
Statistics & Econometrics

for CS|DS@UCU

University of Augsburg
Chair of Statistics
Prof. Dr. Yarema Okhrin

Introduction

Statistics deals with the analysis of processes, which are driven by random factors. For this purpose we collect data on the process. Numerous methods and tools were developed to help us to collect, describe, analyze and draw conclusions from the data (observations).

Wikipedia: Statistics is a mathematical science pertaining to the collection, analysis, interpretation or explanation, and presentation of data.

Examples: number of clicks on a ad-banner, number of orders of a particular product, price of a particular financial assets, number and size of insurance claims, creditability of a particular company, customer churn, etc

Econometrics deals with the modelling of causal dependence between one or several dependent variable and set a set of explanatory variables. Thus it aims to “explain” the relationship. Special tools for forecasting, modelling specific types of data and specific functional relationships.

Examples: impact of expenditures on ad campaigns, training for employees, quality assurance, research, etc. on the sales/profit

Time series analysis deals with modelling and forecasting of time ordered data.

Examples: modelling the dynamics of sales, asset prices, website traffic

Bulding blocks of Statistics

Descriptive Statistics

- Presentation of data using tables and graphs
- Characterizing the data using a few but powerful measures

Probability Theory

- The concept of probability, conditional probability
- random variables, distribution and density function, characterization of RV's

Inferential Statistics

- Inference about the population on the basis of a sample
- Testing statistical hypothesis, building confidence intervals, measuring reliability of tests

Additional advanced components

- Theory of point estimation
- Nonparametric statistics
- Large sample theory
- Bayesian statistics

Chapter 1

Descriptive Statistics

Descriptive Statistics

Basic statistical concepts

- real world problem \rightsquigarrow statistical analysis
- The complete set of the objects, which are subject of the analysis is called **population** and is usually denote by Ω . We denote the elements of Ω by ω .
- **Note:** we are not interested in the population itself, but more in the properties of the population measured by one or several quantities of interest X (**characteristics/attributes**).

$X: \Omega \rightarrow S$, where S is the space of possible values of X .

$x = X(\omega)$ is called a **realization** or an **observation**.

Example: public appeal of a new movie

Ω = the set of all audience members,

X = (assessment of the movie, age, gender, occupation)

	assessment	age	gender	occupation
1	good	23	m	student
2	very good	14	m	pupil
3	good	19	f	shop assistant
4	satisfactory	35	m	worker
5	adequate	29	f	school teacher

Data sampling

- **complete survey**: we collect and analyze all elements of Ω (for example, population census).

Disadvantage: too expensive, too costly, not always feasible in practice (for example, life expectancy of bulbs)

- **partial survey**: we collect only a small part of the elements of the population.
- The set of the considered elements is called **sample**.

Classification of variables I

- **nominal scale:**

Let x and y denote two realizations of an attribute. If the attribute is nominal, then we can only conclude that either

$$x = y \text{ (equality) or } x \neq y \text{ (inequality)}$$

Example: marital status, gender, occupation

- **ordinal scale:** the realizations can be naturally ordered, i.e. statements with „smaller/less “ and „larger/more “ have clear interpretation. This implies that for all realizations x and y

$$x = y \quad \text{or} \quad x > y \quad \text{or} \quad x < y.$$

Examples: grades, rankings

Classification of variables II

- **interval scale**: if the differences between two realizations of an ordinaly scaled attribute has natural meaning.

Example: temperature values in Celsius, year of birth

- **ratio scale**: additionally to definition of the interval scale we require that there is a meaningful non-arbitrary zero in the set of realizations.

Examples: income, price, turnover, age

- **absolute scale**: in addition to the interval scale we have a natural, scale-independent unit.

Examples: quantity, number of students enrolled at a university

Classification of variables III

- An attribute is called **qualitative** , if it has a finite set of possible realizations and is at most ordinally scaled. The realizations reflect the difference/strength, but not the magnitude (e.g. gender, colour).
- If, however, the realizations reflect both the difference and the magnitude, then we speak about **quantitative attributes** (for example, age, income, price).
- We observe an increasing informational content by moving from nominal to interval scale, but the observations may suffer from assessment errors.

Classification of variables IV

A variable/attribute is **discrete**, if the set of possible realizations is a countable set. The attribute/variable is **continuous**, if it has uncountably many possible realizations.

Examples: height, speed, time, grade, quality

Note:

- Despite of the fact that many variables are continuous by nature, it is **not** possible to measure them with an arbitrary precision.
- Often a discrete attribute has very many realizations (for example, prices, income). In this case it is reasonable to treat them as continuous attributes.

Characteristics of univariate data sets

Starting point: the quantity of interest X

- the sample x_1, \dots, x_n with $x_i \in \mathbb{R}$ (univariate);
- let a_1, \dots, a_k denote all possible but different realizations

absolute frequency of a_i :

$n(a_i)$ = frequency of the occurrence of the realization a_i in the sample

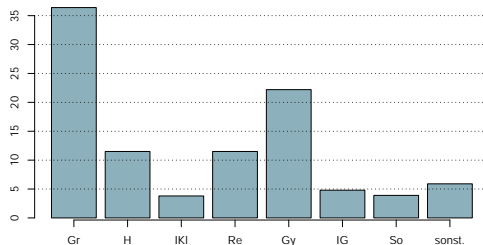
relative frequency of a_i : $h(a_i) = n(a_i)/n$

Graphical presentation of the frequencies I

bar plot: for each realization plot bars/sticks. The height of the bars equals the absolute OR relative frequency.

Example:

Out of 9 558 455 pupils in Germany (in 1993) 36.4% went to elementary school, 11.5% to secondary modern, 3.8% to integrated secondary and junior high school, 11.5% to junior high school, 22.2% to “Gymnasium”, 4.8% to integrated school, 3.9% to special schools and 5.9% to other types of schools.



Empirical cumulative distribution function

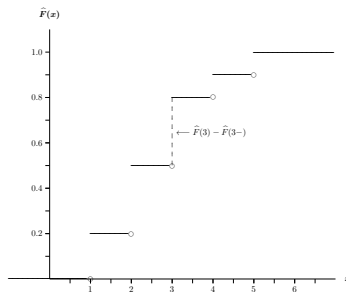
Requirement: at least the ordinal scale

empirical cumulative distribution function (ECDF):

$\hat{F}(x)$ = relative number of observations equal to or less than x

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

Example: public appeal of a movie (grades: 1, 1, 2, 2, 2, 3, 3, 3, 4, 5)



Properties of the ECDF:

- a) $\hat{F}(x) = 0$ for $x < x_{(1)}$, $\hat{F}(x) = 1$ for $x \geq x_{(n)}$
- b) $\hat{F}(x)$ is increasing
- c) $\hat{F}(x)$ is continuous from the right
- d) $\hat{F}(x_j) - \hat{F}(x_j-) = \text{relative frequency of } x_j$

Note: The ECDF contains all the information about the sample in an aggregated form.

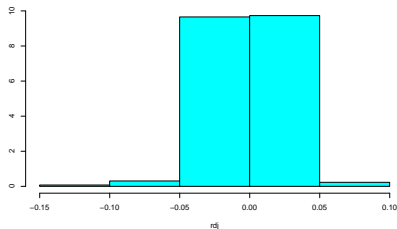
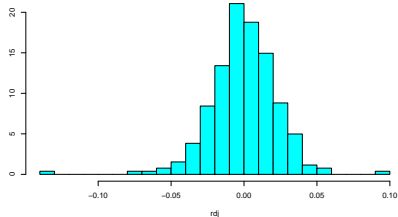
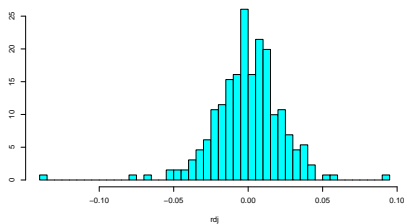
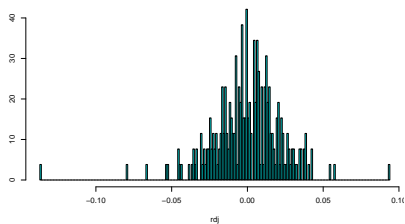
Histogram

- (a) Let $K_j : [x_0 + (j - 1)h, x_0 + jh)$, $j \in \mathbb{Z}$ be the classes of possible values with starting point y_0 and bandwidth h ;
- (b) count the observations in each K_j (class frequency $n(K_j)$);
- (c) calculate the relative class frequency $h(K_j) = n(K_j)/n$, where T is the sample size;
- (d) normalise to 1: $f_j = \frac{n(K_j)}{nh}$ (relative class frequency divided by h);
- (e) plot rectangles of height f_j for each class K_j .

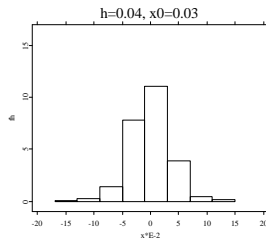
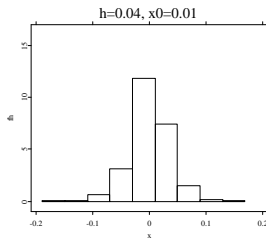
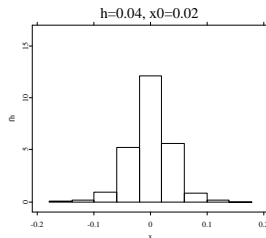
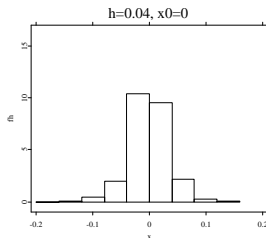
Histogram

$$\hat{f}_h(x) = h(K_j)/h \quad \text{for} \quad x \in K_j$$

Here: Dow Jones index returns with the bandwidth
 $h = 0.001, 0.005, 0.01, 0.05$



Four histograms for the same data with different starting points:
 $x_0 = 0$, $x_0 = 0.01$, $x_0 = 0.02$, $x_0 = 0.03$; bandwidth $h = 0.04$



conditions on the classes:

- disjunct classes
- each realization falls in one of the classes
- **desirable:** all classes have equal width
- the square above the class K_i : $h(K_i)/|K_i| \cdot |K_i| = h(K_i)$,
i. e. the key information about the histogram is revealed by the squares of the rectangles!

•

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \sum_{i=1}^k h(K_i) = 1$$

- special method are required to determine the “best” bandwidth

Characteristics/Parameters

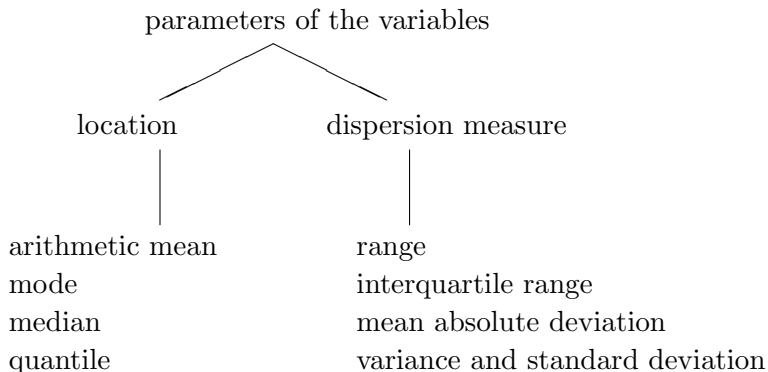
Parameters are measures, that quantify important characteristics of the empirical distribution function.

Important parameters are e.g.:

Location parameter: Gives insights into the central tendency of the the data.

Dispersion measure: Contains information about the variability of the data.

Overview



Location measure

Mean characterizes the average location of the data.

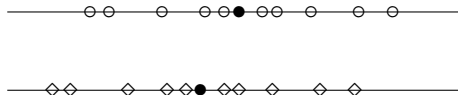
Example: Monthly personal income of elves and orcs in €

Elves: 1000, 1200, 1750, 2200, 2400, 2800, 2950, 3300, 3800, 4150 (◊)

$\bar{x}_{elf} = 2555 \text{ €}$ (●)

Orcs: 600, 800, 1350, 1800, 2000, 2400, 2550, 2900, 3400, 3750 (◊)

$\bar{x}_{orc} = 2155 \text{ €}$ (●)



i) Mean (arithmetic mean, average)

Mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^k n(a_i) a_i = \sum_{i=1}^k h(a_i) a_i$$

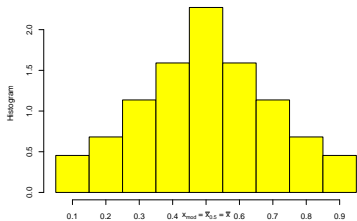
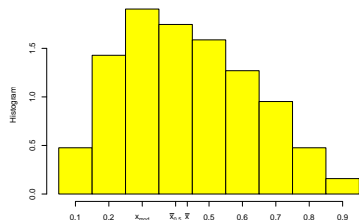
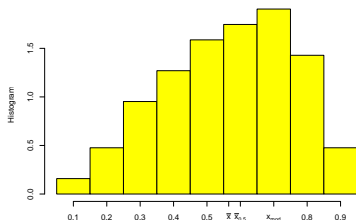
Properties:

- The mean is the value with the smallest possible mean-squared deviation, i. e. it holds for all $a \in \mathbb{R}$

$$\sum_{i=1}^n (x_i - \bar{x})^2 \leq \sum_{i=1}^n (x_i - a)^2.$$

- The mean is very sensitive to **outliers** (for example, monthly income of 1000.0, 1000.0, 1000.0, 10000.0 returns $\bar{x} = 3250$).
- **Note:** the mean is meaningful **only** for symmetric data. Otherwise it is difficult to draw conclusions.

Symmetric and nonsymmetric distributions



ii) α -trimmed mean \bar{x}_α

$x_{(i)}$ is the i -th order statistics, if $x_{(i)}$ is on the i -th position in the ordered sample.

 α -trimmed mean

$$\bar{x}_\alpha = \frac{1}{n - 2[n\alpha]} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} x_{(i)}$$

with $\alpha \in [0, 0.5)$, $[z]$ denotes the largest natural number, which is smaller than z

Example: grades 2.7, 3.0, 3.0, 3.0, 3.3, 3.3, 3.3, 3.7, 4.0, 6.0

It holds that $\bar{x} = 3.53$, but $\bar{x}_{0.1} = 26.6/8 = 3.325$.

Note: it is much more robust to outliers compared to the simple mean

iii) p -quantile \tilde{x}_p p -quantile

$$\tilde{x}_p = \begin{cases} x_{([np]+1)} & \text{for } np \notin \mathbb{Z} \\ (x_{(np)} + x_{(np+1)}) / 2 & \text{for } np \in \mathbb{Z} \end{cases}, \quad p \in (0, 1]$$

$\tilde{x}_{0.25}$ is called **the lower quartile**, $\tilde{x}_{0.5}$ is the **median** and $\tilde{x}_{0.75}$ is the **the upper quartile**

- The arithmetic mean is not robust to outliers.
- The median is, however, **robust**, as it is determined by the ranks of the observations and not by the exact values.

Sample quantiles correspond to $\hat{F}^{-1}(p)$ (in some sense)

Properties:

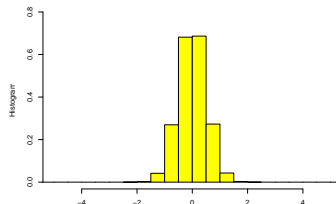
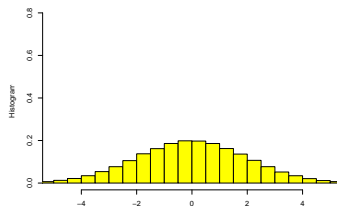
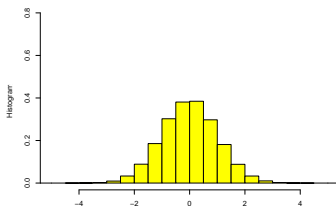
- The number of observations, which are smaller than \tilde{x}_p or equal to \tilde{x}_p , is larger or equal to $[np]$.
- It holds that $x_{([np])} \leq \tilde{x}_p \leq x_{([np]+1)}$.
- It holds for $a \in \mathbb{R}$ that

$$\sum_{i=1}^n |x_i - med| \leq \sum_{i=1}^n |x_i - a|$$

i.e. the median minimizes the mean absolute deviation to all data points.

- The median can also be used to characterize asymmetric data.

Dispersion/Volatility/Variability



Volatility measures

Problem: the location measures do not characterize the data sufficiently

Aim: statements about the variation of the data around the center (a location measure)

i) range

$$\tilde{R} = x_{(n)} - x_{(1)}$$

Note: the range is *extremely* sensitive to the data/outliers.

ii) interquartile range

$$QA = \tilde{x}_{0.75} - \tilde{x}_{0.25}$$

Properties:

a) the interquartile range is robust to outliers.

b) There are at least $[n/2]$ of all observations in the interval $[\tilde{x}_{0.25}, \tilde{x}_{0.75}]$

iii) empirical variance

$$\tilde{s}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^k n(a_i) (a_i - \bar{x})^2 = \sum_{i=1}^k h(a_i) (a_i - \bar{x})^2$$

\tilde{s}^2 is the average squared deviation of the observations from the mean.

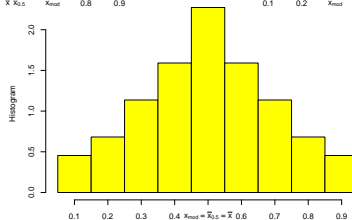
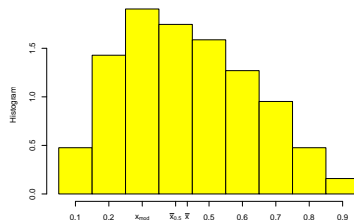
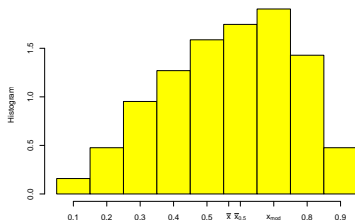
iv) sample variance

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

- s is the **sample standard deviation**. $\tilde{s} = \sqrt{\tilde{s}^2}$ is the **empirical standard deviation**.
- The empirical/sample variance/standard deviation is very sensitive to outliers .
- The empirical/sample variance/standard deviation is only reasonable for symmetric data.

Measures of skewness

Symmetric and nonsymmetric distributions



Aim: statements about the asymmetry of a sample

Note: it is reasonable only for unimodal distributions.

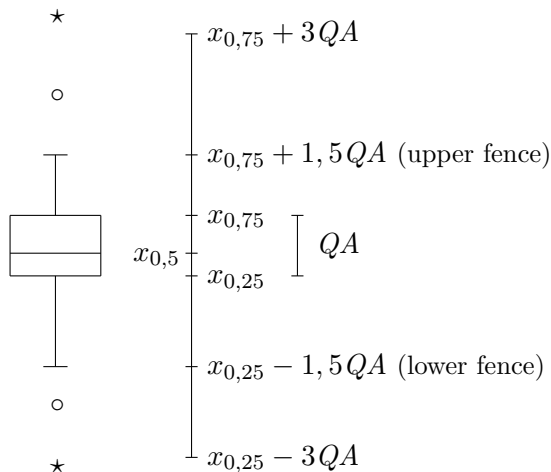
A distribution is **right-skewed**, if the peak is located at the left part of the distribution. Otherwise the distribution is **left-skewed**.

Sample skewness (empirical skewness)

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\tilde{s}} \right)^3$$

If it is larger (smaller) than zero, then we conclude that the distribution is right-skewed (left-skewed).

Boxplot - Graphical representation of some measures of location and variation



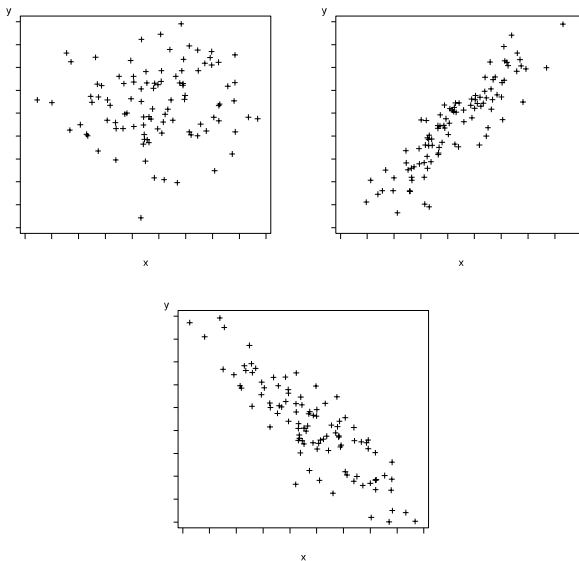
Characteristics of bivariate data sets

now: 2 variables/attributes X, Y , sample: $(x_1, y_1), \dots, (x_n, y_n)$

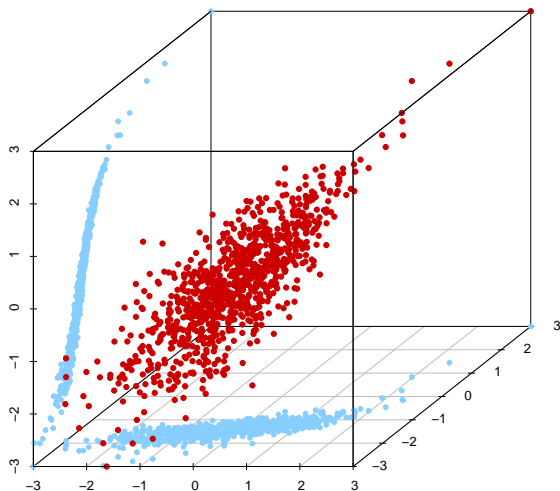
But: for each of the variables we can determine the individual measures of location and volatility as for univariate data sets.

For bivariate data sets we are particularly interested in the relationship between X and Y . This is the subject of the following discussion.

Scatterplots



3D-Scatterplot



Correlation measures for interval-scaled variables

Requirement: X and Y have interval scale

Aim: measure of correlation

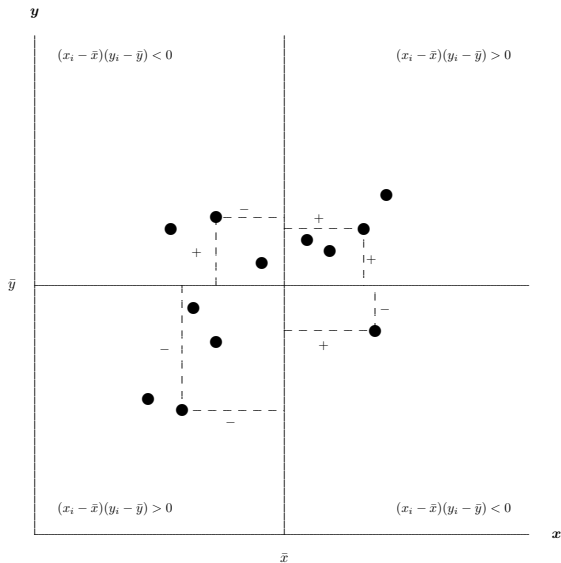
positiv relationship: large (small) values of X with large (small) values of Y

negativ relationship: inverse tendency

empirical covariance

$$\tilde{s}_{XY} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

An alternativ measure is the **sample covariance** $s_{XY} = \frac{n}{n-1} \tilde{s}_{XY}$.



Properties:

- $\tilde{s}_{XY} = \tilde{s}_{YX}$
- Invariant to shifts in the location, i. e. for $x_i^* = a x_i + b$ and $y_i^* = c y_i + d$ it holds that $\tilde{s}_{X^*Y^*} = a c \tilde{s}_{XY}$.
- $|\tilde{s}_{XY}| \leq \tilde{s}_X \tilde{s}_Y$
- It is sensitive to outliers.

Disadvantage: the empirical variance is not normalized and, therefore, depends on the scale

Sample correlation coefficient of Pearson:

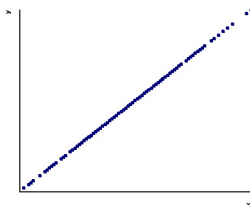
$$r_{XY} = \frac{s_{XY}}{s_X s_Y} = \frac{\tilde{s}_{XY}}{\tilde{s}_X \tilde{s}_Y}$$

Properties:

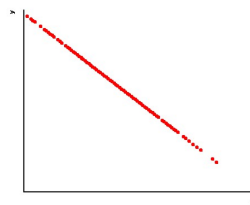
- $r_{XY} = r_{YX}$
- Invariant with respect to shifts in the location **and** in the scale
- $|r_{XY}| \leq 1$.
- If $r_{XY} = 1$ (or -1), then all observations (x_i, y_i) , $i = 1, \dots, n$ lie on a single straight line with positive (negative) slope.
- The empirical correlation coefficient is a measure of **linear** dependence between two variables.
- We cannot conclude about causality of the relationship!

Perfect correlation

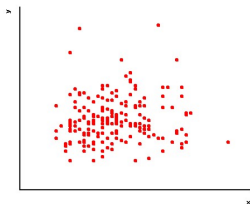
$$r = +1$$



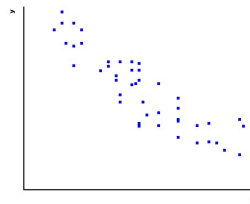
$$r = -1$$



weak correlation



strong correlation



Correlation measures for ordinal data

Requirement : X and Y are ordinal

Example: the relationship between the exam results (X , grade: $1, \dots, 5$) and the participation in tutorials (Y , seldom, regularly, always)

Idea of the ranks: assign to each observation of the sample x_1, \dots, x_n its position in the ordered sample $x_{(1)}, \dots, x_{(n)}$:

$$R(x_j) = v \quad \Leftrightarrow \quad x_j = x_{(v)}$$

$R(x_j)$ is the **rank** of the observation x_j .

Example: $x_1 = 2, x_2 = 5, x_3 = 1, x_4 = 3$. ordered sample:
 $x_3 < x_1 < x_4 < x_2$. Thus $R(x_1) = 2, R(x_2) = 4, R(x_3) = 1, R(x_4) = 3$.

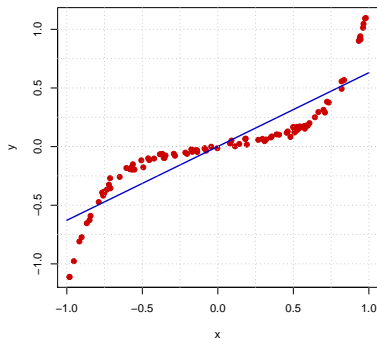
Given: sample $(x_1, y_1), \dots, (x_n, y_n)$; assign to x_1, \dots, x_n the ranks $R(x_1), \dots, R(x_n)$ and to y_1, \dots, y_n the ranks $R(y_1), \dots, R(y_n)$.

Rank correlation coefficient of Spearman

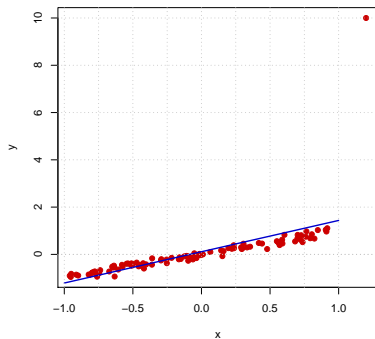
$$R_{XY} = r_{R(X), R(Y)} = \frac{\sum_{i=1}^n (R(x_i) - \bar{R})(R(y_i) - \bar{R})}{\sqrt{\sum_{i=1}^n (R(x_i) - \bar{R})^2 \sum_{i=1}^n (R(y_i) - \bar{R})^2}}$$

with $\bar{R} = (n + 1)/2$.

$$\begin{aligned}\delta &= 0.892, \\ \rho &= 0.996\end{aligned}$$



$$\begin{aligned}\delta &= 0.659, \\ \rho &= 0.982\end{aligned}$$



Correlation measures for nominal variables

Now: 2 nominal variables with realizations a_1, \dots, a_k for X and b_1, \dots, b_l for Y

Example: 156 graduates, 93 boys, 63 girls. 9 boys and 2 girls failed the exam.

Contingency table of absolute frequencies:

X	Y		Σ
	passed	failed	
B	84	9	93
G	61	2	63
Σ	145	11	156

Contingency table of relative frequencies :

X	Y		Σ
	passed	failed	
B	0.538	0.058	0.596
G	0.391	0.013	0.404
Σ	0.929	0.071	1.0

Bivariate frequency table

- absolute frequency for (a_i, b_j) :

$n_{ij} = n(X = a_i, Y = b_j)$ = the number of cases, where the pair (a_i, b_j) is observed in the sample

- absolute marginal frequency of a_i :

$n_{i\cdot}$ = the number of cases, where the realization a_i is observed in x_1, \dots, x_n

- the relative frequencies are $h_{ij} = n_{ij}/n$ and $h_{i\cdot} = n_{i\cdot}/n$ respectively.

on analogy: $n_{\cdot j}$, for Y

Example:

- relative marginal frequencies for gender: (0.596, 0.404)
- relative marginal frequencies for exam results: (0.929, 0.071)

contingency table for absolute frequencies

X	Y				Σ
	b_1	b_2	\cdots	b_l	
a_1	n_{11}	n_{12}	\cdots	n_{1l}	$n_{1\cdot}$
a_2	n_{21}	n_{22}	\cdots	n_{2l}	$n_{2\cdot}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
a_k	n_{k1}	n_{k2}	\cdots	n_{kl}	$n_{k\cdot}$
Σ	$n_{\cdot 1}$	$n_{\cdot 2}$	\cdots	$n_{\cdot l}$	n

Aim: a measure of dependency

Idea: weak dependency, if for all i, j

$$n_{ij} \approx \frac{n_{i.} \cdot n_{.j}}{n}$$

$$\rightsquigarrow \chi^2 = \sum_{i=1}^k \sum_{j=1}^l \frac{(n_{ij} - n_{i.} \cdot n_{.j} / n)^2}{n_{i.} \cdot n_{.j} / n}$$

χ^2 „large“ $\rightsquigarrow X$ and Y are dependent.

Since χ^2 increases with n , we consider

The contingency coefficient of Pearson

$$C = \sqrt{\chi^2/(\chi^2 + n)}, \text{ with } C_{max} = \sqrt{\frac{\min\{k, l\} - 1}{\min\{k, l\}}}$$

Thus

Corrected contingency coefficient of Pearson

$$C_{Corr} = C/C_{max} \in [0, 1]$$

The smaller is C_{Corr} , the „weaker “is the dependence. $C_{Corr} = 0$ only if X and Y are independent .

Chapter 2

Elements of Probability Theory

Probability of events

Origins of probability theory: Jakob Bernoulli (1655-1705),
Pierre-Simon de Laplace (1749-1827)

The probability theory originated from the analysis of games of chance (gambling).

Aim: statements about probabilities of random events

- Subsets consisting of a single element of Ω are called **elementary events**: $\{\omega\} \in \Omega$
- Any subset of Ω is called an **event**: $A = \{\omega_1, \dots\} \in \Omega$.

Laplace probability

Starting point: All elementary events have the same probability!

If Ω is finite, then it holds

$$P(A) = \frac{\text{the number of for } A \text{ „favourable cases“}}{\text{the number of all possible cases}} = \frac{|A|}{|\Omega|},$$

where $|A|$ denotes the number of elements in A and similarly for $|\Omega|$.

Example: roulette game ($\Omega = \{0, \dots, 36\}$)

- A = the set of numbers divisible by 3
- B = the even numbers

It holds $P(\{0\}) = P(\{1\}) = \dots = P(\{36\}) = 1/37$, i.e. it is a Laplace experiment. Then

$$P(A) = \frac{|A|}{|\Omega|} = \frac{12}{37}.$$

The probability, that we observe a number of pips, which is divisible by 3, but not divisible by 2, is

$$P(A \cap \bar{B}) = \frac{|\{3, 9, 15, 21, 27, 33\}|}{37} = \frac{6}{37}.$$

Statistical probability

Let $A \subset \Omega$. The experiment is repeated n times. $h(A)$ denotes the relative frequency of A .

Example: roulette ($\Omega = \{0, 1, \dots, 36\}$)

Let A be the event “*we observe a number from the first dozen*”, i.e. $A = \{1, 2, \dots, 12\}$.

16 replications produce the sample

23	34	13	11	28	9	8	21
16	33	31	15	3	13	23	32

Then $h(A) = \frac{4}{16} = 0.25$.

Example: We throw a coin n times. We obtain

n	$n(K)$	$h_n(K)$
10	7	0.700
20	11	0.550
100	47	0.470
400	204	0.510
1000	492	0.492
2000	1010	0.505

The coin is symmetric. Therefore the relative frequencies converge to the true probability of 0.5.

Richard von Mises (1931)

The probability of observing A :

$$P(A) := \lim_{n \rightarrow \infty} h_n(A)$$

Disadvantages: difficult to implement in practice

Axioms of the probability theory

Both the Laplace probability and the statistical probability have their pros and cons. A general approach to probability was suggested by Kolmogorov (1933).

A. N. Kolmogorov (1933)

The **probability measure** P is mapping, which assigns a number to (almost all) events $A \subseteq \Omega$ (namely $P(A)$) and fulfills the following properties:

- $0 \leq P(A) \leq 1$
- $P(\Omega) = 1$
- $P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$ for all $A_i \subset \Omega$ with $A_v \cap A_j = \emptyset$ for $v \neq j$.

$P(A)$ is the **probability of event A** .

Rules for the probabilities

Let P be a probability measure on Ω . Then it holds:

- $P(\bar{A}) = 1 - P(A)$
- $P(\emptyset) = 0$
- $P(A) = P(A \cap B) + P(A \cap \bar{B})$
- If $B \subseteq A$, then $P(B) \leq P(A)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- If Ω is finite, then it holds for $A \subseteq \Omega$ that:

$$P(A) = \sum_{a \in A} P(\{a\}).$$

Note: Both the Laplace probability and the statistical probability are probability measures.

Conditional probability and independence

Now: conditional probability of event A under the condition B (of A , if B is given or observed)

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} \quad \text{for } P(B) > 0$$

Note: $P(A \mid B) + P(\bar{A} \mid B) = 1$

Law of total probability

Let A_1, \dots, A_k be events, which are disjoint in pairs, with $A_1 \cup \dots \cup A_k = \Omega$. Then for an arbitrary event B it holds

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

Bayes' rule (1702 – 1761)

Let A_1, \dots, A_k be events, which are disjoint in pairs with $A_1 \cup \dots \cup A_k = \Omega$. Furthermore, let B be an arbitrary event. Then it holds for $i \in \{1, \dots, k\}$

$$P(A_i \mid B) = \frac{P(B \mid A_i) \cdot P(A_i)}{\sum_{j=1}^k P(B \mid A_j) \cdot P(A_j)}.$$

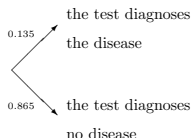
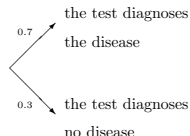
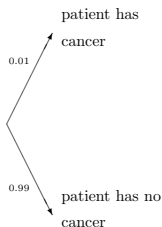
Bayes rule

Example: Extensive studies have shown that appr. 1.0% (**a-priori probability**) of all men between 40 and 50 have the cancer of prostate. A simple diagnostic test is the PSA test.

The PSA test has the property, that it makes the correct diagnosis with probability 0.7 for healthy patients (**sensitivity**) and with probability of 0.865 for ill patients.

What is the probability that a patient with negative (positive) test results is truly healthy (ill) (**posteriori probability**)?

Bayes rule II



Aim: $P(\text{patient is ill} \mid \text{test diagnoses the disease})$

$$\begin{aligned}
 &= \frac{P(\text{patient is ill AND the test diagnoses the disease})}{P(\text{test diagnoses the disease})} \\
 &= \frac{P(\text{test diagnoses the disease, if the patient is ill}) P(\text{patient is ill})}{P(\text{test diagnoses the disease})} \\
 &= \frac{0.7 \cdot 0.01}{P(\text{test diagnoses the disease})}
 \end{aligned}$$

Bayes rule III

Using the rule of total probability we obtain

$$\begin{aligned} P(\text{the test diagnoses the disease}) \\ = 0.7 \cdot 0.01 + 0.135 \cdot 0.99 = 0.14065. \end{aligned}$$

The probability that the patient is really ill, even if the test diagnosed it, equals

$$\frac{0.7 \cdot 0.01}{0.14065} \approx 0.0498.$$

Independent events

Two events $A, B \subseteq \Omega$ are (stochastically) independent, if

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

Note: If A and B are independent, then it holds that $P(B \mid A) = P(B)$ and $P(A \mid B) = P(A)$, since

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

If two events are not (stochastically) independent, then we say, that they are (stochastically) dependent.

Random variables and distribution functions

A **random variable** (attribute) X is an appropriate mapping of the population Ω into the set S . In general $S \subset \mathbb{R}$.

Thus

$$X(\omega) = x,$$

where $\omega \in \Omega$ is a “state of the world” which causes the particular outcome $x \in S$ of the RV X .

If $S \subset \mathbb{R}^n$, then X is an n -dimensional random variable or a **random vector**.

The **distribution function** F_X of a random variable X is defined as

$$F_X(x) = P\left(\{\omega \in \Omega : X(\omega) \leq x\}\right), \quad x \in \mathbb{R}.$$

Usually a short-hand form is used $F(x) = P(X \leq x)$ or $X \sim F$

- The distribution function is a mapping from a set of real numbers into the interval $[0, 1]$.
- The distribution function assigns to each event $\{X \leq x\}$ the corresponding probability.

Properties of distribution functions

Def: The **distribution function** F of a random variable X is a function with the following properties:

- $0 \leq F(x) \leq 1$ for all x
 - $F(\infty) = \lim_{x \rightarrow \infty} F(x) = 1, \quad F(-\infty) = \lim_{x \rightarrow -\infty} F(x) = 0$
 - $F(x)$ is monotone increasing in x
 - F is right-side continuous
-
- Each function F which satisfies the above conditions is a distribution function.
 - If there is a function, which satisfies the above properties, then we can construct a random variable and a probability measure, such that the distribution function of the random variable coincides with the given function.

Computation of the probabilities

The distribution function contains all the information relevant to a statistician. Using the distribution function we can compute all the probabilities related to the random variable.

Assuming $a < b$ it holds:

- $P(a < X \leq b) = F(b) - F(a)$
- $P(a \leq X \leq b) = F(b) - F(a - 0)$
- $P(X > a) = 1 - P(X \leq a) = 1 - F(a)$
- $P(X \geq a) = 1 - P(X < a) = 1 - F(a - 0).$

where $F(a - 0)$ denotes the left-sided limit of F at a , i.e.

$$F(a - 0) = \lim_{\varepsilon \rightarrow 0} F(a - \varepsilon), \text{ with } \varepsilon > 0.$$

Example: we toss a die till the first “6”. Let X denote the number of tosses. Thus $\Omega = \mathbb{N}$.

Then it holds

$$f(i) = P(X = i) = \frac{1}{6} \left(\frac{5}{6}\right)^{i-1}.$$

$F(x) = P(\emptyset) = 0$ for $x < 1$.

For $n \in \mathbb{N}$, we obtain

$$\begin{aligned} F(n) &= P(X \leq n) = \sum_{i=1}^n f(i) \\ &= \frac{1}{6} \sum_{i=0}^{n-1} \left(\frac{5}{6}\right)^i = 1 - \left(\frac{5}{6}\right)^n. \end{aligned}$$

For $n \leq x < n+1$ we obtain $F(x) = F(n)$.

- Probability of more than 10 tosses :

$$P(X > 10) = 1 - F(10) = \left(\frac{5}{6}\right)^{10} \approx 0.16$$

- Probability of more than 3 but less than 8 tosses:

$$\begin{aligned} P(3 < X < 8) &= P(3 < X \leq 7) \\ &= F(7) - F(3) = \left(\frac{5}{6}\right)^3 - \left(\frac{5}{6}\right)^7 \approx 0.3 \end{aligned}$$

Discrete random variables and discrete distribution functions

If X has a countable set of possible values, then X is a **discrete random variable** and F_X is a **discrete distribution function**.

- Let X take the values x_1, x_2, \dots and $p_i = P(X = x_i)$. Then

$$f(x) = \begin{cases} p_i & \text{if } x = x_i \\ 0 & \text{if } x \neq x_i \quad \forall i \end{cases}$$

is the **probability function of X** .

- Let $x_1 < x_2 < \dots$. If $x_i \leq x < x_{i+1}$, then

$$F(x) = \sum_{v=1}^i f(x_v) = P(X = x_1) + \dots + P(X = x_i).$$

Particularly $F(x) = 0$ for $x < x_1$, $F(x) = 1$ for $x > x_n$.

Examples for discrete distribution functions

a) Binomial distribution

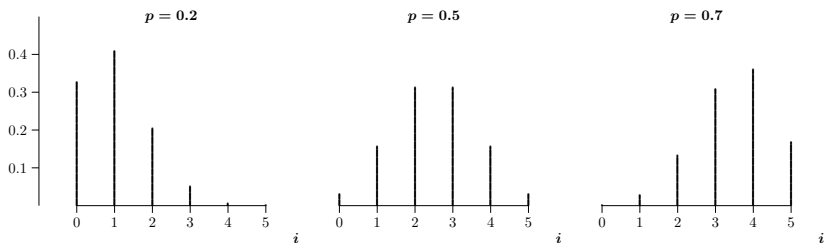
We repeat an experiment independently n times. The probability of observing the event A is $p = P(A)$. Then the probability that we observe A exactly k times is given by

$$\binom{n}{k} p^k (1 - p)^{n-k} = b_{n,p}(k).$$

$b_{n,p}$ is the probability function. The respective distribution function – the binomial distribution function – is denoted by $B_{n,p}$. If $n = 1$, then we call this distribution Bernoulli distribution.

Example: A drug has the probability of success equal to 0.9. What is the probability that exactly 48 out of 50 patients would be healed up?

$$\begin{aligned}P(X = 48) &= \binom{50}{48} \cdot 0.9^{48} \cdot 0.1^2 \\&= 49 \cdot 25 \cdot 0.9^{48} \cdot 0.1^2 \\&\approx 0.078\end{aligned}$$

Probability function of the binomial distribution for $n = 5$ 

b) Hypergeometric distribution

We consider a box with n balls. r of them are red, the rest are white. We draw k balls without replacement. Let the random variable X denote the number of drawn red balls.

$$P(X = i) = \frac{\binom{r}{i} \binom{n-r}{k-i}}{\binom{n}{k}}$$

This is the probability function of the [hypergeometric distribution](#).

Example: from experience we know that the production of particular devices results in 20% of defective products. On a given day we produce 100 devices and randomly select arbitrary 10 of them. What is the probability, that the sample contains exactly 2 flaw products?

$$P(X = 2) = \frac{\binom{20}{2} \binom{80}{8}}{\binom{100}{10}} \approx 0.3181.$$

c) Poisson distribution

Let $X \sim B(n, p)$, i. e. $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$.

It is often the case that for the Binomial distribution n is large and p is small. Let p be a function of n , i.e. $p = p(n)$. If $\lim_{n \rightarrow \infty} np(n) = \lambda > 0$, then

$$\lim_{n \rightarrow \infty} b(n, p(n))(k) = \frac{\lambda^k}{k!} e^{-\lambda}.$$

We denote the limiting distribution by **Poisson distribution** and write $P(\lambda)$.

Example: A large insurance company computes the price of a vehicle insurance contract. On the basis of historical data the company assumes that the number of accidents X in a particular period follows the Poisson distribution with $\lambda = 3$.

Then

$$P(X = 2) = \exp(-3) \frac{3^2}{2!} \approx 0.224,$$

$$P(X \geq 2) = 1 - P(X = 0) - P(X = 1) = 1 - \exp(-3) - 3 \exp(-3) \approx 0.8009.$$

Note: The assumption of Poisson distribution is suitable here, because there are very many contracts and relatively few accidents.

Discrete distributions

	probability- function $f(m)$	Parameter space	Expected value $\mu = E(X)$	Variance $\sigma^2 = E([X - \mu]^2)$
Binomial $B(n, p)$	$\binom{n}{m} p^m (1-p)^{n-m}$	$0 < p < 1$ $n \in \{1, 2, \dots\}$	$n p$	$n p (1-p)$
	$m \in \{0, 1, \dots, n\}$			
Hyper- geometric $H(N, M, n)$	$\frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}$	$N \in \{1, 2, \dots\},$ $M \in \{0, 1, \dots, N\},$ $n \in \{1, 2, \dots, N\}$	$n \frac{M}{N}$	$n \frac{M}{N} \frac{N-M}{N} \frac{N-n}{N-1}$ (for $N > 1$)
	$m \in \{m_{\min}, m_{\min} + 1, \dots, m_{\max}\},$ $m_{\min} := \max\{0, n - (N - M)\},$ $m_{\max} := \min\{n, M\}$			
Poisson $P(\lambda)$	$\frac{\lambda^m}{m!} e^{-\lambda}$	$\lambda > 0$	λ	λ
	$m \in \{0, 1, \dots\}$			
Geometric $G(p)$	$p (1-p)^{m-1}$	$0 < p < 1$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
	$m \in \{1, 2, \dots\}$			

Continuous random variables

X is a **continuous** random variable, if there exists a non-negative function f , such that:

$$F(x) = \int_{-\infty}^x f(t) dt \quad \text{for all } x \in \mathbb{R}.$$

The function f is called the **density (probability density) function** of X .

Properties:

- $P(a < X \leq b) = \int_a^b f(t) dt$
- It holds $P(X = x) = 0$ for all $x \rightsquigarrow P(a < X < b) = P(a \leq X \leq b)$.
- If F is a continuous function, then $F' = f$.
- $\int_{-\infty}^{\infty} f(t) dt = 1$.
- The inverse CDF $F^{-1}(\beta)$ is called **the quantile function**.

$$F^{-1}(\beta) = \inf\{x : F(x) > \beta\} \rightsquigarrow P(X \leq F^{-1}(\beta)) \geq \beta$$

Continuous distributions I

	Density f	Parameter space	Expected value $\mu = E(X)$	Variance $\sigma^2 = E(X - \mu)^2$
Uniform $U(a, b)$	$\frac{1}{b - a}, x \in [a, b]$	$-\infty < a < b < \infty$	$\frac{a + b}{2}$	$\frac{(b - a)^2}{12}$
Normal $N(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}$	$\mu \in \mathbb{R}, \sigma > 0$	μ	σ^2
Exponential $E(\lambda)$	$f(x) = \lambda e^{-\lambda x}, x \geq 0$	$\lambda > 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
χ_n^2 χ_n^2	$\frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, x > 0, n \in \mathbb{N}$		n	$2n$
t-distr. (Student) t_n	$\frac{\left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}}{B(n/2, 1/2)\sqrt{n}}, x \in \mathbb{R}, n \in \mathbb{N}$		$0 \quad (n > 1)$	$\frac{n}{n-2} \quad (n > 2)$

Continuous distributions II

	Dichte f	Parameter space	Expected value $\mu = E(X)$	Variance $\sigma^2 = E(X - \mu)^2$
F -distr. $F_{m,n}$	$\frac{(m/n)^{m/2}}{B(\frac{m}{2}, \frac{n}{2})} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n} x\right)^{-\frac{m+n}{2}},$ $x \geq 0, m, n \in \mathbb{N}$		$\frac{n}{n-2}$ $(n > 2)$	$\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}$ $(n > 4)$
Gamma-distr.	$\frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x}, x \geq 0$	$\lambda > 0, r > 0$	$\frac{r}{\lambda}$	$\frac{r}{\lambda^2}$
Cauchy distr.	$\frac{1}{\pi \beta \{1 + [(x - \alpha)/\beta]^2\}}, x \in \mathbb{R}$	$\beta > 0, \alpha \in \mathbb{R}$	-	-

with $\Gamma(x) = \int_0^\infty \exp(-t)t^{x-1}dt$ and $B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$.

Normal (Gaussian) distribution

The Normal distribution is the most important continuous distribution. It depends on 2 parameters, $\mu \in \mathbb{R}$ and $\sigma > 0$. Its density is given by

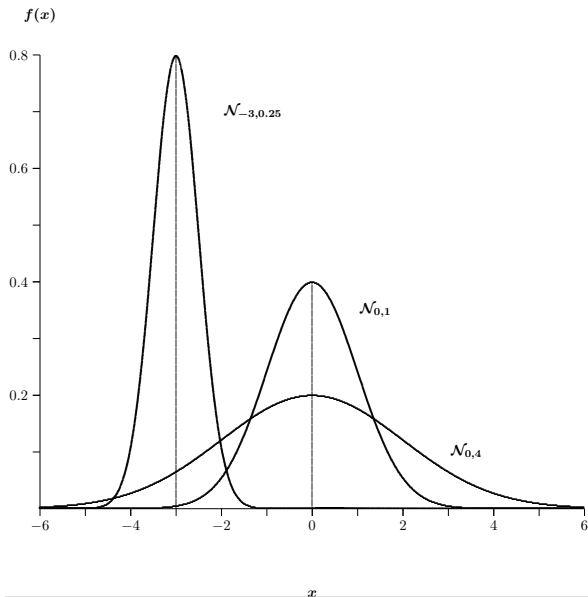
$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left\{ -\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right\}, \quad x \in \mathbb{R}.$$

For the distribution function of normal distribution we use the symbol $N(\mu, \sigma^2)$ or N_{μ, σ^2} .

Properties:

- f symmetric w.r.t. $x = \mu$, i.e. it holds $f(\mu + x) = f(\mu - x)$ for all x .
- The maximum of f is attained at μ .
- f has two turning points at $\mu \pm \sigma$.

Density functions of normal distribution for different parameters



Standard normal distribution

By **standard normal distribution** we denote the normal distribution with $\mu = 0$ and $\sigma = 1$. We write Φ for the distribution function and ϕ for the density.

Properties:

- Since $\phi(x) = \phi(-x)$, it follows that $\Phi(x) = 1 - \Phi(-x)$.
- If $X \sim N(\mu, \sigma^2)$, then $\frac{X - \mu}{\sigma} \sim \Phi$. This implies

$$F_X(x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right).$$

- If $X \sim \Phi$, then $\mu + \sigma X \sim N(\mu, \sigma^2)$.
- If $X \sim N(\mu, \sigma^2)$, then $aX + b \sim N(a\mu + b, a^2\sigma^2)$.

Further properties

- Probability for deviation from the mean for at most c :

$$\begin{aligned}
 P(\mu - c \leq X \leq \mu + c) &= F(\mu + c) - F(\mu - c) \\
 &= \Phi\left(\frac{\mu + c - \mu}{\sigma}\right) - \Phi\left(\frac{\mu - c - \mu}{\sigma}\right) \\
 &= \Phi\left(\frac{c}{\sigma}\right) - \Phi\left(-\frac{c}{\sigma}\right) \\
 &= \Phi\left(\frac{c}{\sigma}\right) - [1 - \Phi\left(\frac{c}{\sigma}\right)] \\
 &= 2 \cdot \Phi\left(\frac{c}{\sigma}\right) - 1
 \end{aligned}$$

$k\sigma$ -intervals $[\mu - k\sigma, \mu + k\sigma]$:

$$P(\mu - k\sigma \leq X \leq \mu + k\sigma) = 2\Phi(k) - 1 = \begin{cases} 0,683, & \text{for } k = 1 \\ 0,954, & \text{for } k = 2 \\ 0,997, & \text{for } k = 3 \end{cases}$$

Exponential distribution

Exponential distribution arises in the analysis of life expectancy. Its density is given by

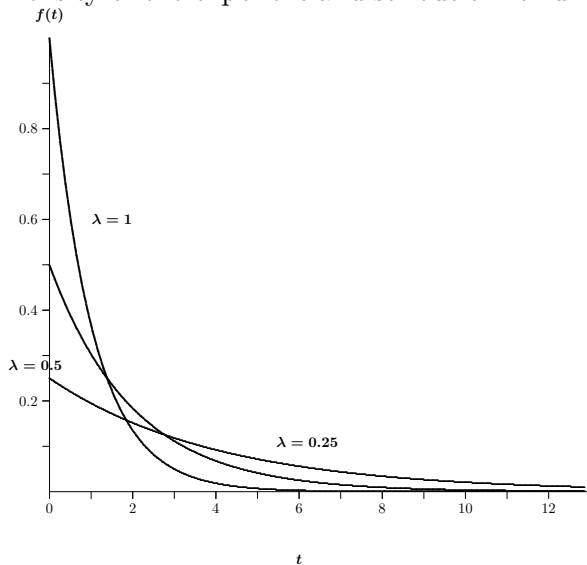
$$f(x) = \begin{cases} \lambda \exp(-\lambda x) & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$$

with $\lambda > 0$. Therefore $F(x) = 1 - \exp(-\lambda x)$. We write $E(\lambda)$.

Example: The life-span of TV-sets follows exponential distribution with $\lambda = 0.08$. What is the probability that the TV-set would have a life-span of more than 10 years?

It holds

$$P(X > 10) = 1 - F(10) = \exp(-0.08 \cdot 10) = \exp(-0.8) \approx \dots$$

Density of the exponential distribution for different parameters λ 

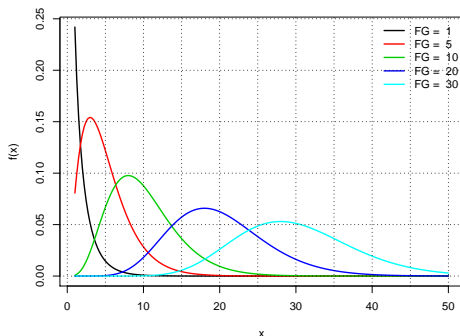
Chi-Square-Distribution (χ_f^2)

Assume that n RV's Z_1, \dots, Z_n

- are independent and
- follow standard normal distribution $Z_i \sim N(0; 1)$ for $i = 1, \dots, n$

Then the sum of squares follows χ^2 distribution with n degrees of freedom

$$Z_1^2 + \dots + Z_n^2 \sim \chi_n^2$$



t -distribution (Student-Distribution)

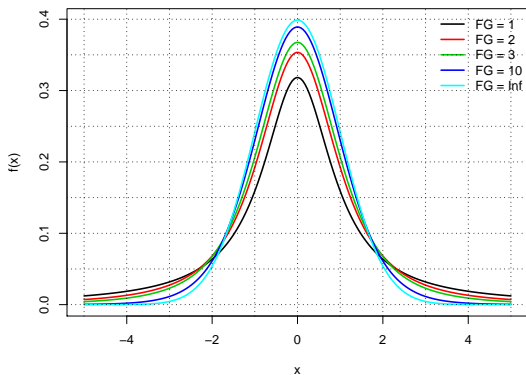
- Z follows the standard normal distribution: $Z \sim N(0, 1)$
- Y is independent from Z and follows the chi-square distributed with df d : $Y \sim \chi_d^2$

Then the random variable

$$T = \frac{Z}{\sqrt{Y/d}}$$

follows the t distribution with degrees of freedom d .

- the density of the t -distribution is a symmetric bell-shaped curve
- the density of the t -distribution has heavier tails compared to the density of the normal distribution
- as $d \rightarrow \infty$ the density function of the t_d -distribution converges to the density of the standard normal distribution.



F-distribution

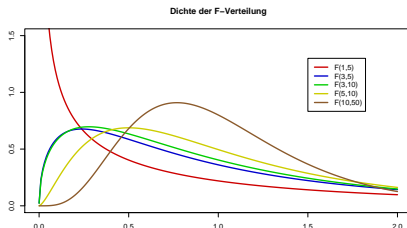
- Having two independent random variables Y_1 and Y_2 , both following the chi-square-distributions with f_1 and f_2 df respectively:

$$Y_1 \sim \chi^2(d_1) \quad Y_2 \sim \chi^2(d_2)$$

Then the distribution of the random variable

$$F = \frac{Y_1/d_1}{Y_2/d_2}$$

is called **F-distribution** with parameters d_1 and d_2



Characteristics of random variables

- In the descriptive statistics we discussed the location and dispersion measures of random samples.
- Here we discuss the measures of location and dispersion for random variables.
- The aim of the discussion is make statements about the center (central tendency) of the distribution.

Let X be a discrete RV und take values x_1, x_2, \dots . Then the expectation of X (or equivalently of F) is given by

$$E(X) = \sum_i x_i P(X = x_i).$$

Examples:

- You win 4 Euro, if you throw “6” on a die and loose 1 Euro if you throw another number of pips. Then

$$E(X) = -1 \cdot \frac{5}{6} + 4 \cdot \frac{1}{6} = -\frac{1}{6}.$$

- Poisson distribution, i.e. $P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}$ for $k \geq 0$

$$\begin{aligned} E(X) &= \sum_{k=0}^{\infty} k P(X = k) = \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} \\ &= e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^k}{(k-1)!} = \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} \\ &= \lambda e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = \lambda e^{-\lambda} e^{\lambda} = \lambda. \end{aligned}$$

Let X be a continuous RV with the density function f . Then the expectation of X (or of F) is given by

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

The integral exists if $\int_{-\infty}^{\infty} |x|f(x)dx < \infty$.

Note:

- If f is symmetric with respect to m , i.e.

$$f(m+x) = f(m-x) \quad \text{for all } x$$

then $E(X) = m$, if it exists.

- This implies that for $X \sim N(\mu, \sigma^2)$ it holds that $E(X) = \mu$.
- The expectation of the Cauchy distribution does not exist.

Rules for computing the expectations

Aim: computation of the expectation of $Y = g(X)$

If X is discrete, then it holds

$$E(Y) = \sum_i g(x_i) P(X = x_i).$$

If X is continuous, then it holds

$$E(Y) = \int_{-\infty}^{\infty} g(x) f(x) dx$$

(if the integral exists)

Examples:

- It holds that $E(aX + b) = aE(X) + b$, since (only continuous case)

Sums and products of random variables

- Let X_1, \dots, X_n be random variables with existing expectations. Then it holds that

$$E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i).$$

- If the RVs X_1, \dots, X_n are additionally independent, then it holds that

$$E\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n E(X_i).$$

- Consider the portfolio consisting of n assets and its return R . Let P_t denote the price of an asset at time point t . The simple **return** of the asset at time point t is given by

$$R_t = 100 (P_t - P_{t-1}) / P_{t-1}.$$

We consider now the returns of n assets at a given time point t . We denote them by R_1, \dots, R_n . Let the relative fraction of the i th asset in the portfolio be given by w_i . This implies $\sum_{i=1}^n w_i = 1$. Then the portfolio return equals $R = \sum_{i=1}^n w_i R_i$. Thus it follows:

$$E(R) = E\left(\sum_{i=1}^n w_i R_i\right) = \sum_{i=1}^n E(w_i R_i) = \sum_{i=1}^n w_i E(R_i).$$

If $E(R_i) = \mu$ for all $i = 1, \dots, n$, then $E(R) = \mu$ too.

Dispersion measures of distribution functions

The **dispersion (variability) measures** for the distribution function measure the concentration of the probability around the center of symmetry.

The most popular dispersion measure is the **variance**. It is measured as the expected quadratic deviation from the expectation $\mu = E(X)$:

$$\text{Var}(X) = E([X - \mu]^2).$$

The variance exists if $E(X^2) < \infty$. Often it is denoted by $\sigma^2 = \text{Var}(X)$.

The quantity σ is called the **standard deviation**.

Let X be a discrete RV with the realizations x_1, x_2, \dots . Then it holds that

$$\text{Var}(X) = \sum_i (x_i - \mu)^2 P(X = x_i).$$

If X is a continuous RV with the density function f , then

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

Note:

- If $\text{Var}(X) = 0$, then $X = E(X)$. For continuous RVs it holds “almost everywhere”.
- For all $a, b \in \mathbb{R}$ it holds that

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

- If $X \sim N(\mu, \sigma^2)$, then X has the same distribution as $\mu + \sigma Y$ with $Y \sim \Phi$. This implies

$$\text{Var}(X) = \text{Var}(\mu + \sigma Y) = \sigma^2 \text{Var}(Y) = \sigma^2.$$

Note: the parameter σ^2 of the normal distribution equals the variance!

- If the RVs X_1, \dots, X_n are independent (!) and the respective variances exist, then

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i).$$

Characteristics of 2D distributions

The most popular measures of comovement are the **covariance** and the **correlation**.

- The **covariance** between X and Y is given by:

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

The covariance exists if $E(|XY|) < \infty$.

- If $\text{Var}(X) > 0$ and $\text{Var}(Y) > 0$, then

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

is called the **correlation coefficient of Pearson**.

- X und Y are **uncorrelated** if $\text{Corr}(X, Y) = 0$.

If X and Y are discrete random variables with realizations $x_1, x_2, \dots, y_1, y_2, \dots$, then

$$\begin{aligned} \text{Cov}(X, Y) &= \sum_i \sum_j (x_i - E(X)) (y_j - E(Y)) \cdot P(X = x_i, Y = y_j) \\ &= \sum_i \sum_j x_i y_j P(X = x_i, Y = y_j) - E(X) E(Y). \end{aligned}$$

If (X, Y) is a continuous random vector with the density function f , then

$$\begin{aligned} \text{Cov}(X, Y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E(X)) (y - E(Y)) \cdot f(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x y f(x, y) dx dy - E(X) E(Y) = E(XY) - E(X)E(Y). \end{aligned}$$

Rules for covariances and correlations:

- $\text{Corr}(aX + b, cY + d) = \text{Corr}(X, Y)$ (if a and c have the same sign)
(invariance w.r.t. to location and scale shifts)
- $|\text{Corr}(X, Y)| \leq 1$,
 $|\text{Corr}(X, Y)| = 1$, if X and Y lie on a straight line, i. e.
 $Y = \alpha + \beta X$.
- If X and Y are independent, then $\text{Cov}(X, Y) = 0$.
The inverse statement is **not correct** in general!!!
- $\text{Var}(aX + bY) = a^2\text{Var}(X) + b^2\text{Var}(Y) + 2ab\text{Cov}(X, Y)$, since

$$\begin{aligned}\text{Var}(aX + bY) &= E[(aX + bY - E(aX + bY))^2] \\ &= E(a(X - E(X)) + b(Y - E(Y)))^2 \\ &= a^2\text{Var}(X) + b^2\text{Var}(Y) + 2ab\text{Cov}(X, Y)\end{aligned}$$

Two dimensional distribution functions

Let $X = (X_1, X_2)'$ (for example, the returns of Daimler and BMW, exchange rates Euro/\$ and Euro/CHF). Then

$$F_X(x_1, x_2) = P\left(\{\omega \in \Omega : X_1(\omega) \leq x_1, X_2(\omega) \leq x_2\}\right), \quad x_1, x_2 \in \mathbb{R}$$

is a (2-dimensional) distribution function of the random vector X . The short-hand notation is $F(x_1, x_2) = P(X_1 \leq x_1, X_2 \leq x_2)$.

$F_X(x_1, \infty)$ is the marginal distribution of X_1 and

$F_X(\infty, x_2)$ is the marginal distribution X_2 .

Note: it holds

$$F_X(x_1, \infty) = P(X_1 \leq x_1) =: F_1(x_1) \quad \text{and}$$

$$F_X(\infty, x_2) = P(X_2 \leq x_2) =: F_2(x_2).$$

Discrete and continuous random vectors

If the set of possible values of X is countable, then X is **discrete** and

$$f(x_1, x_2) = P(X_1 = x_1, X_2 = x_2)$$

is the **(joint) probability function** of (X_1, X_2) .

If X is **continuous**, then the distribution function F of X is given by

$$F(x_1, x_2) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} f(t_1, t_2) dt_2 dt_1, \quad x_1, x_2 \in \mathbb{R}$$

with $f(t_1, t_2) \geq 0$ for all t_1, t_2 . The function f is a **(2-dimensional) probability density function (pdf)** of (X_1, X_2) .

Note: If f is given, then the density function of f_1 (f_2) of X_1 (X_2) can be obtained in the following way

$$f_1(x_1) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_2, \quad f_2(x_2) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_1.$$

Multivariate normal distribution

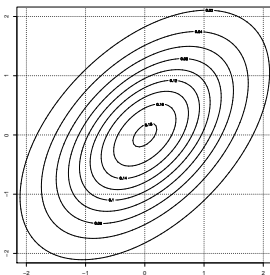
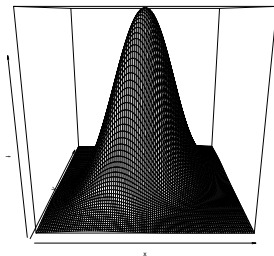
Def: The random vector \mathbf{X} follows a *p*-dimensional multivariate normal distribution ($\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$), if its density is given by

$$f(\mathbf{x}) = (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right].$$

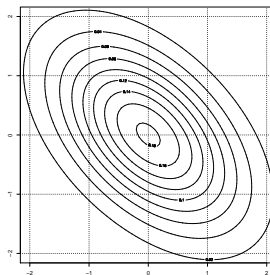
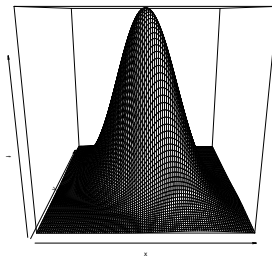
Other multivariate distributions known in explicit form: *t*, Laplace, Wishart, and very few others.

Example (2-dimensional normal distribution)

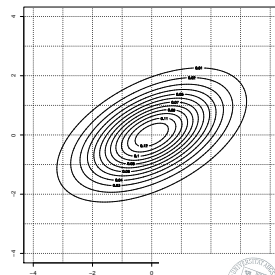
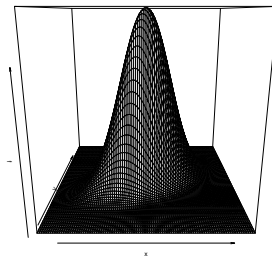
$$\sigma_1^2 = 1, \sigma_2^2 = 1, \rho = 0.5$$



$$\sigma_1^2 = 1, \sigma_2^2 = 1, \rho = -0.5$$



$$\sigma_1^2 = 2, \sigma_2^2 = 1, \rho = 0.5$$



Multivariate RV

Def: \mathbf{X} is a p -dimensional **random vector**, if the components X_1, \dots, X_p are scalar RVs.

The joint CDF is given by

$$F(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_p \leq x_p)$$

For a continuous random vector \mathbf{X} it holds:

$$F(x_1, \dots, x_{i-1}, -\infty, x_{i+1}, \dots, x_p) = 0$$

$$F(+\infty, \dots, +\infty) = 1$$

$$F(\mathbf{x}) = \int_{-\infty}^{x_p} \dots \int_{-\infty}^{x_1} f(\mathbf{u}) d\mathbf{u}$$

Expectation and covariance matrix

Def: For a random vector \mathbf{X} the **expectation** is defined by

$$E(\mathbf{X}) = \boldsymbol{\mu} = (\mu_1, \dots, \mu_p)' = (EX_1, \dots, EX_p)'$$

and the **covariance matrix** by

$$\begin{aligned} \text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma} &= \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_p^2 \end{pmatrix} \\ &= E(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var}(X_p) \end{pmatrix} \end{aligned}$$

The **correlation matrix** is given by $\mathbf{R} = (\rho_{ij})_{i,j=1,\dots,p}$ with $\rho_{ii} = 1$ and $\rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$.

Rules

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y})$$

$$E(a\mathbf{X} + b) = aE(\mathbf{X}) + b$$

$$\text{Cov}(\mathbf{X}) = E(\mathbf{X}\mathbf{X}') - \boldsymbol{\mu}\boldsymbol{\mu}'$$

$$\text{Var}(\mathbf{a}'\mathbf{X}) = \mathbf{a}'\text{Cov}(\mathbf{X})\mathbf{a} = \sum_{i,j=1}^p a_i a_j \sigma_{ij}$$

$$\text{Cov}(\mathbf{A}\mathbf{X} + \mathbf{b}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}'$$

$\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ and \mathbf{R} is symmetric and positive semidefinite.

Let $\mathbf{Z} = (\mathbf{X}', \mathbf{Y}')'$, where \mathbf{X} and \mathbf{Y} are p and q -dim. Then it holds

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{Z}} &= (\boldsymbol{\mu}'_{\mathbf{X}}, \boldsymbol{\mu}'_{\mathbf{Y}})' \\ \boldsymbol{\Sigma}_{\mathbf{ZZ}} &= \begin{pmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{xx} \end{pmatrix} = \begin{pmatrix} \text{Cov}(\mathbf{X}) & \text{Cov}(\mathbf{X}, \mathbf{Y}) \\ \text{Cov}(\mathbf{Y}, \mathbf{X}) & \text{Cov}(\mathbf{Y}) \end{pmatrix}. \end{aligned}$$

Note: $\boldsymbol{\Sigma}_{xy} = \boldsymbol{\Sigma}_{yx}$.

Independent random vectors

up to now: independence of events

Recall: two events A_1 and A_2 are independent, if $P(A_1 \cap A_2) = P(A_1) \cdot P(A_2)$. Then it holds $P(A_1 | A_2) = P(A_1)$.

Example: A_1 = „success of a therapy“, A_2 = „a drug was given“.

X_1, \dots, X_n are (stochastically) independent, if it holds for all $x_1, \dots, x_n \in \mathbb{R}$

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = \prod_{i=1}^n P(X_i \leq x_i).$$

Note:

- If X_1, \dots, X_n are independent and g_1, \dots, g_n are function, then $g_1(X_1), \dots, g_n(X_n)$ are also independent.
- Let f be the probability function (density) of (X_1, \dots, X_n) and let f_i denote the probability function (density) of X_i .

X_1, \dots, X_n are independent if and only if

$$f(x_1, \dots, x_n) = \prod_{i=1}^n f_i(x_i)$$

holds for all $x_1, \dots, x_n \in \mathbb{R}$

Example: toss two symmetric dice: X_1 = number on the first die, X_2 = number on the second die

$$P(X_1 = i, X_2 = j) = 1/36 = P(X_1 = i) P(X_2 = j)$$

The random variables X_1 and X_2 are independent .

Marginal distributions

Let a $p + q$ -dim. vector \mathbf{Z} be partitioned into $\mathbf{Z} = (\mathbf{X}', \mathbf{Y}')'$, such that \mathbf{X} and \mathbf{Y} are p and q dim. respectively.

$$F_{\mathbf{X}}(\mathbf{x}) = P(\mathbf{X} \leq \mathbf{x}) = F_{\mathbf{Z}}(x_1, \dots, x_p, +\infty, \dots, +\infty)$$

$$f_{\mathbf{X}}(\mathbf{x}) = \int_{-\infty}^{+\infty} f(\mathbf{x}, \mathbf{y}) d\mathbf{y}$$

Independency

Def: \mathbf{X} and \mathbf{Y} are independent iff

$$F_{\mathbf{Z}}(\mathbf{x}, \mathbf{y}) = F_{\mathbf{X}}(\mathbf{x})F_{\mathbf{Y}}(\mathbf{y}) \quad \text{or} \quad f_{\mathbf{Z}}(\mathbf{x}, \mathbf{y}) = f_{\mathbf{X}}(\mathbf{x})f_{\mathbf{Y}}(\mathbf{y}).$$

Conditional distributions

We consider the distribution of the explained variables \mathbf{y} conditional on a set of explanatory variables \mathbf{x} .

$$f(\mathbf{y}|\mathbf{x}) = \frac{f(\mathbf{x}, \mathbf{y})}{f_{\mathbf{X}}(\mathbf{x})}$$

↪ The conditional expectation plays a key role in econometrics and a large portion of research is aimed to estimate it.

$$E(\mathbf{y}|\mathbf{x}) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \mathbf{y} f(\mathbf{y}|\mathbf{x}) d\mathbf{y}.$$

For $\mathbf{x} = f(\mathbf{w})$ it holds that:

$$E(\mathbf{y}|\mathbf{x}) = E[E(\mathbf{y}|\mathbf{w})|\mathbf{x}]$$

$$E(\mathbf{y}|\mathbf{x}) = E[E(\mathbf{y}|\mathbf{x})|\mathbf{w}]$$

$$E(\mathbf{y}|\mathbf{x}) = E[E(\mathbf{y}|\mathbf{x}, \mathbf{z})|\mathbf{x}]$$

$$E[E(\mathbf{y}|\mathbf{x})] = E(\mathbf{y})$$

Transformation of random variables

Requirement: X_1 and X_2 are independent.

- If X_1 and X_2 are discrete, then

$$\begin{aligned}
 P(X_1 + X_2 = x) &= \sum_{\substack{u, t \\ u + t = x}} P(X_1 = u, X_2 = t) \\
 &\stackrel{\text{indep.}}{=} \sum_t P(X_1 = x - t) P(X_2 = t).
 \end{aligned}$$

- Let f_1 and f_2 be the densities of X_1 and X_2 . Then the density of $X_1 + X_2$ is given by

$$f_{X_1+X_2}(x) \stackrel{\text{indep.}}{=} \int_{-\infty}^{\infty} f_1(x-t) f_2(t) dt.$$

Implications:

If X_1, \dots, X_n are independent with

- $X_i \sim N(\mu_i, \sigma_i^2)$, then

$$\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right).$$

- $X_i \sim N(\mu, \sigma^2)$, then

$$\bar{X} \sim N(\mu, \sigma^2/n).$$

- $X_i \sim B(1, p)$, then

$$\sum_{i=1}^n X_i \sim B(n, p).$$

Lemma 1: Let $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$, where \mathbf{A} is a $q \times p$ -matrix with $rg(\mathbf{A}) = q \leq p$. Then $\mathbf{Y} \sim \mathcal{N}_q(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$.

Lemma 2: Let $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\mathbf{Y} = \boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$, where $\boldsymbol{\Sigma}^{-1/2}$ is the Cholesky decomposition of matrix $\boldsymbol{\Sigma}$. Then $\mathbf{Y} \sim \mathcal{N}_p(\mathbf{0}, \mathbf{I})$.