100} \approx 15\%$

d) switch to $p = 0,8 \quad \binom{4}{3} (0,8)^3 0,2 + \binom{4}{4} 0,8^4 \cdot 0,2^0 \approx 81\%$

Continuous Random Variables

The probability of the realization of X within some specific interval $[a,b]$ is the area below the density function $f(x)$:

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



Properties of the density function

- Values of $f(x)$ cannot be interpreted!
- If the length of the intervall reaches zero ($a = b$), also the area reaches zero:

$$\rightarrow \boxed{P(X = x) = 0}$$

- The whole area below the density function equals 1 (certain event):

$$\boxed{\int_{-\infty}^{+\infty} f(x) dx = 1}$$

Example Density Function

Student Mike is always late in the statistics lecture. His delay is a continuous random variable with the following density function

$$f(x) = \begin{cases} a - \frac{1}{8}x & \text{for } 0 \leq x \leq 4 \\ 0 & \text{otherwise} \end{cases}$$

1. What is a?
2. Sketch the density graphically
3. Calculate the probability that Mike is between 1 and 2 minutes late at some specific day.
4. Sketch the cumulative distribution function graphically
5. Calculate the Expected value and Variance

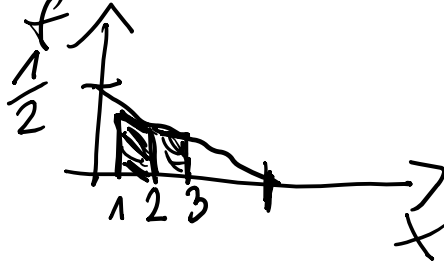
$$\int_0^4 (a - \frac{1}{8}x) dx = \left[ax - \frac{1}{16}x^2 \right]_0^4 = 4a - \frac{1}{16} \cdot 16 = 4a - 1 = 1 \Rightarrow a = \frac{1}{2}$$

$$4a - \frac{1}{16} \cdot 16 = 1 \Rightarrow 4a - 1 = 1 \Rightarrow a = \frac{1}{2}$$

a [1,2] density, cumulative distribution, E(x), Var(x)

Example Density Function

$$f(x) = \begin{cases} a - \frac{1}{8}x & \text{for } 0 \leq x \leq 4 \\ 0 & \text{otherwise} \end{cases} \Rightarrow f(4) = \frac{1}{2} - \frac{1}{8} \cdot 4 \Rightarrow \frac{1}{2} = \frac{1}{8} \cdot 4 \Rightarrow x = 4$$



$$3.) \int_1^2 \left(\frac{1}{2} - \frac{1}{8}x \right) dx$$

$$= \left. \frac{1}{2}x \right|_1^2 - \left. \frac{x^2}{16} \right|_1^2 = \left(1 - \frac{1}{2} \right) - \left[\frac{4}{16} - \frac{1}{16} \right]$$

$$= \frac{1}{2} - \frac{3}{16} = \frac{8-3}{16} = \frac{5}{16}$$

```

solve( int(a-(1/8)x dx; 0..4)=1,a)
plot(1/2-x/8,0..4)
int(1/2-x/8),0..y
plot(-1/16 (-8 + y) y,0..4)
int(x(1/2-x/8)),0..4
int((x-4/3)^2(1/2-x/8)),0..4
    
```

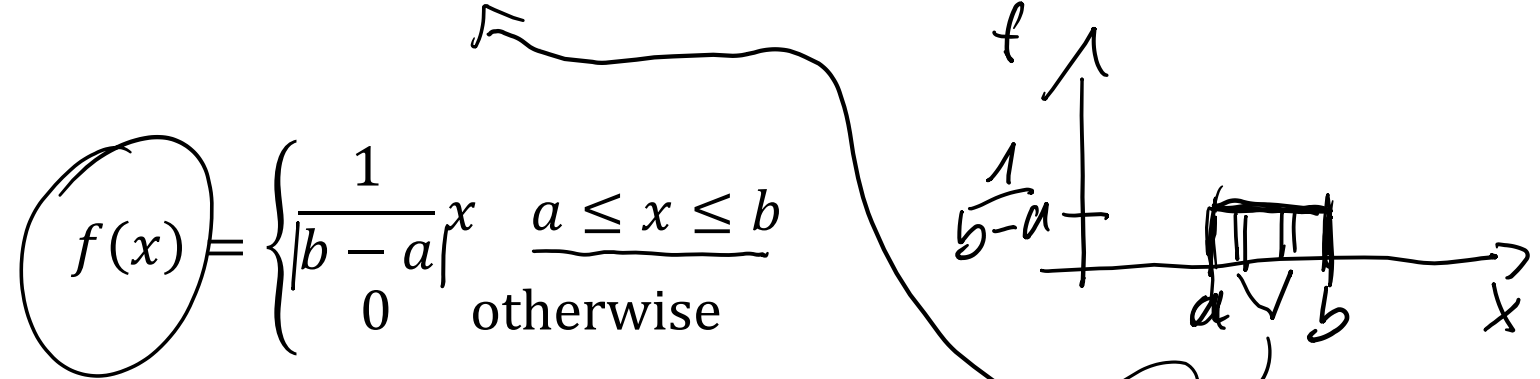
$$4.) F(x) = \int_0^x f(y) dy$$

<https://www.wolframalpha.com/>

$$5.) E(x) = \int_0^4 x f(x) dx \quad \text{Var}(x) = \int_0^4 (x - E(x))^2 f(x) dx$$

Uniform distribution

The density is constant over a given Interval $[a,b]$. This means, that every subinterval with same length have same probability

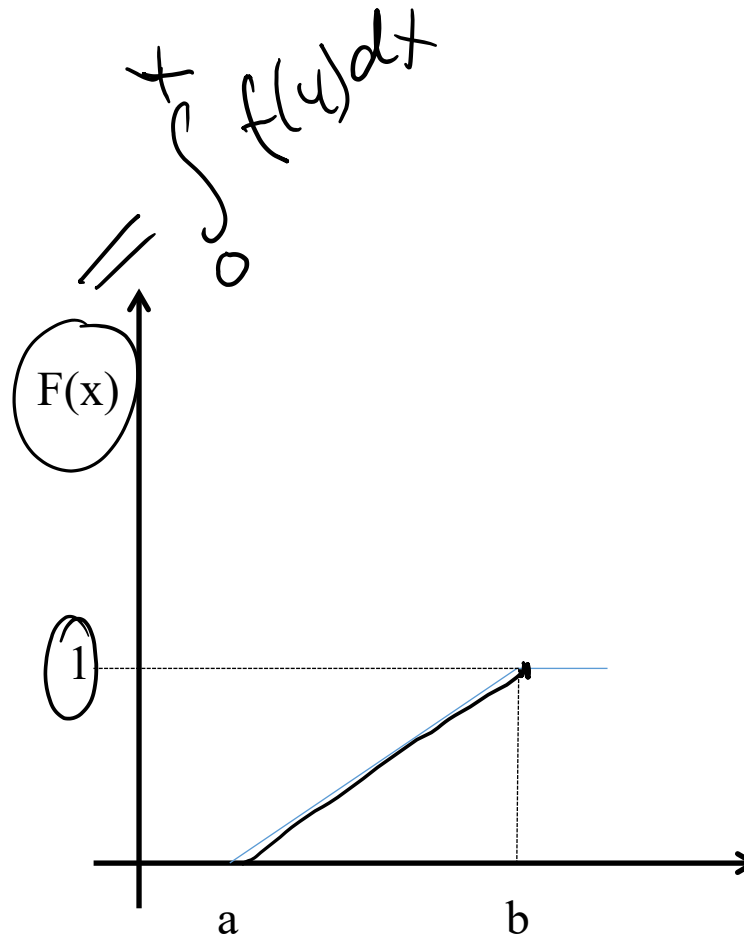
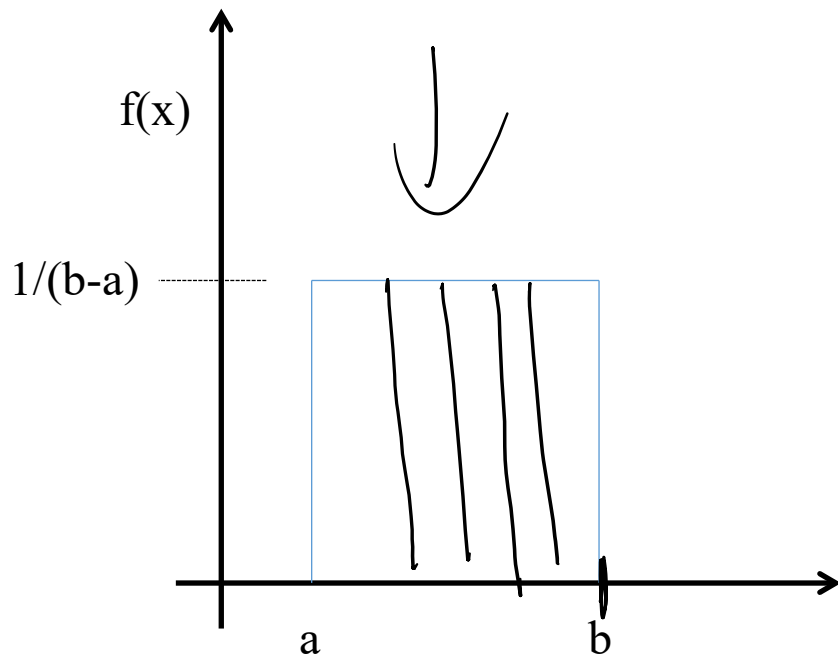


$$E(X) = \int x f(x)$$

$$E(X) = \frac{a+b}{2}$$

$$\int (x - E(X))^2 f(x) dx = \text{Var}(X) = \frac{1}{12} (b-a)^2$$

Uniform distribution

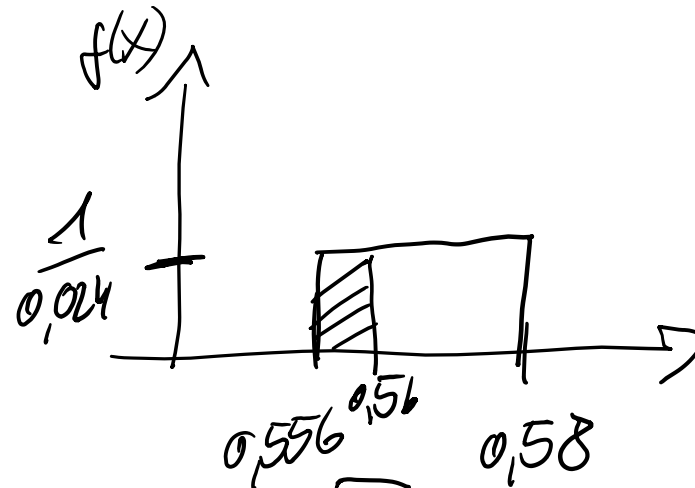


Example

Suppose you are the COP (Chief of Production) within a brewery. Furthermore the machine shows a filling quantity of the bottles is 0,568L. But you suppose, that real filling quantity is uniformly distributed around 0,556L and 0,58L.

What is the probability that the filling quantity is less than 0,56L?

$$f(x) = \frac{1}{0,58 - 0,556} \times \begin{cases} 0,556 \leq x \leq 0,58 \\ 0 \text{ otherwise} \end{cases}$$
$$= \frac{1}{0,024} \times$$



$$\Rightarrow P(x \leq 0,56L) = \frac{0,004}{0,024} = \frac{1}{6} \approx 16,6\%$$

Normal Distribution

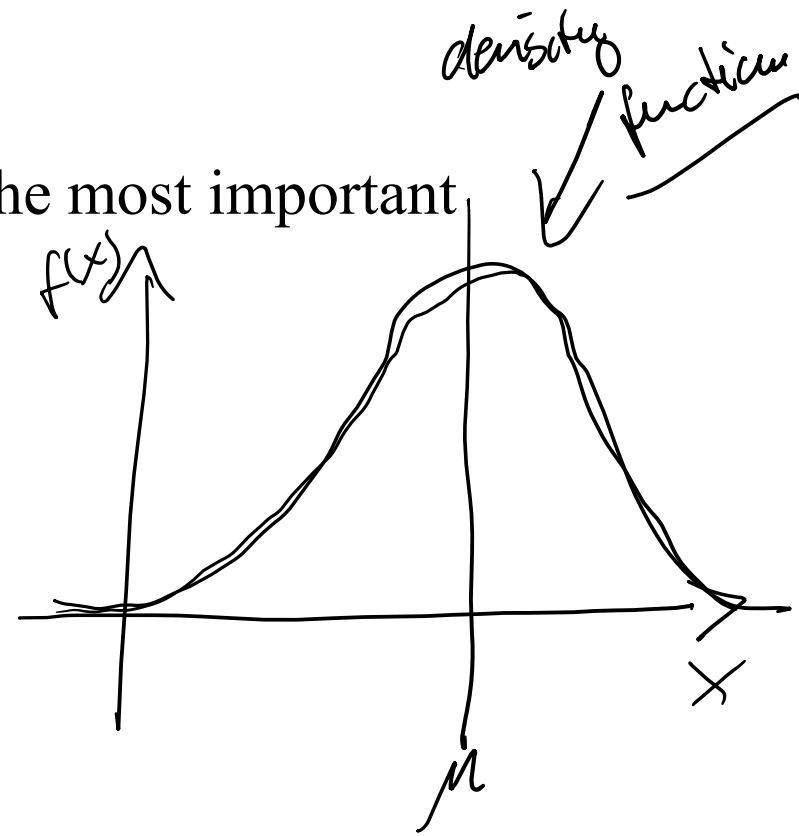
The normal distribution or gaussian distribution is the most important distribution.

$$X \sim N(\mu; \sigma^2)$$

We write:

The density function is given by:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



The cumulative distribution function by:

$$F(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_0^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = \frac{1}{2} \left(1 + \operatorname{erf} \left(\frac{x-\mu}{\sqrt{2\pi}\sigma} \right) \right)$$

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy$$

(Error function)

$$E(X) = \mu$$

$$\operatorname{Var}(X) = \sigma^2$$

The importance of the normal distribution

- Many empirical distributions follow at least approximately the normal distribution
- Many discrete distributions can be approximated by the normal distribution. I.e. Binomial distribution
- The distribution of the sample independent with respect to the underlying true distribution is approximated by the normal distribution for large sample size N (central limit theorem)
- The normal distribution is the basis theoretical models (i.e. white noise)

$$f(x) + \epsilon$$

often assumed to be normally distributed

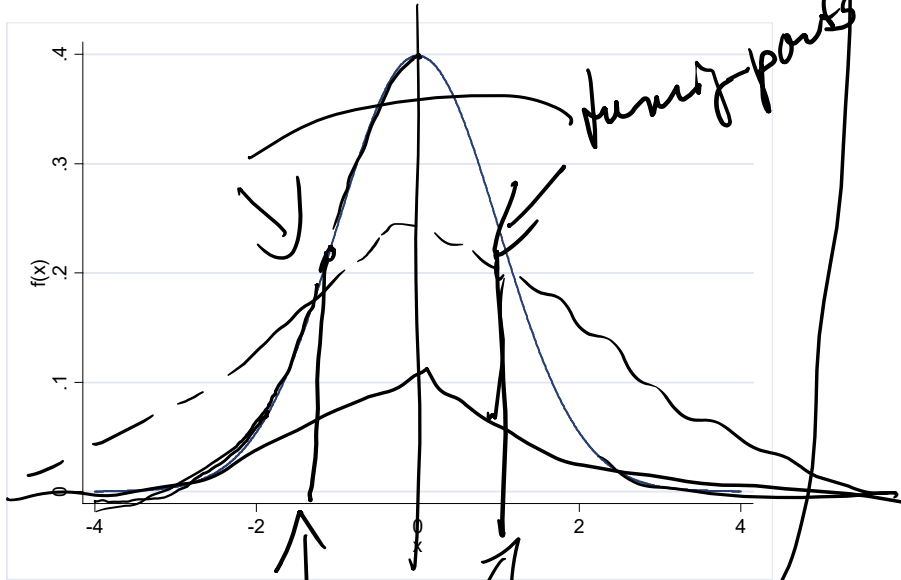
Normal distribution

$(x-\mu)^2 \in$ symmetric function



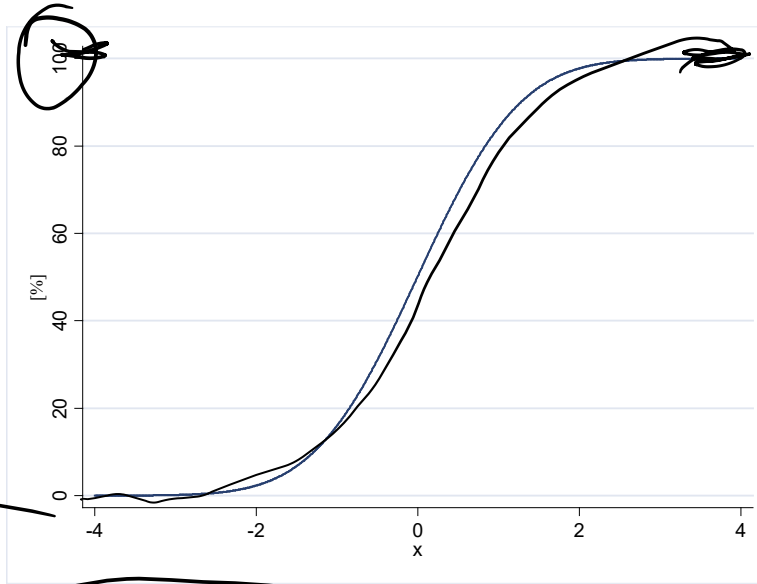
$$f^{(\mu, \sigma^2)}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$F^{(\mu, \sigma^2)}(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt$$



Density function $\mu = 0$ und $\sigma = 1$

$\mu - 2$ $\mu + 2$



Cumulative distribution function $\mu = 0$ und $\sigma = 1$

Properties of the normal distribution

- The normal distribution is a 2-parameter-distribution in μ and σ
- $N(\mu, \sigma^2)$ is symmetrically around $x = \mu$
- The density function has turning points at $x = \mu + \sigma$ and $x = \mu - \sigma$
- The density function flattens if the variance goes up.
- The density function reaches the x-axes asymptotically at $+\infty$ and $-\infty$.

The standard normal distribution

- Normal distribution with $\mu=0$ $\sigma=1$

$$Z \sim N(0;1)$$

- Maximum at $z = 0$

- Turning points at $z = -1$ and $z = 1$

- In order to calculate probabilities the cumulative distribution function is tabulated and implemented in every standard spreadsheet programm like excel (MS office) or calc (libreoffice)

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \int_0^z e^{-\frac{1}{2}x^2} dx .$$

- Every normally distributed random variable can be transformed in a standard normally distributed random variable (Standardisierung)

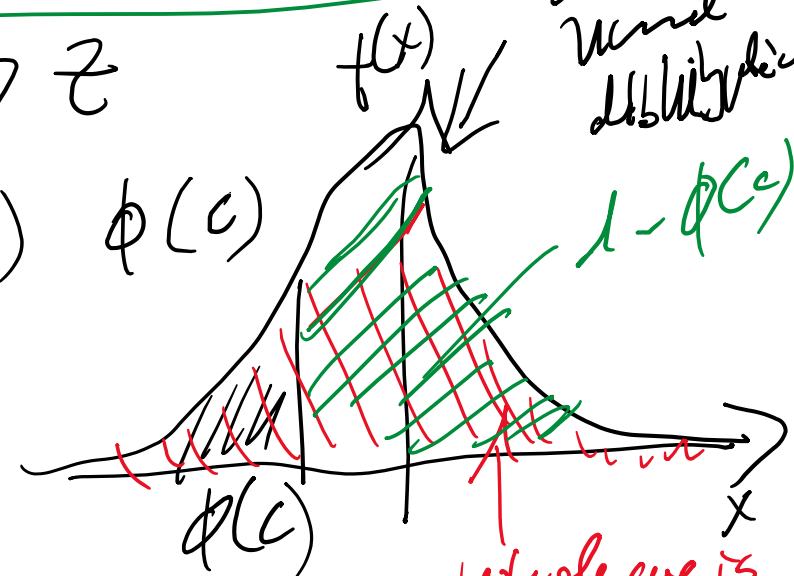
Normal distribution and standard normal distribution

- Suppose $X \sim N(\mu; \sigma^2)$ then

$$X \rightsquigarrow Z$$

$$Z = \frac{x - \mu}{\sigma} \sim N(0; 1)$$

$$\phi(z) \quad \phi(c)$$



- We have $\Phi(-c) = 1 - \Phi(c)$

- Generally, we obtain:

$$x = z\sigma + \mu$$

$$P(X \leq x) = P(\sigma Z + \mu \leq x) = P\left(Z \leq \frac{x - \mu}{\sigma}\right)$$

or

$$\begin{aligned} \sigma z &\leq x - \mu \\ z &\leq \frac{x - \mu}{\sigma} \end{aligned}$$

$$F^{(\mu, \sigma^2)}(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

whole area is just 1

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Normal distribution and standard normal distribution

- Suppose $X \sim N(\mu; \sigma^2)$ then

$$Z = \frac{x - \mu}{\sigma} \sim N(0;1)$$

- We have $\Phi(-c) = 1 - \Phi(c)$
- Generally, we obtain:

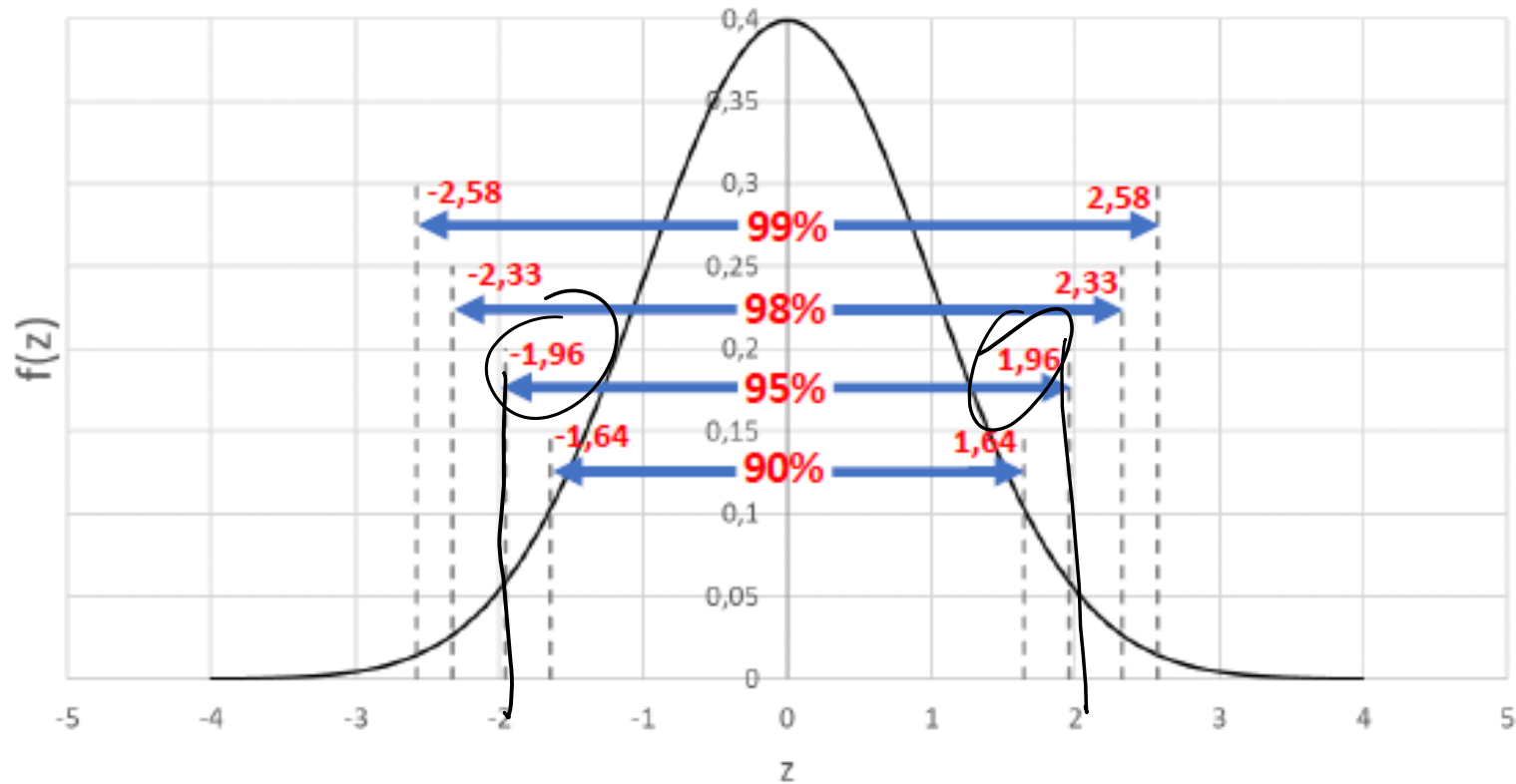
$$P(X \leq x) = P(\sigma Z + \mu \leq x) = P(Z \leq \frac{x - \mu}{\sigma})$$

or

$$F^{(\mu, \sigma^2)}(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

Relevant Values of the standard normal distribution

z	F(z)
-2.58	0.5 %
-2.33	1 %
-1.96	2.5 %
-1.64	5 %
-1.23	10 %



Point estimator for the mean μ of a parent distribution

General random variable

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

specific sample with concrete values

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

X_i : Capital letters \rightarrow random variables.

x_i : Lowercase letters \rightarrow concrete values of the random variables within a specific sample

\Rightarrow The arithmetic mean of a sample is the estimator of the true mean of the parent population.

Point estimator for the proportion π of a parent distribution

In analogy to the estimation of the unknown mean via the arithmetic mean, the estimator for the unknown proportion π of a parent population is given by:

$$\hat{\pi} = \frac{k}{n} \quad k: \text{ number of successes out of sample size } n$$

$\implies \hat{\pi}$ is the estimator for the true proportion π of the parent population.

Point estimator for the variance σ^2 of a parent distribution

The estimator $\hat{\sigma}^2$ of the unknown variance σ^2 of a parent population is given by:

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

\Rightarrow The sample variance $\hat{\sigma}^2$ is the (unbiased) estimator of the true variance σ^2 of the parent population.

Classical criteria of optimal estimators

Unbiased estimator:

An estimator $\hat{\theta}$ of a true parameter θ is called unbiased, if the expected value of the estimator equals the true value of the estimated parameter θ :

$$E(\hat{\theta}) = \theta \quad \forall \theta \in \Theta \quad (\Theta : \text{parameter space})$$

Classical criteria of optimal estimators

Asymptotically unbiased estimator:

An estimator $\hat{\theta}$ of a true parameter θ is called asymptotically unbiased, if a sequence of the expected value of the estimator $E(\hat{\theta}_n)$ reaches in the limit the true value of the estimated parameter θ :

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta \quad \forall \theta \in \Theta \quad (\Theta : \text{parameter space})$$

Classical criteria of optimal estimators

Efficiency:

Efficiency describes the accuracy of an unbiased estimator. The accuracy is generally measured via the variance. The estimator with the minimal variance is called efficient.

\implies An estimator $\hat{\theta}^*$ is called efficient if

$$E(\hat{\theta}^*) = \theta$$

$$\text{Var}(\hat{\theta}^*) \leq \text{Var}(\hat{\theta})$$

for all unbiased estimators $\hat{\theta}$.

(desired) properties of estimators

Consistency:

An estimator, which approaches the true parameter, if the sample size is raised is called consistent. A sequence $(\hat{\theta}_n)$ of estimators is called consistent if for all $\epsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| < \epsilon) = 1$$

this means $\hat{\theta}_n$ is converging stochastically to θ and is located for a given sample size n within an epsilon-neighbourhood $U_\epsilon(\theta)$ of the true value.

(desired) properties of estimators

- ▶ **consistency in mean square**

(MSE: Mean squared error):

A sequence $(\hat{\theta}_n)$ of estimators is called consistent in mean square if:

$$\lim_{n \rightarrow \infty} E([\hat{\theta}_n - \theta]^2) = 0$$

with

$$MSE := E([\hat{\theta}_n - \theta]^2) = Var(\hat{\theta}_n) - (E(\hat{\theta}_n) - \theta)^2$$

→ with raising sample size n the difference between variance and the displacement vanishes.

(desired) properties of estimators

Sufficiency:

An estimator is called sufficient, if all information about the true parameter is used.

Idea:

Suppose you have a data set \mathbf{S} of n independently distributed random variables. Then we look for some mapping $T(\mathbf{S})$ whose values contain all information of original data set. Thus, efficiency is a concept in order to reduce high-dimensional data-vectors to a usable size.

Exercise

1. $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
2. $\frac{K}{\bar{X}}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
3. $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
4. $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
5. $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
6. Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

Exercise

7. Suppose $X \sim B(n, p)$: There is no unbiased estimator of the standard deviation of the binomial distribution.
8. Show $MSE := E([\hat{\theta}_n - \theta]^2) = Var(\hat{\theta}_n) - (E(\hat{\theta}_n) - \theta)^2$
9. Suppose X_1, \dots, X_n are independently Bernoulli-distributed, then $T(X) = X_1 + \dots + X_n$ is a sufficient estimator for p (T is the total number of successes).
- * Suppose X_1, \dots, X_n are independently normally distributed and σ^2 is known. Then \bar{X} is a sufficient estimator for μ .

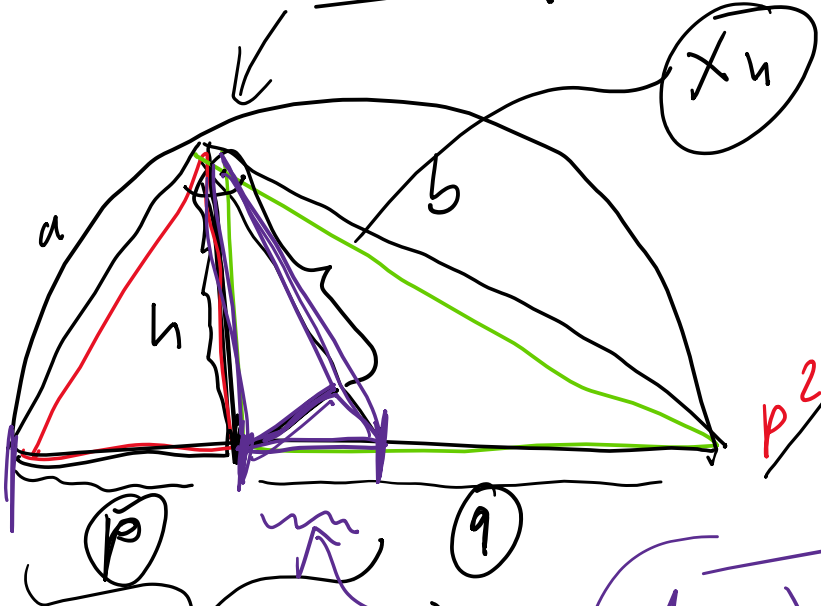
$x_g = \sqrt{x_h \cdot x_a}$ | $x_g^2 = x_h \cdot x_a \Rightarrow p \cdot q = \frac{2}{\frac{1}{p} + \frac{1}{q}} \cdot \frac{1}{2}(p+q) = \frac{1}{\frac{p+q}{p \cdot q}} \cdot (p+q)$ two numbers
 $= \frac{1}{\frac{1}{p \cdot q}} = p \cdot q = x_g^2$ q.e.d. ①
 theorem of Thales

$h = \sqrt{p \cdot q}$

$x_a = \frac{1}{2}(p+q)$
 $x_h = \frac{2}{\frac{1}{p} + \frac{1}{q}}$
 $x_g = \sqrt{p \cdot q}$

$a^2 + b^2 = c^2$

$p^2 + h^2 + q^2 + h^2 = (p+q)^2 = p^2 + q^2 + 2pq$
 $\Rightarrow 2h^2 = 2pq \Rightarrow h = \sqrt{pq} = x_g$



$x_a = \frac{1}{2}(p+q)$

$\frac{1}{2}(p+q) - p$

$c = p+q$

geometrischer Mittel
 Theorem

$$a) \binom{k}{n} p^k (1-p)^{n-k} = P(k=3) \quad 3)$$

$$= \binom{3}{12} \left(\frac{3}{4}\right)^3 \left(\frac{1}{4}\right)^{12-3} = \frac{12!}{3!(12-3)!} \cdot \frac{3^3 \cdot 1}{4^{12}}$$

$$= \frac{1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 \cdot 9 \cdot 10 \cdot 11 \cdot 12}{1 \cdot 2 \cdot 3 \cdot 1 \cdot \cancel{2 \cdot 3} \cdot 4 \cdot \dots \cdot 9} = \frac{1 \cdot 5 \cdot 11 \cdot 12 \cdot 7 \cdot 3 \cdot 9}{7 \cdot 3 \cdot 4^{12}} = \frac{5 \cdot 11 \cdot 3}{4^{11}} = \frac{15 \cdot 11}{4^{11}} = \frac{165}{4^{11}}$$

$$b) P(k \geq 8) = \sum_{k=8}^{12} \binom{k}{12} p^k (1-p)^{12-k} \rightarrow \text{do in excel} = 0,004\% \\ = \binom{8}{12} (0,75)^8 (0,25)^4 + \binom{9}{12} (0,75)^9 (0,25)^3 + \dots \quad k=12$$

$$c) E(x) = np = 12 \cdot \frac{3}{4} = 9$$

$$\text{Var}(x) = np(1-p) = 9 \cdot \frac{1}{4} = \frac{9}{4}$$

$$\text{Stdev} = \sqrt{\text{Var}} = \sqrt{np(1-p)} = \frac{3}{2}$$

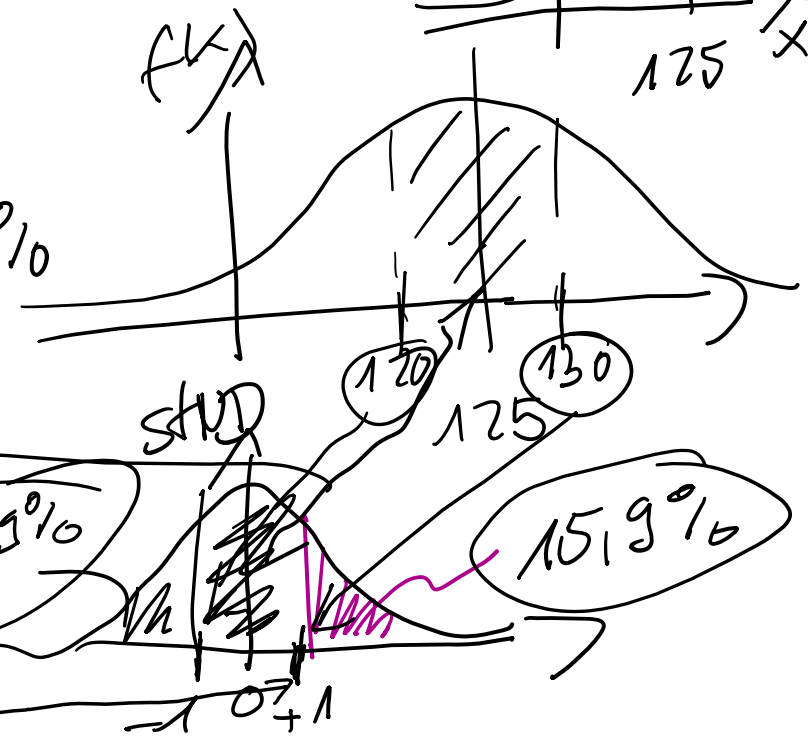
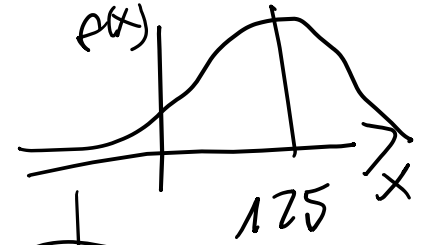
Within a manufacturing process on average 75% of the tools are correct.

- Calculate the probability, that within a sample of $n = 12$ you have exactly 3 correct tools.
- Calculate the probability, that within a sample of $n = 9$ you have at least 8 correct tools.
- Calculate the expected value, variance and standard variation of a sample of $n=25$.

4) A producer of cocoa knows from experience, that the weight of the 125g-packs is normally distributed with $\mu = 125$ g and variance of $\sigma^2 = 25$ g.

- a) What is the probability that the weight of a pack is exactly 125 g (argue)?
- b) What is the probability, that the weight of a pack is within 120 g and 130 g?
- c) What is the probability, that the weight of a pack is less than 110 g?
- d) What is the probability, that the weight of a pack is more than 140 g?
- e) Calculate the symmetric interval around the expected value, such that with a probability of 95% the true weight of a pack is within this interval.
- f) Sketch your results graphically with the given distribution and the standard normal distribution.

$P(x=125) = 0$
 $P(120 \leq x \leq 130)$



$$z = \frac{x - \mu}{\sigma}$$

$$z_d = \frac{120 - 125}{5} = -1$$

$$\phi(1) = 84.1\%$$

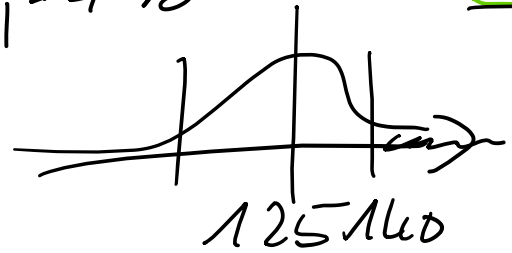
$$P(120 \leq x \leq 130) = P(-1 \leq z \leq 1)$$

$$= 100\% - 2 \cdot 15.3\% = 68.2\%$$

c) $P(x \leq 110) = P(z \leq -3) = 1 - \phi(3)$

$$= 1 - 99.86\% = 0.14\%$$

d) $P(x \geq 140) = \phi(3) = 99.86\%$



e) 95%

$$e) 95\% = P(\underline{x}_d \leq x \leq \overline{x}_u)$$

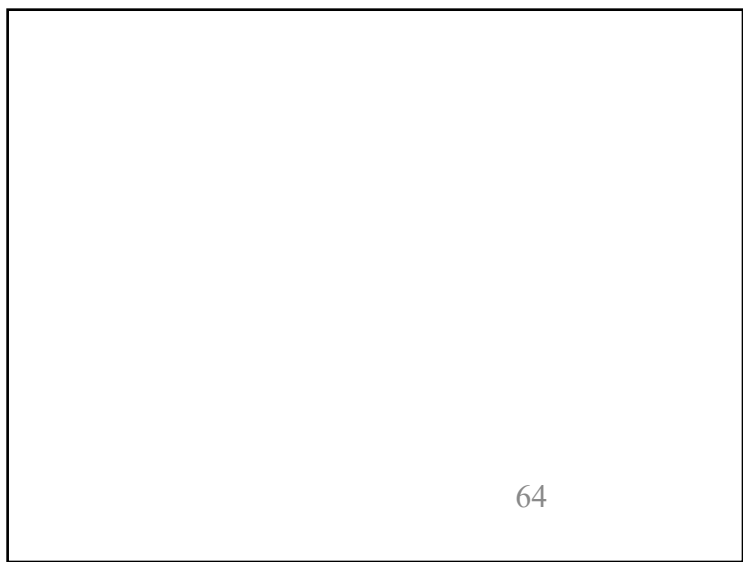
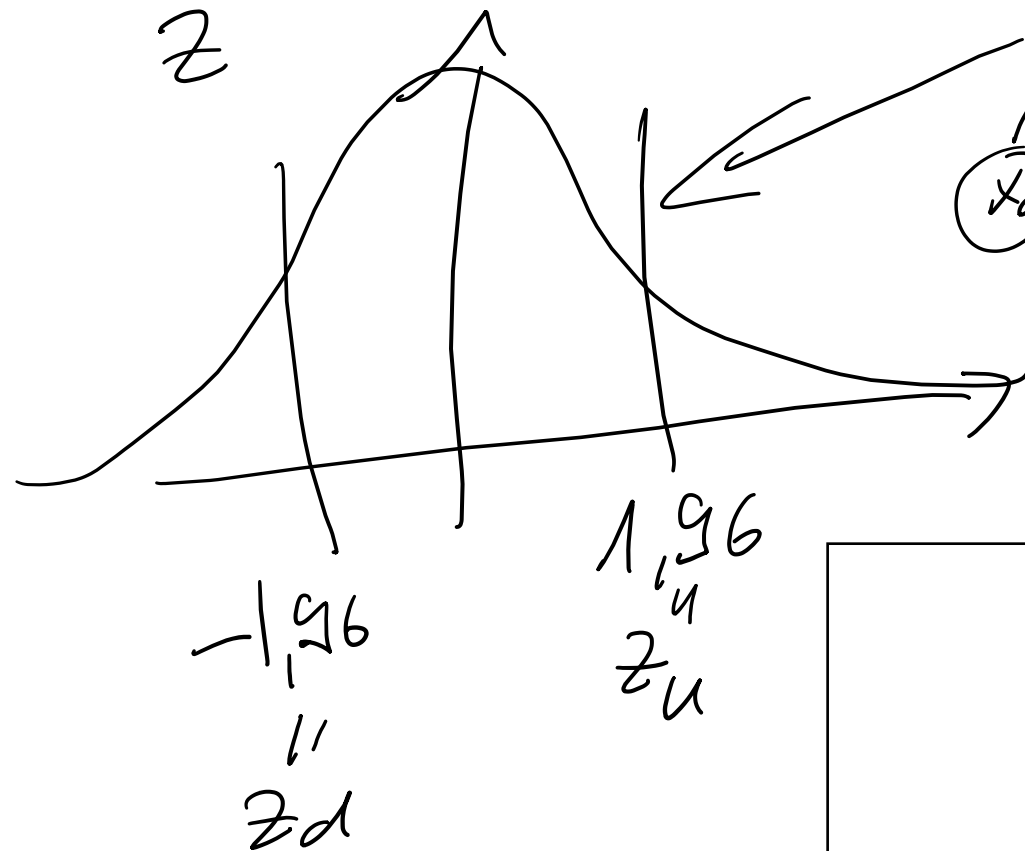
$$-1,96 = \frac{x_d - 125}{5}$$

$$\approx 5(-2) + 125 = x_d$$

$$x_d \approx 115,2$$

$$x_u \approx 134,8$$

$$z = \frac{x - \mu}{\sigma}$$



The annual yield [%] of stock investment can be approximated with a normally distributed random variable with $\mu=10$ and $\sigma=2$.

- What is the probability that the yield is within 8% und 14% liegt?
- Assume that the yields of two different years are statistically independent.
 - What is the probability that the yields in two following years is at least 8%?
 - What is the probability that the yields in the next 10 years will be exactly three times less than 11%?
- Which yield can be maximally expected with a probability of 99%?

$$\Rightarrow P(8 \leq X \leq 12)$$

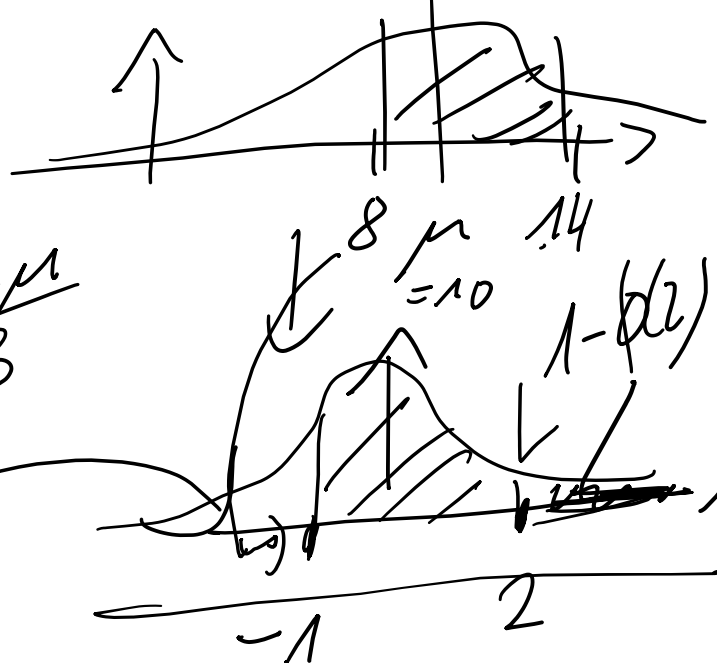
$$= 1 - [1 - \Phi(1) + 1 - \Phi(2)]$$

$$= 1 - (1 - \Phi(1) - \Phi(2))$$

$$= \Phi(1) + \Phi(2) - 1 \approx \underline{\underline{81,8\%}}$$

a)

$$z = \frac{x - \mu}{\sigma}$$



$$\Phi(1) = 84,1\% \quad \Phi(2) = 97,7\%$$

Bivariate
distribute

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Normal distribution and standard normal distribution

- Suppose $X \sim N(\mu; \sigma^2)$ then

$$Z = \frac{x - \mu}{\sigma} \sim N(0;1)$$

- We have $\Phi(-c) = 1 - \Phi(c)$
- Generally, we obtain:

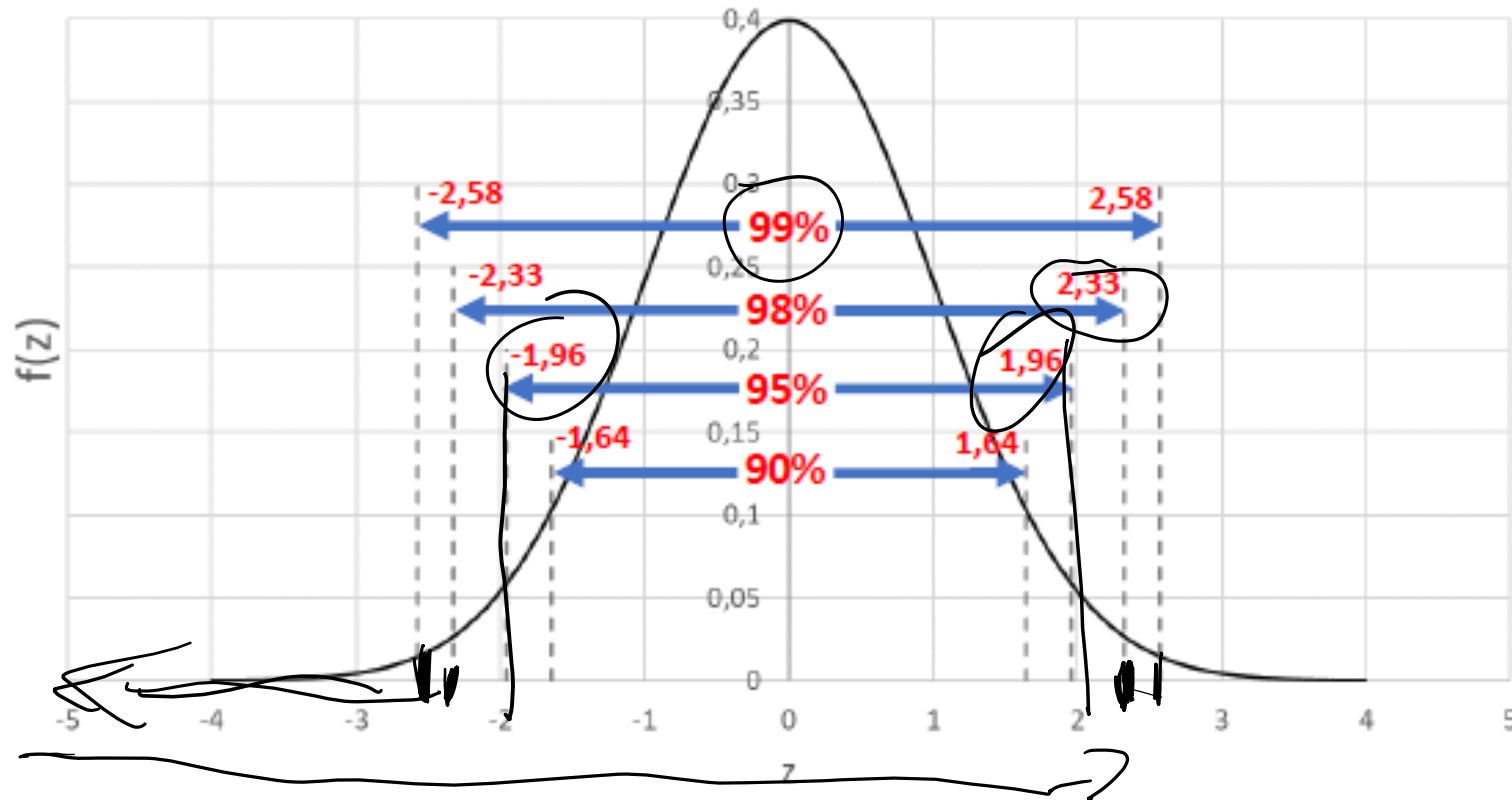
$$P(X \leq x) = P(\sigma Z + \mu \leq x) = P\left(Z \leq \frac{x - \mu}{\sigma}\right)$$

or

$$F^{(\mu, \sigma^2)}(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

Relevant Values of the standard normal distribution

z	$F(z)$
-2.58	0.5 %
-2.33	1 %
-1.96	2.5 %
-1.64	5 %
-1.23	10 %



Point estimator for the mean μ of a parent distribution

General random variable

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

specific sample with concrete values

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

X_i : Capital letters \rightarrow random variables.

x_i : Lowercase letters \rightarrow concrete values of the random variables within a specific sample

\Rightarrow The arithmetic mean of a sample is the estimator of the true mean of the parent population.

Point estimator for the proportion π of a parent distribution

In analogy to the estimation of the unknown mean via the arithmetic mean, the estimator for the unknown proportion π of a parent population is given by:

$$\hat{\pi} = \frac{k}{n} \quad k: \text{ number of successes out of sample size } n$$

$\implies \hat{\pi}$ is the estimator for the true proportion π of the parent population.

Point estimator for the variance σ^2 of a parent distribution

The estimator $\hat{\sigma}^2$ of the unknown variance σ^2 of a parent population is given by:

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

\Rightarrow The sample variance $\hat{\sigma}^2$ is the (unbiased) estimator of the true variance σ^2 of the parent population.

Classical criteria of optimal estimators

Unbiased estimator:

An estimator $\hat{\theta}$ of a true parameter θ is called unbiased, if the expected value of the estimator equals the true value of the estimated parameter θ :

$$E(\hat{\theta}) = \theta \quad \forall \theta \in \Theta \quad (\Theta : \text{parameter space})$$

Classical criteria of optimal estimators

Asymptotically unbiased estimator:

An estimator $\hat{\theta}$ of a true parameter θ is called asymptotically unbiased, if a sequence of the expected value of the estimator $E(\hat{\theta}_n)$ reaches in the limit the true value of the estimated parameter θ :

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta \quad \forall \theta \in \Theta \quad (\Theta : \text{parameter space})$$

Classical criteria of optimal estimators

Efficiency:

Efficiency describes the accuracy of an unbiased estimator. The accuracy is generally measured via the variance. The estimator with the minimal variance is called efficient.

\implies An estimator $\hat{\theta}^*$ is called efficient if

$$E(\hat{\theta}^*) = \theta$$

$$\text{Var}(\hat{\theta}^*) \leq \text{Var}(\hat{\theta})$$

for all unbiased estimators $\hat{\theta}$.

(desired) properties of estimators

Consistency:

An estimator, which approaches the true parameter, if the sample size is raised is called consistent. A sequence $(\hat{\theta}_n)$ of estimators is called consistent if for all $\epsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| < \epsilon) = 1$$

this means $\hat{\theta}_n$ is converging stochastically to θ and is located for a given sample size n within an epsilon-neighbourhood $U_\epsilon(\theta)$ of the true value.

(desired) properties of estimators

- ▶ **consistency in mean square**

(MSE: Mean squared error):

A sequence $(\hat{\theta}_n)$ of estimators is called consistent in mean square if:

$$\lim_{n \rightarrow \infty} E([\hat{\theta}_n - \theta]^2) = 0$$

with

$$MSE := E([\hat{\theta}_n - \theta]^2) = Var(\hat{\theta}_n) - (E(\hat{\theta}_n) - \theta)^2$$

→ with raising sample size n the difference between variance and the displacement vanishes.

(desired) properties of estimators

Sufficiency:

An estimator is called sufficient, if all information about the true parameter is used.

Idea:

Suppose you have a data set \mathbf{S} of n independently distributed random variables. Then we look for some mapping $T(\mathbf{S})$ whose values contain all information of original data set. Thus, efficiency is a concept in order to reduce high-dimensional data-vectors to a usable size.

Exercise

1. $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
2. $\frac{K}{\bar{X}}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
3. $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
4. $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
5. $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
6. Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

Exercise

7. Suppose $X \sim B(n, p)$: There is no unbiased estimator of the standard deviation of the binomial distribution.
8. Show $MSE := E([\hat{\theta}_n - \theta]^2) = Var(\hat{\theta}_n) - (E(\hat{\theta}_n) - \theta)^2$
9. Suppose X_1, \dots, X_n are independently Bernoulli-distributed, then $T(X) = X_1 + \dots + X_n$ is a sufficient estimator for p (T is the total number of successes).
- * Suppose X_1, \dots, X_n are independently normally distributed and σ^2 is known. Then \bar{X} is a sufficient estimator for μ .

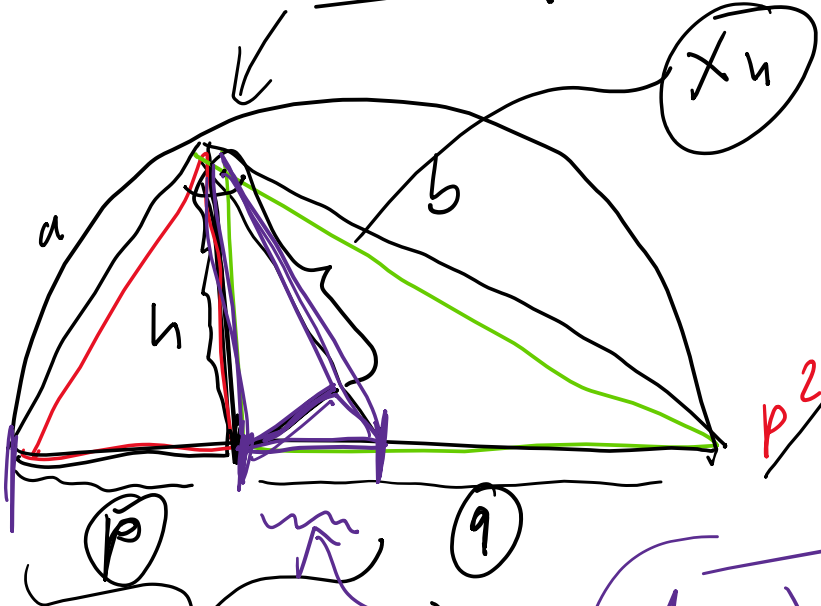
$x_g = \sqrt{x_h \cdot x_a}$ | $x_g^2 = x_h \cdot x_a \Rightarrow p \cdot q = \frac{2}{\frac{1}{p} + \frac{1}{q}} \cdot \frac{1}{2}(p+q) = \frac{1}{\frac{p+q}{p \cdot q}} \cdot (p+q)$ two numbers
 $= \frac{1}{\frac{1}{p \cdot q}} = p \cdot q = x_g^2$ q.e.d. ①
theorem of Thales

$h = \sqrt{p \cdot q}$

$a^2 + b^2 = c^2$

$p^2 + h^2 + q^2 + h^2 = (p+q)^2 = p^2 + q^2 + 2pq$
 $\Rightarrow 2h^2 = 2pq \Rightarrow h = \sqrt{pq} = x_g$

$x_a = \frac{1}{2}(p+q)$
 $x_h = \frac{2}{\frac{1}{p} + \frac{1}{q}}$
 $x_g = \sqrt{p \cdot q}$



$x_a = \frac{1}{2}(p+q)$

$\frac{1}{2}(p+q) - p$

$c = p+q$

geometric mean theorem

$$a) \binom{k}{n} p^k (1-p)^{n-k} = P(k=3) \quad 3)$$

$$= \binom{3}{12} \left(\frac{3}{4}\right)^3 \left(\frac{1}{4}\right)^{12-3} = \frac{12!}{3!(12-3)!} \cdot \frac{3^3 \cdot 1}{4^{12}}$$

$$= \frac{1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 \cdot 9 \cdot 10 \cdot 11 \cdot 12}{1 \cdot 2 \cdot 3 \cdot 1 \cdot \cancel{2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 \cdot 9}} = \frac{1 \cdot 5 \cdot 11 \cdot 12}{7 \cdot 3 \cdot 4^{12}} = \frac{5 \cdot 11 \cdot 3}{4^{11}} = \frac{15 \cdot 11}{4^{11}} = \frac{165}{4^{11}}$$

$$b) P(k \geq 8) = \sum_{k=8}^{12} \binom{k}{12} p^k (1-p)^{12-k} \rightarrow \text{do in excel} = 0,004\%$$

$$= \binom{8}{12} (0,75)^8 (0,25)^4 + \binom{9}{12} (0,75)^9 (0,25)^3 + \dots \quad k=12$$

$$c) E(x) = np = 12 \cdot \frac{3}{4} = 9$$

$$\text{Var}(x) = np(1-p) = 9 \cdot \frac{1}{4} = \frac{9}{4}$$

$$\text{Stdev} = \sqrt{\text{Var}} = \sqrt{np(1-p)} = \frac{3}{2}$$

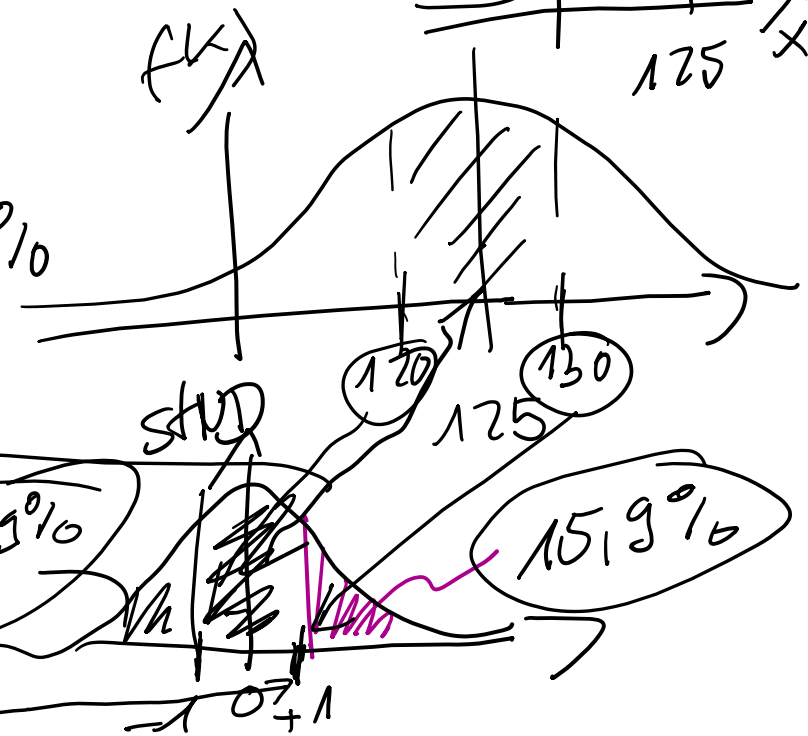
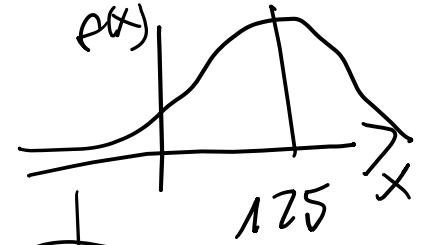
Within a manufacturing process on average 75% of the tools are correct.

- Calculate the probability, that within a sample of $n = 12$ you have exactly 3 correct tools.
- Calculate the probability, that within a sample of $n = 9$ you have at least 8 correct tools.
- Calculate the expected value, variance and standard variation of a sample of $n=25$.

4) A producer of cocoa knows from experience, that the weight of the 125g-packs is normally distributed with $\mu = 125$ g and variance of $\sigma^2 = 25$ g.

- a) What is the probability that the weight of a pack is exactly 125 g (argue)?
- b) What is the probability, that the weight of a pack is within 120 g and 130 g?
- c) What is the probability, that the weight of a pack is less than 110 g?
- d) What is the probability, that the weight of a pack is more than 140 g?
- e) Calculate the symmetric interval around the expected value, such that with a probability of 95% the true weight of a pack is within this interval.
- f) Sketch your results graphically with the given distribution and the standard normal distribution.

$P(x=125) = 0$
 $P(120 \leq x \leq 130)$



$$z = \frac{x - \mu}{\sigma}$$

$$z_d = \frac{120 - 125}{5} = -1$$

$$\Phi(1) = 84,1\%$$

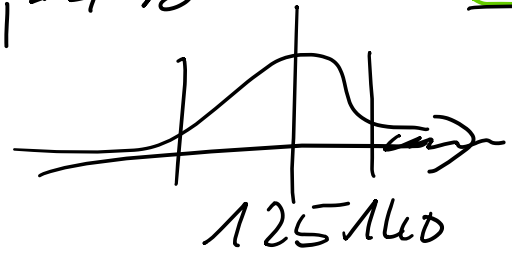
$$P(120 \leq x \leq 130) = P(-1 \leq z \leq 1)$$

$$= 100\% - 2 \cdot 15,9\% = 68,2\%$$

c) $P(x \leq 110) = P(z \leq -3) = 1 - \Phi(3)$

$$= 1 - 99,86\% = 0,14\%$$

d) $P(x \geq 140) = \Phi(3) = 99,86\%$



e) 95%

$$e) 95\% = P(\underline{x}_d \leq x \leq \underline{x}_u)$$

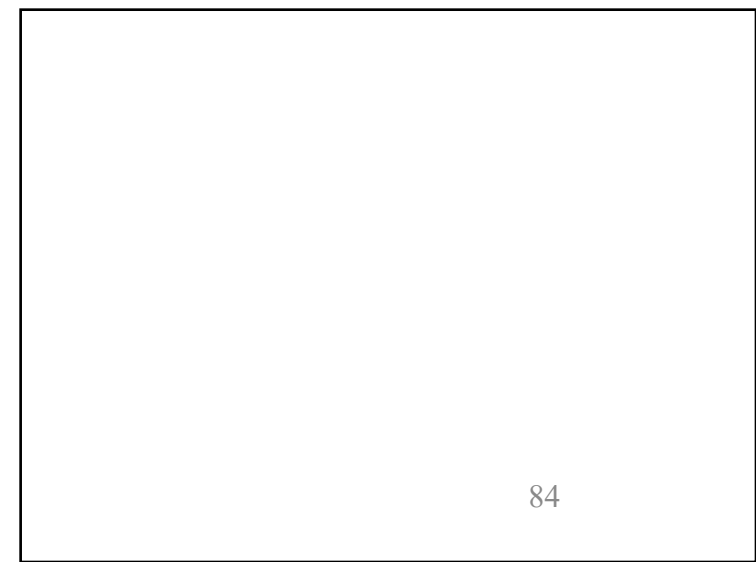
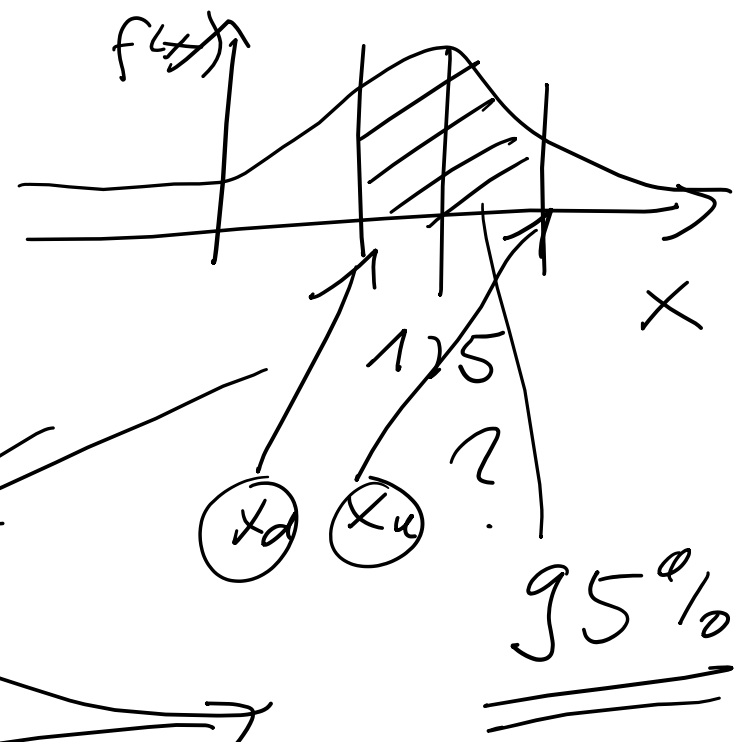
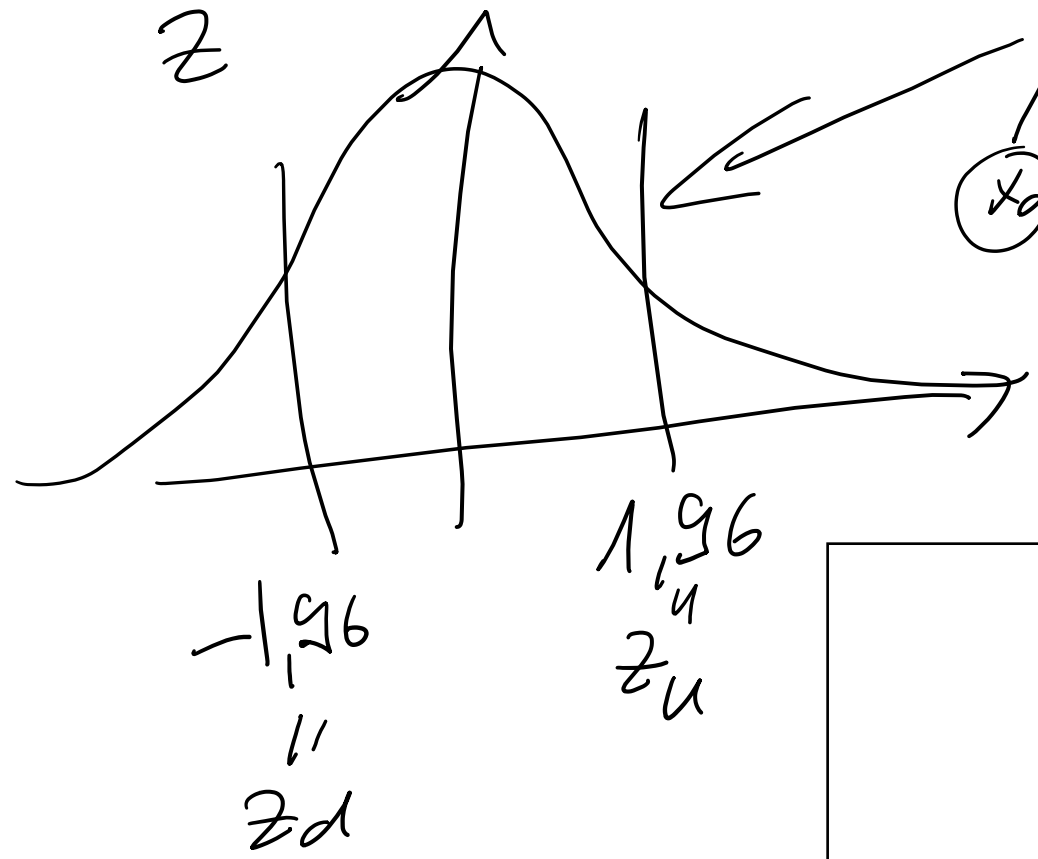
$$-1,96 = \frac{x_d - 125}{5}$$

$$\approx 5(-2) + 125 = x_d$$

$$x_d \approx 115,2$$

$$x_u \approx 134,8$$

$$z = \frac{x - \mu}{\sigma}$$



The annual yield [%] of stock investment can be approximated with a normally distributed random variable with $\mu=10$ and $\sigma=2$.

- What is the probability that the yield is within 8% und 14% liegt?
- Assume that the yields of two different years are statistically independent.
 - What is the probability that the yields in two following years is at least 8%?
 - What is the probability that the yields in the next 10 years will be exactly three times less than 11%?
- Which yield can be maximally expected with a probability of 99%?

$$\Rightarrow P(8 \leq X \leq 12)$$

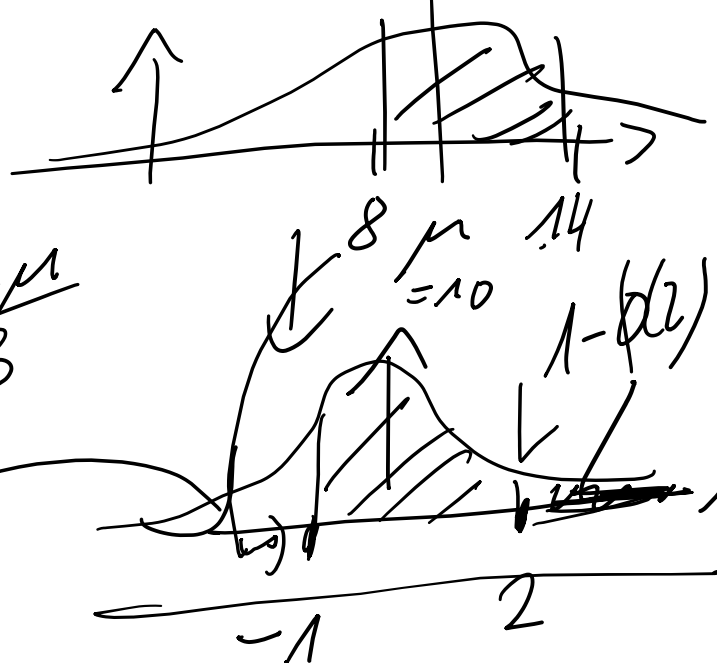
$$= 1 - [1 - \Phi(1) + 1 - \Phi(2)]$$

$$= 1 - (1 - \Phi(1) - \Phi(2))$$

$$= \Phi(1) + \Phi(2) - 1 \approx \underline{\underline{81,8\%}}$$

a)

$$z = \frac{x - \mu}{\sigma}$$



$1 - \Phi(1)$

Bivariate
distribute

$$\Phi(1) = 84,1\% \quad \Phi(2) = 97,7\%$$

5) The annual yield [%] of stock investment can be approximated with a normally distributed random variable with $\mu=10$ and $\sigma=2$.

- a) What is the probability that the yield is within 8% und 14% liegt?
- b) Assume that the yields of two different years are statistically independent.
 - i. What is the probability that the yields in two following years is at least 8%?
 - ii. What is the probability that the yields in the next 10 years will be exactly three times less than 11%?
- c) Which yield can be maximally expected with a probability of 99%?

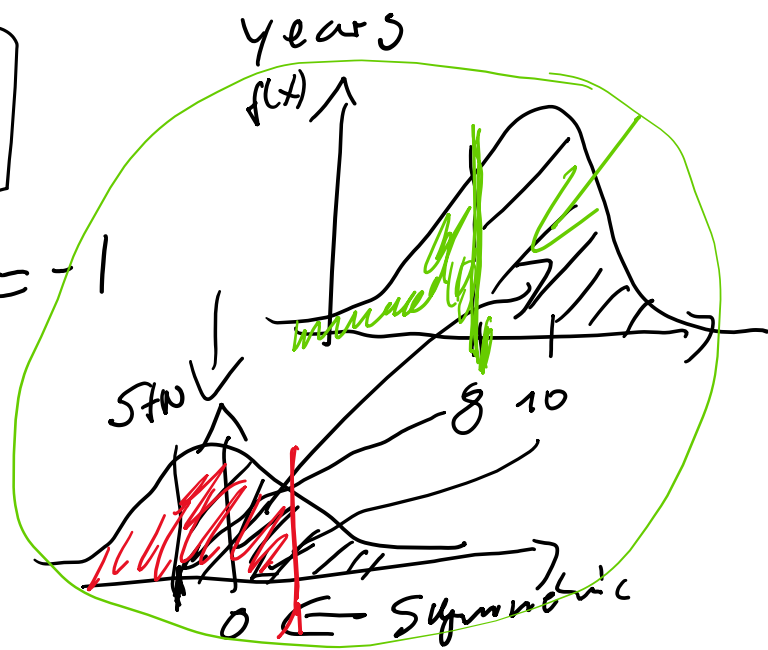
$$X \sim N(10, 4) \quad N(\mu, \sigma^2)$$

b(i) $P(X \geq 8)$ for two following

$$z = \frac{x - \mu}{\sigma}$$

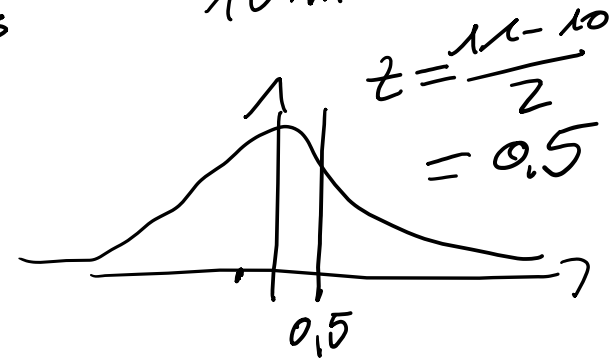
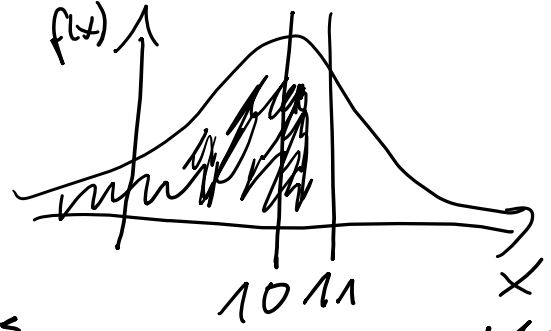
$$z = \frac{8 - 10}{2} = -1$$

$$z = 11$$



$$\Rightarrow \phi(1) = 84\% = P(X \geq 8)$$

$\Rightarrow [P(X \geq 8)]^2$ because the events are statistically independent = 0,707 \approx 71%



$$(ii) P(X \leq 11) = P(Z \leq 0,5) = 69\%$$

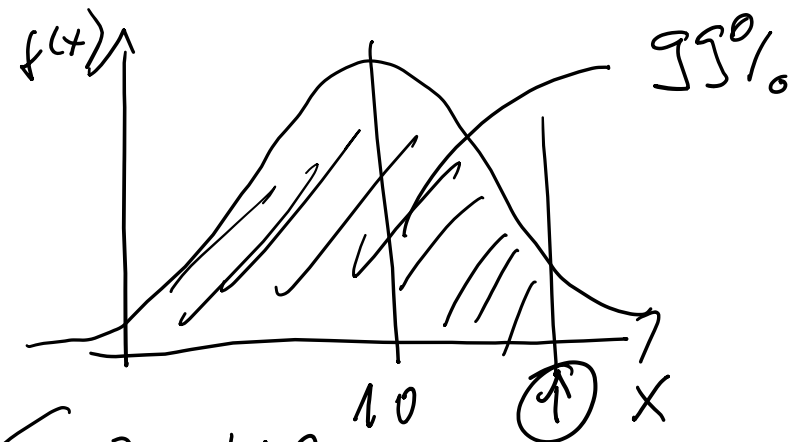
\Rightarrow there $\binom{10}{3}$ possibilities that 3 times the yield is $X \leq 11\%$ out ten events

$$\frac{10!}{7! \cdot 3!} [P(X \leq 11)]^3 \cdot [1 - P(X \leq 11)]^7$$

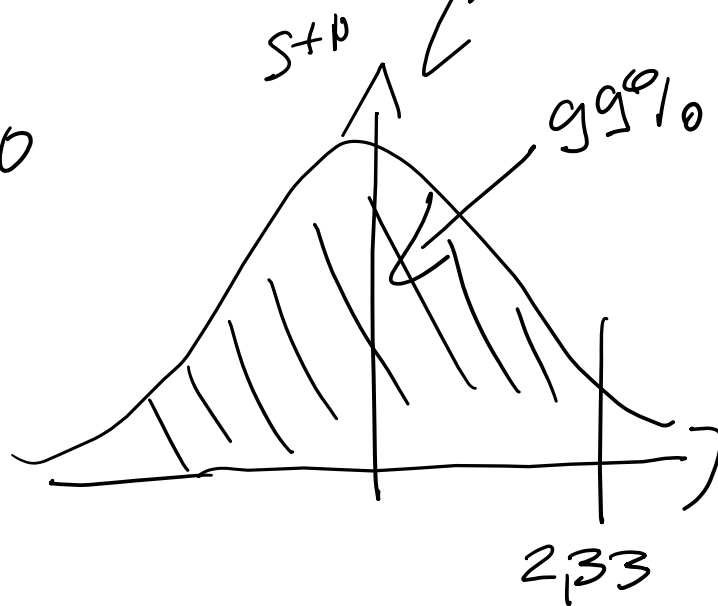
$$\approx 120 (0,69)^3 (0,31)^7 \approx 0,0105 \approx 1,05\%$$

5) The annual yield [%] of stock investment can be approximated with a normally distributed random variable with $\mu=10$ and $\sigma=2$.

- What is the probability that the yield is within 8% und 14% liegt?
- Assume that the yields of two different years are statistically independent.
 - What is the probability that the yields in two following years is at least 8%?
 - What is the probability that the yields in the next 10 years will be exactly three times less than 11%?
- Which yield can be maximally expected with a probability of 99%?



\Rightarrow which X coincides with the area below the density function of 0,99



$$P(Z \leq 2,33) = 0,99$$

$$\Rightarrow X = ? \quad Z = \frac{x - \mu}{\sigma}$$

$$= \sigma Z + \mu = 2 \cdot 2,33 + 10$$

$$= 14,66$$

maximum yield is

$$\underline{14,66\%}$$

Exercise

1. $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
2. $\frac{k}{n}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
3. $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
4. $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
5. $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
6. Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

weighted average

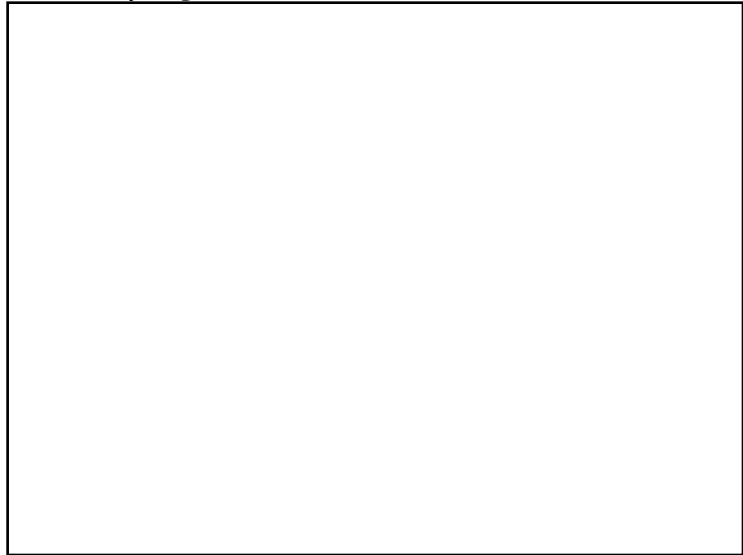
$$\begin{aligned}
 1) \quad E\left(\sum_{i=1}^n a_i X_i\right) &= \sum_{i=1}^n E(a_i X_i) \\
 &= \sum_{i=1}^n a_i E(X_i) \stackrel{!}{=} \sum_{i=1}^n a_i \cdot \mu \\
 &= \mu \left(\sum_{i=1}^n a_i\right) = \mu
 \end{aligned}$$

$$\begin{aligned}
 E(x+y) &\stackrel{!}{=} E(x) + E(y) \\
 E(ax) &= a E(x)
 \end{aligned}$$

no index i anymore

2) $\left(\frac{k}{n}\right)$ is unbiased Est. for p for binomial distribution
 number of successes \rightarrow $\left(\frac{k}{n}\right)$ is just in arithmetic mean of X_i binomial distributed
 $X_i \rightarrow 1$
 $X_i \rightarrow 0$
 number of events

but $\frac{k}{n}$ is an estimator out of the class of estimators of question 1
 $\Rightarrow \frac{k}{n}$ is also unbiased!



Exercise

- $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
- $\frac{K}{\bar{X}}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
- $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
- $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
- $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
- Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

$$(a-b)^2 = a^2 - 2ab + b^2$$

$$\underbrace{1 + 1 + 1 + 1 + \dots + 1}_{10 \text{ times}} =$$

$$\underbrace{\mu + \mu + \dots + \mu}_{10 \text{ times}}$$

\Rightarrow we know the variance of the arithmetic mean
 scale with $\frac{1}{n} = \frac{1}{n-1} [n\sigma^2 - \sigma^2] = \frac{1}{n-1} [(n-1)\sigma^2] = \sigma^2$

$$\begin{aligned} 3) E(\hat{\sigma}^2) &= E\left(\frac{1}{n-1} \sum (X_i - \bar{X})^2\right) = \frac{1}{n-1} E\left(\sum (X_i - \bar{X})^2\right) \\ &= \frac{1}{n-1} E\left(\sum (X_i - \mu + \mu - \bar{X})^2\right) = \frac{1}{n-1} E\left(\sum (X_i - \mu - (\bar{X} - \mu))^2\right) \\ &= \frac{1}{n-1} E\left(\sum (X_i - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2\right) \\ &= \frac{1}{n-1} \left[\sum E((X_i - \mu)^2) - 2 \sum E((X_i - \mu)(\bar{X} - \mu)) + \sum E((\bar{X} - \mu)^2) \right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n E(X_i - \mu)^2 - 2(\bar{X} - \mu) \left[\sum_{i=1}^n X_i - \sum_{i=1}^n \mu \right] + (\bar{X} - \mu)^2 \sum_{i=1}^n 1 \right] \\ &= \frac{1}{n-1} \left[\sum E(X_i - \mu)^2 - 2n(\bar{X} - \mu)^2 + n(\bar{X} - \mu)^2 \right] \\ &= \frac{1}{n-1} \left[\sum E(X_i - \mu)^2 - n(\bar{X} - \mu)^2 \right] \end{aligned}$$

Exercise

1. $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
2. $\frac{K}{\bar{X}}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
3. $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
4. $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
5. $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
6. Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

$$T = \sum a_i X_i; \quad \boxed{\sum a_i = 1} \quad a_i = \frac{1}{n} \text{ for efficiency}$$

$$\text{Var}(T) = \text{Var}(\sum a_i X_i) = \sum \text{Var}(a_i X_i) = \sum a_i^2 \text{Var}(X_i)$$

$$= \sum a_i^2 \sigma^2$$

$$\Rightarrow \min \sum a_i^2 \sigma^2 \text{ s.t. } \sum a_i = 1$$

$$L = \sum a_i^2 \sigma^2 + \lambda (1 - \sum a_i) = L(a_1, a_2, a_3, \dots, \lambda)$$

$$\frac{\partial L}{\partial a_i} = 2a_i \sigma^2 - \lambda = 0$$

$$= a_i = \frac{\lambda}{2\sigma^2}$$

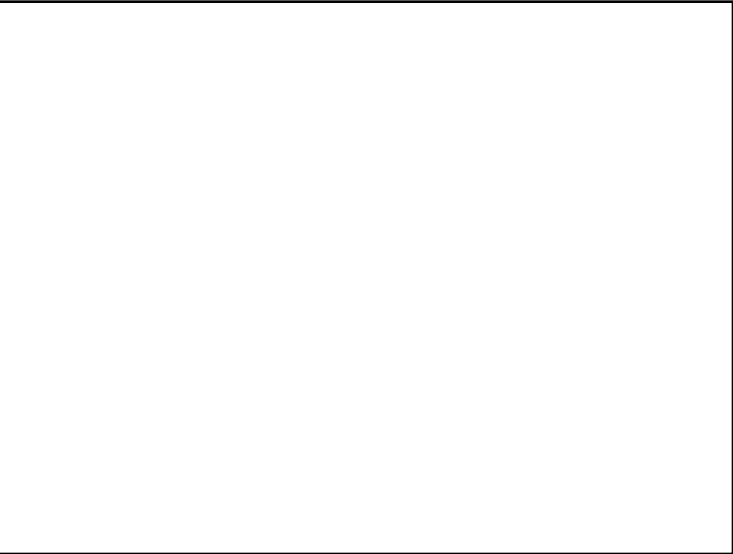
for all a_i

$$\Rightarrow a_i = a_j \text{ for } i \neq j$$

\Rightarrow all parameters have to be the same
 $\Rightarrow a_i = \frac{1}{n}$ because $\sum_{i=1}^n \frac{1}{n} = n \cdot \frac{1}{n} = 1$

X_i are statistically identically distributed and independent





Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Exercise

1. $T(X_1, \dots, X_n) = \sum_{i=1}^n a_i X_i$ with $\sum_{i=1}^n a_i = 1$ ($a_i > 0$) is an unbiased estimator for μ .
2. $\frac{K}{X}$ is an unbiased estimator for p of the binomial distribution $B(n, p)$.
3. $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator for the variance, if μ is unknown.
4. $\hat{\sigma}^2(\mu) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an unbiased estimator for the variance, if μ is known.
5. $\hat{\sigma}^2(\bar{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is an asymptotically unbiased estimator for the variance, if μ is unknown.
6. Within the linear unbiased estimators for μ , the arithmetic mean $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is an efficient estimator.

unbiased just
by definition
↓
 $\hat{\sigma}^2 = \frac{1}{n} \sum (X_i - \mu)^2$
= $E(\hat{\sigma}^2)$
of variance

$$\begin{aligned}
 5) E(\underbrace{S^2(X)}_{(3)}) &= E\left(\frac{1}{n} \sum (x_i - \bar{x})^2\right) = \frac{1}{n} \underbrace{E\left(\sum (x_i - \bar{x})^2\right)}_{(n-1)\sigma^2} \\
 &\stackrel{\text{result from (3)}}{=} \frac{n-1}{n} \sigma^2 = \lim_{n \rightarrow \infty} \left(\frac{n-1}{n}\right) (\sigma^2) = \sigma^2 \lim_{n \rightarrow \infty} \left(\frac{n-1}{n}\right) \\
 &= \sigma^2 \lim_{n \rightarrow \infty} \left(1 - \left(\frac{1}{n}\right)\right) = \sigma^2 (1 - 0) = \sigma^2 \quad \checkmark
 \end{aligned}$$

$$\begin{aligned}
 E(aX) \\
 &= a E(X)
 \end{aligned}$$

Exercise

7. Suppose $X \sim B(n, p)$: There is no unbiased estimator of the standard deviation of the binomial distribution.
8. Show $MSE := E([\hat{\theta}_n - \theta]^2) = Var(\hat{\theta}_n) - (E(\hat{\theta}_n) - \theta)^2$
9. Suppose X_1, \dots, X_n are independently Bernoulli-distributed, then $T(X) = X_1 + \dots + X_n$ is a sufficient estimator for p (T is the total number of successes).
- * Suppose X_1, \dots, X_n are independently normally distributed and σ^2 is known. Then \bar{X} is a sufficient estimator for μ .

7) Suppose there is an unbiased estimator $\hat{\theta}$ is a Polynomial in $p \rightarrow$ This is differentiable $p \geq 0$

$$E(\hat{\theta}) = \sum_{k=0}^n \theta(k) \binom{n}{k} p^k (1-p)^{n-k}$$

g) $E((\hat{\theta} - \theta)^2) = E(\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2) = E(\hat{\theta}^2) - 2E(\hat{\theta}\theta) + E(\theta^2)$

$Var(\hat{\theta}) = E((\hat{\theta} - E(\hat{\theta}))^2) = E(\hat{\theta}^2 - \underline{E(\hat{\theta})E(\hat{\theta})} + \underline{(E(\hat{\theta}))^2}) = E(\hat{\theta}^2) - (E(\hat{\theta}))^2$
 \rightarrow next lecture



Confidence level

The given probability of the **confidence interval** is called **confidence level**

confidence interval: $\hat{\theta}_a \leq \theta \leq \hat{\theta}_b$

with the lower level a and the upper level b

limits are also random variables

Estimating the interval: $P(\hat{\theta}_a \leq \theta \leq \hat{\theta}_b) = 1 - \alpha$

90%
95%
99%

most common probabilities

confidence level $1 - \alpha$: With probability $1 - \alpha$ the true parameter θ will be covered by the interval with the limits $[\hat{\theta}_a, \hat{\theta}_b]$. With probability α the true value is not covered.

Note the limits $\hat{\theta}_a, \hat{\theta}_b$ are also random variables!

↓ Keep in mind because in many statistics book the definition of confidence levels or confidence intervals do not take this into account

One- and two-sided confidence intervals

One-sided confidence interval

Either the lower limit is $-\infty$ or the upper limit is $+\infty$. The Parameter θ has with probability $1 - \alpha$ at most the value $\hat{\theta}_b$ or at least $\hat{\theta}_a$.

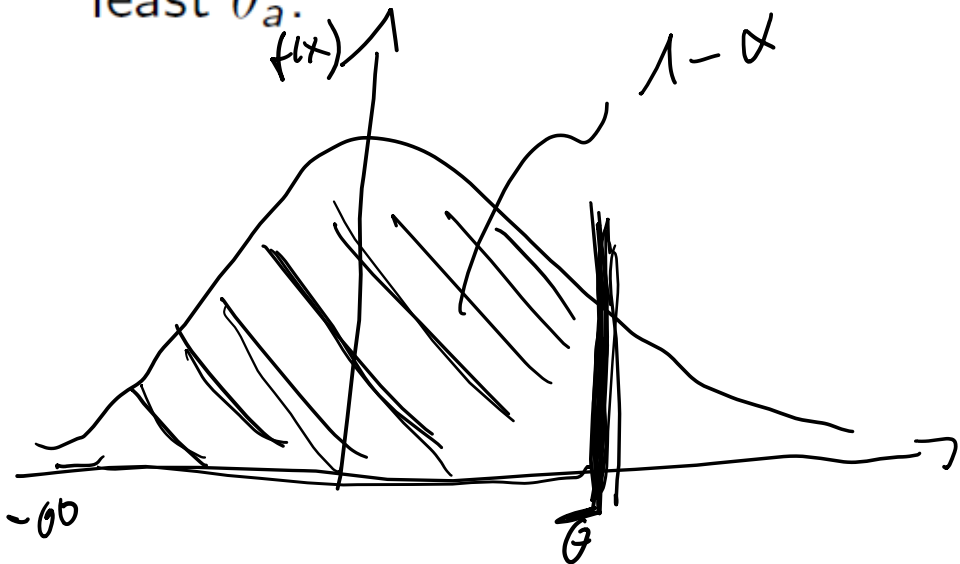
Two-sided confidence interval

We calculate the upper and lower limit ($\hat{\theta}_b$ und $\hat{\theta}_a$) of the interval, which should cover the unknown true parameter.

One- and two-sided confidence intervals

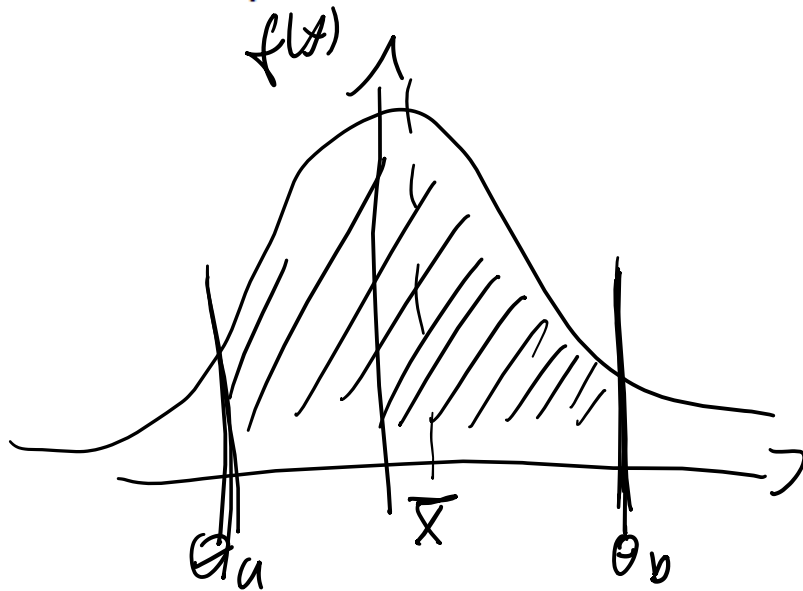
One-sided confidence interval

Either the lower limit is $-\infty$ or the upper limit is $+\infty$. The Parameter θ has with probability $1 - \alpha$ at most the value $\hat{\theta}_b$ or at least $\hat{\theta}_a$.



Two-sided confidence interval

We calculate the upper and lower limit ($\hat{\theta}_b$ und $\hat{\theta}_a$) of the interval, which should cover the unknown true parameter.



Calculating confidence intervals

In order to calculate the confidence interval we need a knowledge about the distribution and the variance of the variables. If we do not know the **variance of the population**, we have to estimate this value via the random sample:

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

- ▶ For large samples ($n \geq 30$) we can use the central limit theorem and after the standardization we can apply the **Standard-Normal distribution**.
- ▶ For small samples ($n < 30$) we use the **Student-t-Vedistribution** if we know, that the population is normally distributed. Doing this, we also standardize the variable.
- ▶ Additionally the sample size n has to be small compared to the size of the population N : $\frac{n}{N} < 5\%$ or $N = \infty$.

Confidence intervals for large samples $n \geq 30$

$$\underbrace{n \geq 30} \Rightarrow P\left(\underbrace{\bar{x} - z \frac{\hat{\sigma}}{\sqrt{n}}}_{\theta_u} \leq \underbrace{\mu}_{\theta_b} \leq \bar{x} + z \frac{\hat{\sigma}}{\sqrt{n}}\right) = \underline{1 - \alpha}$$

- ▶ $1 - \alpha$ confidence level
- ▶ \bar{x} (sample mean) point estimator for μ
- ▶ z Value of the standard normal distribution for a probability of $1 - \frac{\alpha}{2}$
- ▶ $\hat{\sigma}$ estimated standard deviation (if the variance is known, we can use σ for $\hat{\sigma}$)
- ▶ n der Stichprobenumfang
- ▶ Since $n \geq 30$ we can use the normal distribution

(If the population is normally distributed $(X) \sim N(\mu, \sigma^2)$ and the variance is known, we have without the restriction ($n \geq 30$) $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$)

z is just number of the standard normal distribution for the confidence level $1 - \alpha$

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

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

Standard deviation of X

\Rightarrow we have for the standard deviation of \bar{X}

$$\sigma_{\bar{X}}^2 = \frac{1}{n} \sigma_X^2$$

$$\Rightarrow \sigma_{\bar{X}} = \frac{\sigma_X}{\sqrt{n}}$$

Confidence intervals for small samples $n < 30$

$n \geq 30 \Rightarrow Z$

$$n < 30 \implies P\left(\bar{x} - t_{n-1} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{n-1} \frac{\hat{\sigma}}{\sqrt{n}}\right) = 1 - \alpha$$

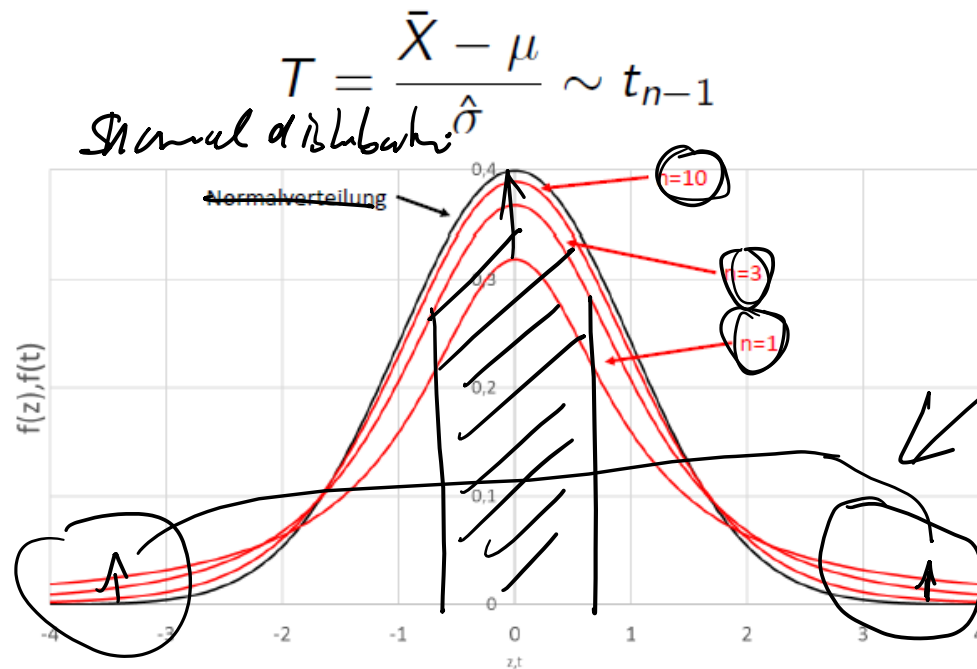
if the population is normally distributed ($X \sim N(\mu, \sigma^2)$) and the variance is not known.

t_{n-1} : value of the Student-t-distribution with $n - 1$ number of degrees of freedom

- ▶ depends of the degrees of freedom
- ▶ converges for large n to the normal distribution
- ▶ can be approximated by the normal distribution, if $n \geq 30$

Student-t-distribution

If the variance of a normally distributed population is not known, we replace within the standardization the variance σ with the sample variance $\hat{\sigma}$. The new variable Student-t distributed



t-distribution has larger tails than the normal-distribution \Rightarrow extreme events are more probable

The t-distribution is similar to the normal distribution but broader. For increasing n the t-distribution reaches the normal distribution.

Student-t-distribution

The standardized estimator of the sample mean of a normally distributed data is not normally distributed but t-distributed

Density function: $f_n(X = x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$

Erwartungswert $E(X) = 0$

Varianz $Var(X) = \frac{n}{n-2}$

Examples:

- ▶ Calculating confidence intervals of expected values of normally distributed data not knowing the variance.
- ▶ Comparison of different random samples.
- ▶ Testing the coefficients within the linear regression assuming that the error term is normally distributed.

$N(0,1) \rightarrow f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2}$

$\Gamma\left(\frac{n+1}{2}\right)$

\rightarrow so called gamma function

$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$

Example

Expected value, unknown variance, normally distributed population $n \geq 30$

The mileage of wheels with a population ($N=100.000$) should be estimated via a random sample with $n = 625$. Within the sample, we obtain an average mileage of 36.000km and $\hat{\sigma}=4.800$ km.

1. Which interval covers the true mileage of the population with probability 95%?
2. What changes, if the sample size increases to 2.500 (confidence level 95%)?
3. What changes, if the confidence level increases to 99% ($n = 625$)?

Example

Expected value, unknown variance, normally distributed population and $n < 30$

What are the average working hours of BA-students of the Jade Hochschule (N=6000) per week in order to finance their studies? In order to answer this question a random sample out of all students is drawn: Within this sample we have:

$$\bar{x} = 18,36h \text{ und } \hat{\sigma} = 3,42.$$

Calculate the confidence interval to a confidence level of 90%— if the size of the random sample in WHV is $n = 16$. Assume, that the population normally distributed

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

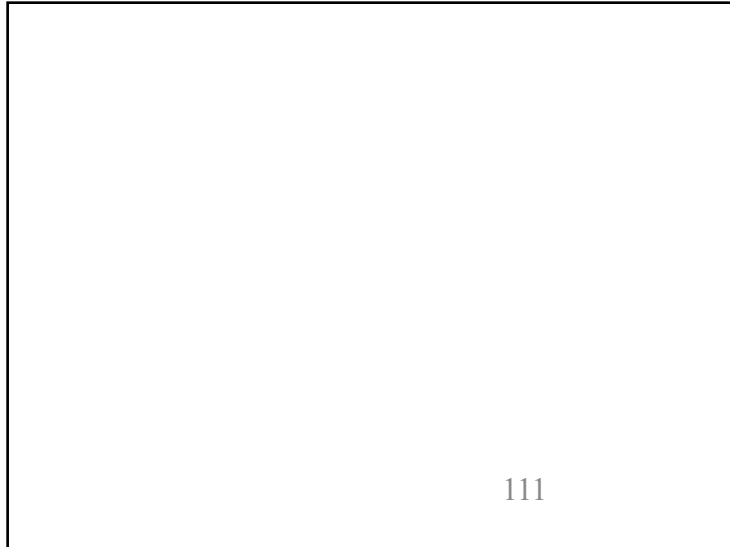
Confidence interval for a proportion

$X_i \rightarrow \begin{matrix} 1 \\ 0 \end{matrix} \Rightarrow$ proportion of successes
 $\hat{\pi} = \bar{x}$
 $X_i \sim B(1, \pi)$
 $\Rightarrow B^2 = \pi(1-\pi)$
 $\hat{\pi}$ is the estimator for the true probability π

$$n\hat{\pi}(1 - \hat{\pi}) \geq 9 \implies P\left(\hat{\pi} - z\sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}} \leq \pi \leq \hat{\pi} + z\sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}\right) = 1 - \alpha$$

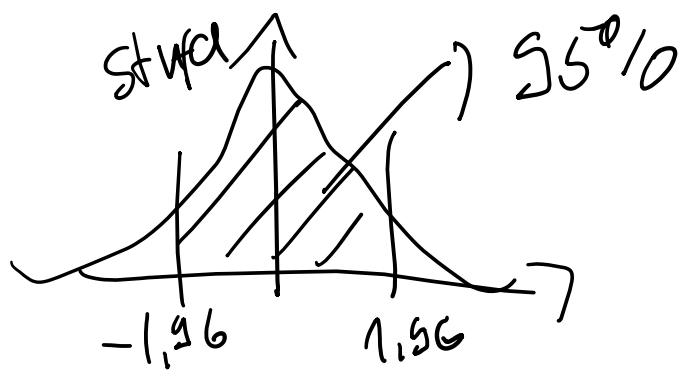
- ▶ Dichotomous attribute with a probability of success of π (i.e. $X = 0, 1$ with $P(X = 1) = \pi$).
- ▶ $\hat{\pi}$ proportion of successes in the random sample
- ▶ For $n\hat{\pi}(1 - \hat{\pi}) \geq 9$ Z is approximately normally distributed:

$$Z = \frac{\hat{\pi} - \pi}{\sqrt{\hat{\pi}(1 - \hat{\pi})}} \sqrt{n} \sim N(0, 1)$$



Example

Proportion



$$P\left(\hat{\pi} - z\sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}} \leq \pi \leq \hat{\pi} + z\sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}\right) = 1-\alpha$$

In advance of the uk poll, we want to estimated the proportion of the liberal party in UK. The random sample (asked people) has a size of $n = 1200$ and 96 person answered to vote for the liberal party. Calculate the confidence interval to the confidence level 95%.

$$\hat{\pi} = \frac{96}{1200} = 8\%$$

$$z = 1.96$$

$$n = 1200$$

$$\Rightarrow 8\% \pm \sqrt{\frac{8\% \cdot 92\%}{1200}}$$

$$\Rightarrow 8\% \pm 1.5\%$$

$$\underline{I_d \approx 6.5\% \quad I_u \approx 9.5\%}$$

Confidence intervals for the Variance

Given the parent population is normally distributed ($X \sim N(\mu, \sigma^2)$), then

$$Y = \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{(n-1)}$$

This is known from Data
Chi² - Distributed

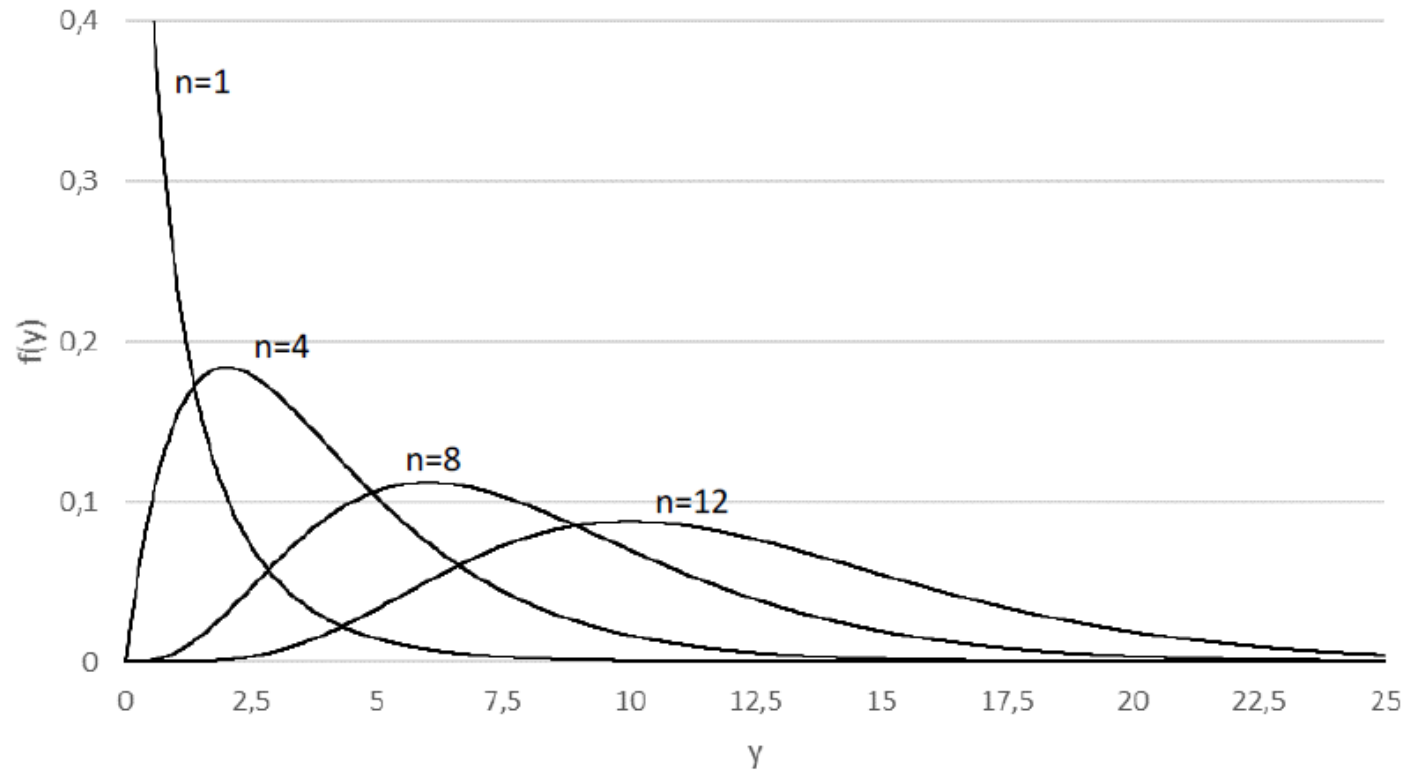
χ^2 -distributed with $n - 1$ degrees of freedom (S^2 random variable corresponding to the variance $\hat{\sigma}^2$ of the random sample). For values $\chi^2_{\frac{\alpha}{2}}$ and $\chi^2_{1-\frac{\alpha}{2}}$ for which the χ^2 -distribution with $n - 1$ degrees of freedom reaches the values $\frac{\alpha}{2}$ and $1 - \frac{\alpha}{2}$, we have:

$$P(\chi^2_{\frac{\alpha}{2}; n-1} \leq Y \leq \chi^2_{1-\frac{\alpha}{2}; n-1}) = 1 - \alpha$$

\implies

$$P\left(\frac{(n-1)S^2}{\chi^2_{1-\frac{\alpha}{2}; n-1}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{\frac{\alpha}{2}; n-1}}\right) = 1 - \alpha$$

χ^2 -Distribution



In opposite to the normal and t-distribution the χ^2 -distribution is not symmetric. This means both limits of a two-sided confidence interval has to be calculated separately.

χ^2 -distribution

Suppose X_1, \dots, X_n are independently identical standard normally distributed, then the distribution of the random variable $X = \sum_{i=1}^n X_i^2$ is called χ^2 -distribution

Density function: $f_n(X = x) = \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{\Gamma(\frac{n}{2}) 2^{\frac{n}{2}}}$, $x > 0, 0$ sonst

Expected value $E(X) = n$

Variance $Var(X) = 2n$

Examples:

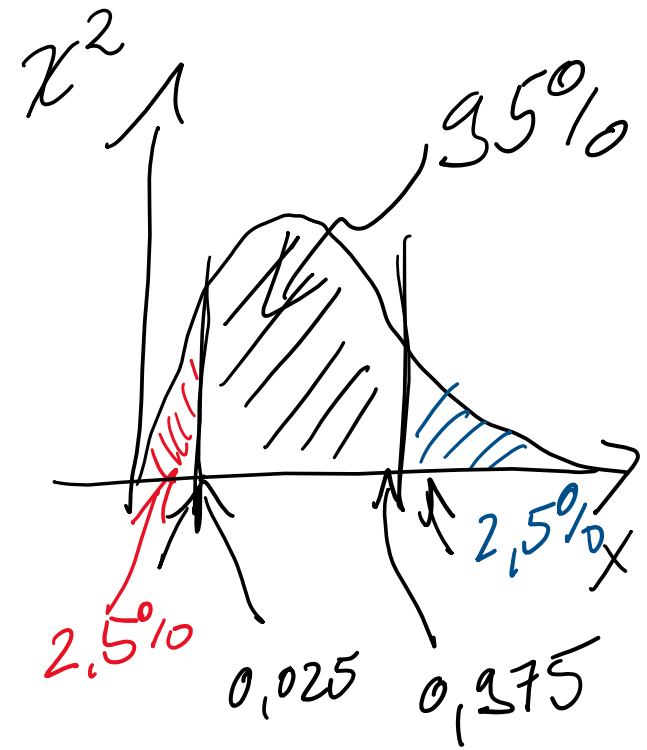
- ▶ Variance of a random sample of size n given the parent distribution is normally distributed
- ▶ χ^2 -goodness of fit test
- ▶ χ^2 -test of independence (contingency table!)

↳ distribution also depends on the sample size

Example

Variance given a normally distributed parent population

A stock of the Oldtech company has within the last 19 years an average profit of $\bar{x} = 15$ US-Dollar. The uncertainty of the profit, estimated with the standard deviation is $\hat{\sigma} = 30$ US-Dollar. We assume that the profits are normally distributed. Given a confidence level of 95% calculate the confidence interval for the unknown variance σ^2 of the parent distribution.



$$\chi_{0,025;18} = 8,231$$

$$\chi_{0,975;18} = 31,53$$

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

$$S^2 \stackrel{!}{=} \hat{\sigma}^2 = 900 \quad n-1 = 18$$

$$I_u = \frac{18 \cdot 900}{8,231}$$

$$\approx 513$$

$$I_d = \frac{18 \cdot 900}{31,53}$$

$$\approx 1968$$

Confidence Interval
for the standard
deviation σ
 $I_d \approx \sqrt{513} \approx 22,7$
 $I_u \approx \sqrt{1968} \approx 44,4$

E2 1) $\bar{x} = 600.000$ $\hat{\sigma} = 90.000$ $n = 225$

(a) Calculate the standard deviation of the average profit?

(b) Calculate a confidence interval of μ for probability of error of $\alpha = 5\%$

(c) Extrapolate for a probability of error $\alpha = 1\%$ a confidence interval for the whole sector.

$$\hat{\sigma}_{\bar{x}} = \sqrt{\text{Var}\left(\frac{1}{n} \sum x_i\right)}$$

$$= \frac{1}{\sqrt{n}} \sqrt{\text{Var}\left(\sum x_i\right)}$$

$$= \frac{1}{\sqrt{n}} \sqrt{\sum \text{Var}(x_i)}$$

$$= \frac{1}{\sqrt{n}} \sqrt{n \cdot \sigma^2}$$

$$= \frac{\sigma}{\sqrt{n}}$$

a) $\hat{\sigma}_{\bar{x}} = \frac{\hat{\sigma}}{\sqrt{n}} = \frac{90.000}{\sqrt{225}} = 6.000$

b) $\bar{x} \pm z \cdot \hat{\sigma}_{\bar{x}} \Rightarrow \bar{x} \pm z \cdot \frac{\hat{\sigma}}{\sqrt{n}}$

$\Rightarrow 600.000 - 1,96 \cdot 6000 = I_d$

$600.000 + 1,96 \cdot 6000 = I_u$

$I_d \approx 588240$

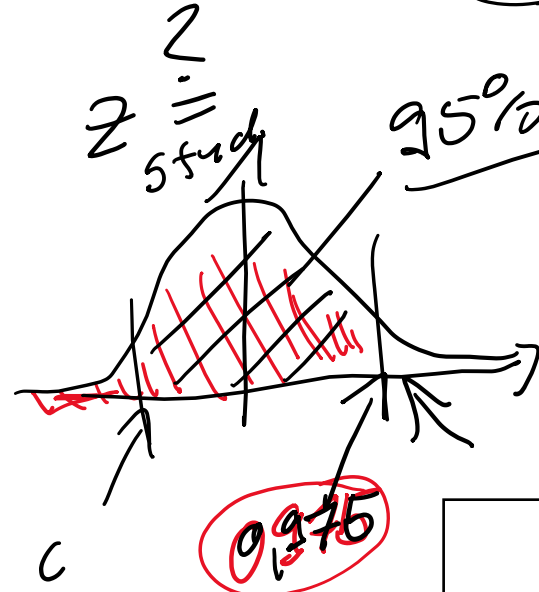
$I_u \approx 611759$

c) \Rightarrow just multiply the limits with 12.100

$I_d = 12.100 \cdot 588240 \approx 7.107 \text{ Mrd. €}$

$I_u = 12.100 \cdot 611759 \approx 7.44 \text{ Mrd. €}$

Bill.



0,975

2. A machine is cutting steel rods of a specific length. Out of the total production of $N = 150$, we take a random sample of $n = 9$. Measuring their length, we obtain:
184,2 mm, 182,6 mm, 185,3 mm, 184,5 mm, 186,2 mm, 183,9 mm, 185,0 mm, 187,1 mm, 184,4 mm.

Due to experience, we know that the parent distribution is normally distributed.

(a) Calculate unbiased estimators for mean and variance of the parent distribution.

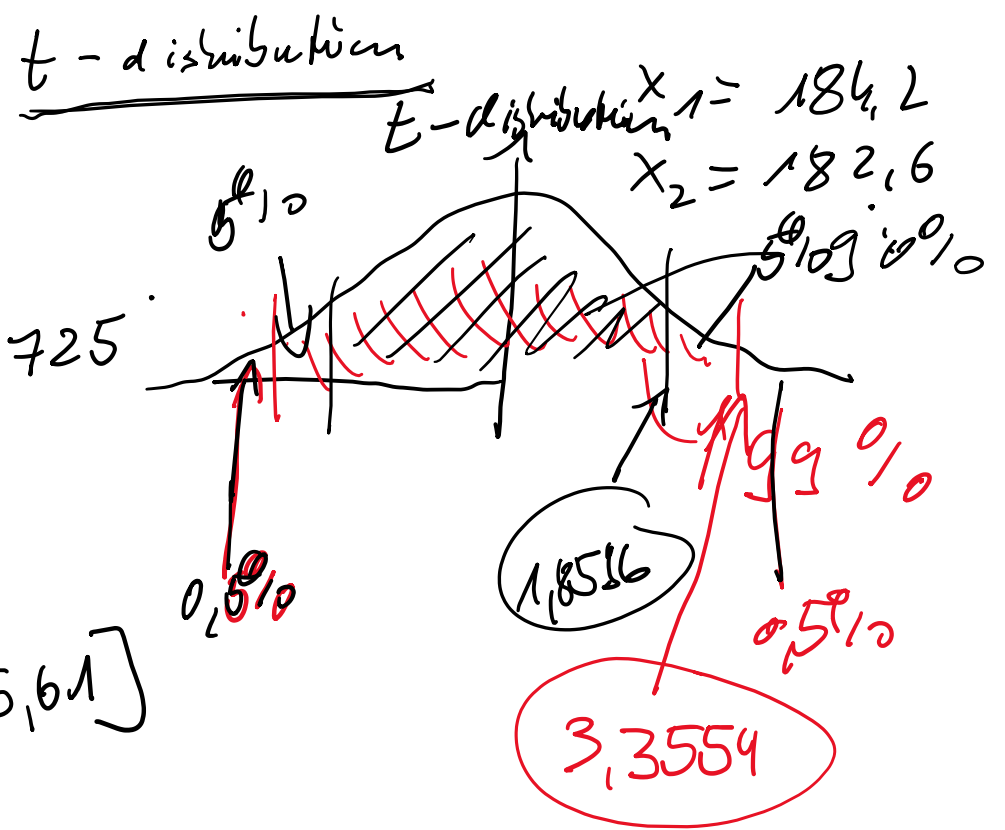
(b) Calculate for μ a confidence interval for the confidence levels 0,9 and 0,99.

Mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 184,8$ $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = 1,725$

b) 90%: $\left(\bar{x}\right) \pm \frac{s}{\sqrt{n}} t_{0,95; n-1}$

$= 184,8 \pm \frac{1,725}{\sqrt{9}} \cdot 1,8536 \Rightarrow [183,99; 185,61]$

99%: $184,8 \pm \frac{1,725}{\sqrt{9}} \cdot 3,3554 \Rightarrow [183,33; 186,27]$



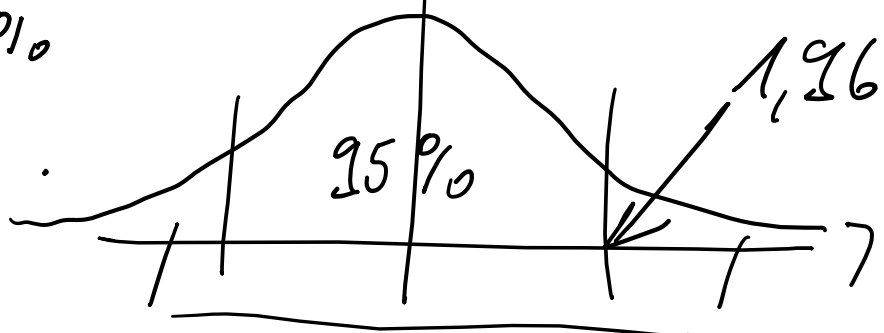
3. Within a random sample of 144 persons in WHV 75% answered, they would like to shop on sundays.

- (a) Calculate a confidence interval for the proportion of inhabitants of WHV, who would prefer a general shop opening on sundays for a probability of error of $\alpha = 0,05$ and $\alpha = 0,1$.
- (b) Calculate the confidence level, if the confidence interval would be represented by $75\% \pm 10\%$.

→ normal distribution sketch

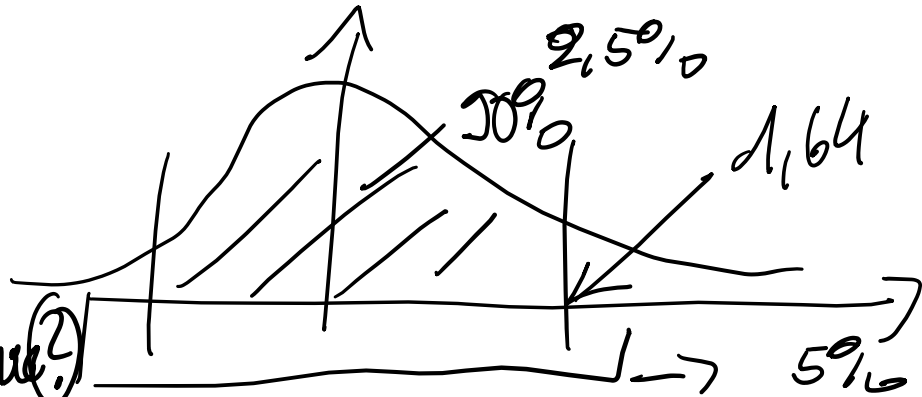
$$\hat{\pi} = 0,75 \quad \hat{\sigma} = \sqrt{\hat{\pi}(1-\hat{\pi})} \quad \text{corresponding to } 97,5\%$$

$$\Rightarrow \left| \frac{\hat{\pi}}{n} \pm \frac{\sqrt{0,75(1-0,75)}}{144} \cdot 1,96 \right| \Rightarrow [0,68; 0,82]$$



$$\alpha = 5\%$$

$$\Rightarrow \left| \frac{\hat{\pi}}{n} \pm \frac{\sqrt{0,75(1-0,75)}}{144} \cdot 1,64 \right| \Rightarrow [0,69; 0,81]$$

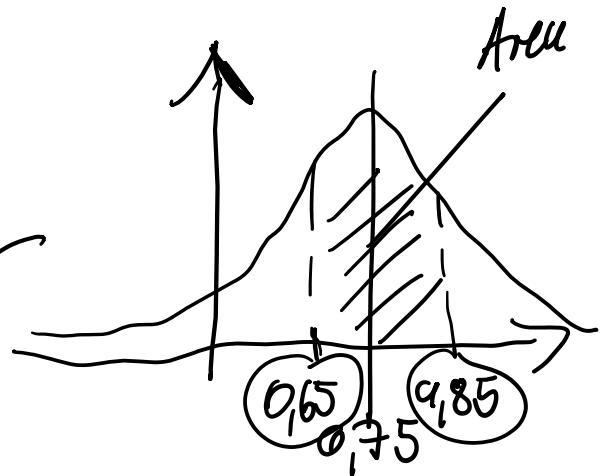
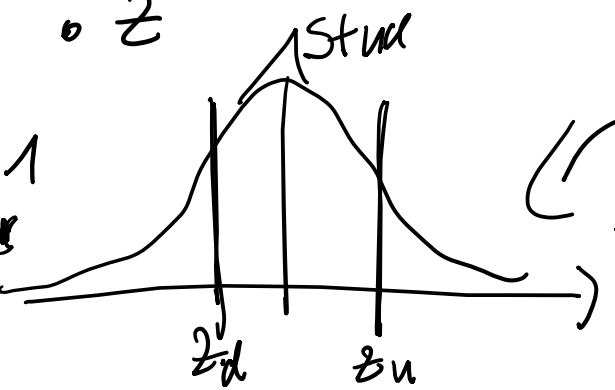


b) $0,75 \pm 0,1$ $[0,65; 0,85] \Rightarrow$ (confidence level?)

$$\hat{\pi} \pm \frac{\sqrt{\hat{\pi}(1-\hat{\pi})}}{n} \cdot z$$

$$\Rightarrow \frac{\sqrt{\hat{\pi}(1-\hat{\pi})}}{n} \cdot z = 0,1$$

$$\frac{\sqrt{0,75(1-0,75)}}{144} \cdot z = 0,1$$

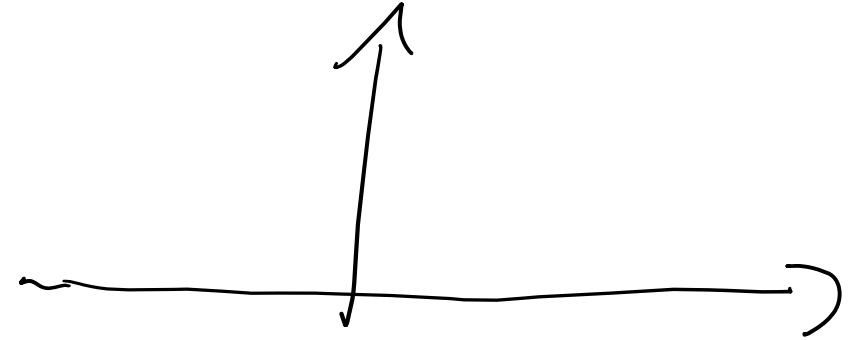


$$* \sqrt{\frac{0,75(1-0,75)}{144}} \cdot z = 0,1 \Rightarrow z = \frac{0,1}{\sqrt{\frac{0,75(1-0,75)}{144}}} = 2,77$$

→ corresponds to 99,72%

$1 - 99,72\% = 0,28\%$
 this we have to multiply by 2 in order
 to get the α -error because we
 a two-sided confidence interval

$$\Rightarrow \boxed{\alpha = 0,56\%}$$



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Hypothesis tests

Assume, you have a random sample and you want to test some hypothesis about for example some parameter of the distribution

Given a specific probability of error, we consider this probability as the probability of making a wrong decision. Thus, given a Hypothesis, the probability of error is the probability that the Hypothesis is not confirmed.

Example:

In your company, you want to know if your marketing activity has increased the awareness of a product. You know, that before the activity, the awareness was roughly 40% (i.e. 40% of the population have known the product). Now you make a survey with $n = 1000$ people asking for the knowledge after the marketing activity. Within this sample 420 People answered, that they know the product.

$$\Rightarrow \frac{420}{1000} = 42\%$$
$$\Rightarrow \frac{420}{1000} \rightarrow 40\%$$

9% more people
are knowing
my product

Question:

Has the awareness of the product increased within the total population??

⇒ Because of random fluctuations an answer via a point estimation is not possible!

Statistical hypothesis testing is asking:

Is it possible to reject the null hypothesis H_0 that the degree of awareness has not changed, because of the result of the survey with the random sample of $n = 1000$? Similar to the confidence intervals one has to specify a degree error probability α , for example 5%.

Hypothesis testing

1. Step:

Create a questioning with two hypotheses contradicting each other:

H_0 : Null hypothesis (incorporates always the equal sign) H_1 : Alternative hypothesis (no equal sign)

Possibilities:

$$\underline{H_0 : \theta \geq c}$$

$$H_1 : \theta < c \quad \underline{\text{left-sided}}$$

$$\underline{H_0 : \theta \leq c}$$

$$H_1 : \underline{\theta > c} \quad \text{right-sided}$$

$$H_0 : \boxed{\theta = c}$$

$$H_1 : \underline{\theta \neq c} \quad \text{two-sided}$$

Error of 1. and 2. kind

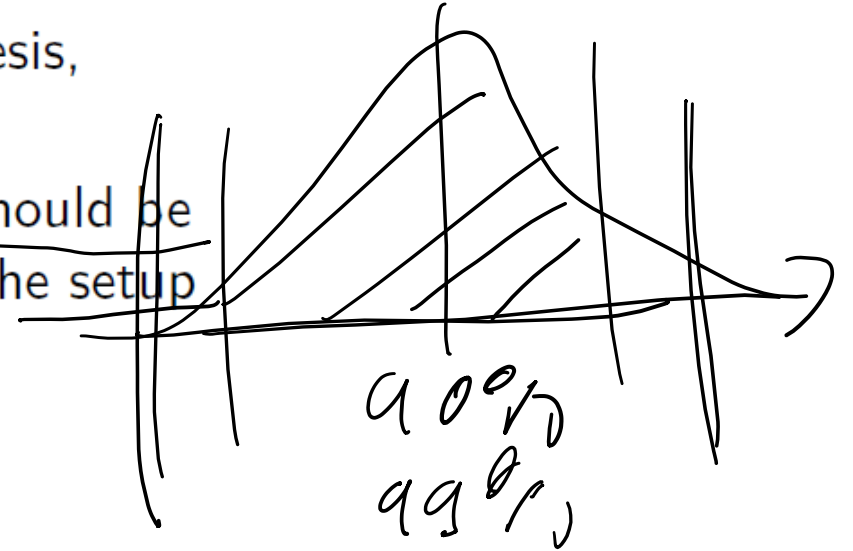
- ▶ Error of 1. kind (or α -error): rejecting the null hypothesis, although it is true.
- ▶ Error of 2. kind (oder β -error): not rejecting the null hypothesis, although it is not true.

Mostly the error of 1. kind is used. Thus the error probability α should be "small" (but still the "right" α) has always to be chosen from the setup of your experiment!):

$$P(\text{H}_0 \text{ rejecting} | \text{H}_0 \text{ true}) = \alpha$$

Problem:

If one choses a small α the probability not rejecting the null hypothesis is increasing, although it is not true. Generally, we choose an α -error of 1%, 5% or 10%.



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

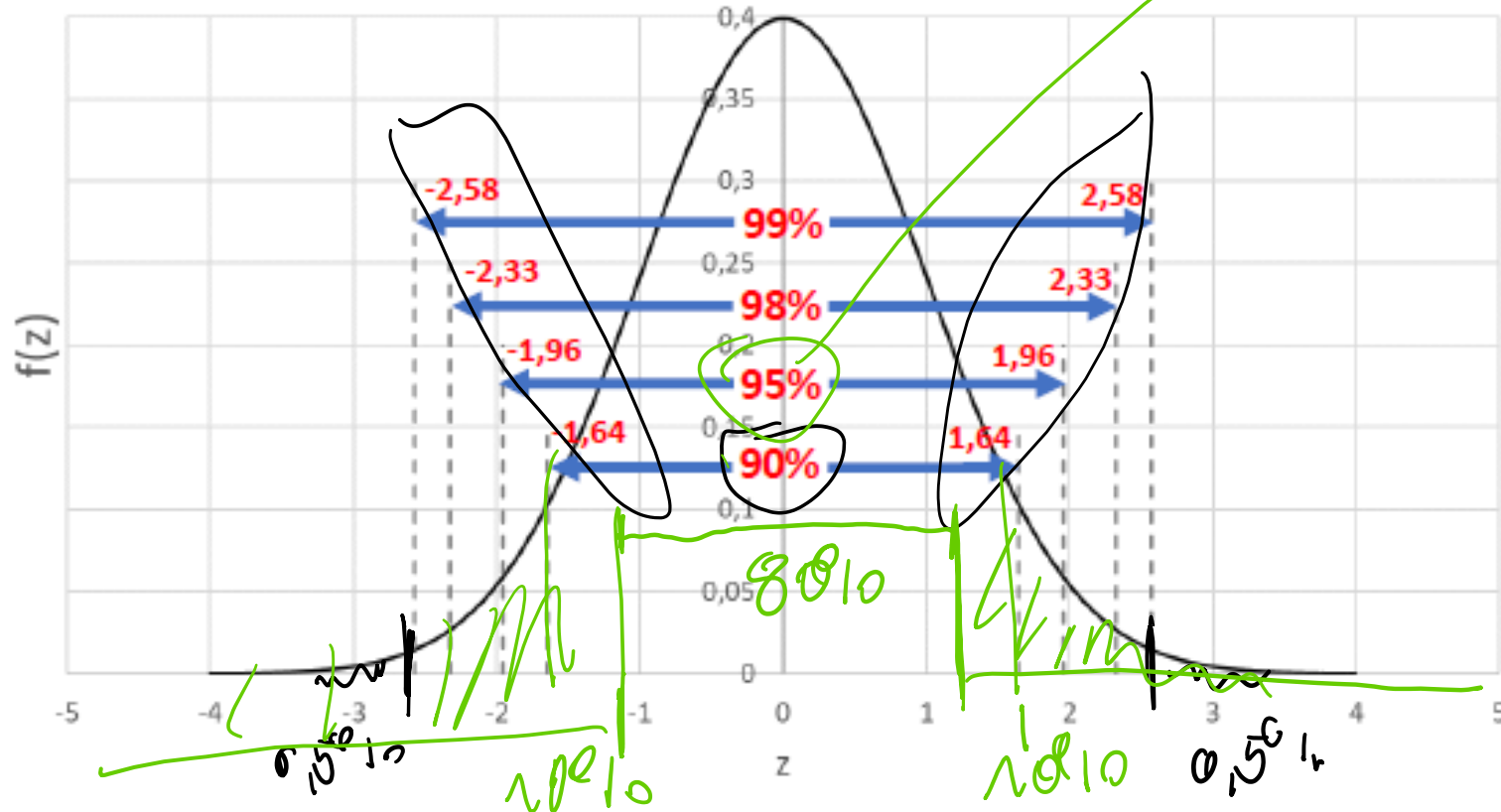
Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Relevant Values of the standard normal distribution

z	F(z)
-2.58	0.5 %
-2.33	1 %
-1.96	2.5 %
-1.64	5 %
-1.23	10 %

$\alpha = 10\%$ but
one-sided



Test-statistic

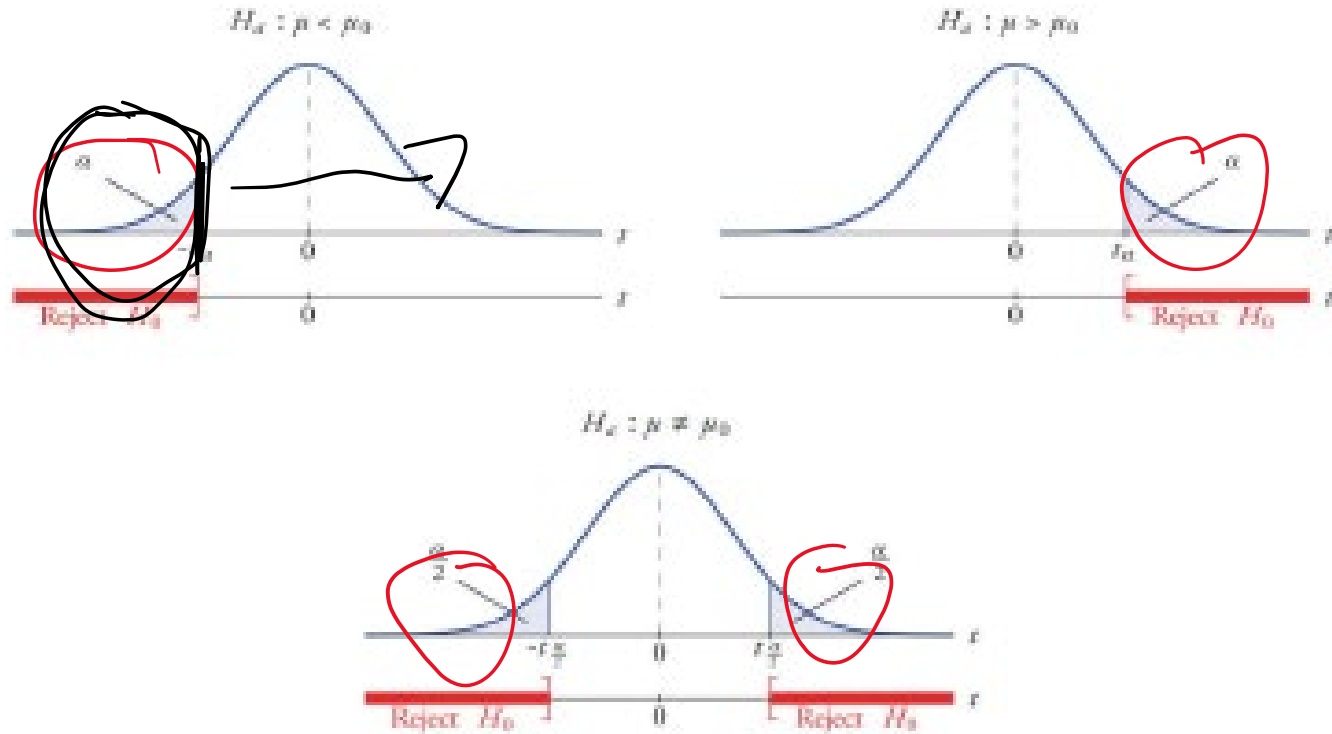
In order to calculate the error probability, we take a random sample out of a probability distribution

⇒ For every test, we need a test-statistic, where we know the corresponding distribution or we make an assumption about the distribution.

Calculating the critical value given the niveau α :

- ▶ **left sided test:** reject H_0 if the calculated test statistic is lower than the critical value CV . *CV is obtained from the assumption α*
- ▶ **right sided test:** reject H_0 if the calculated test statistic is larger than the critical value CV .
- ▶ **two sided test:** reject H_0 if the calculated test statistic is larger than the critical value $+CV$ or lower than the critical value $-CV$.

Example: Test for the mean



Rejecting the null hypothesis
for one- and two sided tests

Hypothesis testing for the mean

Test for μ :

$$H_0: \mu \geq, \leq, = \mu_0$$

$$H_1: \mu <, >, \neq \mu_0$$

(a) If $n \geq 30$ and unknown standard deviation

Test-statistic: $Z = \frac{\bar{X} - \mu}{\hat{\sigma}} \sqrt{n}$ and Z is standard normally distributed

(b) if $n < 30$ and unknown standard deviation

Test-statistic: $T = \frac{\bar{X} - \mu}{\hat{\sigma}} \sqrt{n}$ and T is t-distributed with $n - 1$ degrees of freedom

Given a random sample we want to check some hypothesis.

Given a specific error probability the hypothesis will be rejected or not.

Example 1

Some machine produces boards with a minimum width of 10mm. Boards with less than the minimum width cannot be used. The machine has an accuracy of 0,75mm (= standard deviation). Within a random sample of $n = 100$ we have an average width of 9,85mm. Is this shortfall significant to a level of $\alpha = 10\%$ or only a random fluctuation?

has to be compared
with the critical
value CV

1) $H_0 : \mu \geq \mu_0 \quad \mu \geq 10 \text{ mm}$
 $H_1 : \mu < 10 \text{ mm}$

2) Test-statistic $\Rightarrow Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{9,85 - 10}{0,75} \sqrt{100} \approx -2$

3) CV is obtained from α and the assumption on the distribution
 $\alpha = 10\%$ Normal distribution
 \Rightarrow from a table or statistics software we have
 $CV \approx -1,23$

4) Since $-2 < -1,23$
 $\Rightarrow H_0$ is rejected

Example 2

Within the distribution of some product, the sellers obtain a fixed income and a provision depending on their revenue. On average a provision of 600 Euro is payed. In the current month the framework changed due to higher selling-prices (may be a reason for higher revenue) and the holiday season (may be a reason for lower revenue). Within a random sample of $n = 100$ we obtained an average provision of 585 Euro. Assuming the standard deviation of the provision is 80 and choosing $\alpha = 10\%$, do we have a significant lowering of the provision?

$$\bar{x} = 585$$

$$\hat{\sigma} = 80$$

$$n = 100 \Leftarrow \text{STN-D}$$

$$1) H_0: \mu \geq \mu_0$$

$$H_1: \mu < \mu_0$$

$$2) Z = \frac{585 - 600}{80} \sqrt{100} \approx -1,875$$

$$3) -1,23 < V \Rightarrow 4) -1,875 < -1,23$$

Example 3

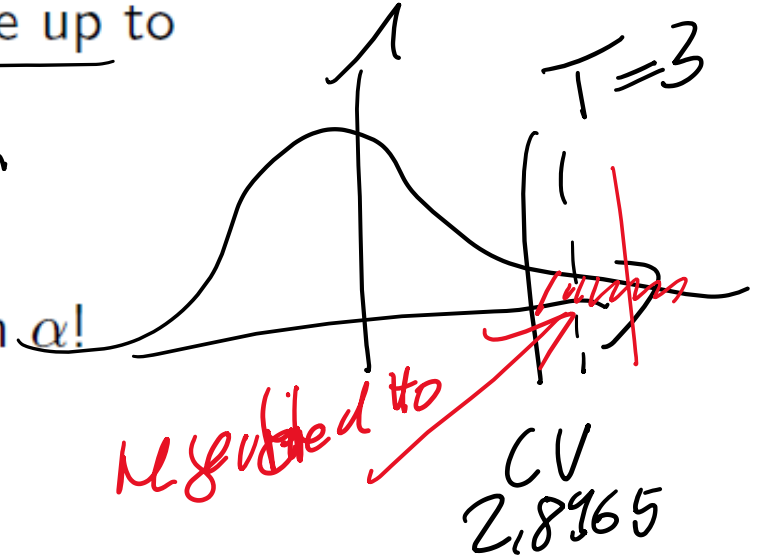
Some new pharmaceutical product should lower the blood pressure up to 20 points. This we know from the advertisement is this true?

small sample size \Rightarrow T-distribution

Assume the parent distribution is normally distributed. Choose

$(\alpha = 0,01; n = 9; \bar{x} = 21,5; \hat{\sigma} = 1,5)$ Discuss the dependance on α !

t-distribution



1) $H_0: \mu \leq 20$

$H_1: \mu > 20$

2) $T = \frac{\bar{x} - \mu}{\hat{\sigma}} \sqrt{n} = \frac{21,5 - 20}{1,5} \sqrt{9} = 3$

3) CV \approx from t-distribution $CV = 2,8965$

4) $3 > 2,8965$

Test comparing two means

Assume we have two random samples with μ_1, μ_2 und $n_1, n_2 \geq 30$ and unknown variances:

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

Test-statistic:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\hat{\sigma}_{\bar{x}_1 - \bar{x}_2}}$$

mit

$$\hat{\sigma}_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}$$

is approximately normally distributed.

Example

Some insurance company has a telephone hotline. The phone calls are evaluated via some new software. Without having the software the phone calls during one month had a connecting on average of 22 minutes within 1500 calls. With the software the connecting time was 18 minutes out of 1200 calls. The unbiased estimation of the standard deviations are given by:

$$\hat{\sigma}_1 = 15 \quad \hat{\sigma}_2 = 12$$

Has the connecting time significantly changed for $\alpha = 5\%$?

4. A provider of mobile games collected the playing time of $n = 120$ randomly drawn users. We suppose, that the playing times x_i are realizations of a normally distributed variable and obtain the following data:

$$\sum_{i=1}^{120} x_i = 21840 \quad \sum_{i=1}^{120} x_i^2 = 4868856$$

(a) Calculate a confidence interval for the variance of the playing time given a confidence level of 99 %

(b) Calculate the confidence interval of the standard deviation for a probability of error of $\alpha = 1\%$

$$\Rightarrow \sigma^2 \approx 7512 \quad \sigma \approx 86,67$$



$$x_1, x_2, \dots, x_{120}$$

$$\sum x_i = 21840 \quad \sum x_i^2 = 4868856$$

$$\sigma^2 = \frac{1}{n-1} \left(\sum_{i=1}^{120} (x_i - \bar{x})^2 \right)$$

$$= \frac{1}{n-1} \left(\sum_{i=1}^{120} (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \right)$$

$$\Rightarrow \frac{1}{n-1} \left(\sum_{i=1}^n (x_i^2) - \frac{1}{n} (\sum x_i)^2 \right)$$

$$Y = \frac{(n-1)\sigma^2}{\sigma^2} \sim \chi^2 - \text{distributed}$$

$$I_d = \frac{(n-1)\hat{\sigma}^2}{\chi^2_{0,99,5}, 119} \approx 5502$$

$$I_u = \frac{(n-1)\hat{\sigma}^2}{\chi^2_{0,5}, 119} \approx 10768$$

$$\chi^2_{0,99,5}, 119 \approx 83$$

$$\chi^2_{0,5}, 119 \approx 162,5$$

just take the rest of the limits
 $\Rightarrow I_d \approx 74,2$
 $I_u \approx 103,8$

from table or statistic software

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Test comparing two means

Assume we have two random samples with μ_1, μ_2 und $n_1, n_2 \geq 30$ and unknown variances:

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

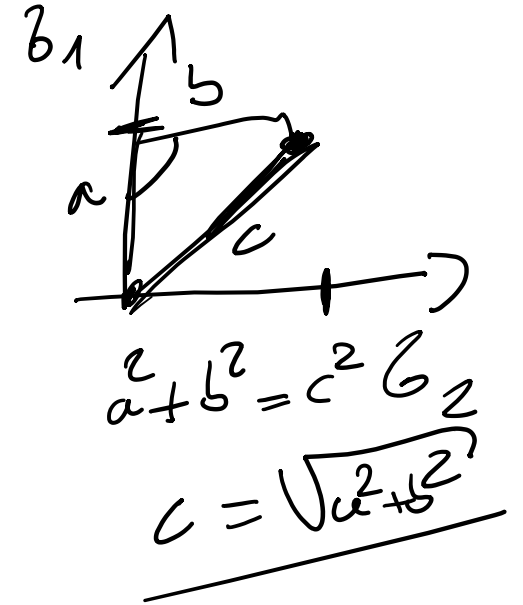
Test-statistic:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\hat{\sigma}_{\bar{X}_1 - \bar{X}_2}}$$

mit

$$\hat{\sigma}_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}$$

is approximately normally distributed.



Example

Some insurance company has a telephone hotline. The phone calls are evaluated via some new software. Without having the software the phone calls during one month had a connecting on average of 22 minutes within 1500 calls. With the software the connecting time was 18 minutes out of 1200 calls. The unbiased estimation of the standard deviations are given by:

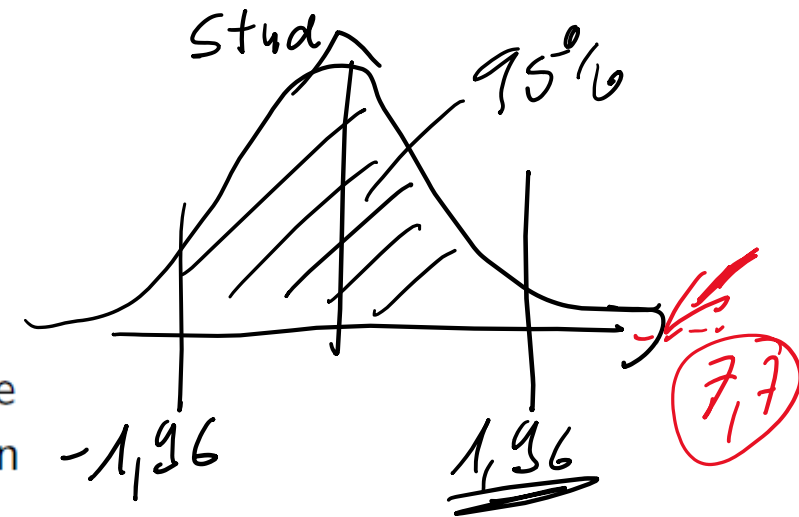
$$\begin{aligned}\bar{x}_1 &= 22 & \bar{x}_2 &= 18 \\ n_1 &= 1500 & n_2 &= 1200 \\ \hat{\sigma}_1 &= 15 & \hat{\sigma}_2 &= 12\end{aligned}$$

Has the connecting time significantly changed for $\alpha = 5\%$?

two sided

$$\begin{aligned}1) & H_0: \mu_1 = \mu_2 \\ & H_1: \mu_1 \neq \mu_2 \\ 2) & z = \frac{\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)}{\hat{\sigma}_{\bar{x}_1 - \bar{x}_2}} \approx \frac{22 - 18 - 0}{0,52} \approx \underline{7,7}\end{aligned}$$

3) $CV \Rightarrow$ we reject the null hypothesis \Downarrow



$$\begin{aligned}\hat{\sigma}_{\bar{x}_1 - \bar{x}_2} &= \sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}} \\ &= \sqrt{\frac{15^2}{1500} + \frac{12^2}{1200}} \approx 0,52\end{aligned}$$

Test for proportion \leftarrow *numbers* this is any or special case for testing for the mean

Given a random sample of X_1, \dots, X_n with X_i are Bernoulli distributed

$$H_0: \pi \geq, \leq, = \pi_0$$

$$H_1: \pi <, >, \neq \pi_0$$

Teststatistic:

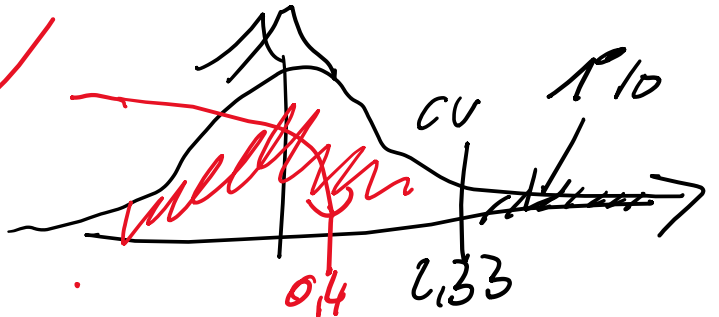
$$Z = \frac{\hat{\pi} - \pi_0}{\hat{\sigma}_\pi} \text{ mit } \hat{\sigma}_\pi = \sqrt{\frac{\pi_0(1 - \pi_0)}{n}}$$

is approximately normally distributed if $n\pi_0(1 - \pi_0) \geq 9$

Example

one-sided

① $400 \cdot 0.5 \cdot 0.5 = 100 \geq 9$ ✓
⇒ we can use the normal distribution



We will test the hypothesis, that a party A has obtained more than 50% of the votes for $\alpha = 0,01$ Within a random sample of $n = 400$ 204 persons, claimed that they have voted for party A.

① $H_0: \pi \leq 0,5 = \pi_0$
② $H_1: \pi > 0,5$

Test statistic
$$z = \frac{\hat{\pi} - \pi_0}{\hat{\sigma}_{\hat{\pi}}} \quad \hat{\sigma}_{\hat{\pi}} = \sqrt{\frac{\pi_0(1-\pi_0)}{n}} = \sqrt{\frac{0,5^2}{400}} = \frac{0,5}{20} = 0,025$$

$$\hat{\pi} = \frac{204}{400} = 0,51 \Rightarrow z = \frac{0,51 - 0,5}{0,025} = 0,4$$

③ CV = 2,33 since $z < CV \Rightarrow H_0$ is not rejected

χ^2 -test of independence

Remember the test for variance also followed the χ^2 -distribution

The χ^2 -test of independence checks whether two nominal scaled variables are likely to be related or not.

- ▶ H_0 : Two variables are independent.
- ▶ H_1 : Two variables are not independent.

The χ^2 -test of independence is based on the contingency table. With the following test statistic, we decide if we reject the Null hypothesis or not

for a given level of α

$$\chi_{emp}^2 = n \sum_{i=1}^k \sum_{j=1}^l \frac{(p_{ij} - p_{i \cdot} p_{\cdot j})^2}{p_{i \cdot} p_{\cdot j}}$$

approximately χ^2 -distributed with $f = (k - 1)(l - 1)$ degrees of freedom:

empirical probabilities are
 theoretical probabilities
 assuming both
 variables are
 independent!

Example

An insurance company wants to increase the number of sold household insurances and is planning a marketing activity. In order to address the right people, it is assumed, that customers who have already a life insurance policy are risk averse and therefore are interested to buy also a household insurance. This is evaluated via the own data of the insurance company. Test the two variables (X: having a household insurance or not/Y: having a life insurance or not) for independence

two-sided \rightarrow
 $H_0: X, Y$ are independent
 $H_1: X, Y$ are not

(2) Test-statistic
 $\chi^2_{emp} \approx 107$

(3) $CV = 7,88$
 H_0 is rejected

$n \rightarrow$

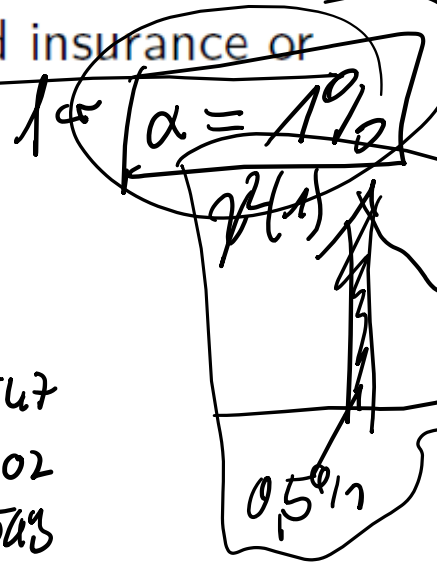
	no household	household	
no Life	436	511	547
Life	121	481	602
	557	992	1549

p-emp	no h	h	
no	28%	33%	61%
l	8%	31%	39%
	36%	64%	100%

p-theo	no h	h	
no	22%	39%	61%
l	14%	25%	39%
	36%	64%	100%

$\Rightarrow \chi^2_{emp} \approx 107$

$f = (k-1)(l-1) = (2-1)(2-1) = 1$
 $\chi(f=1)$



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

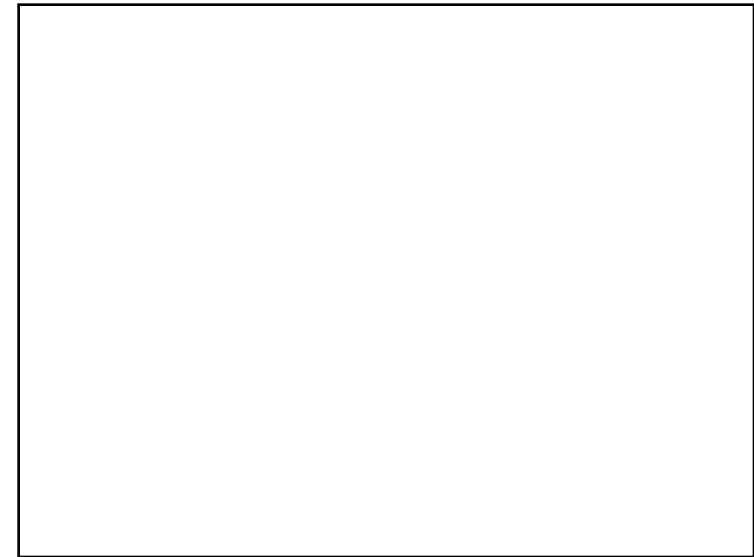
Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Scatterplot

Graphical representation of the attributes of two variables within a two-dimensional coordinate-system.

Number	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
Order	2	4	3	6	2	1	5	3	4	2	2	7	4	3	2	1	5	4	3	3	1	6
Volume [€]	36	44	30	67	21	8	62	35	33	37	22	72	65	45	23	17	44	56	32	37	22	48

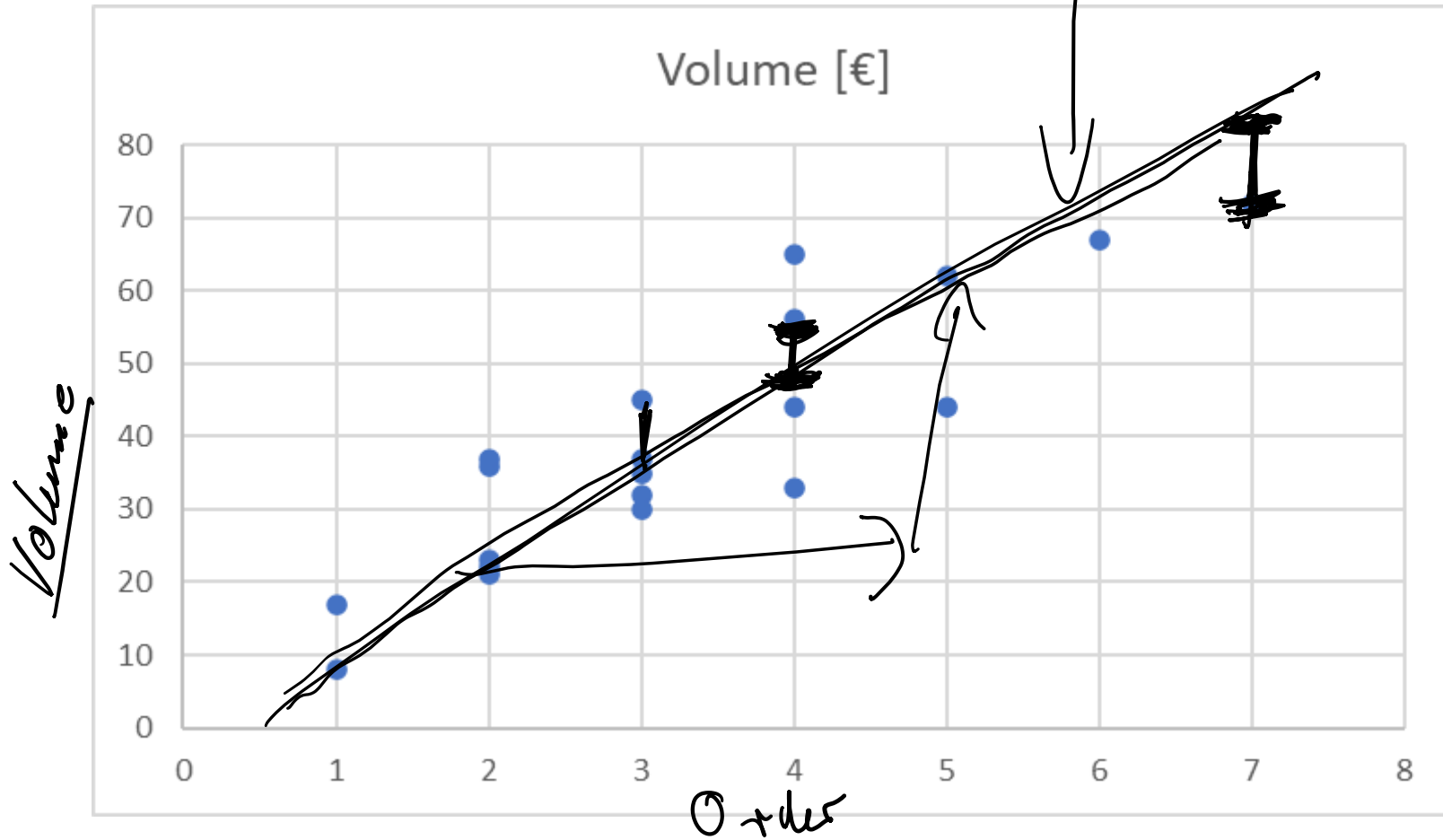


Scatterplot

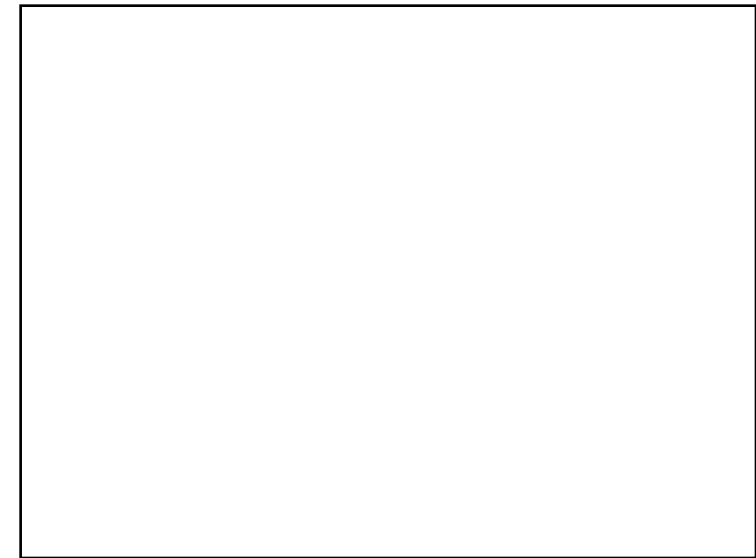
X-coordinate: order

Y-coordinate: Volume [Euro]

How should we fit this line into the data?



Assumption:
The higher the number of orders, the higher is the volume
(4-9)



Linear Regression

Example 2:

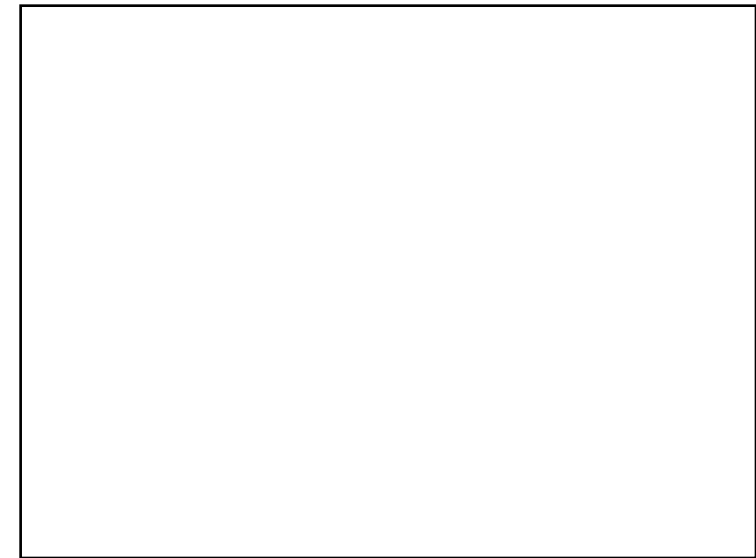
A company wants to know how the sales volume depends on visits of representatives.

The company has collected the following data

	Visits	Sales volume [€]
year	x	y
2008	9	24
2009	11	33
2010	5	10
2011	13	29
2012	20	42
2013	12	24



We also want to answer the question, from a descriptive point of view, which sales volume we can expect, if we increase the visits up to $x=30$



Linear Regression

- Via a linear regression we want to analyse the dependence between two variables

- We distinguish between a dependent variable (y) and an independent variable (x)

- (x) is also called the explanatory variable

- A change of (x) implicates a change (y)

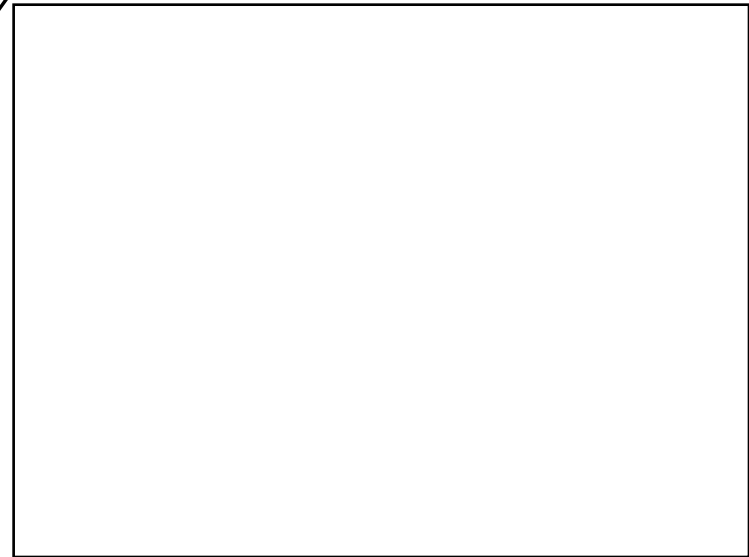
- The linear regression defines an empirical relation via a linear function

i.e. $y = f(x) = 3 + 2x$

or generally

$y = a + bx$

point on the vertical axis
vert
slope of the line



Simple Examples

$19 = a + (\frac{1}{2}) \cdot 2$
 $\Rightarrow a = 20$
 $= a - 1$

$x=2 \quad y=19$

blue data means 3
 $\Delta x = 1 \Rightarrow \Delta y = 3$
 true for all data points

$$y = a + bx$$

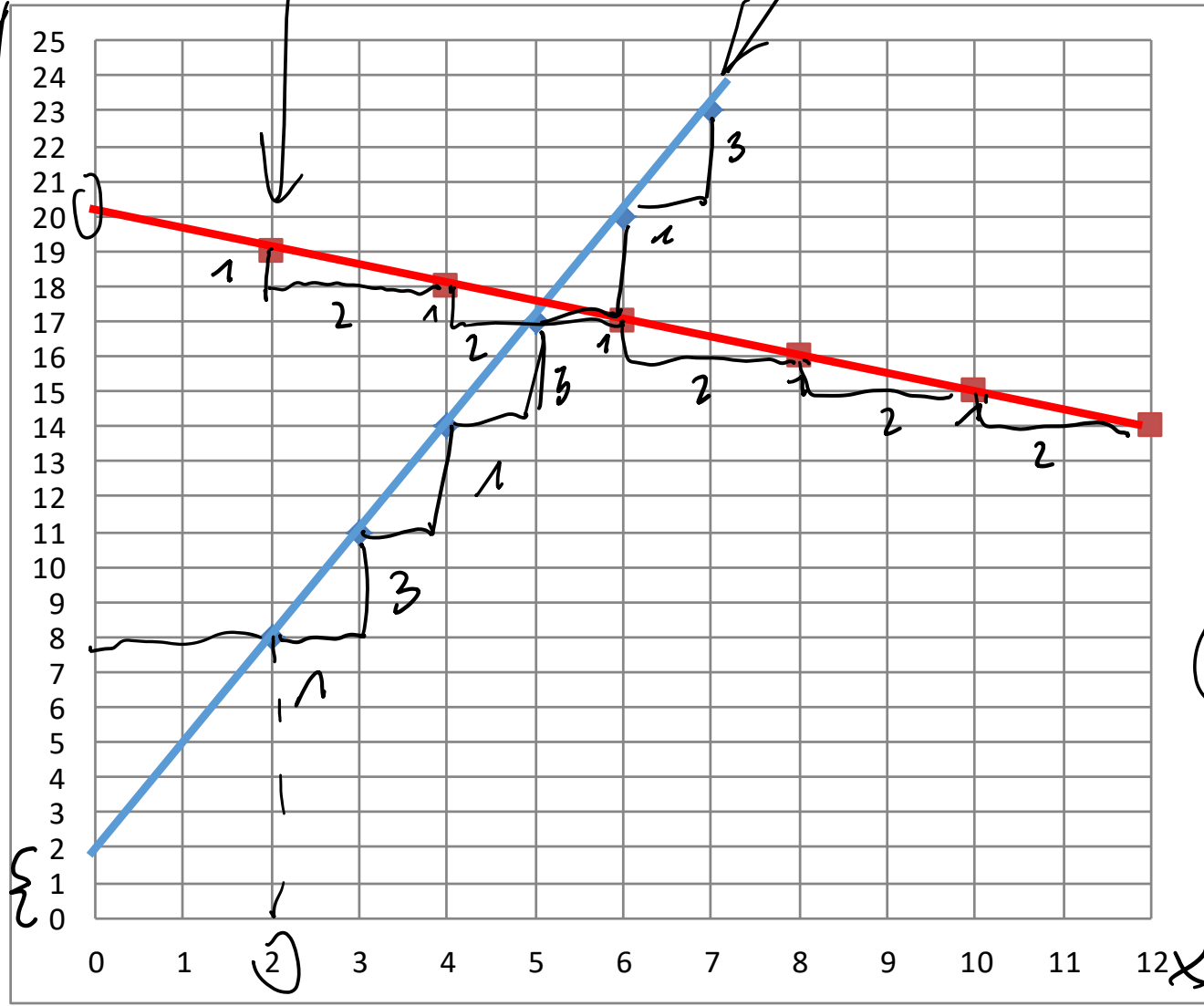
$$\Rightarrow b = \frac{\Delta y}{\Delta x} = 3$$

and $a = 2$
 same data point $(x=2, y=8)$
 $\Rightarrow 8 = a + 3 \cdot 2 \Rightarrow$

$a = 2$

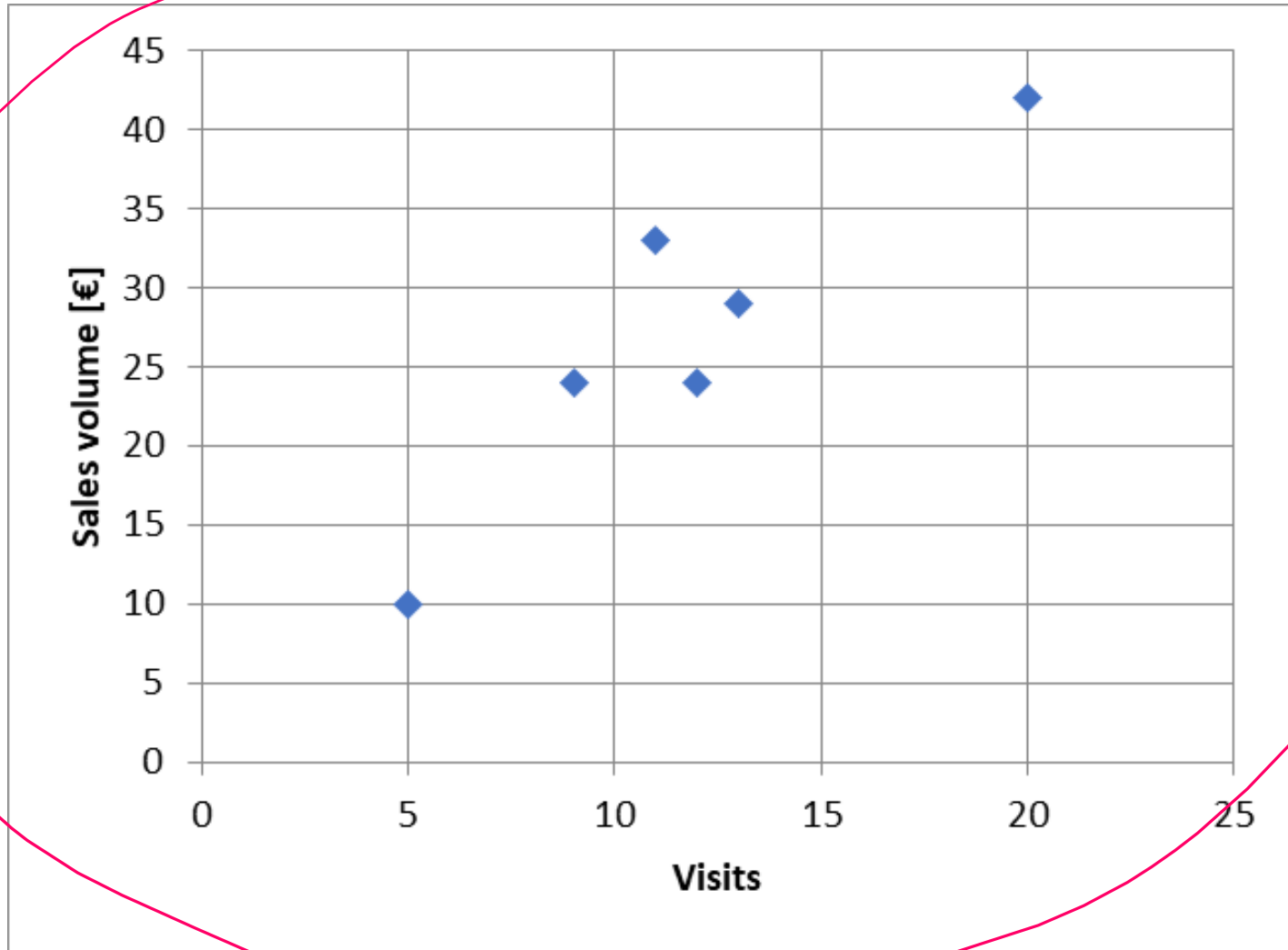
red points
 $a = 20$
 $b = -\frac{1}{2}$

$a = 2$

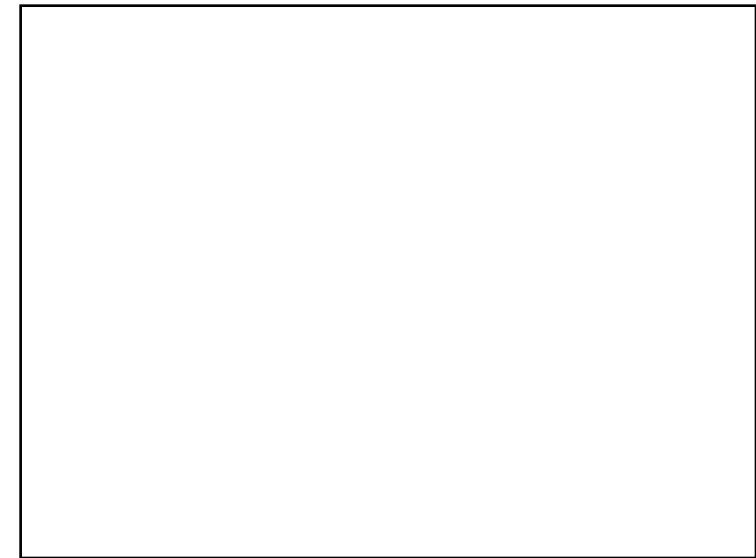


Determine graphically a regression line for the blue and red data points. What kind of dependence do we have?

Example



Try to fit graphically (by eye) a regression line through the data points



Optimizing problem

How to calculate a regression line $\hat{y}_i = a + bx_i$?

The coefficients a and b has be calculated from the collected data.

Value on the regression line
collected data point

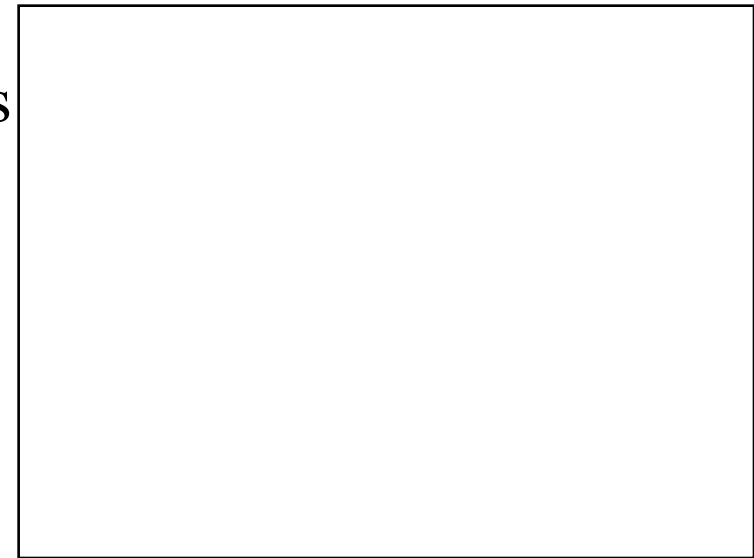
Difference between the data point and the theoretical value on the regression line

$$\hat{y}_i$$
$$y_i$$
$$\varepsilon_i = |y_i - \hat{y}_i|$$

error which we make
by the assumption
 $\hat{y}_i = a + bx_i$

Determine a and b such, that the sum of all quadratic differences ε_i^2 is minimized

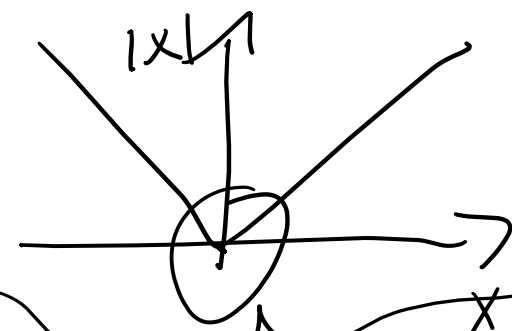
$$\text{min: } \sum_{i=1}^n \varepsilon_i^2 \rightarrow \text{min}$$



Least squared method

The formal optimizing problem

$$F(a, b) = \sum_i^n (y_i - (a + bx_i))^2 = \sum_i^n \epsilon_i^2 \rightarrow \min$$



not differentiable
=> no derivative
=> $f'(x)$

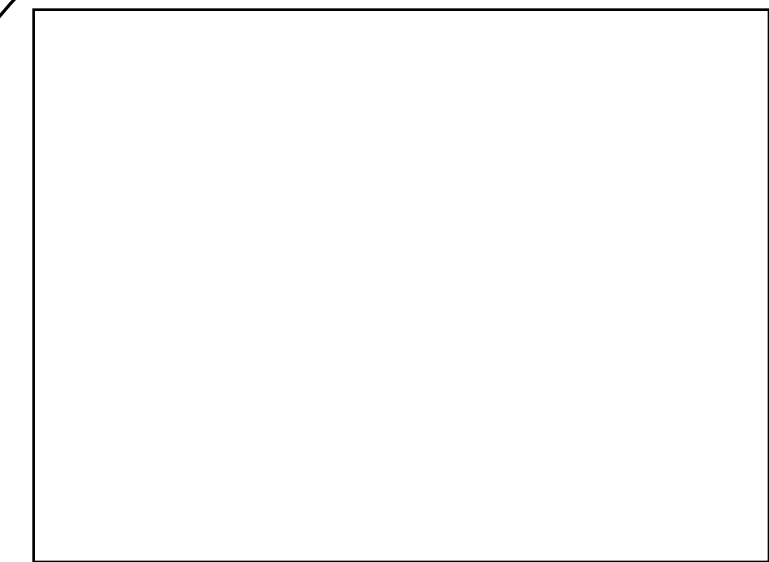
connected
with the
normal
distribution

Why this square

- 1) no problem with $(+)$ is den
- 2) $(4 - 4)^2$ is differentiable
and $|4 - 4|$ is not

$$(4 - 4)$$
$$|4 - 4|$$

3) Normal distribution
 $f(x) = \frac{1}{\sqrt{2\pi} \sigma}$



$f(b) = 5b \Rightarrow f'(b) = 5$ Least squared method

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$f(b) = bx_i \Rightarrow f'(b) = x_i$$

The formal optimizing problem

$$F(a, b) = \sum_i^n (y_i - (a + bx_i))^2 = \sum_i^n \epsilon_i^2 \rightarrow \min$$

$$\frac{\partial F}{\partial a} = \left[\sum_{i=1}^n 2(y_i - (a + bx_i)) \right] (-1) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i - (a + bx_i)) = 0 \quad (1)$$

$$\frac{\partial F}{\partial b} = \left[\sum_{i=1}^n 2(y_i - (a + bx_i)) x_i \right] (-1)$$

$$\Rightarrow \sum_{i=1}^n (y_i - (a + bx_i)) x_i = 0 \quad (2) \quad \text{FOC}$$

$$(1) \Rightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n a + \sum_{i=1}^n bx_i \Rightarrow n\bar{y} = na + nb\bar{x} \Rightarrow \bar{y} = a + b\bar{x}$$

$$\left. \begin{aligned} b \sum_{i=1}^n x_i &= b \cdot n\bar{x} \\ \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i \end{aligned} \right\}$$

\Rightarrow means, the data point (\bar{x}, \bar{y}) is on the optimal regression line



Least squared method

The formal optimizing problem

$$\frac{\partial F}{\partial a} = -2 \left[\sum_i^n (y_i - (a + bx_i)) \right] = 0$$

$$\frac{\partial F}{\partial b} = -2 \left[\sum_i^n (x_i y_i - (a + bx_i)) \right] = 0 \quad (2)$$

$$\Rightarrow \sum_{i=1}^n (x_i y_i - (a x_i + b x_i^2)) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

Since (1) $\bar{y} = a + b\bar{x}$

$$\Rightarrow a = \bar{y} - b\bar{x} \quad (1)$$

Known data

$$\bar{x} \cdot \bar{y} - \bar{y} a \bar{x} - b \bar{x}^2 =$$

$$\Rightarrow \bar{x} \cdot \bar{y} = \bar{x} (\bar{y} - b \bar{x}) - b \bar{x}^2 = 0$$

$$\Rightarrow \bar{x} \cdot \bar{y} - \bar{x} \bar{y} + (b \cdot \bar{x} \cdot \bar{x} - b) \bar{x}^2 = 0$$

$$\bar{x} \cdot \bar{y} - \bar{x} \cdot \bar{y} = b(\bar{x}^2 - \bar{x}^2) \Rightarrow b =$$

$$b = \frac{\bar{x} \cdot \bar{y} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - \bar{x}^2}$$

put this in (1) and we obtain for a: *

$$a = \bar{y} - \frac{\bar{x} \cdot \bar{y} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - \bar{x}^2} \cdot \bar{x}$$

Least squared method

The formal optimizing problem

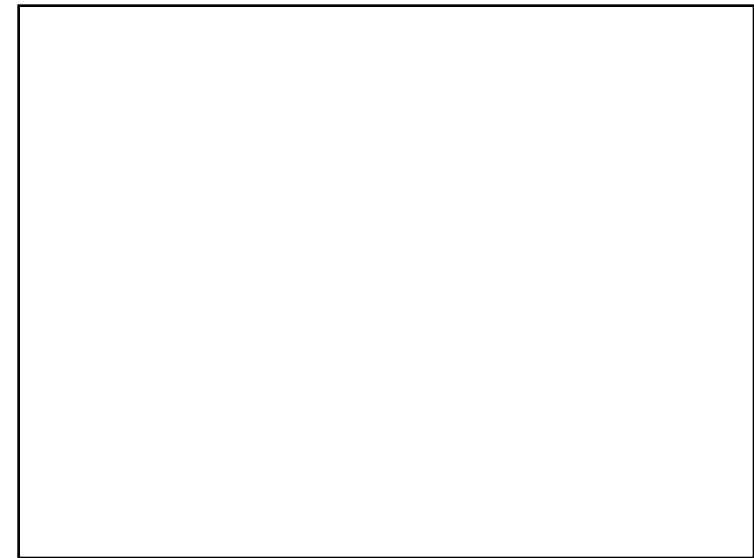
$$F(a, b) = \sum_i^n (y_i - (a + bx_i))^2 = \sum_i^n \epsilon_i^2 \rightarrow \min$$

FOC:

$$\frac{\partial F}{\partial a} = -2 \left[\sum_i^n (y_i - (a + bx_i)) \right] = 0$$

$$\frac{\partial F}{\partial b} = -2 \left[\sum_i^n x_i (y_i - (a + bx_i)) \right] = 0$$

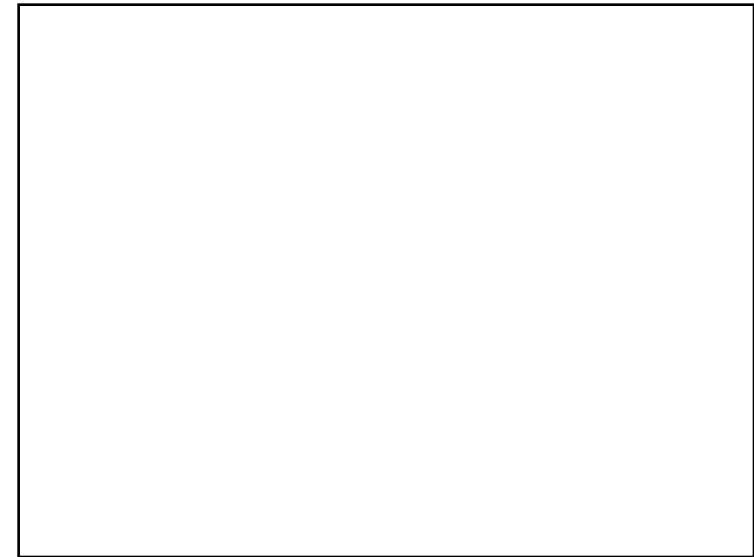
2 dimensional linear system which can directly be solved



$$\frac{\partial F}{\partial a} = -2 \left[\sum_i^n (y_i - (a + bx_i)) \right] = 0$$

Least squared method

$$\frac{\partial F}{\partial b} = -2 \left[\sum_i^n x_i (y_i - (a + bx_i)) \right] = 0$$



Formulas

$$a = \frac{\sum_{i=1}^N x_i^2 \sum_{i=1}^N y_i - \sum_{i=1}^N x_i \sum_{i=1}^N x_i y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i\right)^2} = \bar{y} - b\bar{x}$$

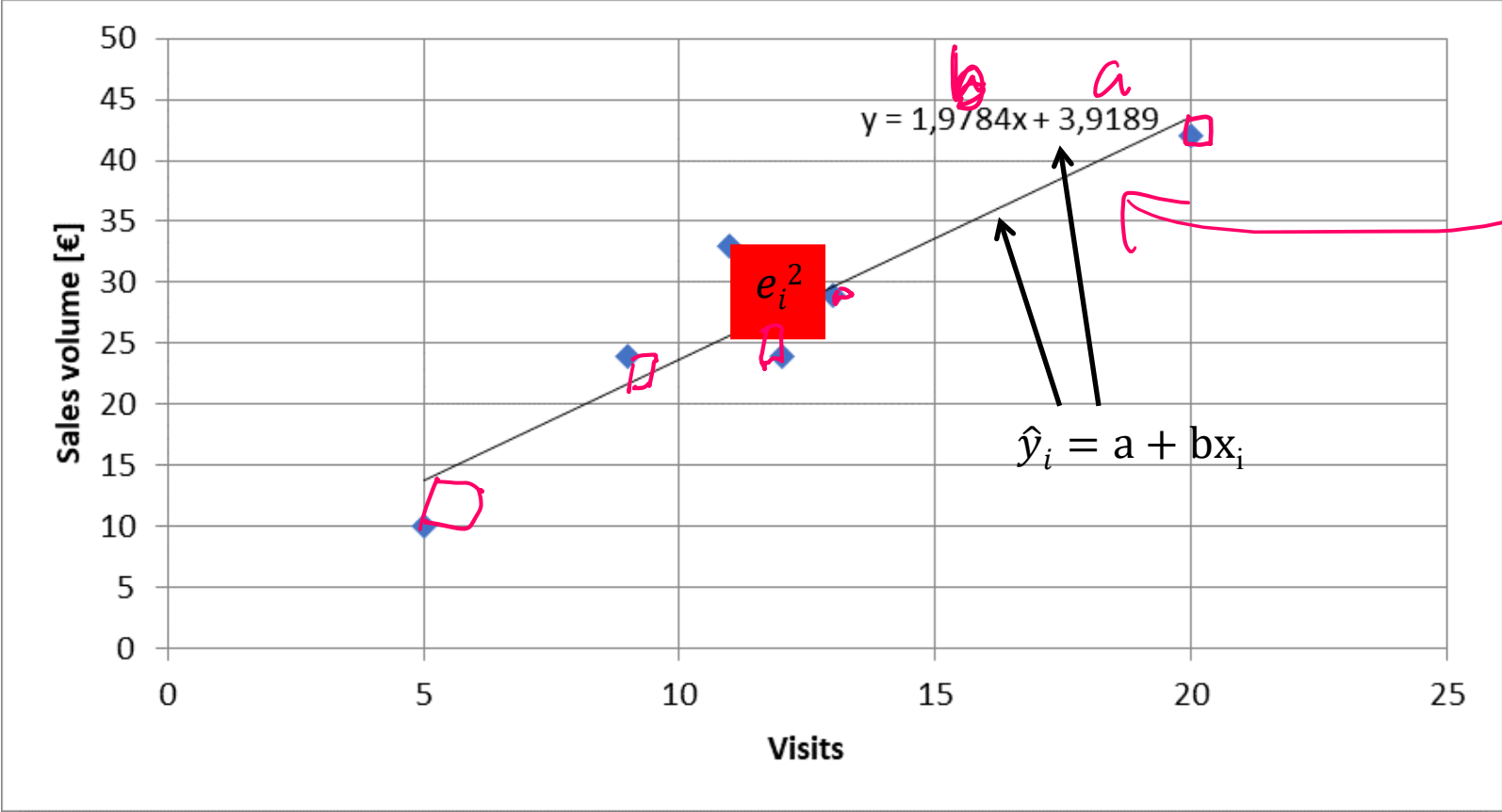
$$b = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i\right)^2} = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

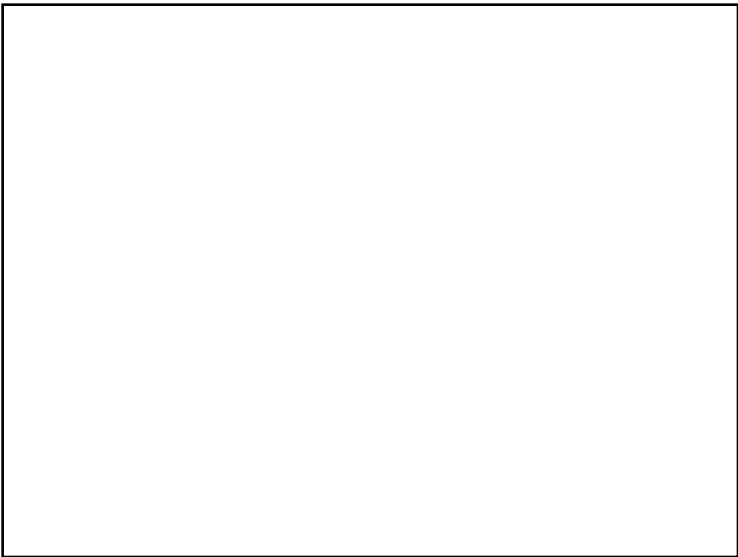
Covariance

$$\begin{array}{l} \bar{x} \quad \bar{y} \\ \overline{x \cdot y} \\ \overline{x^2} \quad \bar{x}^2 \end{array}$$

Calculating the regression line from empirical data



This line
minimizes
 $\sum_{i=1}^n e_i^2$



Lineare Regression

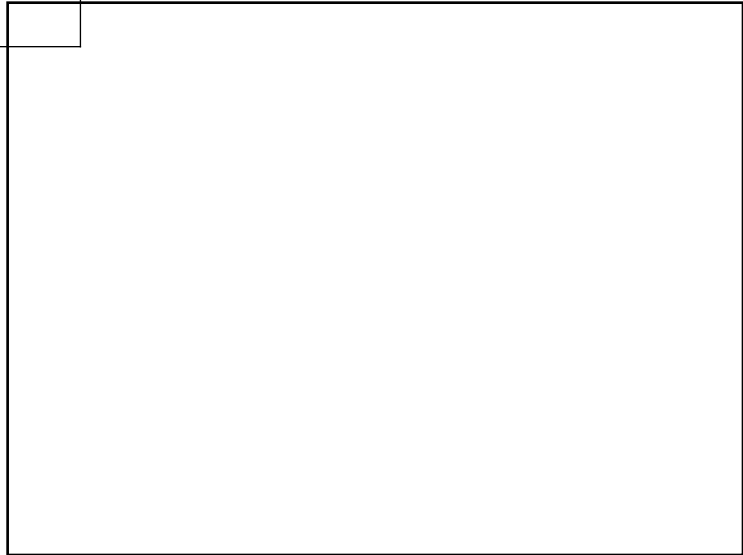
$$\bar{x} \cdot \bar{y} = 315$$

$$\begin{aligned} \bar{x} &= 11,67 \\ \bar{y} &= 27 \\ \bar{x} \cdot \bar{y} &= 315 \\ \sum xy &= 355,67 \\ \bar{x}^2 &= 136,11 \\ \sum x^2 &= 156,67 \end{aligned}$$

	Visits	Sales volume [€]			
year	x	y	xy	x ²	
2008	9	24			
2009	11	33			
2010	5	10			
2011	13	29			
2012	20	42			
2013	12	24			

$$b = \frac{\sum xy - \bar{x} \cdot \bar{y}}{\sum x^2 - \bar{x}^2} = \frac{355,67 - 315}{156,67 - 136,11} \approx 1,97$$

$$a = \bar{y} - b \bar{x} = 27 - 1,97 \cdot 11,67 \approx 3,52$$



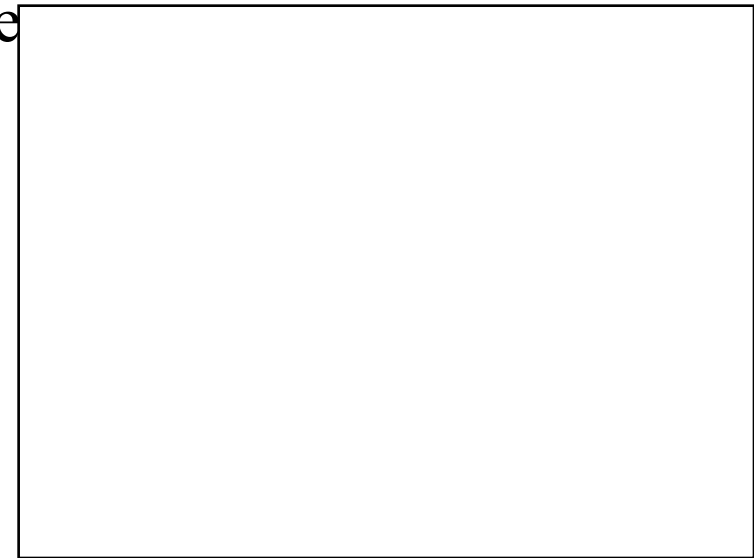
Extrapolation for $x=30$

	Visits	Sales volume [€]			
year	x	y	xy	x^2	yhat
2008	9	24			
2009	11	33			
2010	5	10			
2011	13	29			
2012	20	42			
2013	12	24			

$\hat{y}(30)=$

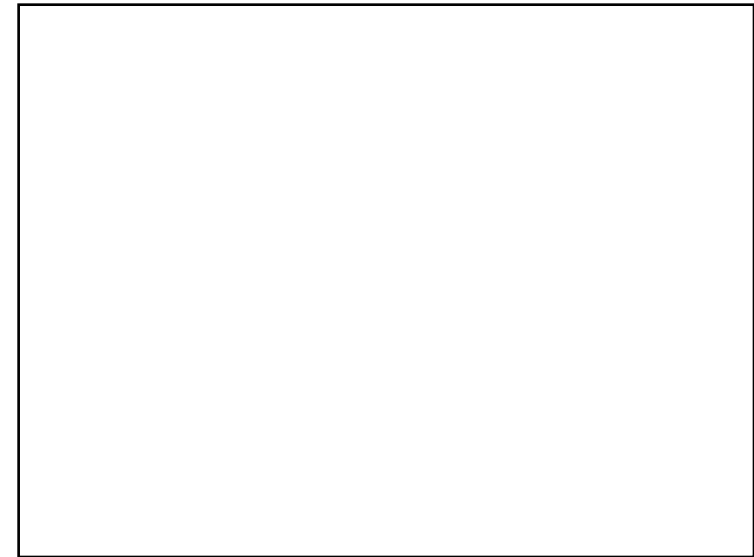
Regressions coefficients

- Regressions coefficient b
 - Slope of the regression line.
 - Determines the marginal effect of a change of one unit of the independent variable x onto the dependent variable y
- Regressions coefficient a
 - The value of the dependent variable if the independent variable $x=0$
 - intercept of the vertical axes



Linear regression

- Measures the linear dependence between two variables
- This dependence can be interpreted, that one variable is quantitatively influencing another variable
- The linear regression is an instrument for forecasting

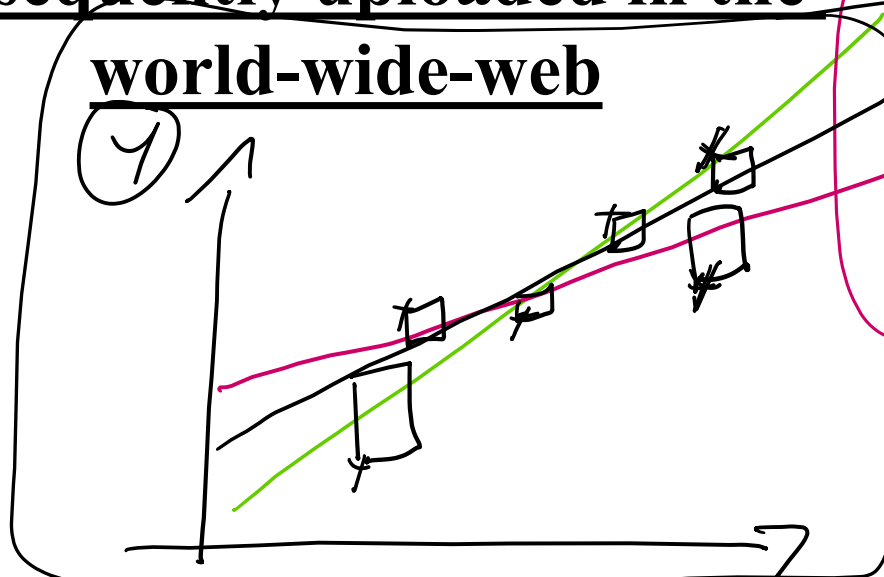


Wilhelmshaven



Least squares Statistics A
was developed by C.F. Gauss, when he
measured the distances in northern Germany

This lecture will be recorded and
Subsequently uploaded in the
world-wide-web



$\hat{y}_1 = a_1 + b_1 x$
 $\hat{y}_0 = a_0 + b_0 x$
 $\hat{y}_2 = a_2 + b_2 x$

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

=> Choose this linear line \hat{y}
for which we have

$\sum \square \rightarrow \min$

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

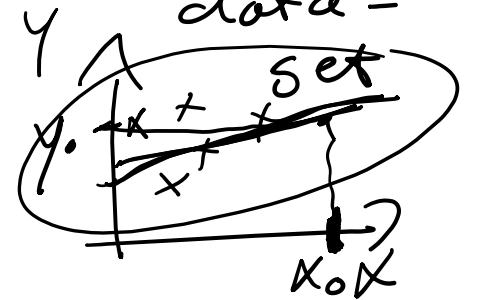
<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Correlation analysis

- In principle for all data sets, we can calculate a regression line.
- But we are also interested in the question meaningful is this calculated dependence?
- For this, we use the correlation analysis, which gives the possibility to measure the strength of the dependence
- For this, we use the correlation coefficient of Bravais-Pearson which is lying between -1 and +1

$$\hat{y}_i = a + b x_i$$

↓
can be
fitted
to every
data -



Covariance

- The covariance is a measure of the joint variability of two random variables

(X, Y) , defined by:

$$\text{cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

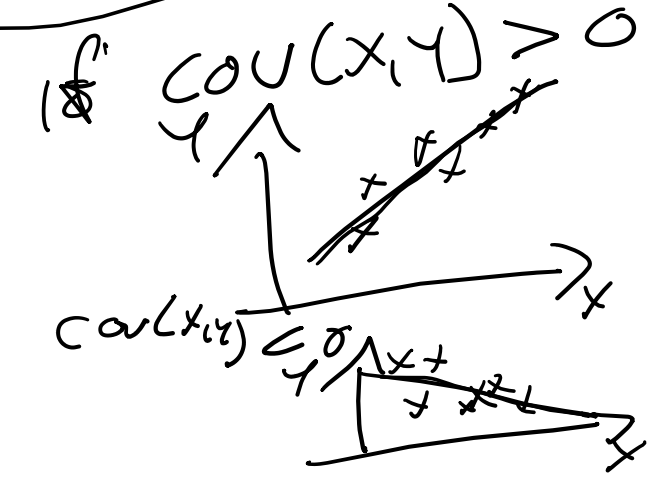
with the unbiased estimator

$$\hat{\sigma}_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$\hat{\sigma}_{xx}$

variance = $E[(x - E(x))^2]$

$\text{cov}(x, x) = \text{var}(x, x)$



- If the greater values of X correspond with the greater values of Y the covariance is positive. In the opposite case, when the greater values of X correspond to the lesser values Y, the covariance is negative. The sign of the covariance shows the tendency in the linear relationship between the variables.

- The magnitude of the covariance can hardly be interpreted, since it is not normalized.

- Therefore, we use the normalized version, called correlation coefficient, showing the strength of the linear relation.

Correlation coefficient

How do we normalize?

$$R = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{\sqrt{\left(N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2 \right) \left(N \sum_{i=1}^N y_i^2 - \left(\sum_{i=1}^N y_i \right)^2 \right)}}$$

$$= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 \sum_{i=1}^N (y_i - \bar{y})^2}}$$

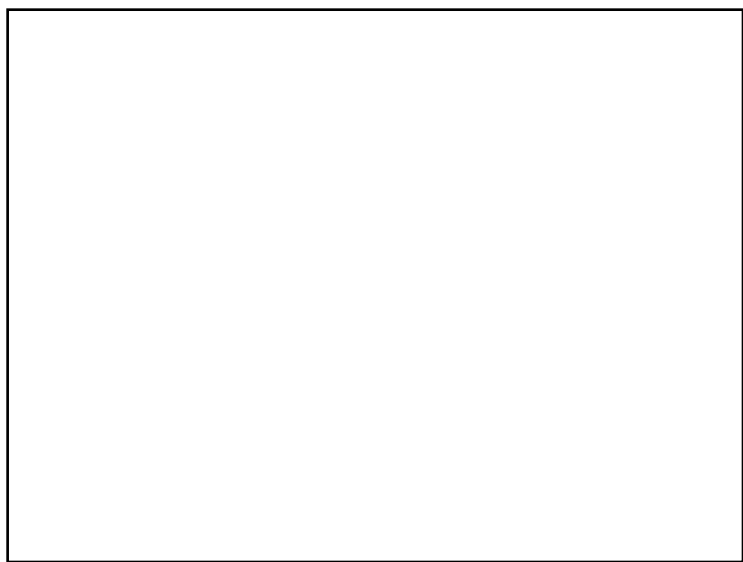
$$\frac{\frac{1}{N-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{N-1} \sum (x_i - \bar{x})^2 \cdot \frac{1}{N-1} \sum (y_i - \bar{y})^2}}$$

in general the covariance is normalized by the product of the standard deviations of X, Y

$$= b \sqrt{\frac{\text{Var}(X)}{\text{Var}(Y)}}$$

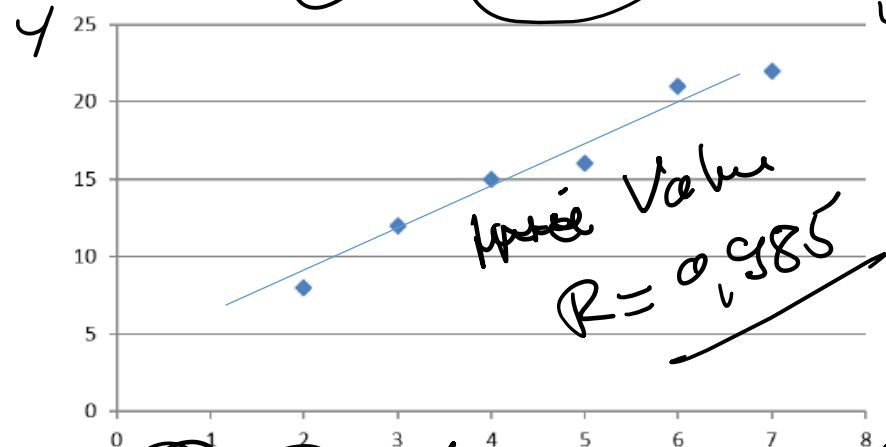
General intervals for $R \in [-1, 1]$

$\pm (0,0;0,2)$	→	almost no dependence
$\pm [0,2;0,4)$	→	low dependence
$\pm [0,4;0,6)$	→	medium dependence
$\pm [0,6;0,8)$	→	high dependence
$\pm [0,8;1,0)$	→	almost full dependence

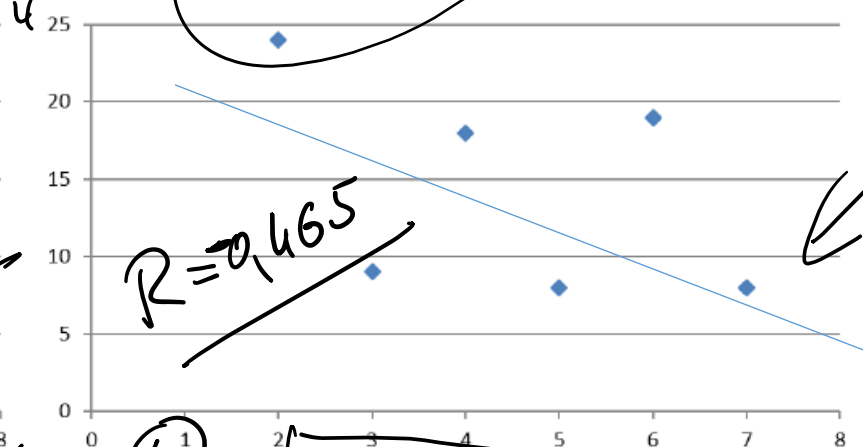


Correlation coefficient examples

① $R \approx 0,8$

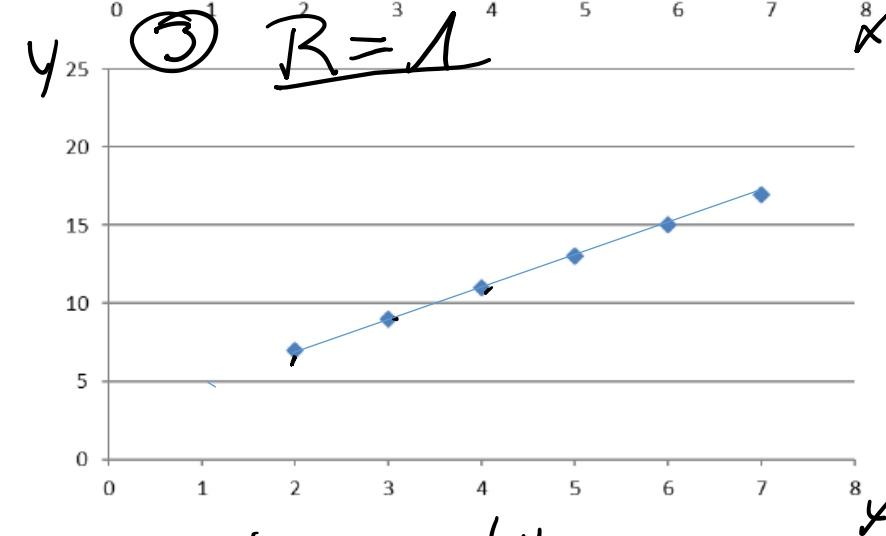


② $R \approx -0,2$



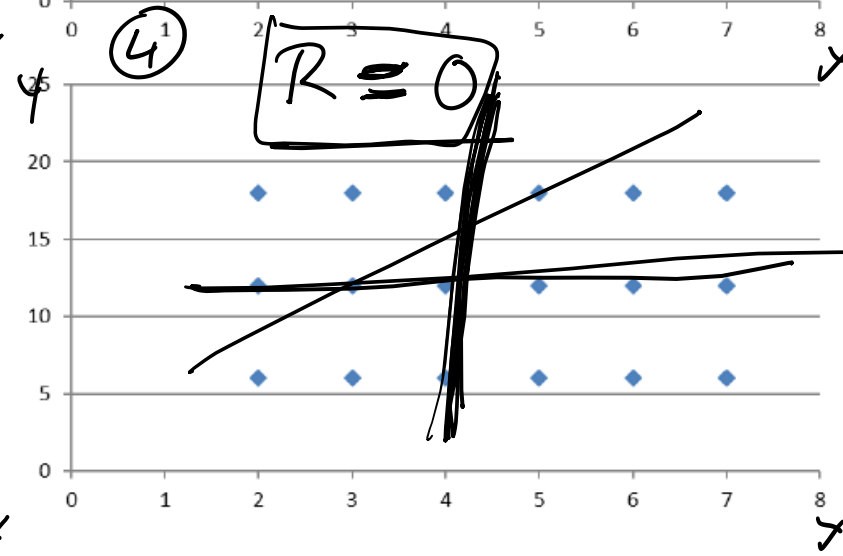
y is more or less decreasing, if x is increasing
 $\Rightarrow R$ should be negative!

③ $R = 1$



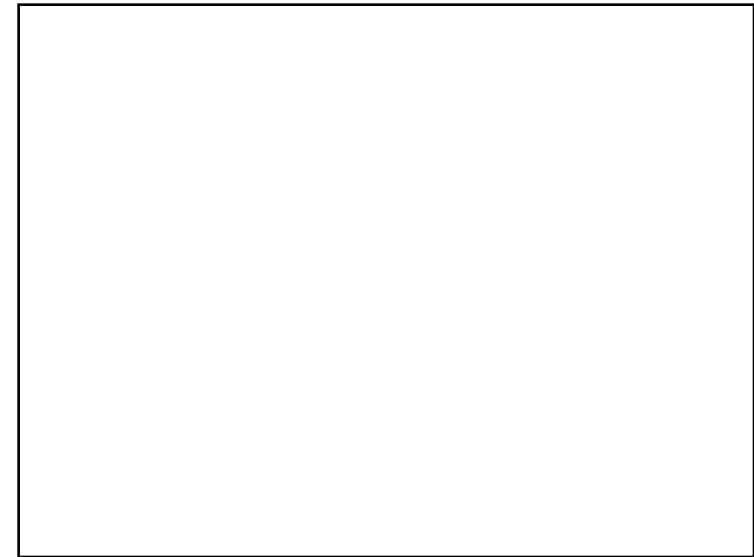
$y = a + bx$
 \Rightarrow all data points are on the line

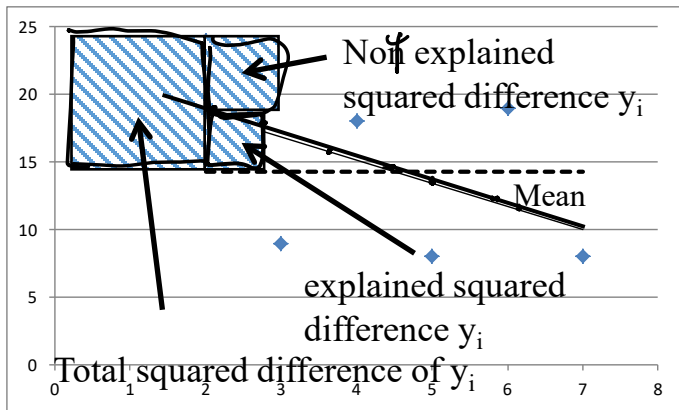
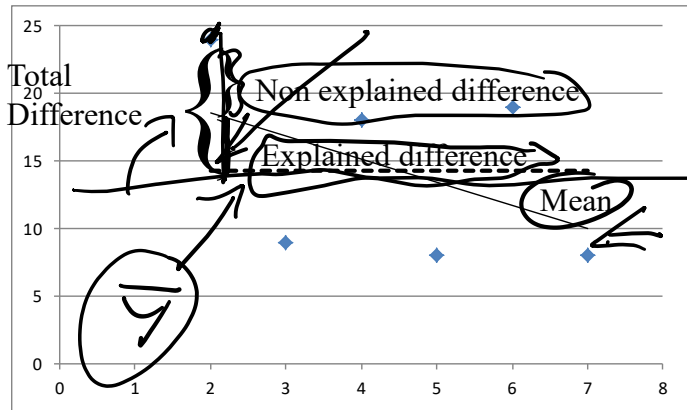
④ $R = 0$



Coefficient of determination R^2

- The squared correlation coefficient R^2 is called coefficient of determination.
- In quantitative analysis, we often use the coefficient of determination R^2 , because it can be intuitively interpreted
- But R^2 does not distinguish between + and - anymore.
- R^2 is in the interval $[0,1]$ since $R \in [-1,1]$
- R^2 graphical interpretation:
 - R^2 equals the proportion of variance explained by the model in relation to the total variance, i.e. how much percent is explained by the regression line
 - $(1-R^2)$ equals the proportion of Variance not explained by the model in relation to the total variance, and has to be explained by other influencing factors





$\hat{y}_i - \bar{y}$ explained difference

$y_i - \hat{y}_i$ not explained difference

$y_i - \bar{y}$ total difference

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

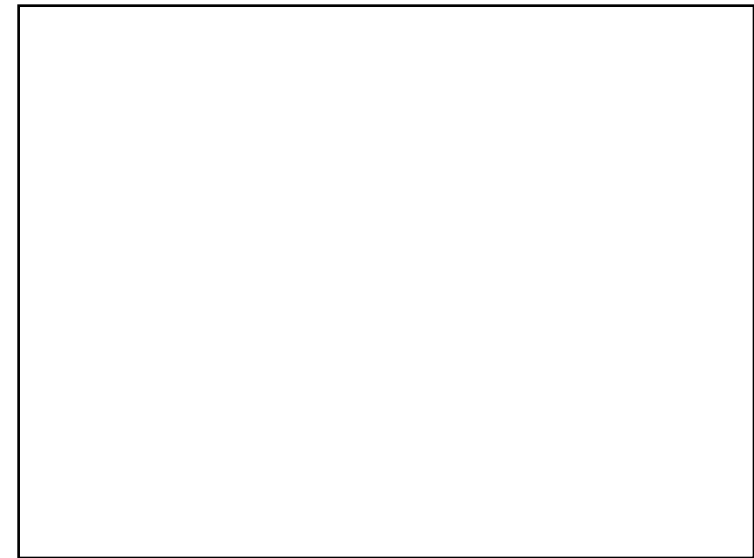
total squared difference

not explained squared difference

explained squared difference

Sum of the squared distances
Sum of the squared total distances

Coefficient of determination $R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \hat{y}_i)^2}$



Correlation Coefficient R^{2C}

$$3,92 + 1,98 \cdot 9$$

$$(21,72 - 27)^2$$

$$\hat{y}_i = 3,92 + 1,98 x_i$$

$$\bar{y} = 27$$

N	x	y	\hat{y}	$(\hat{y} - \bar{y})^2$	$(y - \bar{y})^2$
1	9	24	21,72		
2	11	33			
3	5	10			
4	13	29			
5	20	42			
6	12	24			
Gesamt					

$R^2 =$

$R =$

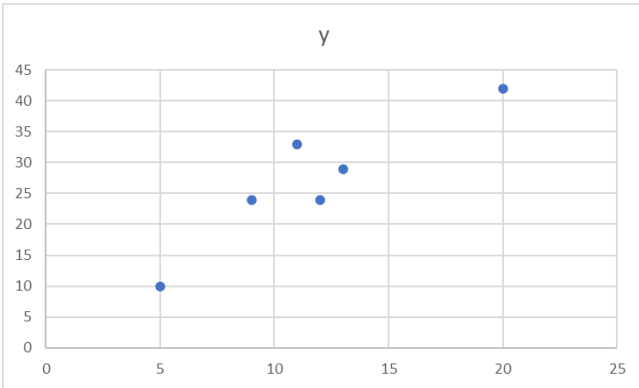


Correlation Coefficient R

$$\hat{y}_i = 3,12 + 1,98 x_i$$

$$\bar{y} = 27$$

	Visits	Sales volume [€]				explained	total
year	x	y	xy	x ²	yhat	(yhat-ybar) ²	(y-ybar) ²
2008	9	24	216	81	21,72	27,83	9,00
2009	11	33	363	121	25,68	1,74	36,00
2010	5	10	50	25	13,81	173,95	289,00
2011	13	29	377	169	29,64	6,96	4,00
2012	20	42	840	400	43,49	271,80	225,00
2013	12	24	288	144	27,66	0,43	9,00
sum	70	162	2134	940		482,72	572,00

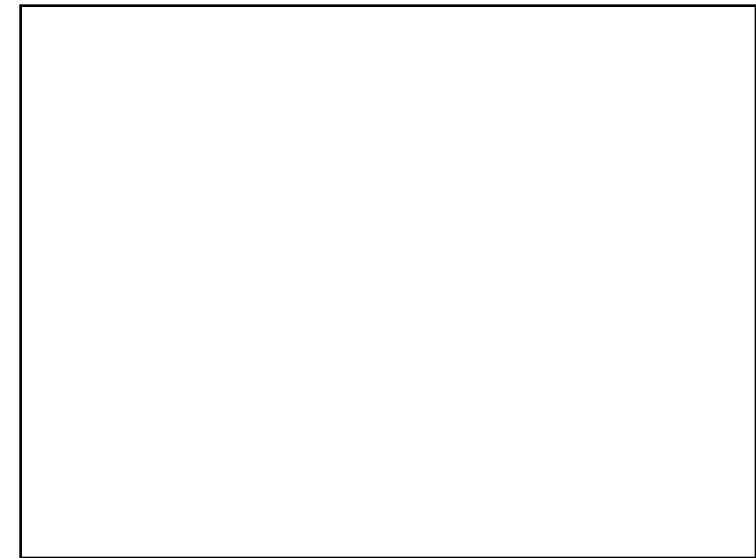


$$R^2 = \frac{487,2}{572}$$

\leftarrow positive dependence
 \Rightarrow taking the root for calculating R we need the st^h sign

$$R^2 = 0,84$$

$$R = +0,92$$



Interpretation of R^2

- R^2 measures the strength of the dependence between two variables X and Y.
- R^2 can be interpreted as the proportion explained by the linear regression.
- $R^2 = 0$, the linear regression model is given just by the constant a and $b=0$
→ The change of the independent variable X has no influence on the dependent variable Y
- $R^2 = 1$, then the regression line fully explains the dependence between X and Y and we have $\hat{y}_i = a + bx_i = y_i$

- Make a forecast for inflation for Germany

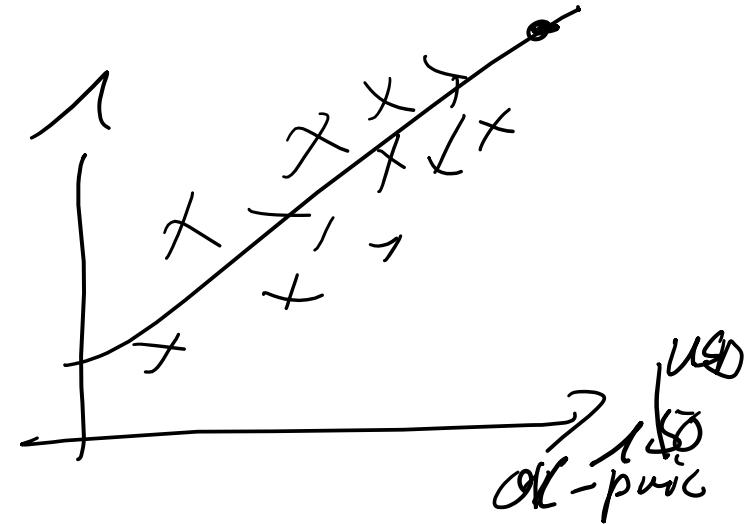
- First take the oil price the last 5 years
(monthly data)

- Exchange rate

- Calculate the dependence between the energy component of the HCPI and the oil price

- And calculate the total inflation coming from the guessed increase of oil prices in the future due to a possible embargo of Russian oil

energy
component



Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

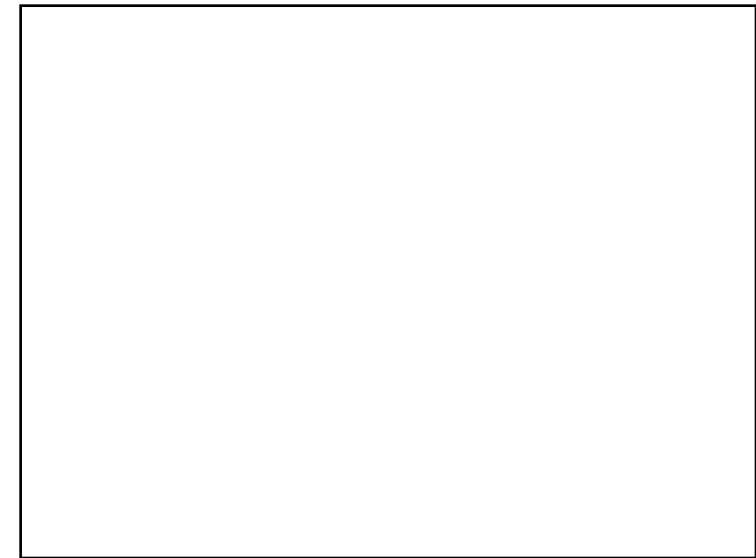
Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Case study I: Inflation

Forecast of the inflation in Germany for 2023 based on the historical development of oil prices

- Data (monthly data):
 - Oil price of the last 5 years
 - Exchange rate Euro/USD
 - HCPI (Harmonized Consumer Price Index) Germany
 - HCPI Energy Component Germany
- Calculate the historical dependence between the energy component of the HCPI and the oil price via a linear regression
- Make an assumption about the future exchange rate and Oil price of 2023
- Assume the non-energy component of HCPI stays constant on the last known level.
- Calculate total inflation for 2023 based on the former assumptions and dependences



Case study I: Inflation

<https://fred.stlouisfed.org/>

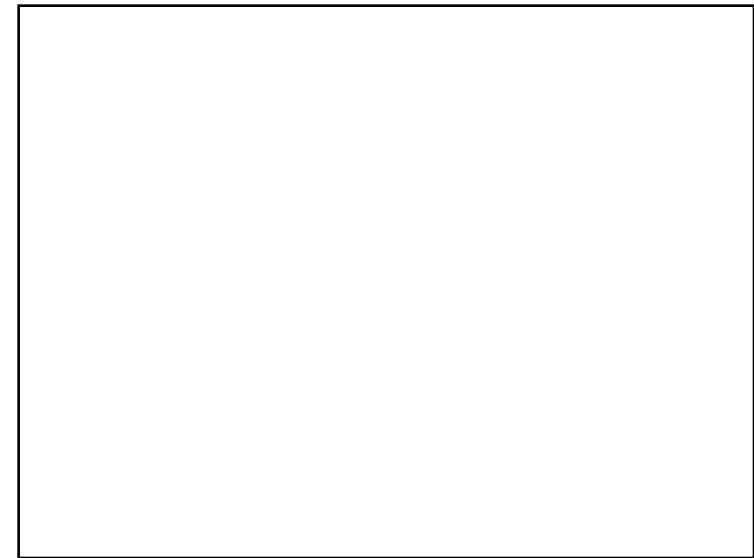
Link for Oil price of sort Brent

https://www.bundesbank.de/dynamic/action/de/statistiken/zeitreihen-datenbanken/zeitreihen-datenbank/759778/759778?listId=www_sdks_b01012_2

Link exchange rate USD-Euro

<https://ec.europa.eu/eurostat/de/>

Link HCPI



Case study I: Inflation

Result:

Dependence between Oil-price and Energy component HCPI

$$\text{HCPI-Energy} = 76,415 + 0,4987 * \text{Oilprice}[\text{Euros}]$$

-> with the assumptions, that Exchangerates USD:Euro will be at 1:1 in Summer and Oilprice will be at 130 USD per Barrel in Summer until December 2023, we obtain, that Oil prices will nove up Inflation by roughly 0,5 Percentagepoints.

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

1. Properties of the linear regression.

(a) Show $\hat{y} = a + b\bar{x}$

(b) Show $b = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - \bar{x}^2}$ and $a = \frac{\bar{x^2}\bar{y} - \bar{x}\bar{xy}}{\bar{x^2} - \bar{x}^2}$

(c) Show $\sum_{i=1}^n \epsilon_i = 0$

(d) Show $\sum_{i=1}^n \epsilon_i x_i = 0$

(e) Show $\sum_{i=1}^N (y_i - \bar{y})^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \sum_{i=1}^N (\hat{y}_i - \bar{y})^2$

This follows from FOC
total variance = explained variance

$\hat{y}_i = a + bx_i$
 $y_i = a + bx_i + \epsilon_i$

$\hat{y}_i - y_i = \epsilon_i$

$\Rightarrow \frac{1}{n} \sum x_i = \bar{x}$ and $\frac{1}{n} \sum y_i = \bar{y}$ are on the regression line

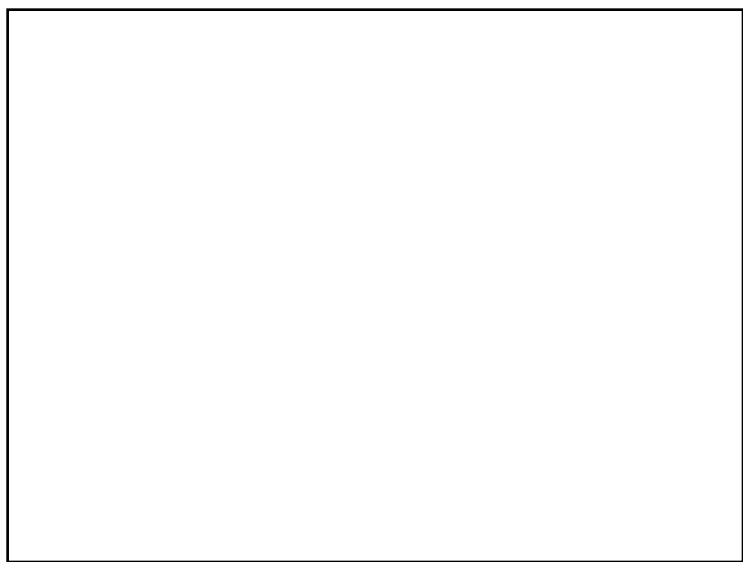
(1) $\frac{\partial F}{\partial a} = -2 \left[\sum_i (y_i - \hat{y}_i) \right] = 0$
(2) $\frac{\partial F}{\partial b} = -2 \left[\sum_i x_i (y_i - \hat{y}_i) \right] = 0$

a) $\Rightarrow \frac{1}{n} \sum y_i - \frac{1}{n} \sum a = \frac{1}{n} \sum bx_i = 0$
 $\bar{y} - a - b\bar{x} = 0 \Rightarrow \bar{y} = a + b\bar{x}$

\Rightarrow thus \bar{y} and \bar{x} are on the regression line

c) Take (1) $\Rightarrow \sum (y_i - \hat{y}_i) = 0 \Rightarrow \sum_{i=1}^n \epsilon_i = 0$ q.e.d.

d) Take (2) $\Rightarrow \sum x_i \epsilon_i = 0$ q.e.d.



1. Properties of the linear regression.

(a) Show $\bar{y} = a + b\bar{x}$

(b) Show $b = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}$ and $a = \frac{\overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy}}{\overline{x^2} - \bar{x}^2}$

(c) Show $\sum_{i=1}^n \epsilon_i = 0$

(d) Show $\sum_{i=1}^n \epsilon_i x_i = 0$

(e) Show $\sum_{i=1}^N (y_i - \bar{y})^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \sum_{i=1}^N (\hat{y}_i - \bar{y})^2$

$$e) \sum (y_i - \bar{y})^2 = \sum y_i^2 - 2\bar{y} \sum y_i + \sum \bar{y}^2$$

$$+ \sum (y_i - \hat{y}_i)^2 = \sum y_i^2 - 2\sum y_i \hat{y}_i + \sum \hat{y}_i^2$$

$$\sum (\hat{y}_i - \bar{y})^2 = \sum \hat{y}_i^2 - 2\bar{y} \sum \hat{y}_i + \sum \bar{y}^2$$

$\Rightarrow 2 \sum \hat{y}_i^2$

$$\Rightarrow -2\bar{y} \sum y_i = \sum \hat{y}_i^2 - 2\sum y_i \hat{y}_i - 2\bar{y} \sum \hat{y}_i + \bar{y} \sum \hat{y}_i$$

by definition $a + bx_i$

$$\bar{y} \sum (\hat{y}_i - y_i) = \sum (\hat{y}_i^2 - y_i \hat{y}_i)$$

$$\underbrace{\epsilon_i}_{\epsilon_i} = \sum \hat{y}_i (\hat{y}_i - y_i)$$

$$\Rightarrow \sum \epsilon_i = \sum (a + bx_i) \epsilon_i = \underbrace{a \sum \epsilon_i}_0 + \underbrace{b \sum x_i \epsilon_i}_0$$

q.e.d

\Rightarrow qualitatively we have
total variance = unexplained variance + explained variance



2. A housing agent sells last April 26 flats with the following data regarding size G in $[m^2]$ and rental M in $[€]$. We know $Cov(G, M) = 5760$ $\sigma_M = 223,61$ $\sigma_G = 21,9$ $\mu_M = 1100$ $\mu_G = 80$

$y \hat{=} M$ $x \hat{=} G$
 $\Rightarrow Var(x) = \sigma_G^2$

- (a) Calculate and draw the regression line.
- (b) Calculate the average and marginal price per m^2 .
- (c) Estimate the rental of flat of size $100 m^2$.

a) $b = \frac{Cov(x, y)}{Var(x)} = \frac{5760}{21,9^2} \approx 12,01$ $a = \bar{y} - b\bar{x} = 1100 - 12,01 \cdot 80 \approx 139,22$
 $\hat{y}_i = a + b x_i$

c) $\hat{M}(100m^2) = 139,22 + 12,01 \cdot 100 \approx 1340,20 €$

3. Within a random sample of 100 data points (X, Y) we have $R = -0,93$. Decide which statement is right or wrong?

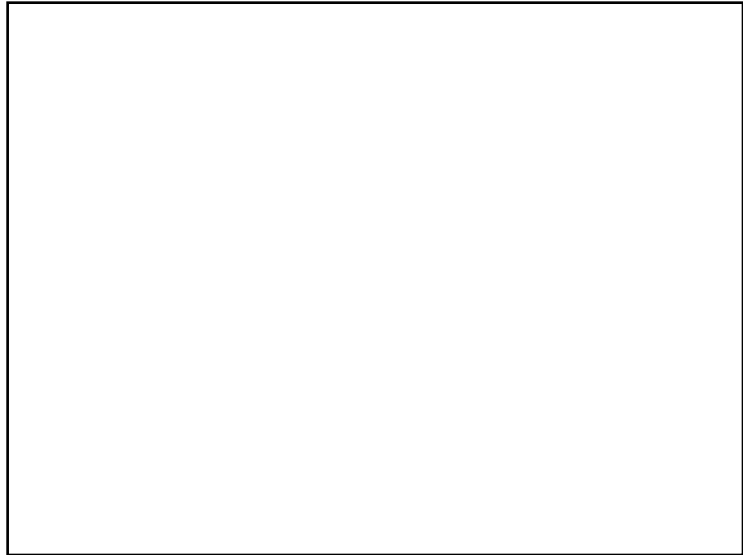
- (1) There is no dependence between X and Y since R is negative.
- (2) The data points scattering very near around the regression line.
- (3) The linear regression has positive slope.
- (4) The dependence of X and Y is inversely proportional.

wrong because no dependence $\Rightarrow R = 0$
 true since R is very near at -1 and $R \in [-1, 1]$
 wrong because R is negative
 true same questions in (3) only that we say y is going down if x is going up

(5) Which proportion of the total variance is explained by the linear regression with a $R = -0,93$?

\Rightarrow coefficient of determination

$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$ \leftarrow explained variation / total variation $= (-0,93)^2 \approx 86,5\%$



4. In 2021 a brewery has the following output and costs:

Month	Output [Hektoliter]	Costs [€]
Jan	600	6500
Feb	680	8200
Mrz	720	7300
Apr	1010	8900
Mai	900	9900
Jun	990	10000
Jul	1270	10300
Aug	1440	12500
Sep	1380	11500
Okt	1010	9200
Nov	830	8200
Dez	1070	9300

we are interested into the cost function of our brewery $\Rightarrow c(x) \rightarrow$ function form!

$$b = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{9627583,3 - 991,6 \cdot 9316,6}{1049650 - (991,6)^2}$$

a) $= 5,865$

- (a) Calculate the Cost function via a linear regression.
- (b) Interpret economically the parameters.
- (c) Calculate the correlation between Output and costs?
- (d) Estimate the total costs and average variable costs of an output of 1100 Hektoliter.

$$a = \bar{y} - b\bar{x} = 9316,6 - 5,865 \cdot 991,6$$

variable costs

$$= 3500,308$$

	x	y	x ²	xy
mean	991,666667	9316,66667	1049650	9627583,33

$\frac{\quad}{\bar{x}} \quad \frac{\quad}{\bar{y}} \quad \frac{\quad}{\overline{x^2}} \quad \frac{\quad}{\overline{xy}}$

\Rightarrow costfunction $c(x) = \underbrace{3500,308}_{\text{fixed costs}} + 5,865 \cdot x$

$c(x) = \alpha + \beta x^2$

b) $a := \text{fix cost}$
 variable costs := $b \cdot x$

$b :=$ variable costs per unit
 because we have a linear dependence

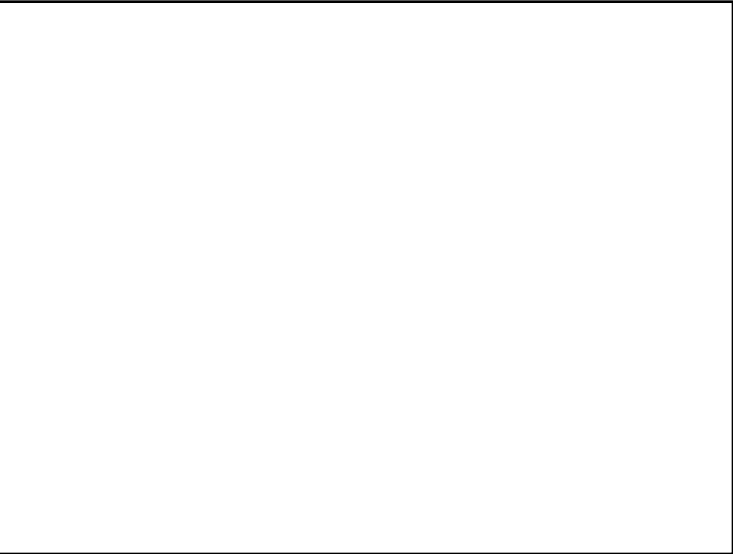
\therefore marginal costs
 $= c'(x) = b$

c) $R =$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2} = 0,87 \Rightarrow R = 0,93$$

$$c(1100) = 3500,308 + 5,865 \cdot 1100 = 9952,667$$





Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

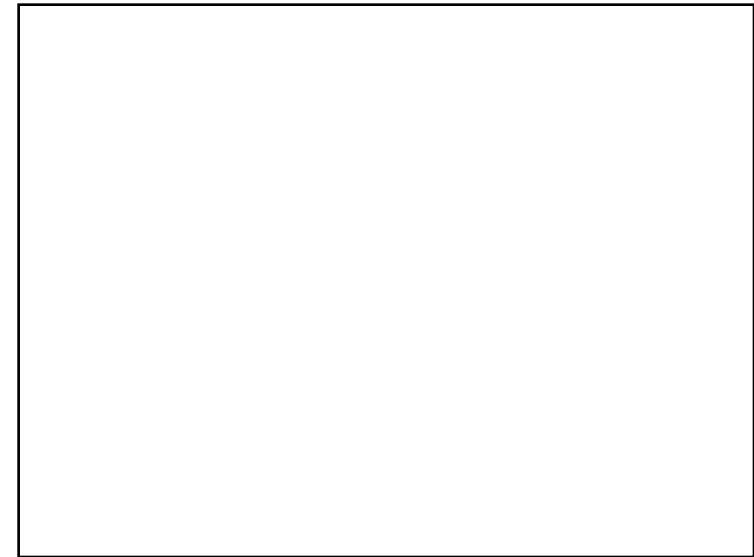
Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Markowitz and CAPM

Bernhard Köster

Summer Term 2022

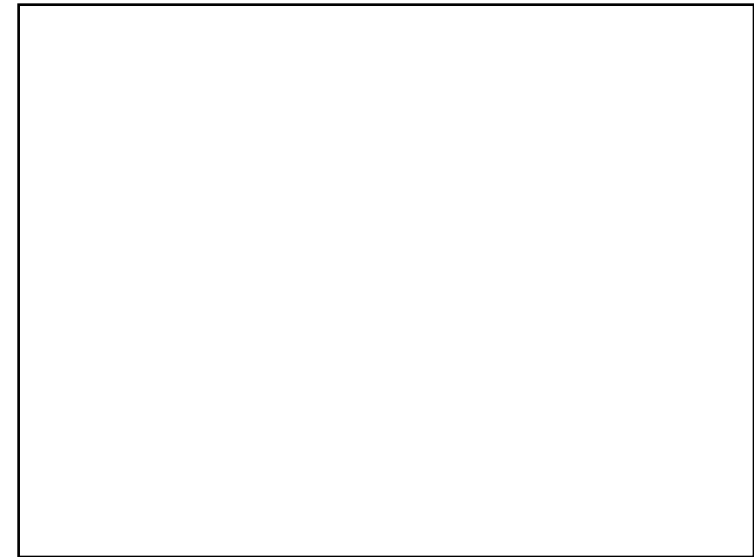


Example: Minimum-Variance-Portfolio

A Portfolio consists of two assets a and b with investment weight $\alpha \in [0, 1]$.

Prices of the assets are random variables with means μ_a, μ_b and variances σ_a^2, σ_b^2 . Calculate the portfolio with the minimal variance, if assets are correlated with $\rho \in (-1, 1)$.

$$\rho = \mathbb{R}$$



Optimizing problems with many variables

Suppose we have $\vec{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$. For continuously twice differentiable functions, we have:

for example $u(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$
 $\frac{\partial f}{\partial x_1} = 0 \quad \frac{\partial f}{\partial x_2} = 0 \dots$

Investmentparaphrase
 $\rightarrow 20\% \text{ in Asset 1}$
 $80\% \text{ in Asset 2}$

$f(\vec{x}) \rightarrow \max \Rightarrow \frac{\partial f}{\partial \vec{x}}(\vec{x}^*) = 0$ and $\text{Hess } f(\vec{x}^*)$ is negative definite
 FOC or

$$\begin{matrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{matrix}$$

$c'(x) = 0$ FOC
 Maximiere $u(x^*) < 0$

$f(\vec{x}) \rightarrow \min \Rightarrow \frac{\partial f}{\partial \vec{x}}(\vec{x}^*) = 0$ and $\text{Hess } f(\vec{x}^*)$ ist positive definite

\Rightarrow there exists a global min/max

Very often in economics, we are confronted not with global optimization but with optimizing problems subject to some restrictions. Then we have

$$f(\vec{x}) \rightarrow \max (\min) \text{ s.t. } g_1(\vec{x}) = 0, g_2(\vec{x}) = 0, \dots, g_m(\vec{x}) = 0 \quad (n > m)$$

microeconomics

$$m - p_1 x_1 - p_2 x_2 = p$$

Budget
restriction

with FOC

Lagrange function i.e. utility

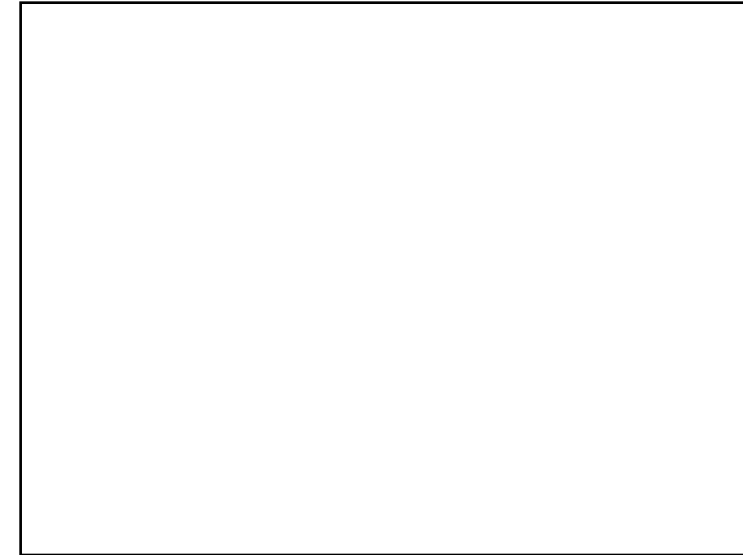
$$\rightarrow \mathcal{L} = f(\vec{x}) + \lambda_1 g_1 + \lambda_2 g_2 + \dots + \lambda_m g_m$$

$$\mathcal{L}(x_1, \dots, x_n, \lambda_1, \dots, \lambda_m)$$

FOC

$$\begin{matrix} \frac{\partial \mathcal{L}}{\partial x_1} = 0 & \dots & \frac{\partial \mathcal{L}}{\partial x_n} = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda_1} = 0 & \dots & \frac{\partial \mathcal{L}}{\partial \lambda_n} = 0 \end{matrix}$$

For SOC including the second order hessian matrix and regularity conditions for the existence of optima (i.e. implicit functions theorem is valid!) are beyond this lecture and we refer to the literature.

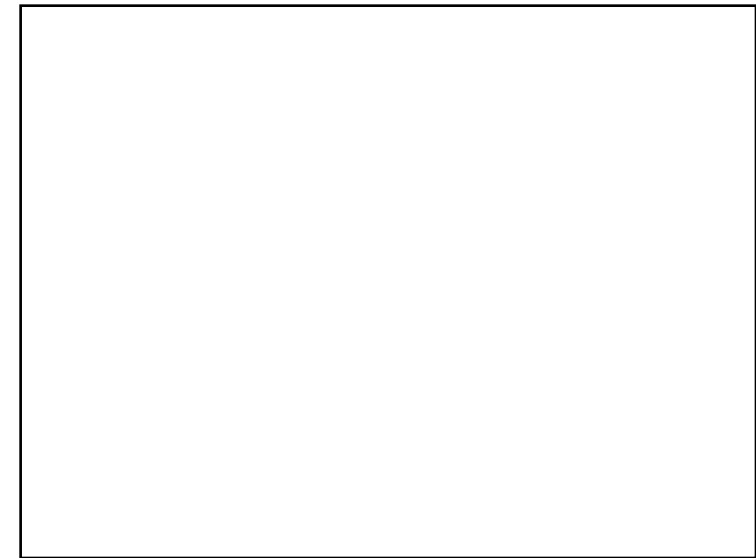


Markowitz Portfolio Modell

Assumptions:

1. An Investor has different risky investment opportunities with finite means and variances.
2. Investors minimize one-period-utility with decreasing marginal utility in wealth and a risk avers structure.
3. Risk is measured via the volatility of the expected yield.
4. Decision variables are only expected yield and variances.
5. Given the Risk investors prefer higher yield to lower yield. Given the yield, investors prefer lower risk to higher risk.

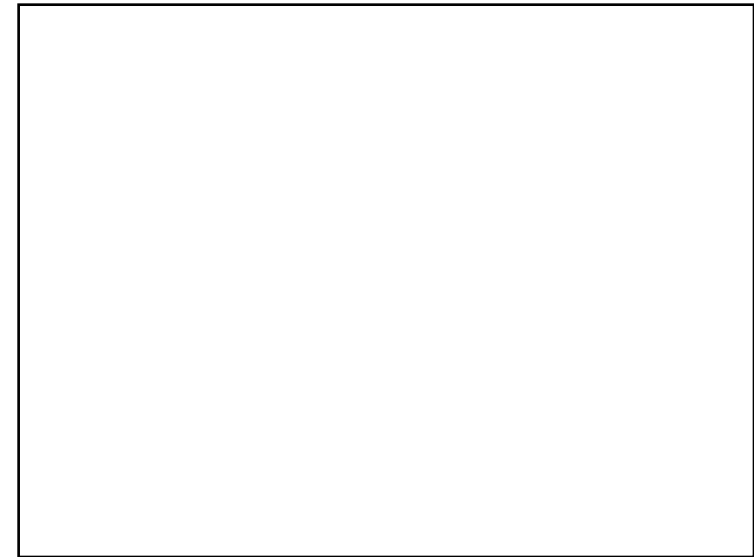
⇒ in general risk avers investors



Notation

random variables

- ▶ R_1, R_2 : Yield asset 1 and 2
- ▶ $\mu_1 = E(R_1), \mu_2 = E(R_2)$: expected yield asset 1 and 2
- ▶ $\sigma_1 = \sigma(R_1), \sigma_2 = \sigma(R_2)$: standard deviation asset 1 and 2
- ▶ $\rho = \rho(R_1, R_2) = \frac{\text{Cov}(R_1, R_2)}{\sigma_1 \sigma_2}$: correlation coefficient of the yields of asset 1 and 2
- ▶ R : yield of a Portfolio consisting of asset 1 and 2
- ▶ $\mu = \mu(R)$: expected yield of a portfolio
- ▶ $\sigma = \sigma(R)$: standard deviation of a portfolios



Expected yield and Variance of a Portfolios

$$\rho = \frac{\text{COV}(R_1, R_2)}{\sigma_1 \sigma_2}$$

$$E(X + \alpha Y)$$

$$= E(X) + \alpha E(Y)$$

X, Y are random variables

$$R = \alpha R_1 + (1 - \alpha) R_2$$

μ f

$$E(R) = \mu = \alpha \mu_1 + (1 - \alpha) \mu_2 \quad (1)$$

$$Var(R) = \sigma^2 = \alpha^2 \sigma_1^2 + (1 - \alpha)^2 \sigma_2^2 + 2\alpha(1 - \alpha) \rho \sigma_1 \sigma_2 \quad (2)$$

(3)

► W.l.o.g. $\mu_1 < \mu_2$ and $\sigma_1 < \sigma_2$ ($\alpha \in [0, 1]$) proportion invested in asset 1 \Rightarrow $(1 - \alpha)$ investment in asset 2

$$Var(\alpha R_1 + (1 - \alpha) R_2) = E([\alpha R_1 + (1 - \alpha) R_2 - E(\alpha R_1 + (1 - \alpha) R_2)]^2)$$

$$= E([\alpha R_1 + (1 - \alpha) R_2 - \alpha \mu_1 - (1 - \alpha) \mu_2]^2)$$

$$= E([\alpha (R_1 - \mu_1) + (1 - \alpha) (R_2 - \mu_2)]^2)$$

$$= E(\alpha^2 (R_1 - \mu_1)^2 + (1 - \alpha)^2 (R_2 - \mu_2)^2 + 2\alpha(1 - \alpha) (R_1 - \mu_1)(R_2 - \mu_2))$$

$$= \alpha^2 E((R_1 - \mu_1)^2) + (1 - \alpha)^2 E((R_2 - \mu_2)^2) + 2\alpha(1 - \alpha) E((R_1 - \mu_1)(R_2 - \mu_2))$$

$\text{COV}(R_1, R_2)$

$$E(R) = \mu = \alpha \mu_1 + (1-\alpha) \mu_2$$

$$\text{Var}(R) = \sigma^2 = \alpha^2 \sigma_1^2 + (1-\alpha)^2 \sigma_2^2 + 2\alpha(1-\alpha) \rho \sigma_1 \sigma_2$$

$\rho = 1$ perfect correlation

$$\Rightarrow \sigma^2 = \alpha^2 \sigma_1^2 + (1-\alpha)^2 \sigma_2^2 + 2\alpha(1-\alpha) \sigma_1 \sigma_2 = x^2 + y^2 + 2xy = (x+y)^2$$

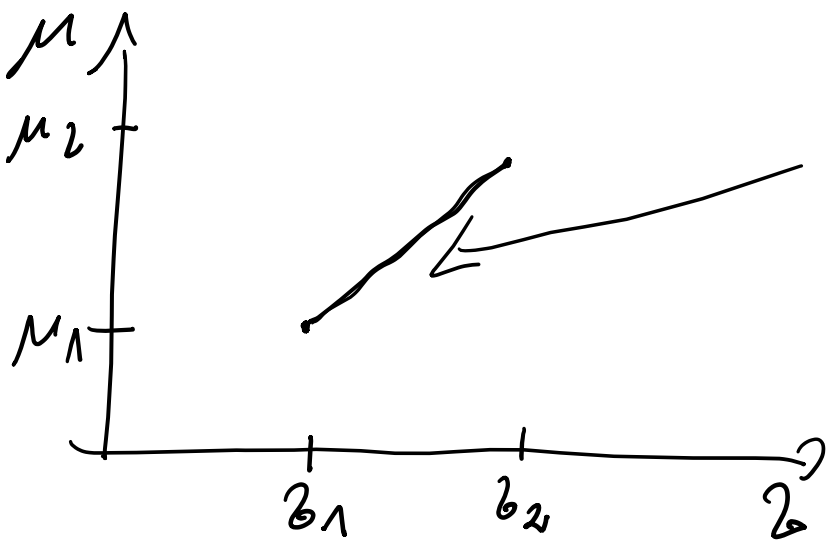
$$= (\alpha \sigma_1 + (1-\alpha) \sigma_2)^2$$

$$\Rightarrow \sigma = \alpha \sigma_1 + \sigma_2 - \alpha \sigma_2 = \sigma_2 + \alpha(\sigma_1 - \sigma_2) \Rightarrow \alpha = \frac{\sigma - \sigma_2}{\sigma_1 - \sigma_2} = \alpha(\sigma)$$

$$\Rightarrow \mu = \alpha(\mu_1 - \mu_2) + \mu_2 = \frac{\sigma - \sigma_2}{\sigma_1 - \sigma_2} (\mu_1 - \mu_2) + \mu_2$$

$\mu(\sigma) \rightarrow$ what kind of dependence?

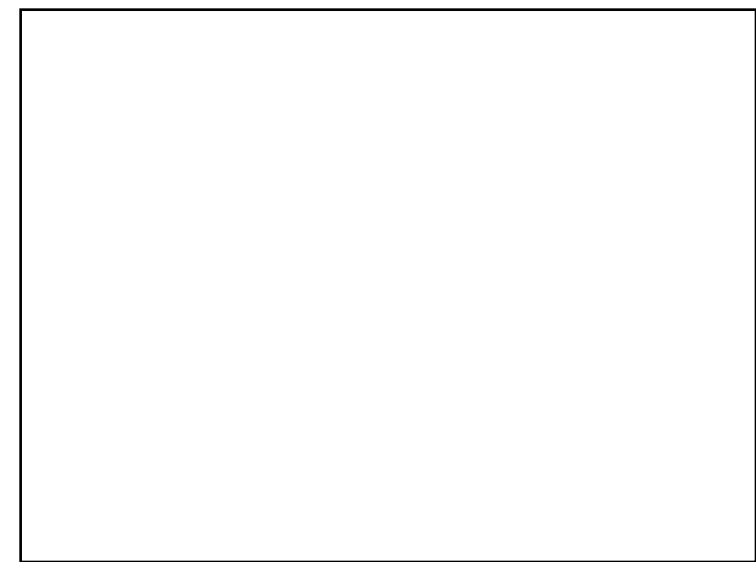
all these combinations of (μ, σ) are accessible



(2) what kind of portfolios are accessible for an investor?
 see $\alpha^2 \sigma_1^2 = (\alpha \sigma_1)^2 \mid \alpha \sigma_1 = x$
 $(1-\alpha) \sigma_2 = y$

are all parameters

this is just a linear function



$$E(R) = \mu = \alpha\mu_1 + (1-\alpha)\mu_2 = \alpha(\mu_1 - \mu_2) + \mu_2 = 0$$

$$\text{Var}(R) = \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 + 2\alpha(1-\alpha)\rho\sigma_1\sigma_2$$

No correlation $\rho = 0$

$$\Rightarrow \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 = \sigma^2(\alpha)$$

$$\sigma^2'(\alpha) = 2\alpha\sigma_1^2 - 2(1-\alpha)\sigma_2^2 = 0$$

$$\Rightarrow \alpha\sigma_1^2 - \sigma_2^2 + \alpha\sigma_2^2 = 0 \Rightarrow \alpha = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

Inserting α into $\sigma^2(\alpha)$

$$\Rightarrow \sigma^2 = \left(\frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}\right)^2 \sigma_1^2 + \left(\frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2}\right)^2 \sigma_2^2 = \frac{\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2} < \sigma_1^2 < \sigma_2^2$$

(2)
(3)

Since we have a quadratic dependence we can ask for the portfolio with the parameter α which creates the minimum risk i.e. minimum variance

$$1 = \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

$$\alpha = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

$$f(x) = (1-x)^2$$

$$f'(x) = 2(1-x)(-1)$$

Investment parameter for the Minimum-Variance-Portfolio

$\sigma_{MVP} = \text{Minimum Variance}$

investing in a portfolio with a proportion of α in asset 1 you can reduce your risk!

What are the accessible portfolios?

from (2) we obtain $\alpha = \frac{\mu - \mu_2}{\mu_1 - \mu_2}$

insert in (3)

$$\Rightarrow \sigma^2 = \left(\frac{\mu - \mu_2}{\mu_1 - \mu_2}\right)^2 \sigma_1^2 + \left(1 - \frac{\mu - \mu_2}{\mu_1 - \mu_2}\right)^2 \sigma_2^2$$



$$\sigma^2(\mu)$$

$$E(R) = \mu = \alpha\mu_1 + (1-\alpha)\mu_2 \quad (2)$$

$$\text{Var}(R) = \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 + 2\alpha(1-\alpha)\rho\sigma_1\sigma_2 \quad (3)$$

$$\sigma^2 = -1 \Rightarrow \sigma^2 = (\alpha\sigma_1)^2 + ((1-\alpha)\sigma_2)^2 - 2\alpha(1-\alpha)\sigma_1\sigma_2 \stackrel{1}{=} x^2 + y^2 - 2xy = (x-y)^2$$

$$\Rightarrow \sigma^2 = (\alpha\sigma_1 - (1-\alpha)\sigma_2)^2$$

$$x = \alpha\sigma_1 \quad y = (1-\alpha)\sigma_2$$

$$\Rightarrow \sigma = |\alpha\sigma_1 + \alpha\sigma_2 - \sigma_2| \Rightarrow \text{can be zero} \Rightarrow \alpha(\sigma_1 + \sigma_2) - \sigma_2 = 0$$

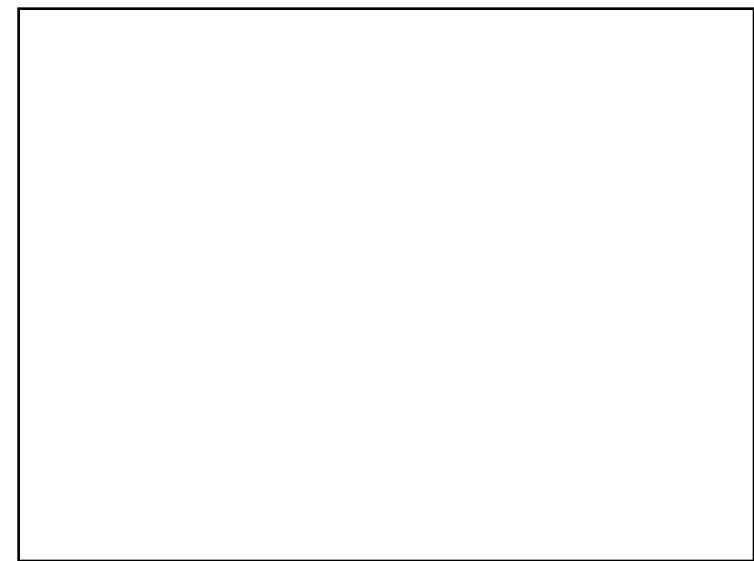
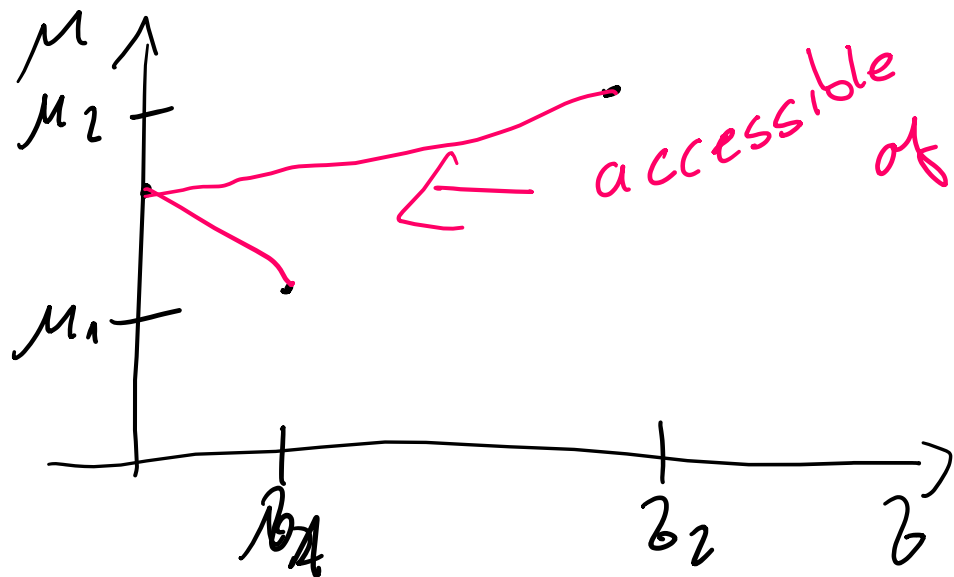
\Rightarrow insert this dependence in μ

$$\Rightarrow \boxed{\alpha^* = \frac{\sigma_2}{\sigma_1 + \sigma_2}}$$

where the portfolio has

\Rightarrow again you obtain
a linear dependence between μ and σ

no risk
 $\sigma^2(\alpha^*) = 0$



Accessible portfolios

We obtain all accessible portfolios (μ, σ) inserting (2) in (3) solving for σ .

$$\sigma^2 = \frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{(\mu_1 - \mu_2)^2} \left[\mu - \frac{\mu_2\sigma_1^2 + \mu_1\sigma_2^2 - (\mu_1 + \mu_2)\rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} \right]^2 \quad (4)$$

$= b^2(\mu)$

$$+ \frac{(1 - \rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} \quad (5)$$

$$\sigma^2 = \text{Var}K^2(R_1, R_2) [\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2 \quad (6)$$

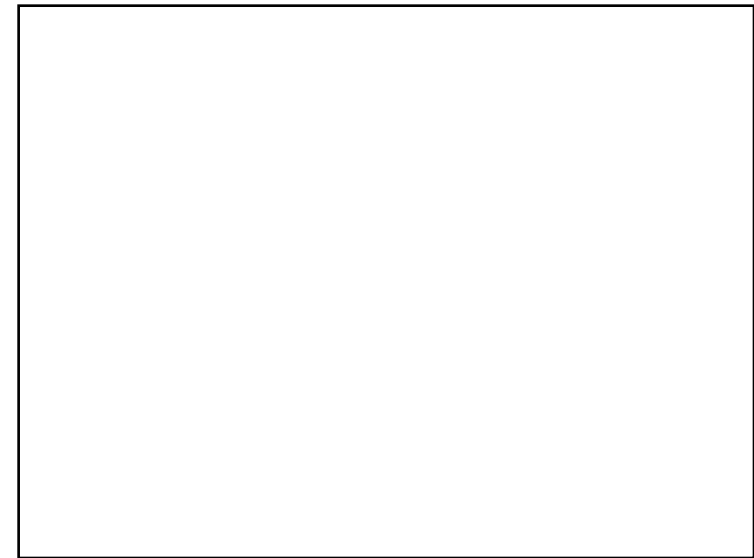


Accessible portfolios

▶ $VarK^2(R_1, R_2) = \frac{Var(R_1.R_2)}{[E(R_1.R_2)]^2} = \frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{(\mu_1 - \mu_2)^2}$ (Coefficient of variation)

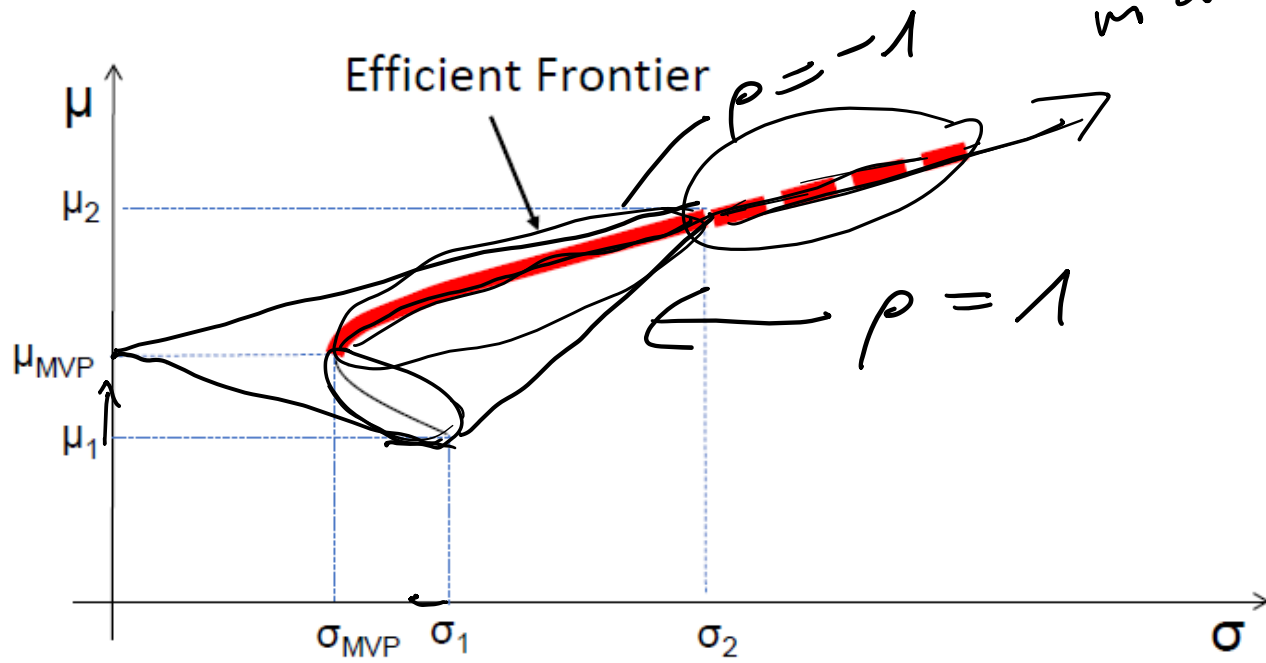
▶ $\mu_{MVP} = \frac{\mu_2\sigma_1^2 + \mu_1\sigma_2^2 - (\mu_1 + \mu_2)\rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$ and $\sigma_{MVP}^2 = \frac{(1 - \rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$

$(\mu_{MVP}, \sigma_{MVP})$: Expected value and standard deviation of the global minimum-variance-portfolio. (6) represents a parabola with the vertex $(\mu_{MVP}, \sigma_{MVP})$. For the extrema $\rho = \pm 1$ we obtain separate solutions.

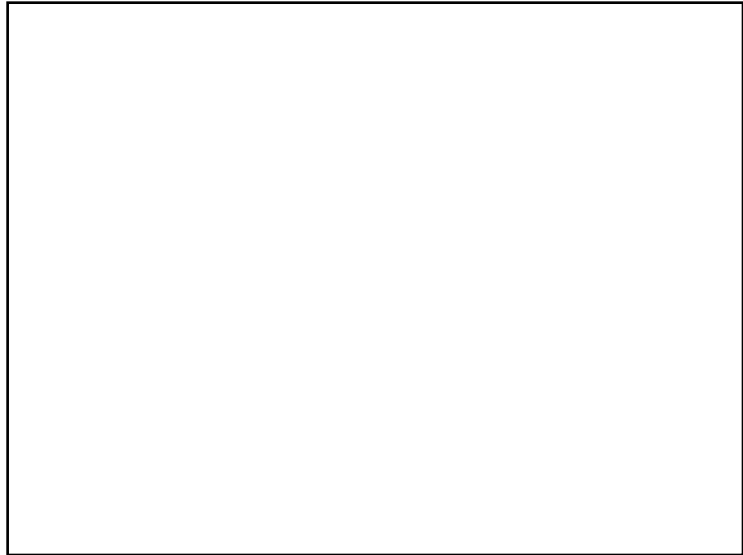


Efficient frontier

Since we assume risk averse investors (higher risk implicates higher yields) only the upper part of the parabola is economically relevant. This is called the *efficient frontier*.



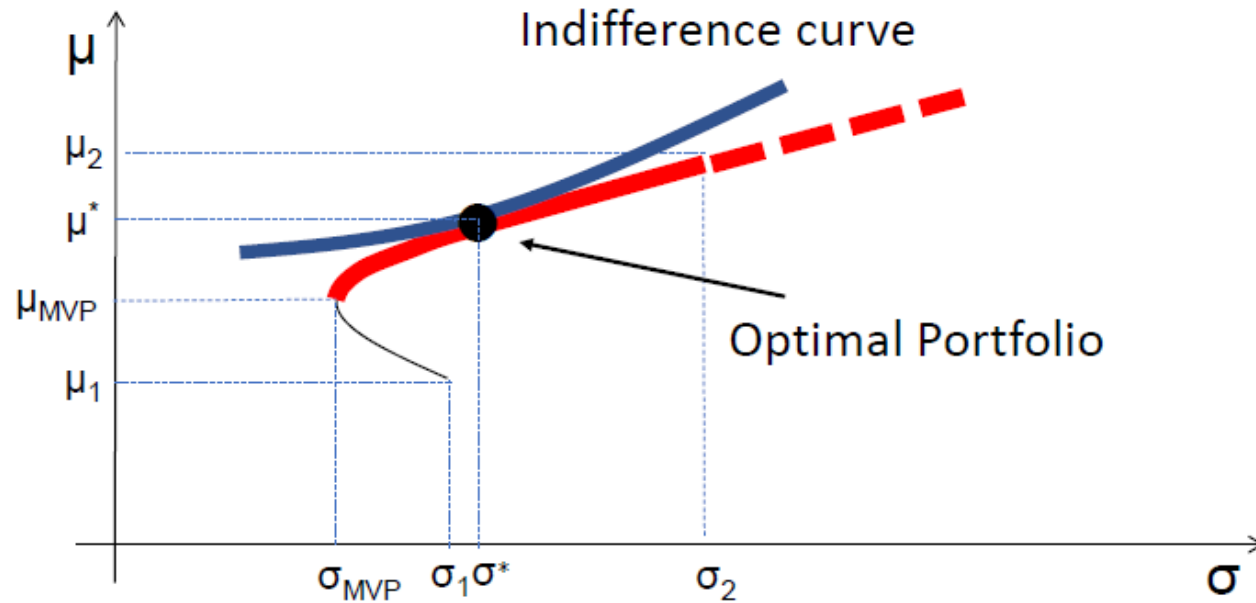
in a region with $\alpha < 0$
 $\alpha > 1$
is also possible
since in financial markets you can short sell an asset



Optimal Portfolio

In general a risk averse utility function is given by $U(\mu, \sigma)$ with $\frac{\partial U}{\partial \mu} > 0$ and $\frac{\partial U}{\partial \sigma} < 0$. The optimizing problem is then given by:

$$\max_{\mu, \sigma} U(\mu, \sigma) \quad \text{N.B. } \sigma^2 = A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2 \quad (A = \text{Var}K^2(R_1, R_2)) \quad (7)$$



Example

The explicit utility is given by $U(\mu, \sigma) = \mu - \frac{a}{2}\sigma^2$ with $(a > 0)$. \implies

$$\mathcal{L}(\mu, \sigma, \lambda) = \mu - \frac{a}{2}\sigma^2 + \frac{\lambda}{2}(\sigma^2 - (A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2))$$

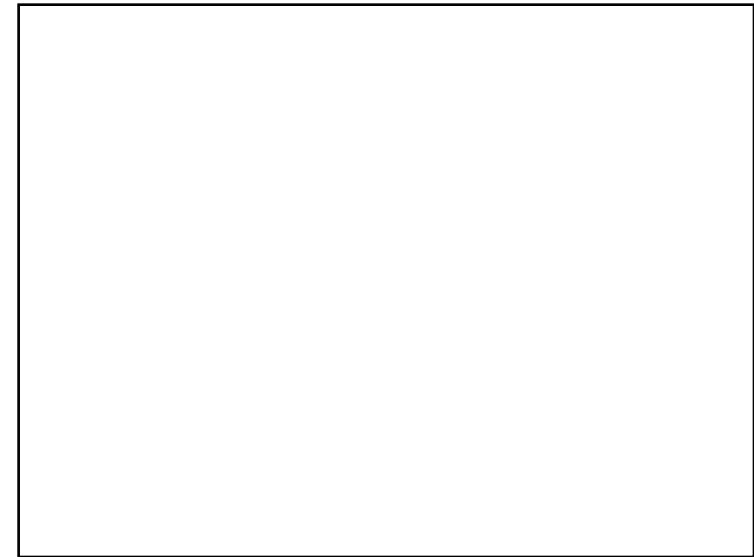
\implies FOC

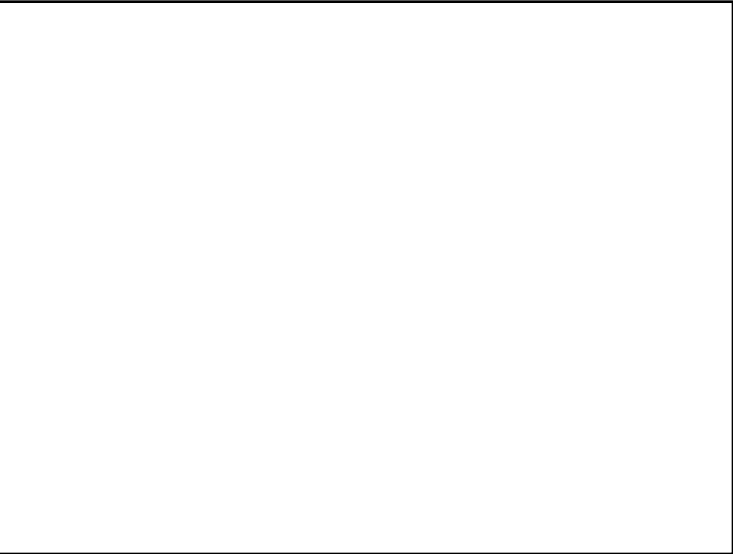
$$\frac{\partial \mathcal{L}}{\partial \mu} = 1 - A\lambda(\mu - \mu_{MVP}) = 0$$

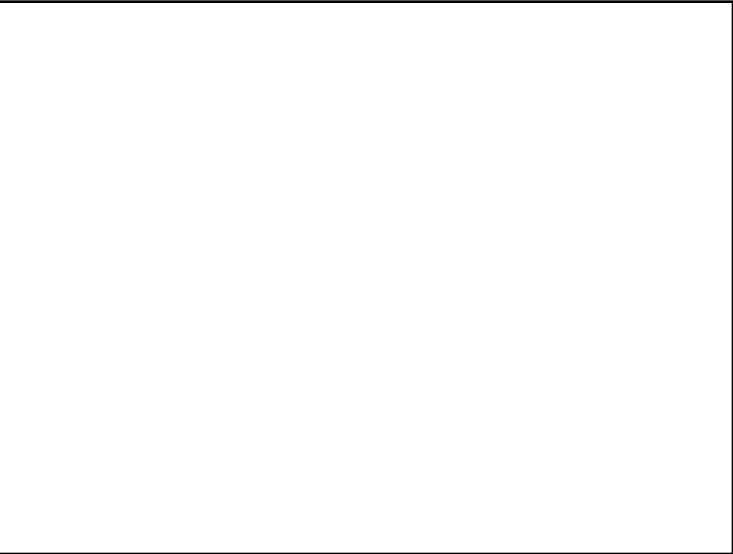
$$\frac{\partial \mathcal{L}}{\partial \sigma} = -a\sigma + \lambda\sigma = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sigma^2 - (A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2) = 0$$

$$\implies \mu^* = \mu_{MVP} + \frac{1}{aA} \quad \sigma^* = \sqrt{\frac{1}{Aa^2} + \sigma_{MVP}^2} \quad \alpha^* = \frac{\mu^* - \mu_2}{\mu_1 - \mu_2}$$







Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web

[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

Markowitz Portfolio Modell

Assumptions:

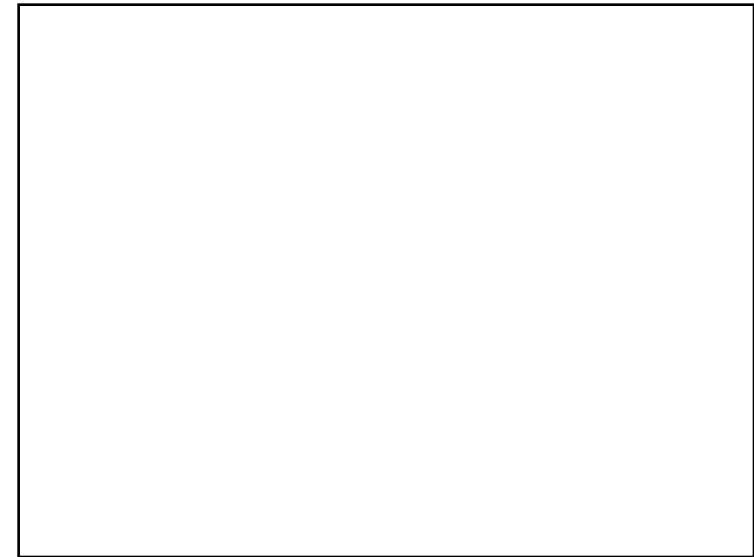
1. An Investor has different risky investment opportunities with finite means and variances.
2. Investors minimize one-period-utility with decreasing marginal utility in wealth and a risk avers structure.
3. Risk is measured via the volatility of the expected yield.
4. Decision variables are only expected yield and variances.
5. Given the Risk investors prefer higher yield to lower yield. Given the yield, investors prefer lower risk to higher risk.

⇒ in general risk avers investors

Notation

random variables

- ▶ R_1, R_2 : Yield asset 1 and 2
- ▶ $\mu_1 = E(R_1), \mu_2 = E(R_2)$: expected yield asset 1 and 2
- ▶ $\sigma_1 = \sigma(R_1), \sigma_2 = \sigma(R_2)$: standard deviation asset 1 and 2
- ▶ $\rho = \rho(R_1, R_2) = \text{Cov}(R_1, R_2) / (\sigma_1 \sigma_2)$: correlation coefficient of the yields of asset 1 and 2
- ▶ R : yield of a Portfolio consisting of asset 1 and 2
- ▶ $\mu = \mu(R)$: expected yield of a portfolio
- ▶ $\sigma = \sigma(R)$: standard deviation of a portfolios



Expected yield and Variance of a Portfolios

$$\text{Var}(\alpha x + \gamma)$$

$$\rho = \frac{\text{COV}(R_1, R_2)}{\sigma_1 \sigma_2}$$

$$\begin{aligned} E(x + \alpha y) \\ = E(x) + \alpha E(y) \end{aligned}$$

x, y are random variables

$$R = \alpha R_1 + (1 - \alpha) R_2$$

(1)

$$E(R) = \mu = \alpha \mu_1 + (1 - \alpha) \mu_2$$

(2)

$$\text{Var}(R) = \sigma^2 = \alpha^2 \sigma_1^2 + (1 - \alpha)^2 \sigma_2^2 + 2\alpha(1 - \alpha) \rho \sigma_1 \sigma_2$$

(3)

$(1 - \alpha)$ investment in asset 2

► W.l.o.g. $\mu_1 < \mu_2$ and $\sigma_1 < \sigma_2$ $\alpha \in [0, 1]$ proportion invested in asset 1

$$\text{Var}(\alpha R_1 + (1 - \alpha) R_2) = E([\alpha R_1 + (1 - \alpha) R_2 - E(\alpha R_1 + (1 - \alpha) R_2)]^2)$$

$$= E([\alpha R_1 + (1 - \alpha) R_2 - \alpha \mu_1 - (1 - \alpha) \mu_2]^2)$$

$$= E([\alpha (R_1 - \mu_1) + (1 - \alpha) (R_2 - \mu_2)]^2)$$

$$= E(\underbrace{\alpha^2 (R_1 - \mu_1)^2}_{\sigma_1^2} + \underbrace{(1 - \alpha)^2 (R_2 - \mu_2)^2}_{\sigma_2^2} + 2\alpha(1 - \alpha) (R_1 - \mu_1) (R_2 - \mu_2))$$

$$= \alpha^2 E((R_1 - \mu_1)^2) + (1 - \alpha)^2 E((R_2 - \mu_2)^2) + 2\alpha(1 - \alpha) \underbrace{E((R_1 - \mu_1) (R_2 - \mu_2))}_{\text{COV}(R_1, R_2)}$$

$$E(R) = \mu = \alpha \mu_1 + (1-\alpha) \mu_2$$

$$\text{Var}(R) = \sigma^2 = \alpha^2 \sigma_1^2 + (1-\alpha)^2 \sigma_2^2 + 2\alpha(1-\alpha) \rho \sigma_1 \sigma_2$$

$\rho = 1$ perfect correlation

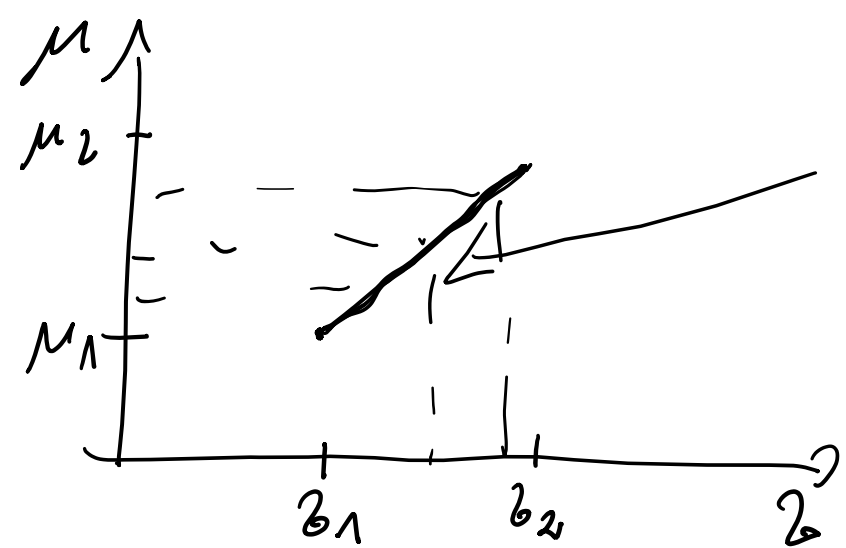
$$\Rightarrow \sigma^2 = \alpha^2 b_1^2 + (1-\alpha)^2 b_2^2 + 2\alpha(1-\alpha) b_1 b_2 = x^2 + y^2 + 2xy = (x+y)^2 = (\alpha b_1 + (1-\alpha) b_2)^2$$

$$\Rightarrow b = \alpha b_1 + b_2 - \alpha b_2 = b_2 + \alpha(b_1 - b_2) \Rightarrow \alpha = \frac{b - b_2}{b_1 - b_2} = \alpha(b)$$

$$\Rightarrow \mu = \alpha(\mu_1 - \mu_2) + \mu_2 = \frac{b - b_2}{b_1 - b_2} (\mu_1 - \mu_2) + \mu_2$$

$\mu(b) \rightarrow$ what kind of dependence?

all parameters are all parameters
 This is just a linear function



all these combinations of (μ, σ) are accessible

(2) what kind of portfolios are accessible for an investor?
 see $\alpha^2 b_1^2 = (\alpha b_1)^2 \mid \alpha b_1 = x$
 $(1-\alpha) b_2 = y$



$$E(R) = \mu = \alpha\mu_1 + (1-\alpha)\mu_2 = \alpha(\mu_1 - \mu_2) + \mu_2 = 0$$

$$\text{Var}(R) = \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 + 2\alpha(1-\alpha)\rho\sigma_1\sigma_2$$

No correlation $\rho = 0$

$$\Rightarrow \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 = \sigma^2(\alpha)$$

$$\sigma^2'(\alpha) = 2\alpha\sigma_1^2 - 2(1-\alpha)\sigma_2^2 = 0$$

$$\Rightarrow \alpha\sigma_1^2 - \sigma_2^2 + \alpha\sigma_2^2 = 0 \Rightarrow \alpha = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

Inserting α into $\sigma^2(\alpha)$

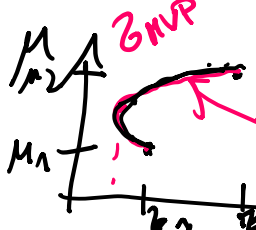
$$\Rightarrow \sigma^2 = \left(\frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}\right)^2 \sigma_1^2 + \left(\frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2}\right)^2 \sigma_2^2 = \frac{\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

investing in a portfolio with a proportion of α in asset 1 you can reduce your risk!

What are the accessible portfolios?

from (2) we obtain $\alpha = \frac{\mu - \mu_2}{\mu_1 - \mu_2}$

insert in (3) $\Rightarrow \sigma^2 = \left(\frac{\mu - \mu_2}{\mu_1 - \mu_2}\right)^2 \sigma_1^2 + \left(1 - \frac{\mu - \mu_2}{\mu_1 - \mu_2}\right)^2 \sigma_2^2 \Rightarrow \mu(\sigma)$



(2) since we have a quadratic dependence we can ask for the portfolio with the parameter α which creates the minimum risk i.e. minimum variance

$$1 = \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

$$\alpha = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}$$

$$f(x) = (1-x)^2$$

$$f'(x) = 2(1-x)(-1)$$

Investment parameter for the Minimum-Variance-Portfolio

$\sigma_{GNVP} = \text{Minimum Variance}$

$$E(R) = \mu = \alpha\mu_1 + (1-\alpha)\mu_2 \quad (2)$$

$$\text{Var}(R) = \sigma^2 = \alpha^2\sigma_1^2 + (1-\alpha)^2\sigma_2^2 + 2\alpha(1-\alpha)\rho\sigma_1\sigma_2 \quad (3)$$

$$x = \alpha z_1 \quad y = (1-\alpha)z_2$$

$$\boxed{z = -1} \Rightarrow z^2 = (\alpha z_1)^2 + ((1-\alpha)z_2)^2 - 2\alpha z_1(1-\alpha)z_2 \stackrel{1}{=} x^2 + y^2 - 2xy = (x-y)^2$$

$$\Rightarrow z^2 = (\alpha z_1 - (1-\alpha)z_2)^2$$

$$\Rightarrow \boxed{z = |\alpha z_1 + \alpha z_2 - z_2|} \Rightarrow \text{can be zero} \Rightarrow \alpha(z_1 + z_2) - z_2 = 0$$

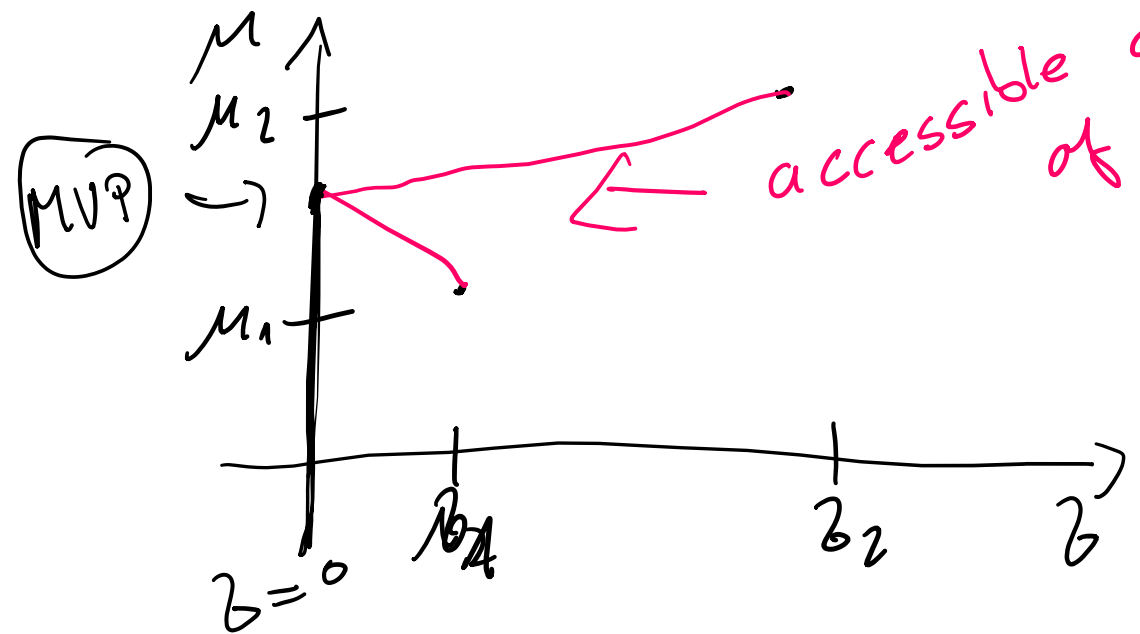
\Rightarrow insert this dependence in μ

$$\Rightarrow \boxed{\alpha^* = \frac{z_2}{z_1 + z_2}}$$

where the portfolio has

\Rightarrow again you obtain
a linear dependence between μ and z

no risk
 $z^2(\alpha^*) = 0$



Accessible portfolios

$$\rho \in (-1, 1)$$

We obtain all accessible portfolios (μ, σ) inserting (2) in (3) solving for σ .

$$\sigma^2 = \frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{(\mu_1 - \mu_2)^2} \left[\mu - \frac{\mu_2\sigma_1^2 + \mu_1\sigma_2^2 - (\mu_1 + \mu_2)\rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} \right]^2 \quad (4)$$

$$= \sigma^2(\mu)$$

$$+ \frac{(1 - \rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} \quad (5)$$

$$\sigma^2 = \text{Var}K^2(R_1, R_2) [\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2 \quad (6)$$

Function for the efficient frontier!

↑
 Difference between the Expected value of the Portfolio $\mu(\rho)$ relative to the MVP

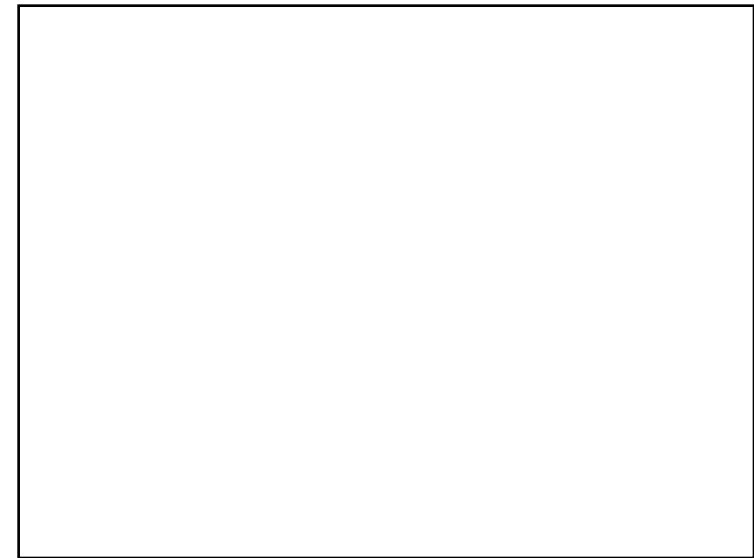


Accessible portfolios

▶ $VarK^2(R_1, R_2) = \frac{Var(R_1.R_2)}{[E(R_1.R_2)]^2} = \frac{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}{(\mu_1 - \mu_2)^2}$ (Coefficient of variation)

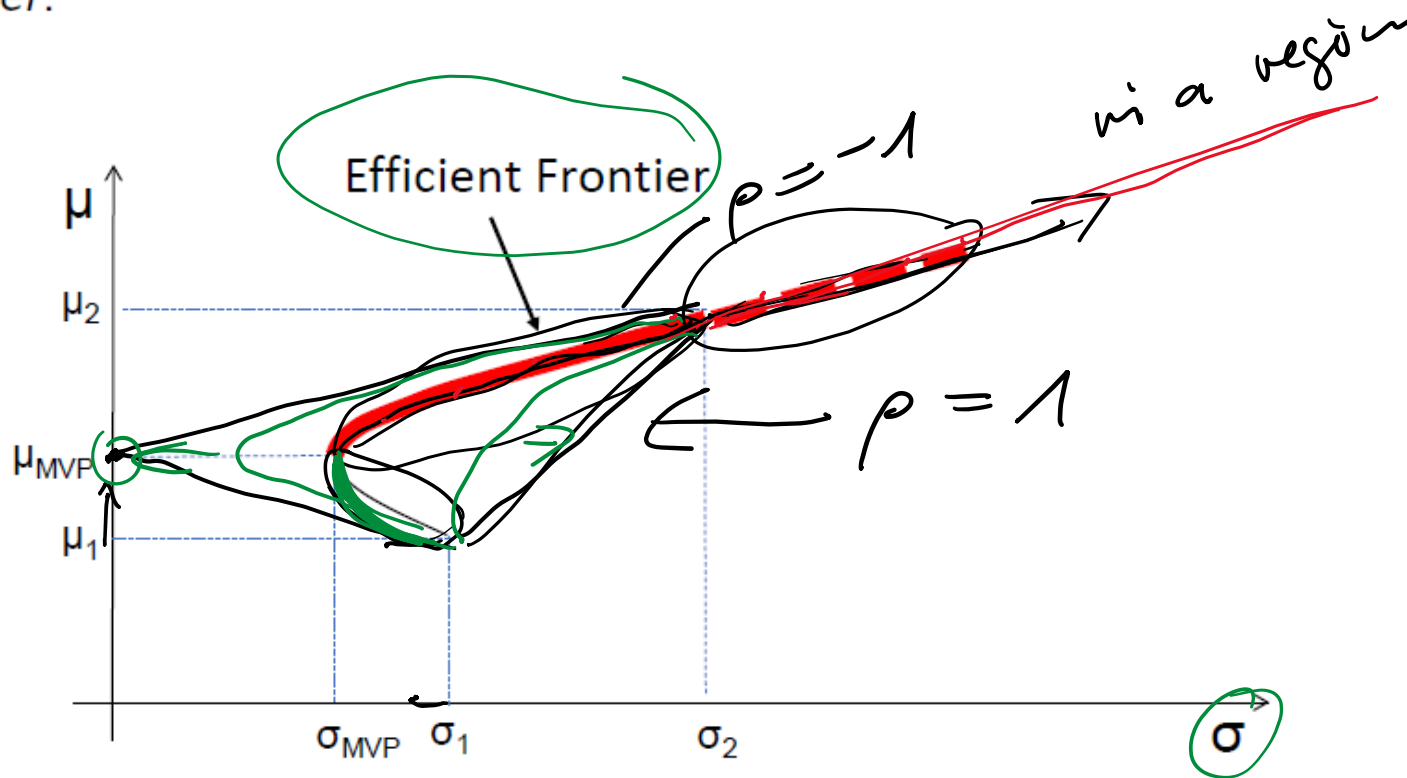
▶ $\mu_{MVP} = \frac{\mu_2\sigma_1^2 + \mu_1\sigma_2^2 - (\mu_1 + \mu_2)\rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$ and $\sigma_{MVP}^2 = \frac{(1 - \rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$

$(\mu_{MVP}, \sigma_{MVP})$: Expected value and standard deviation of the global minimum-variance-portfolio. (6) represents a parabola with the vertex $(\mu_{MVP}, \sigma_{MVP})$. For the extrema $\rho = \pm 1$ we obtain separate solutions.

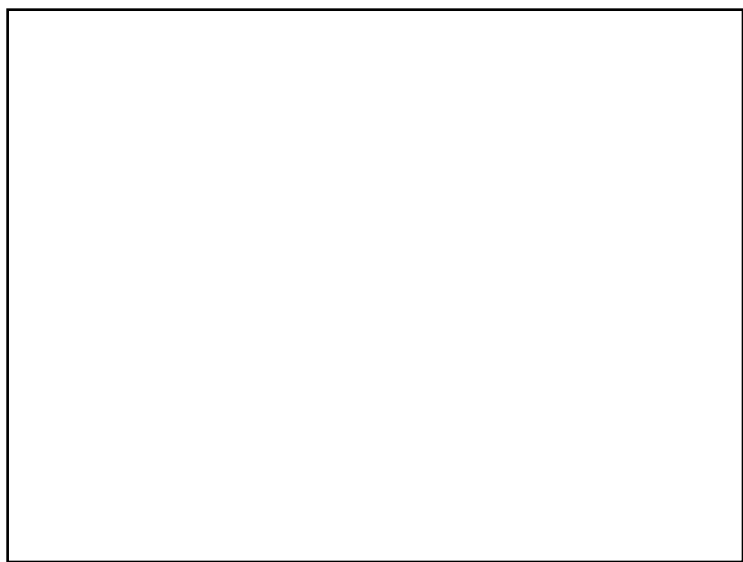


Efficient frontier

Since we assume risk averse investors (higher risk implicates higher yields) only the upper part of the parabola is economically relevant. This is called the *efficient frontier*.



$\alpha \in [0, 1]$
 $\alpha < 0$
 $\alpha > 1$
is also possible
since in financial markets you can short sell an asset
 α could be negative or positive



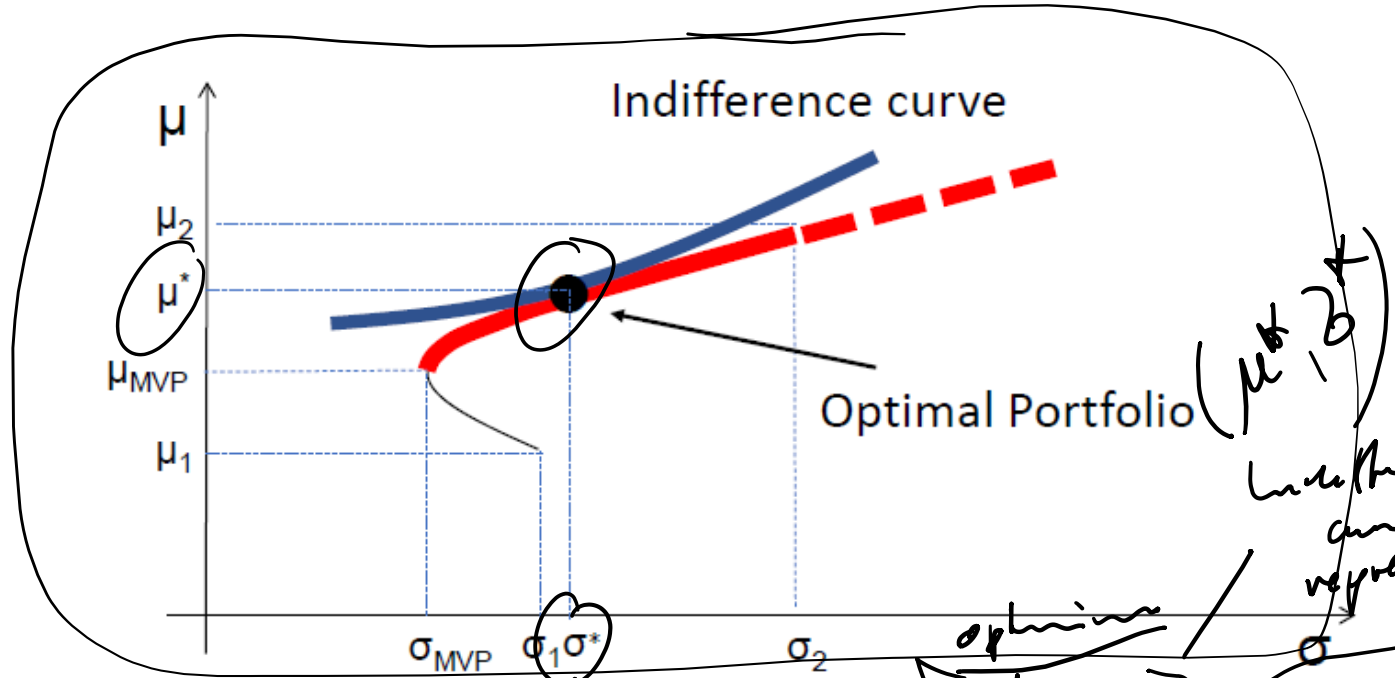
Optimal Portfolio

In general a risk averse utility function is given by $U(\mu, \sigma)$ with $\frac{\partial U}{\partial \mu} > 0$ and $\frac{\partial U}{\partial \sigma} < 0$. The optimizing problem is then given by:

$\max_{\mu, \sigma} U(\mu, \sigma)$ N.B. $\sigma^2 = A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2$ ($A = \text{Var}K^2(R_1, R_2)$) (7)

u good = bad u
 $\begin{matrix} + \\ - \end{matrix}$
 $\begin{matrix} \mu \\ \sigma \end{matrix}$

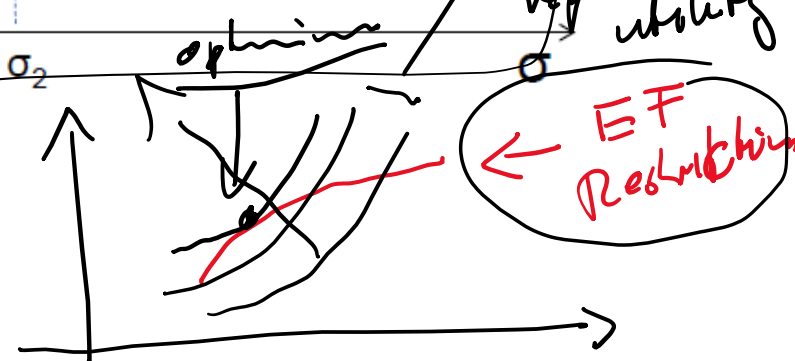
$\max(U(x_1, x_2))$ ($u = \gamma x_1 + \beta x_2$)



Indifference curves with a good μ and a bad σ

Indifference curves representing utility

Restriction is just the efficient frontier



Example

we get $a\mu$ quadratic dependence of σ

micro $u = \sqrt{x_1 x_2}$

$L = \sqrt{x_1 x_2} + \lambda (u - p_1 x_1 - p_2 x_2)$

The explicit utility is given by $U(\mu, \sigma) = \mu - \frac{a}{2}\sigma^2$ with $(a > 0)$. \Rightarrow

$\mathcal{L}(\mu, \sigma, \lambda) = \mu - \frac{a}{2}\sigma^2 + \lambda(\sigma^2 - (A[\mu - \mu_{MVP}] + \sigma_{MVP}^2))$

Utility Restriction

\Rightarrow FOC

$\frac{\partial \mathcal{L}}{\partial \mu}$
 $\frac{\partial \mathcal{L}}{\partial \sigma}$
 $\frac{\partial \mathcal{L}}{\partial \lambda}$

$1 - A\lambda(\mu - \mu_{MVP}) = 0 \Rightarrow 1 = A(\mu - \mu_{MVP}) \Rightarrow \frac{1}{A} = \mu - \mu_{MVP}$

$-a\sigma + \lambda\sigma = 0 \Rightarrow a\sigma = \lambda\sigma \Rightarrow a = \lambda$

$\frac{1}{2}(\sigma^2 - (A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2)) = 0 \Rightarrow$ solve for σ from u or calculate μ

$\mu^* = \mu_{MVP} + \frac{1}{aA}$ $\sigma^* = \sqrt{\frac{1}{Aa^2} + \sigma_{MVP}^2}$ $\alpha^* = \frac{\mu^* - \mu_2}{\mu_1 - \mu_2}$

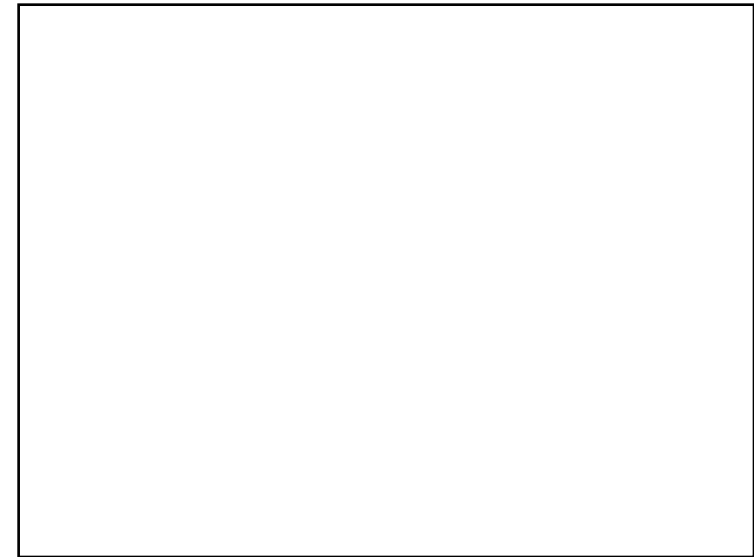
$\mu^* = \alpha \mu_1 + (1-\alpha) \mu_2$

my investment-plan

Portfolio with an additional riskless asset

Assumptions:

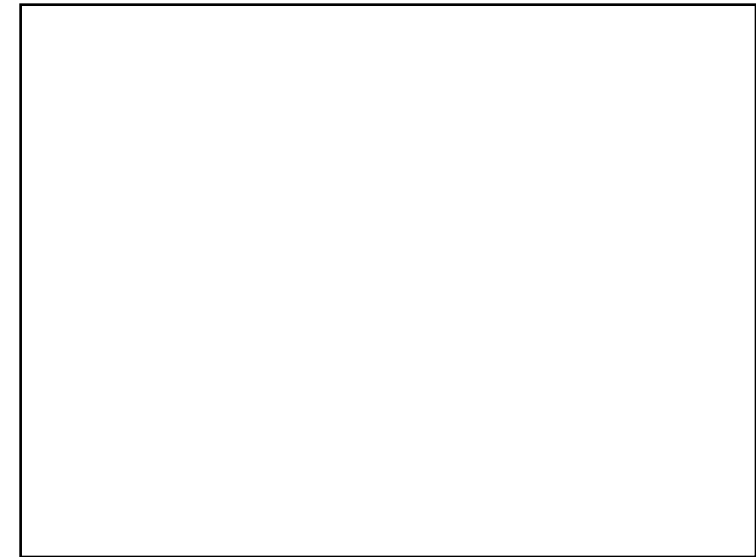
1. It is possible to invest any amount in the riskless asset with yield R_0 .
2. Any investment can be financed via a riskless credit with interest rate R_0 .



Notation

in general non zero standard deviation of the risky investment

- ▶ R_p, R_0 : Yield of the risky portfolio, riskless interest rate
- ▶ $\mu_p = E(R_p), \mu_0 = R_0$: Expected yield of the risky portfolio, expected yield of the riskless investment
- ▶ $\sigma_p = \sigma(R_p), \sigma_0 = 0$: Standard deviation of the risky portfolio and the riskless investment
- ▶ R : yield of the total investment consisting of the risky portfolio and the riskless investment
- ▶ $\mu = \mu(R)$: Expected yield of the total investment
- ▶ $\sigma = \sigma(R)$: Standard deviation of the total investment



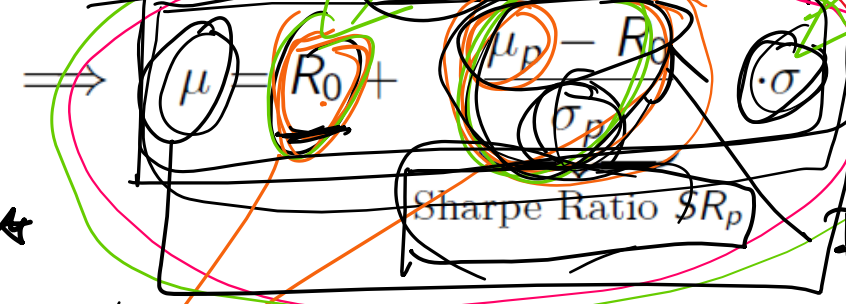
Expected yield and standard deviation of the total investment

risky part
riskless part
put on horizontal axes

$$R = \beta R_p + (1 - \beta) R_0$$

$$\mu(\beta) = E(R) = \mu = \beta \mu_p + (1 - \beta) R_0 = R_0 + \beta(\mu_p - R_0)$$

$$\sigma(\beta) = \text{Var}(R) = \sigma^2 = \beta^2 \sigma_p^2$$



therefore
 β drops as σ

find from data and statistical estimations

are all accessible risk profiles (μ, σ) and riskless investment opportunities
 given investment opportunities
 $\nearrow A_1 \leftarrow$ risky
 $\nearrow A_2$
 $\searrow R_0 \leftarrow$ riskless

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

how does this factor σ^2 from books like numbers

(8)

(9)

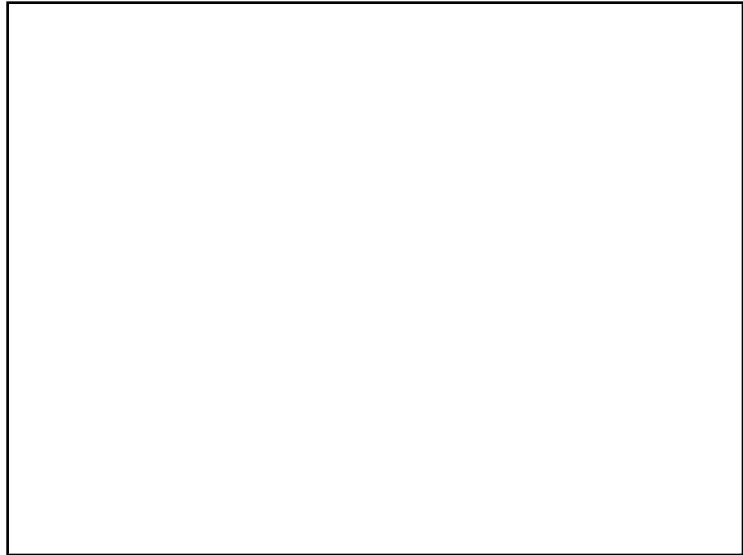
(10)

(11)

$$\mu = A + B\sigma$$

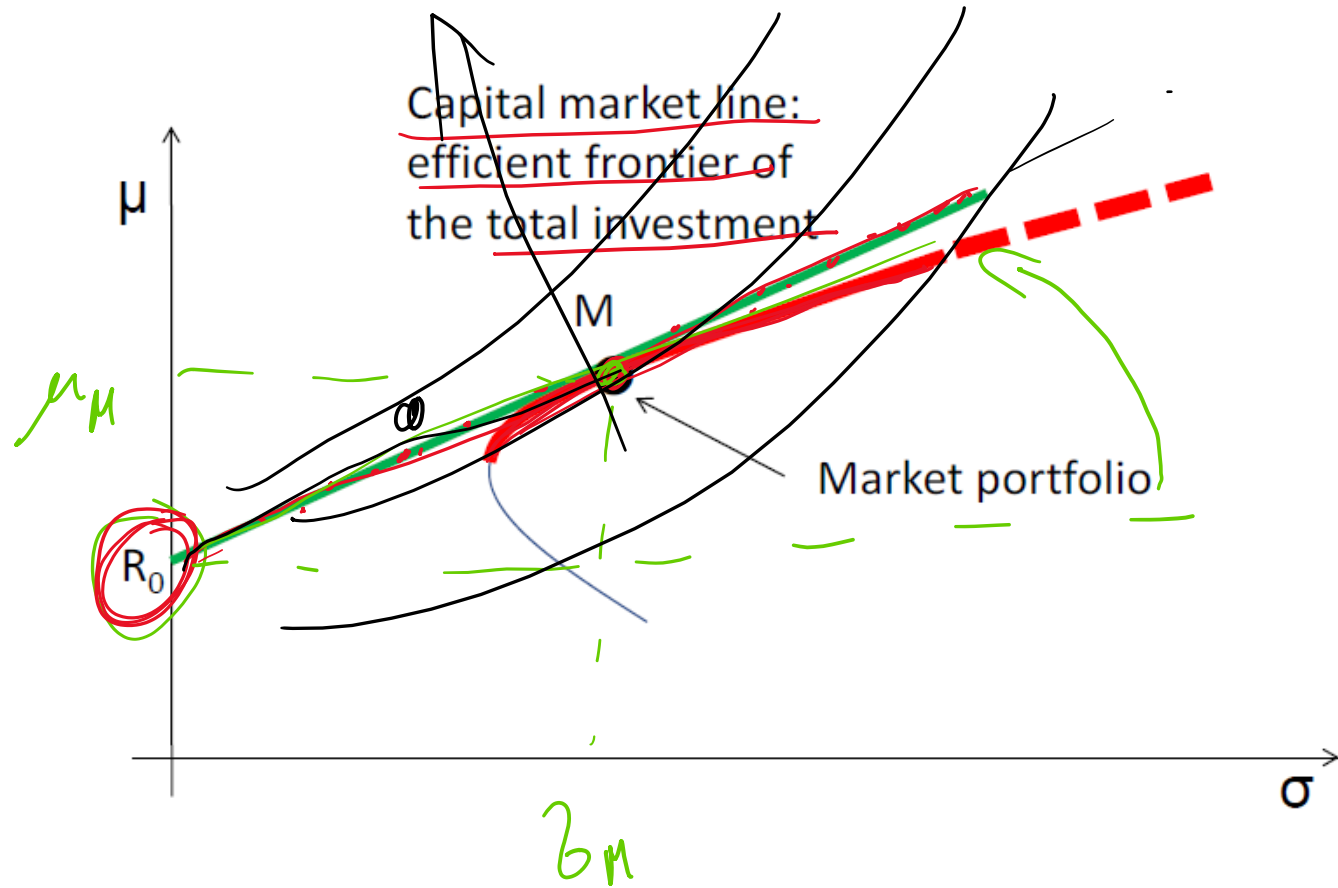
Ex: $y = 3 + 6x$

we have β_0 of the riskless investment is $\beta_0 = 0$



Capital market line

The accessible portfolios of the total investment are obtained via equation (11). The efficient total investments are calculated via the tangent including the point $(0, R_0)$ touching the efficient frontier of the risky portfolio:



Capital market line and the market portfolio

In general we obtain the capital market line via

$$\mu = R_0 + B\sigma \quad \text{oder} \quad \mu = R_0 + \frac{\sigma}{\sqrt{\tilde{B}}} \quad \text{mit} \quad \tilde{B} := \frac{1}{B^2} \quad (12)$$

$$\Rightarrow \sigma^2 = \tilde{B}(\mu - R_0)^2 = g(\tilde{B}, \mu) \quad (13)$$

← solve for \tilde{B}^2

μ must be optimized

For the efficient frontier of the risky portfolio we have:

$$\sigma^2 = A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2 = h(\mu) \quad (14)$$

← Restriktion

The capital market line is (equation (12)) the line with maximum slope B touching the efficient frontier of risky portfolio (equation (14)). Thus, we have to minimize the parameter \tilde{B} subject to the capital market line and efficient frontier are touching each other.



Capital market line and the market portfolio

$\min_{\mu, \tilde{B}} \tilde{B}$ N.B. $g(\tilde{B}, \mu) = h(\mu)$ (15)

$\Rightarrow \mathcal{L}(\tilde{B}, \mu, \lambda) = \tilde{B} + \lambda(\tilde{B}(\mu - R_0)^2 - (A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2))$ (16)

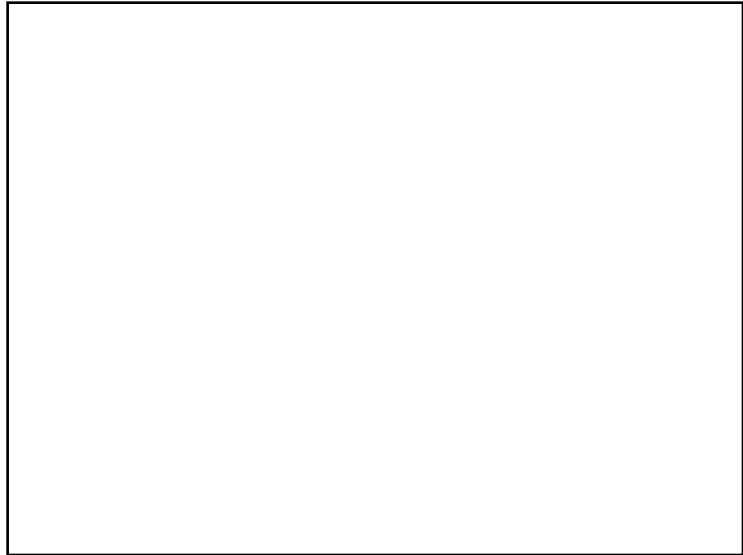
FOC:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{B}} &= 1 + \lambda(\mu - R_0)^2 = 0 \\ \frac{\partial \mathcal{L}}{\partial \mu} &= 2\lambda(\tilde{B}(\mu - R_0) - A(\mu - \mu_{MVP})) = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \tilde{B}(\mu - R_0)^2 - (A[\mu - \mu_{MVP}]^2 + \sigma_{MVP}^2) = 0 \end{aligned}$$

Restriktion

$\lambda = -\frac{1}{(\mu - R_0)^2}$
 this can be solved for μ (\tilde{B})

$\tilde{B}(\mu)$ \leadsto we can solve for \tilde{B}



Capital market line and the market portfolio

Solving FOC, we obtain:

$$\underline{B_{KML}} = \sqrt{\frac{1}{A} + \left[\frac{\mu_{MVP} - R_0}{\sigma_{MVP}} \right]^2} \quad (17)$$

$$\mu_M = \mu_{MVP} + \frac{\sigma_{MVP}^2}{A(\mu_{MVP} - R_0)} \quad (18)$$

$$\sigma_M^2 = \sigma_{MVP}^2 \left[\frac{\sigma_{MVP}^2}{A(\mu_{MVP} - R_0)^2} + 1 \right] \quad (19)$$

with B_{KML} : slope of the capital market line (μ_M, σ_M) : Expected yield and standard deviation of the market portfolio. The market portfolio consists of

$\alpha_M = \frac{\mu_M - \mu_2}{\mu_1 - \mu_2}$ of asset 1 and $1 - \alpha_M = \frac{\mu_M - \mu_1}{\mu_2 - \mu_1}$ of asset 2.

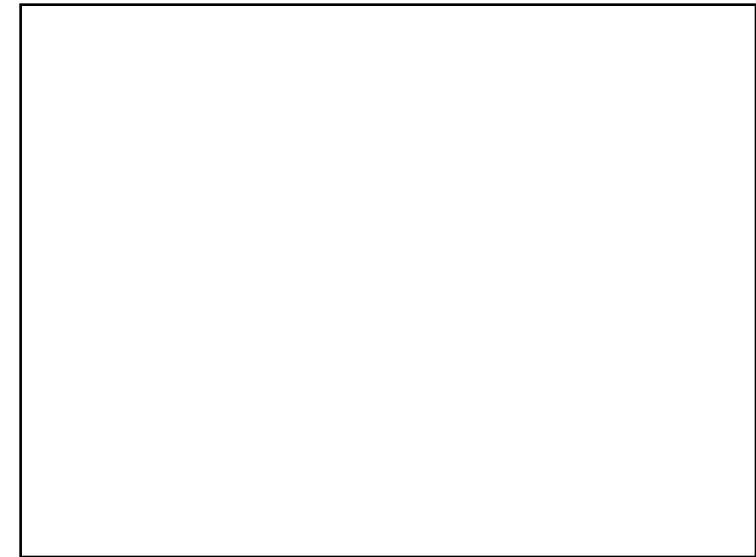
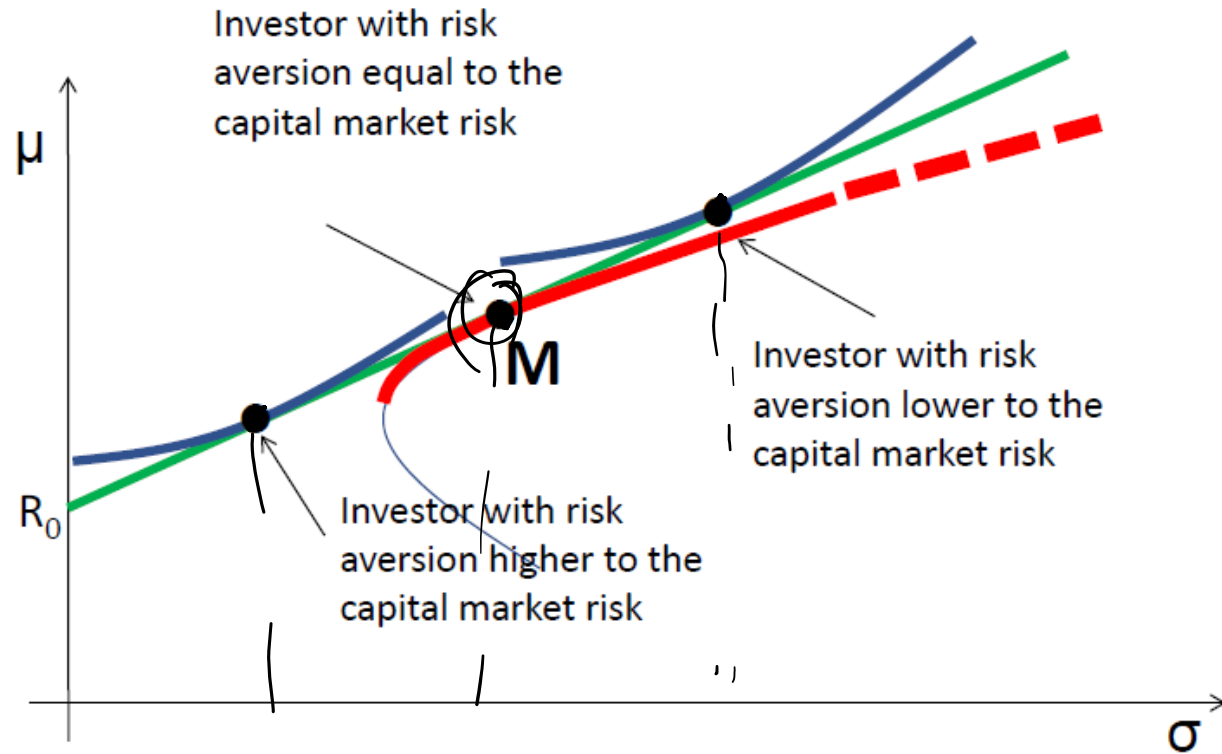
\implies It is possible to replicate every risk profile via investing in the riskless asset and the risky market portfolio, starting with 100% investment in the riskless asset $(\mu, \sigma) = (R_0, 0)$.

*represents in general
the risk of risk
average risk
of the capital
market including
the risky investment
possibilities*

Optimal Investment of the riskless asset

The condition for the optimal investment is now given by the tangent point of capital market line and the indifference curve.

$$\max_{\mu, \sigma} U(\mu, \sigma) \quad \text{N.B. } \mu = R_0 + \sigma B_{KML} \quad (20)$$



Example

Again utility is given by $U(\mu, \sigma) = \mu - \frac{a}{2}\sigma^2$ mit ($a > 0$). \implies

$$\mathcal{L}(\mu, \sigma, \lambda) = \mu - \frac{a}{2}\sigma^2 + \lambda(R_0 + \sigma B_{KML} - \mu)$$

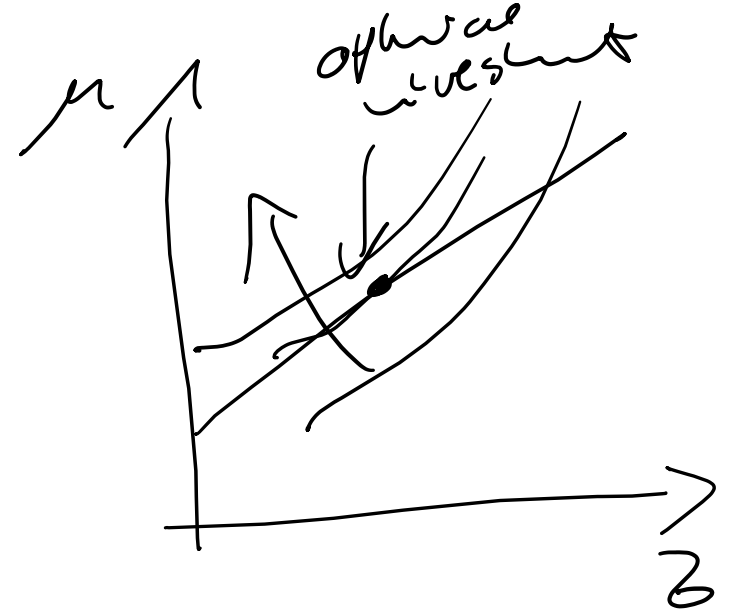
\implies FOC

$$\frac{\partial \mathcal{L}}{\partial \mu} = 1 - \lambda = 0$$

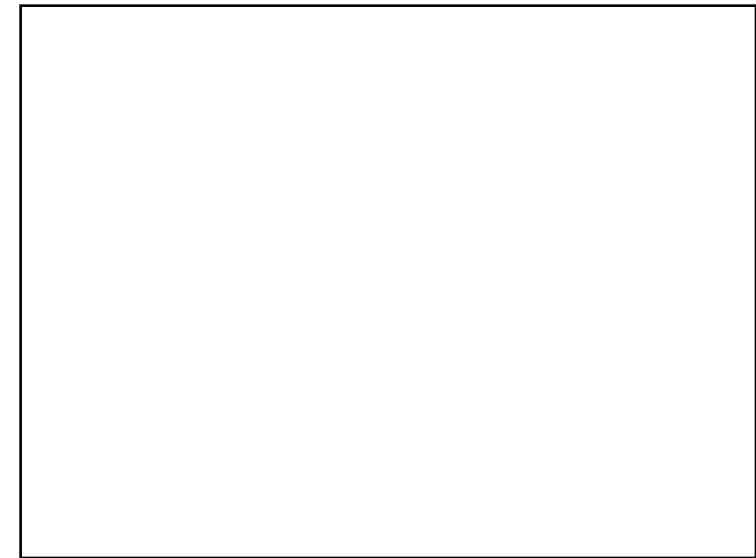
$$\frac{\partial \mathcal{L}}{\partial \sigma} = -a\sigma + \lambda B_{KML} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = R_0 + \sigma B_{KML} - \mu = 0$$

$$\implies \mu^* = R_0 + \frac{B_{KML}^2}{a} \quad \sigma^* = \frac{B_{KML}}{a} \quad \beta^* = \frac{\mu^* - R_0}{\mu_M - R_0}$$



here work



CAPM and linear regression

Over or under performance of security i in relation to the market portfolio M is given by (due to the assumptions of the Markowitz model):

$$\mu_i - R_0 = \frac{\text{COV}(R_i, R_M)}{\sigma_M^2} (\mu_M - R_0) \quad \text{CAPM-equation}$$

β_i

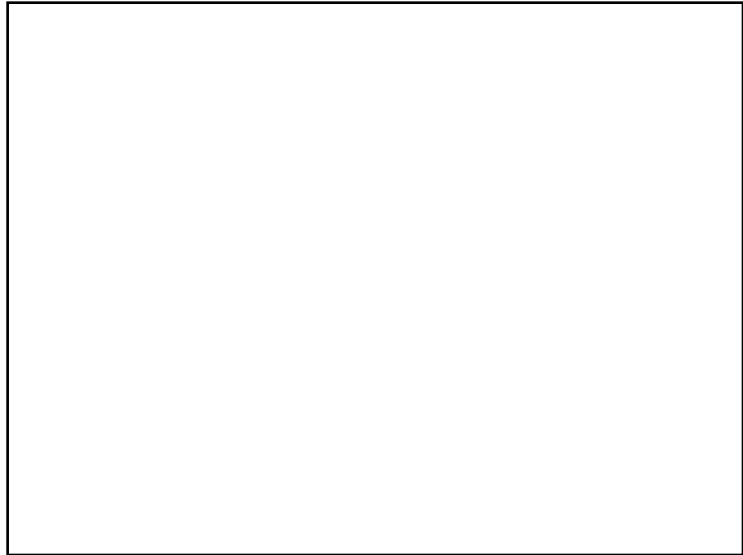
famous β -factors

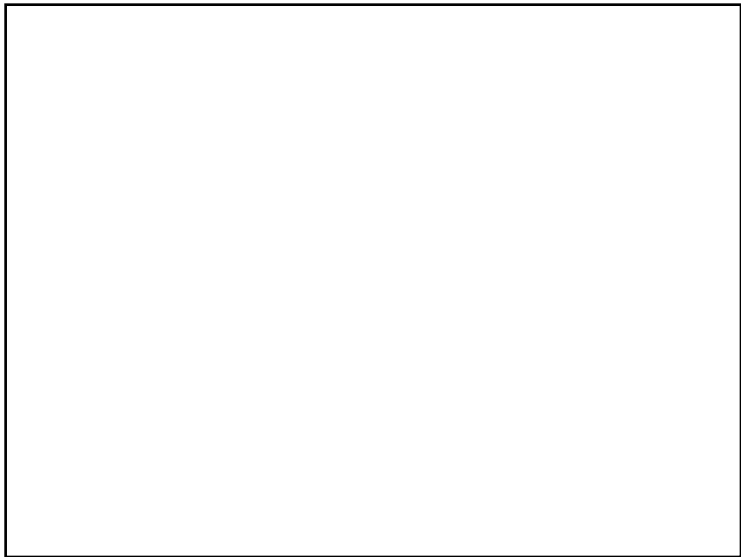
In analogy to the capital market line this condition related to a specific security is named security market line. *Compare with the capital market line*

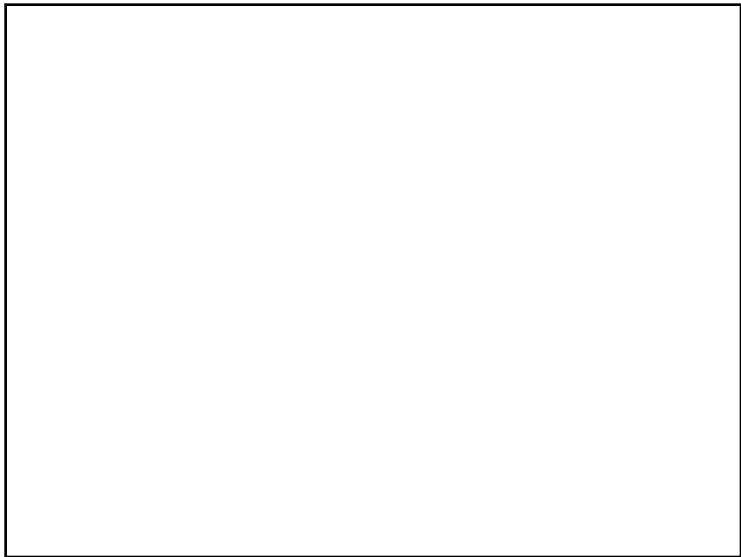
$\mu_i - R_0$ is the risk premium of security i related to the riskless interest rate R_0 equaling β times the risk premium of the market portfolio related to R_0 . Therefore β can be interpreted as the individual risk of security i relative to the market risk.

all DAX securities
DAX itself

estimation via linear regression
→ google finance } → direct
 yahoo finance } puzii
 ↳ into excel





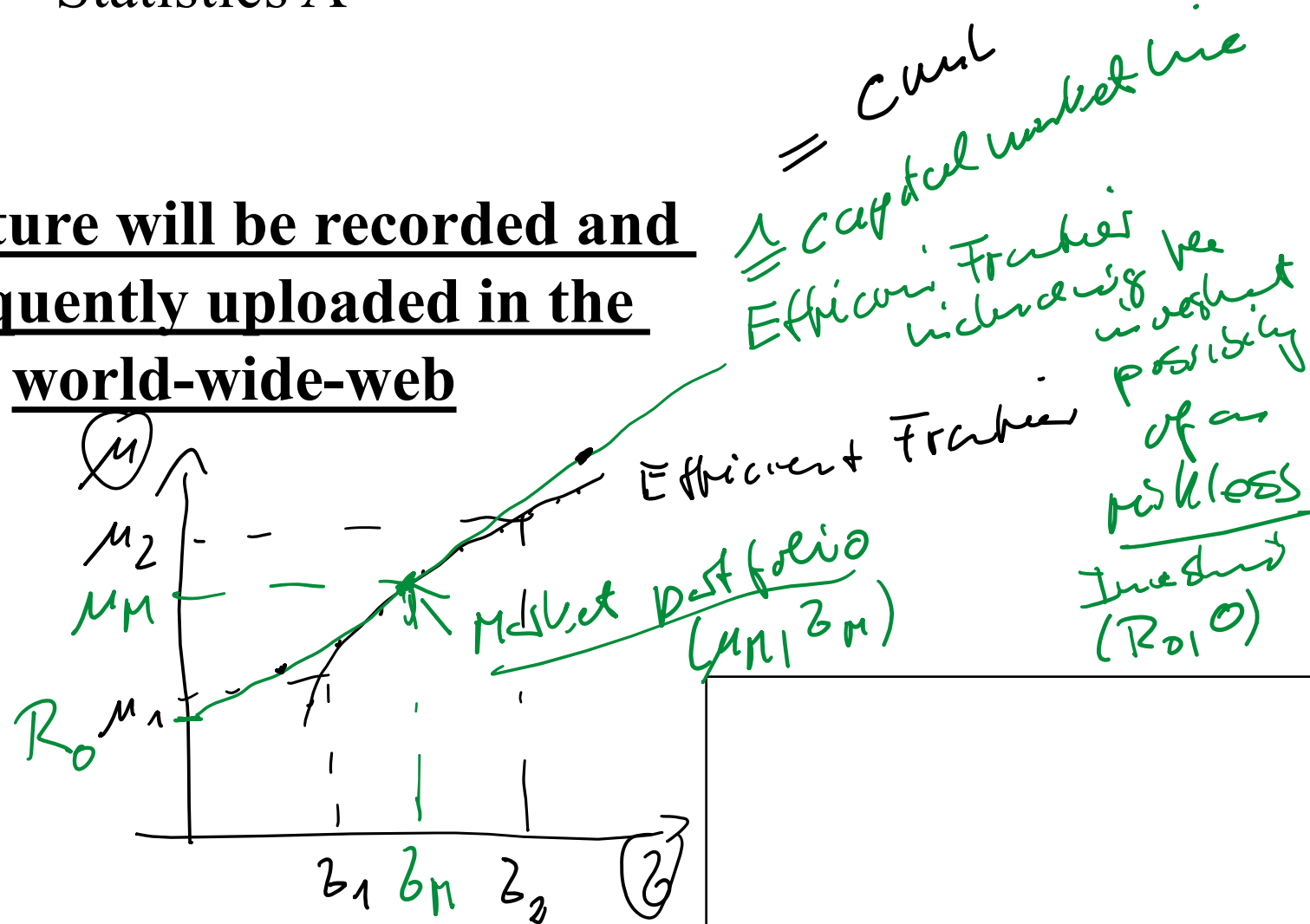


Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web



[Function translator \(webpage\)](#)

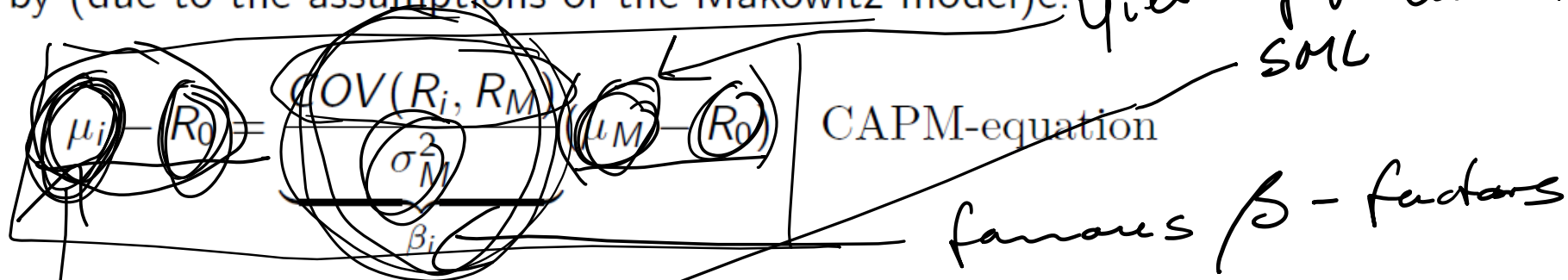
[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

CAPM and linear regression → Capital Asset Pricing - Modell

Over or under performance of security i in relation to the market portfolio M ⇒ Equilibrium-Modell
 is given by (due to the assumptions of the Markowitz model): yield of the market portfolio



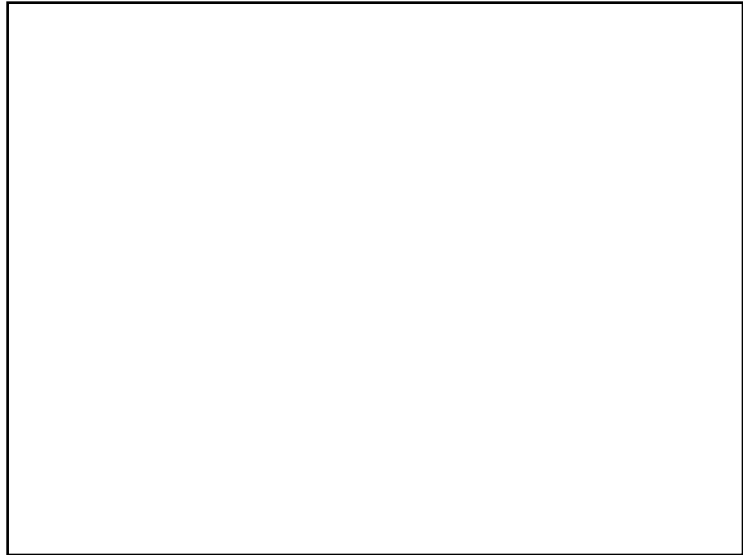
In analogy to the capital market line this condition related to a specific security is named security market line. Compare with the capital market line

$\mu_i - R_0$ is the risk premium of security i related to the riskless interest rate R_0 equaling β times the risk premium of the market portfolio related to R_0 . Therefore β can be interpreted as the individual risk of security i relative to the market risk.

all DAX securities
 DAX itself

estimation via linear regression

→ google finance } → direct
 yahoo finance } praxis
 into excel



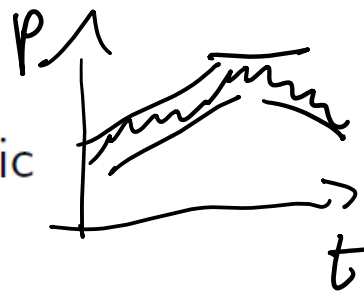
CAPM and linear regression

Over or under performance of security i in relation to the market portfolio M is given by (due to the assumptions of the Markowitz model):

$$\mu_i - R_0 = \underbrace{\frac{\text{COV}(R_i, R_M)}{\sigma_M^2}}_{\beta_i} (\mu_M - R_0) \quad \text{CAPM-equation}$$

$$\mu_M \hat{=} 1$$

large stock-market index
 S&P 500
 MSX-Index
DAX
 Nikkei
 Dow Jones



In analogy to the capital market line this condition related to a specific security is named security market line.

$\mu_i - R_0$ is the risk premium of security i related to the riskless interest rate R_0 equaling β times the risk premium of the market portfolio related to R_0 .

Therefore β can be interpreted as the individual risk of security i relative to the market risk.

What data should we take in order to estimate the β -factors?

- $\mu_i \hat{=} \text{historical yields of single securities} \leftarrow \text{individual DAX securities}$
- $R_0 \hat{=} \text{interest rates of government's bonds or we can take the reference interest rates of the financial markets}$
- Libor oder Euribar \leftarrow 3-month



Which Period and which frequency should we take

Why not taking ^{daily} data which is accessible?

→ but Daily can be in general highly correlated

→ the general calculated risk of a surge security is mainly driven by a short subperiod

Why not taking yearly data?

→ but the die sample size goes down

→ we should increase the whole period, we are looking at

→ i.e. we got back the last 30 years

→ dot com ~ 2001

• financial crisis 2008/09

• Corona crisis 2020

→ but this

means

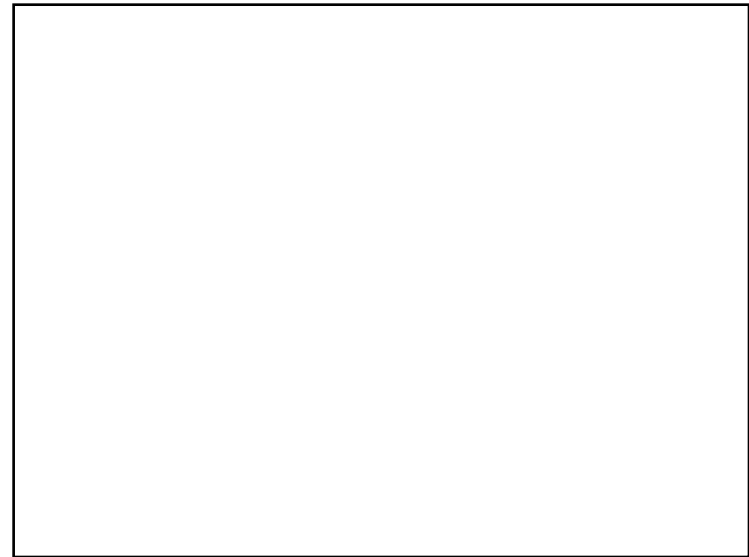
that the

founder

of financial

markets was

highly changed in this period



→ measurable

choice of

frequency

Period length

5-10 years

Where do I get the data

→ -DAX

<https://de.finance.yahoo.com/quote/ADS.DE/history?period1=1262304000&period2=1653264000&interval=1mo&filter=history&frequency=1mo&includeAdjustedClose=true>

-DAX-tiles

-Euribor

<https://sdw.ecb.europa.eu/>
<https://fred.stlouisfed.org/>
you have also a excel add-in

https://www.bundesbank.de/dynamic/action/de/statistiken/zeitreihen-datenbanken/zeitreihen-datenbank/759778/759778?listId=www_sgeldmkt_mb03_neu

<https://query1.finance.yahoo.com/v7/finance/download/NFLX?period1=1022112000&period2=1589241600&interval=1d&events=history>

Advis 1, Mai 160
1. April 150

$$R_i = \frac{160}{150} - 1$$

p.a.
=> $\frac{\text{Euribor}}{12}$

Dax-Performance

$$\rightarrow \Delta R_i = R_{it} - R_{ot} = \beta (R_{Mt} - R_{ot})$$

=> linear regression with no absolute Parameter

$$y = a + bx$$

=> a = 0 in our case

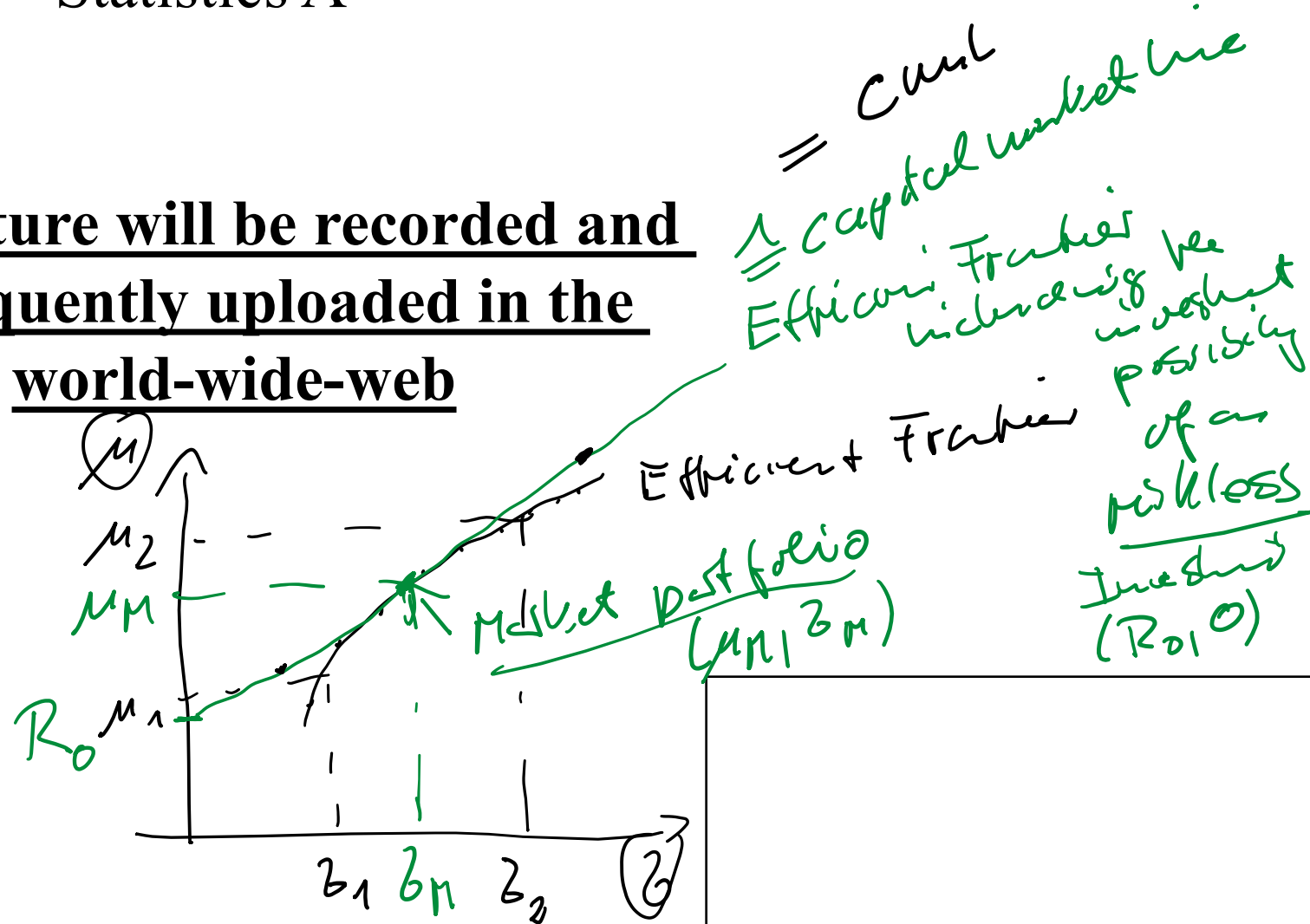
$$\sim R_{it} - R_{ot} = \beta (R_{Mt} - R_{ot})$$

Statistics A

Wilhelmshaven



This lecture will be recorded and
Subsequently uploaded in the
world-wide-web



[Function translator \(webpage\)](#)

[Function translator Excel 1 \(add in\)](#)

Prof. Dr. Bernhard Köster
Jade-Hochschule Wilhelmshaven

<http://www.bernhardkoester.de/vorlesungen/inhalt.html>

CAPM and linear regression

Over or under performance of security i in relation to the market portfolio M is given by (due to the assumptions of the Markowitz model):

$$\mu_i - R_0 = \frac{\text{COV}(R_i, R_M)}{\sigma_M^2} (\mu_M - R_0) + \alpha$$

CAPM-equation

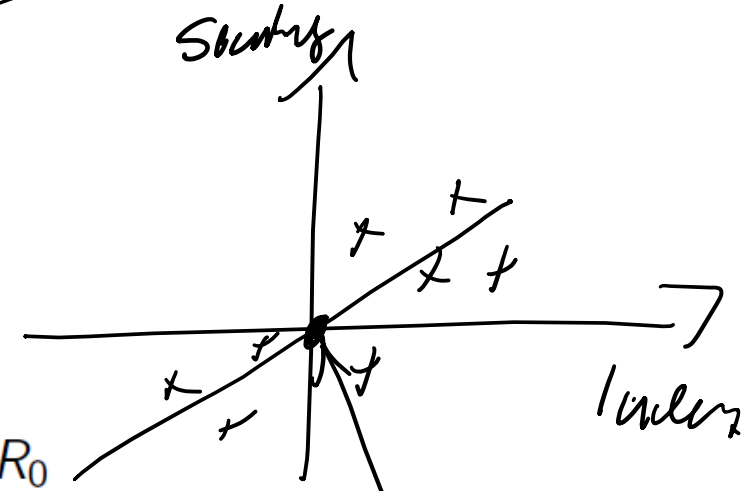
$\alpha = 0$

Handwritten annotations:
 - $\mu_i - R_0$ is labeled R_i and "gross Interest".
 - $\text{COV}(R_i, R_M)$ is labeled β_i .
 - σ_M^2 is labeled σ_M^2 .
 - $\mu_M - R_0$ is labeled "DAX" and "gross interest rate".

SML
 CML (Markowitz)

In analogy to the capital market line this condition related to a specific security is named security market line.

$\mu_i - R_0$ is the risk premium of security i related to the riskless interest rate R_0 equaling β times the risk premium of the market portfolio related to R_0 . Therefore β can be interpreted as the individual risk of security i relative to the market risk.



Handwritten note:
 linear regression -
 line has to pass
 through the origin

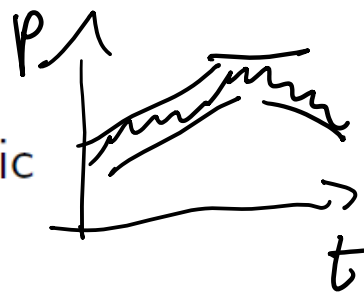
CAPM and linear regression

Over or under performance of security i in relation to the market portfolio M is given by (due to the assumptions of the Markowitz model):

$$\mu_i - R_0 = \underbrace{\frac{\text{COV}(R_i, R_M)}{\sigma_M^2}}_{\beta_i} (\mu_M - R_0) \quad \text{CAPM-equation}$$

$$\mu_M \hat{=} 1$$

large stock-market index
 S&P 500
 MSI-wake
DAX
 Nikkei
 Dow Jones



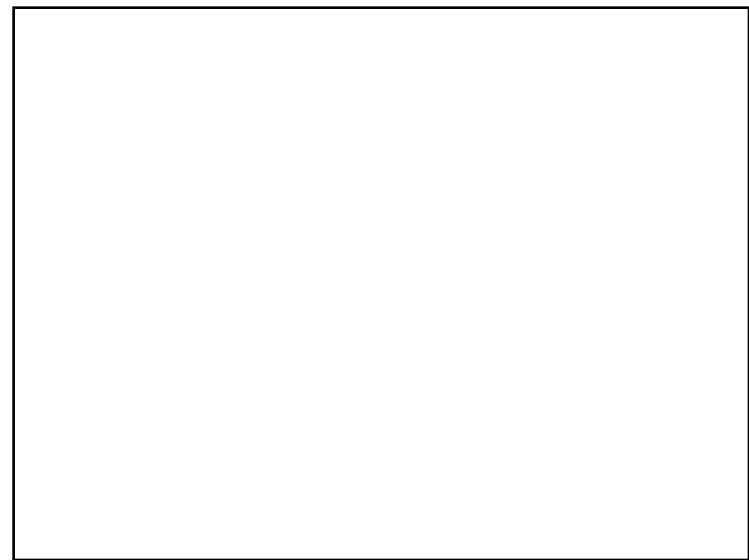
In analogy to the capital market line this condition related to a specific security is named security market line.

$\mu_i - R_0$ is the risk premium of security i related to the riskless interest rate R_0 equaling β times the risk premium of the market portfolio related to R_0 .

Therefore β can be interpreted as the individual risk of security i relative to the market risk.

What data should we take in order to estimate the β -factors?

- $\mu_i \hat{=} \text{historical yields of single securities} \leftarrow \text{individual DAX securities}$
- $R_0 \hat{=} \text{interest rates of government's bonds or we can take the reference interest rates of the financial markets}$
- Libor oder Euribar \leftarrow 3-month



Which Period and which frequency should we take

Why not taking ^{daily} data which is accessible?

→ but Daily can be in general highly correlated

→ the general calculated risk of a surge security is mainly driven by a short subperiod

Why not taking yearly data?

→ but the die sample size goes down

→ we should increase the whole period, we are looking at

→ i.e. we got back the last 30 years

→ dot com ~ 2001

• financial crisis 2008/09

• Corona crisis 2020

→ but this

means

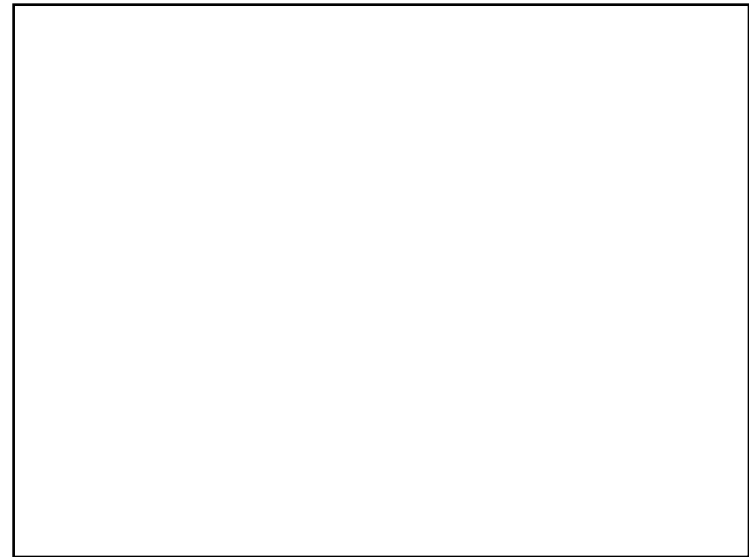
that the

founder

of financial

markets was

highly changed in this period



→ measurable

choice of

frequency

Period length

5 - 10 years

Where do I get the data

→ -DAX

<https://de.finance.yahoo.com/quote/ADS.DE/history?period1=1262304000&period2=1653264000&interval=1mo&filter=history&frequency=1mo&includeAdjustedClose=true>

-DAX-tiles

-Euribor

<https://sdw.ecb.europa.eu/>
<https://fred.stlouisfed.org/>
you have also a excel add-in

https://www.bundesbank.de/dynamic/action/de/statistiken/zeitreihen-datenbanken/zeitreihen-datenbank/759778/759778?listId=www_sgeldmkt_mb03_neu

<https://query1.finance.yahoo.com/v7/finance/download/NFLX?period1=1022112000&period2=1589241600&interval=1d&events=history>

Advis 1, Mai 160
1. April 150

$$R_i = \frac{160}{150} - 1$$

p.a.
=> $\frac{\text{Euribor}}{12}$
Dax-Reference

→ $\Delta R_i = R_{it} - R_{ot} = \beta (R_{Mt} - R_{ot})$
=> linear regression with no absolute Parameter
 $y = a + bx$
~) $R_{it} + R_{ot} = (\beta)(R_{Mt} - R_{ot}) + R_{ot}$

