

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Repetition of statistics

Julius Fenn¹

¹University of Freiburg

PhD student at the Institute of Experimental Psychology Freiburg (Cognition, Action, and Sustainability)

14.12.2021

Any questions?

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature



Motivation statistics

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Central areas of life science, among others, deal with numbers.
- Empirical phenomena are usually not deterministic, but stochastic
- an intuitive approach is not sufficient for the description of a variety of empirical phenomena, so that a mathematically sound treatment is necessary
- Statistics is a key qualification that is increasingly demanded by the labour market
- ...

Motivation statistics: Investment in the education system

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

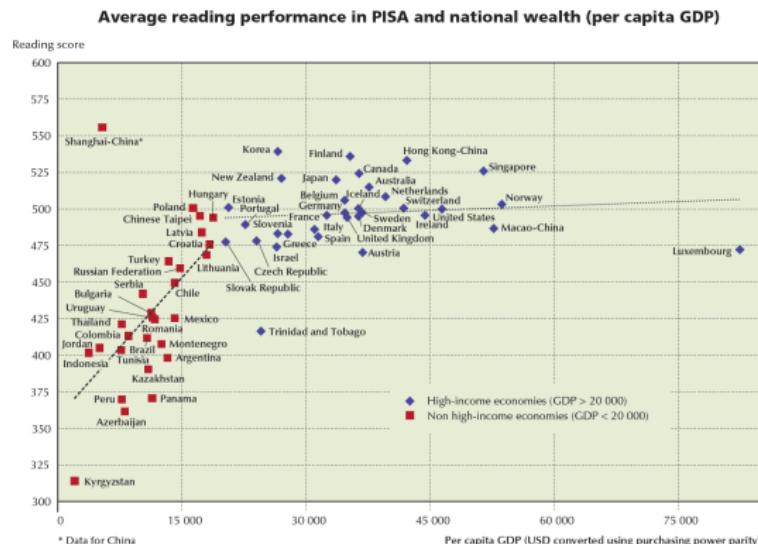
Distributions

discrete

continuous

Literature

data → statistic → political decisions



PISA in Focus 13

O'Connell (2019)

Motivation statistics: unstoppable (?) disruptive technologies

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Table 1. Five successive technological revolutions, 1770s to 2000s

<i>Technological revolution</i>	<i>Popular name for the period</i>	<i>Core country or countries</i>	<i>Big-bang initiating the revolution</i>	<i>Year</i>
First	The "Industrial Revolution"	Britain	Arkwright's mill opens in Cromford	1771
Second	Age of steam and railways	Britain (spreading to Continent and USA)	Test of the "Rocket" steam engine for the Liverpool-Manchester railway	1829
Third	Age of steel, electricity, and heavy engineering	USA and Germany forging ahead and overtaking Britain	The Carnegie Bessemer steel plant opens in Pittsburgh, Pennsylvania	1875
Fourth	Age of oil, the automobile, and mass production	USA (with Germany at first vying for world leadership), later spreading to Europe	First Model-T comes out of the Ford plant in Detroit, Michigan	1908
Fifth	Age of information and telecommunications	USA (spreading to Europe and Asia)	The Intel microprocessor is announced in Santa Clara, California	1971

Kaldor (2018)

10 most expensive companies worldwide

Course structure

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- ① fundamentals: basic terms, data collection, excursus on measurement, scale levels
- ② descriptive, explorative statistics
 - univariate: graphs, measures, symmetry of distributions
 - multivariate: graphs, correlation; excursus term theory
- ③ probabilities, randomness
- ④ distributions: discrete, continuous
- ⑤ (inferential statistical like hypothesis testing, linear regression, ...)

central learning platform: ILIAS

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

You can find all literature recommendations, slides, ... on ILIAS „**WS 21 - Fundamentals of Questionnaire Design and Analysis using lab.js and R**“

central terms: Describe, Search, Conclude

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Statistics represents the theoretical basis of all empirical research and is a subfield of mathematics¹

- **descriptive statistics:** descriptive and graphical processing and compression of data, as well as for the detection of errors in the data set (data validation).
- **explorative statistics:** search for structures and peculiarities in data in order to arrive at new questions and the choice of suitable statistical models
- **inductive statistics:** using probability theory / stochastics to reach more general conclusions beyond the surveyed sample

¹Separation according to the method controversy into two camps using purely quantitative / qualitative methods is (from my point of view) extremely harmful, as the methods complement and overlap each other, see e.g. Wolf (2008), Kelle (2008)

central terms

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **statistical units:** Objects over which quantities (characteristics) of interest are observed
- **characteristics:** Observed quantities of the statistical units, e.g.
 - characteristic X_{gender} with the (characteristic) expressions x_{female} , x_{male} , x_{diverse} (also observation).
- **population:** Set of all statistical units about which one wants to arrive at statements
- **sample:** subset of the population under investigation
 - Sample should be as representative as possible (see explanations from slide 12 onwards)

central terms

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **variable:** characteristics which can take at least two different values
- **constant:** characteristic with only one value

ID	IDZP	GR	GES	OG	VG	CC	KU	IQ	KG	IQ.cat
1	1	1	1 2	1913.88	1005	6.08	54.7	96	57.607	unterdurchschnittlich
2	2	1	2 2	1684.89	963	5.73	54.2	89	58.968	unterdurchschnittlich
3	3	2	1 2	1902.36	1035	6.22	53	87	64.184	unterdurchschnittlich
4	4	2	2 2	1860.24	1027	5.8	52.9	87	58.514	unterdurchschnittlich
5	5	3	1 2	2264.25	1281	7.99	57.8	101	63.958	mittel
6	6	3	2 2	2216.4	1272	8.42	56.9	103	61.690	mittel
7	7	4	1 2	1866.99	1051	7.44	56.6	103	133.358	mittel
8	8	4	2 2	1850.64	1079	6.84	55.3	96	107.503	unterdurchschnittlich
9	9	5	1 2	1743.04	1034	6.48	53.1	127	62.143	überdurchschnittlich
10	10	5	2 2	1709.3	1070	6.43	54.8	126	83.009	überdurchschnittlich
11	11	6	2 1	1689.6	1173	7.99	57.2	101	61.236	mittel
12	12	6	1 1	1806.31	1079	8.76	57.2	96	61.236	unterdurchschnittlich
13	13	7	2 1	2136.37	1067	6.32	57.2	93	83.916	unterdurchschnittlich
14	14	7	1 1	2018.92	1104	6.32	57.2	88	79.380	unterdurchschnittlich
15	15	8	2 1	1966.81	1347	7.6	55.8	94	97.524	unterdurchschnittlich
16	16	8	1 1	2154.67	1439	7.62	57.2	85	99.792	unterdurchschnittlich
17	17	9	1 1	1767.56	1029	6.03	57.2	97	81.648	unterdurchschnittlich
18	18	9	2 1	1827.92	1100	6.59	56.5	114	88.452	überdurchschnittlich
19	19	10	2 1	1773.83	1204	7.52	59.2	113	79.380	mittel
20	20	10	1 1	1971.63	1160	7.67	58.5	124	72.576	überdurchschnittlich

Study: Tramo et al. 1998 Brain size, head size, and intelligence quotient in monozygotic twins

Study of identical twins

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Sample size $N = 20$

Variables:

- **ID:** Identifier of the individual person
- **IDZP:** Identifier of the twin pair
- **GR:** Birth order
- **GES:** Sex (1 = male, 2 = female)
- **OG:** Surface area of the cerebral cortex in cm^2
- **VG:** Volume of the forebrain in cm^3
- **CC:** Area of the corpus callosum in cm^2
- **KU:** Circumference of head in cm
- **IQ:** Intelligence quotient
- **KG:** Body weight in kg
- additionally **IQ.cat:** grouped intelligence quotient in 3 groups

Data collection: basics

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

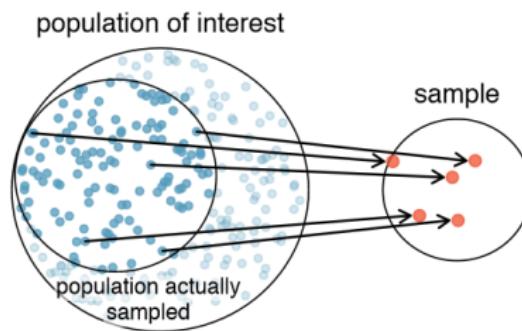
discrete

continuous

Literature

Sampling theory (*sampling*) is a separate field of research with the aim of exploring methods of sampling without making systematic errors

Example of systematic error - causes e.g. *non-response bias*, incomplete *frame population* (see slide 17):



Graphics taken from: [OpenIntro Statistics \(Diez et al 2019\)](#); Zucchini et al. (2009), 27ff.

Data collection: sampling methods I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

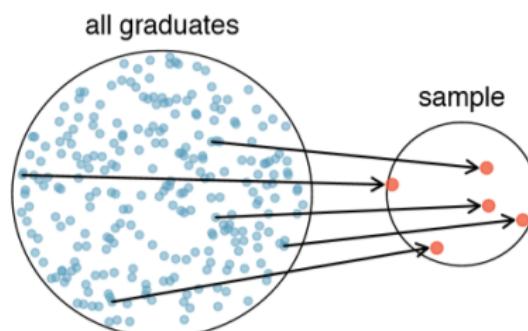
discrete

continuous

Literature

- **simple random sample:** random selection of statistical units; each object has the same probability of being included in the sample.
 - Advantage: systematic errors are avoided; can be representative

Example simple random sample:



Data collection: sampling methods II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

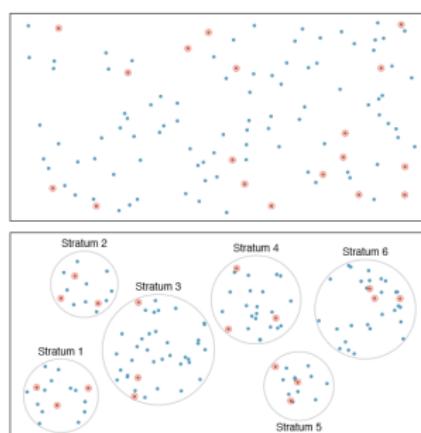
discrete

continuous

Literature

- **simple random sample**
- **stratified random sample.** Population is divided into groups (*strata*) with the aim of finding strata that are similar in their characteristics - „divide-and-conquer sampling strategy“

Example (1) simple random sample, (2) stratified random sample:



Data collection: sampling methods III

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

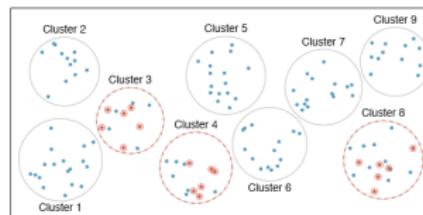
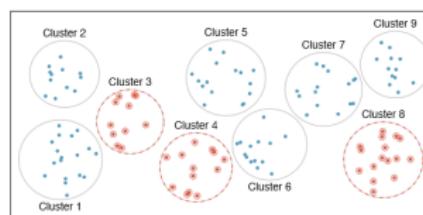
discrete

continuous

Literature

- **cluster sampling:** also groups (mostly according to natural concentrations of study units) formed; no systematic difference in their characteristics and thus small images of the population

Example cluster sampling (*multistage sampling*):



Data collection: sampling methods IV

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

other procedures (but these are methodologically criticised !):

- **quota selection:** certain quotas of samples are fixed from the outset (e.g. 90% of all education science students are male).
- **selection of typical cases:** statistical objects are selected as typical representatives of the population according to subjective criteria
- **selection according to convenience (*convenience sample*):** Statistical objects that are easy to reach have a higher probability of being included in the sample
- ...

Framework data collection

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

The population from which the sample is drawn determines the general validity of the results and conclusions!



→ Diagram can be rewritten as $T \subseteq F \subseteq R \subseteq S$

- Presence of non-response: $S \subset R$
- List of statistical objects to be drawn (*frame population*) not complete, *target population* not covered: $F \subset T$

Take Home Messages: central terms, data collection

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Statistics can be divided into descriptive, exploratory and inductive statistics.
- Basic concepts allow the description of data sets (row are the statistical units and columns are the variables; the observed value of a variable represents a characteristic)
- Knowledge about methods of data collection allow (already within the experimental design) to critically question the generalisability of the results

Level of measurement: briefly measurement

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

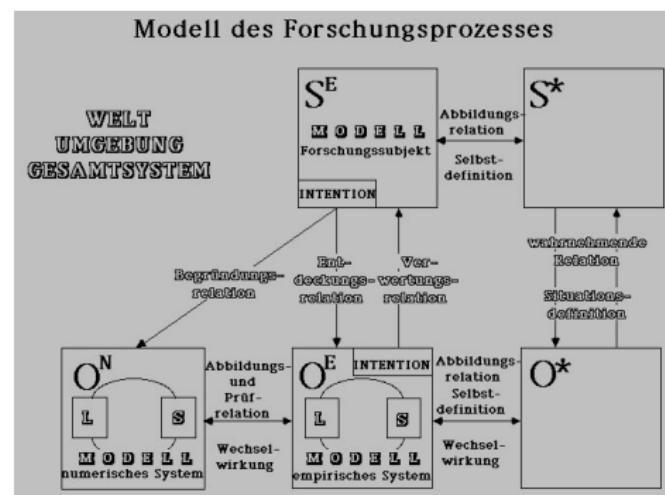
discrete

continuous

Literature

Level of measurement indicates the information contained in a **measurement**

Definition of measurement: the entire process of observing empirical facts (research as a total system) up to the allocation of numbers to the observed phenomena (mapping relation).



Level of measurement: nominal

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Expressions are names or categories assigned to statistical units
- Examples:
 - characteristic X_{gender} with the expressions x_{female} , x_{male} , x_{diverse}
 - Characteristic $Y_{\text{Hair colour}}$ with the expressions y_{brown} , y_{blond} , y_{other}
- lowest level of information
- nominal scaled variables can be assigned any labels and symbols / numbers in the dataset

allowed are injective transformations:

a, b, x, y Values of a nominal scale where $a = b$ and $x \neq y$
 $\rightarrow f(a) = f(b)$ und $f(x) \neq f(y)$

Erklärvideos: [studyflix Skalenniveau](#)

Level of measurement: ordinal

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Expressions can be ordered, but intervals are not interpretable
- Example:
 - school grades² 1 better 2 but the grade interval between 1.2 and 4.5 cannot be compared / interpreted
 - characteristic values are often assigned to ranks, these fulfil the ordering property of numbers (numbers can be put in an order with regard to their height)

allowed are strictly monotonic transformations:

the less-than relation is maintained if the following applies $x < y \rightarrow f(x) < f(y)$

²School grades are often assumed to be interval-scaled, but their measurement-theoretical goodness is critical, see e.g.. Neumann (2007), Lintorf (2012)

Level of measurement: interval

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- distances between numbers can be interpreted in a meaningful way
- there is no natural zero (not useful to form quotients)
- Examples:
 - Temperature in Celsius (set 0° at freezing point, 100° at boiling point water).
 - Results of an intelligence test (often set at IQ: 100 ± 15).

allowed are positive affine transformations:

$$f(x) = a + b * x \text{ mit } b > 0$$

Note: Quadratic transformations are not allowed

Level of measurement: ratio

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- distances between numbers can be interpreted meaningfully and statements about quotients (e.g. person 1 is twice as tall as person 2) are possible
- the ratio scale has a natural zero point and a defined unit of measurement.
- Examples:
 - temperature in Fahrenheit
 - Body height in mm, cm, m
- highest level of information

allowed are formal transformations:

$$f(x) = b * x \text{ mit } b > 0$$

Level of measurement: summary

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Skalenart	count	order	form differences	form ratios
nominal	yes	no	no	no
ordinal	yes	yes	no	no
interval	yes	yes	yes	no
ratio	yes	yes	yes	yes

Tabelle: Meaningful calculations for the different scale levels

- Nominal scaled variables are also called **qualitative variables**
- Interval scaled and ratio scaled variables are also called **quantitative variables**
- Nominal scaled and ordinal scaled variables represent **discrete variables**.
- Interval-scaled and ratio-scaled variables represent **metric variables**

Level of measurement: classification of variables

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **qualitative:** Variables can reflect different categories of a characteristic
- **quantitative:** Variables can reflect a different extent or intensity of a characteristic
 - ordinal scaled variables are in a grey area here, because values can be ordered, but distances between the characteristic expressions (quantitative aspect) cannot be interpreted.
- **discrete:** set of expressions of a variable can be represented by the first n natural numbers ($\mathbb{N} = \{1, 2, \dots, n\}$) and is thus countable
- **continuous:** between two values there can always be (theoretically) infinitely many other values, range of values are all real numbers ($\mathbb{R} = \{-\infty, \infty\}$)

Take Home Messages: level of measurement

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

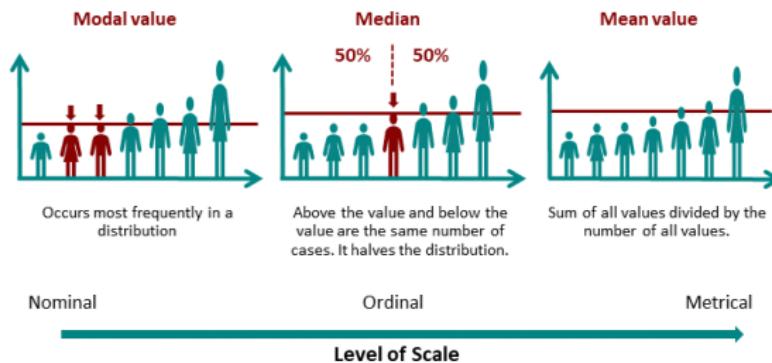
Distributions

discrete

continuous

Literature

- Scale level of a variable indicates the content of information of a measurement.
- the information content of the scales is as follows: ratio > interval > ordinal > nominal
- scale levels have a **downward compatibility** any procedure that we may apply to a low scale level we can also apply to a higher one



R-commands section level of measurement

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

R is an object-oriented language and it must override the „object“ - i.e. the variable of the record - with the following syntax structure:

```
variable <- command(variable)
```

- nominal scale: `factor(x = daten$variable)`

- ordinal scale:

```
factor(x = daten$variable, ordered = TRUE)
```

- metric scales:

```
as.numeric(as.character(x = daten$variable))
```

→ the distinction between interval and ratio scales is not relevant for
R (as well as e.g. SPSS)

note: furthermore, **R** distinguishes between integer (\mathbb{N}) and continuous (\mathbb{R}) variables by the commands: `as.integer(x =)` resp.
`as.numeric(x =)`

Data generation process

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Data generation process: „When [...] data set has been created, a good model has to be constructed for the data generation process (DGP). This is the stage at which the actual [...] analysis begins... When an [...] model has been constructed for the DGP, it should **only be used for the analysis if it reflects the ongoing in the system of interest properly**“ – Lütkepohl und Krätsig (2004)

→ a model should make a parsimonious, plausible and substantially meaningful contribution to the data generation process – Hoyle (2012)

→ in statistics, we assume that we can represent the DGP by discrete and continuous distributions

Distributions I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

a random variable is defined by assigning numbers to an event of a random process (realisation), for this purpose **distributions** are necessary:

- Random processes can only be described by functional relations $f(X)$, types of representation:
 - for discrete random variables → probability function
 - for continuous random variables → density function
- distributions are determined by functions and their parameters

A random variable $f(X)$ is a function that assigns a real number to each possible outcome of a random experiment / process.

further information see: [Wikipedia random variables](#); Fahrmeir et al. (2016); Zucchini et al. (2009)

Distributions II: example I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Random variables can only take on finitely many values or countably infinitely many values

→ each possible event can be assigned a natural number

Example: rolling 6 or 30 fair dices several times at once and we want to know $A = \{\text{result is 1}\} = \{1\} \rightarrow P(\{1\}) = 1/6$ (binomial distributed):

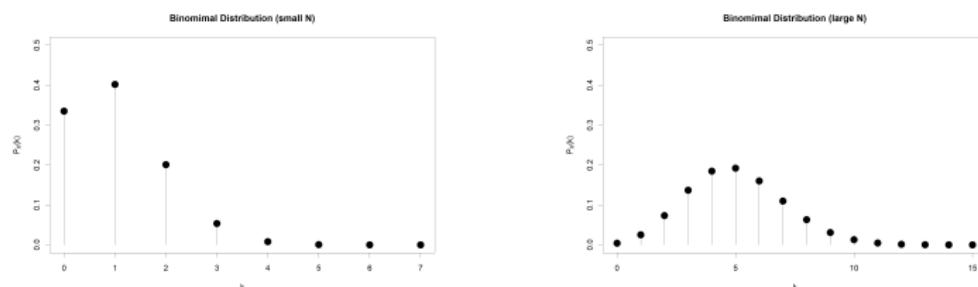


Abbildung: left rolling 6 dices, right 30 dices

Is rolling a dice a Gauss distribution?

Distributions II: example II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

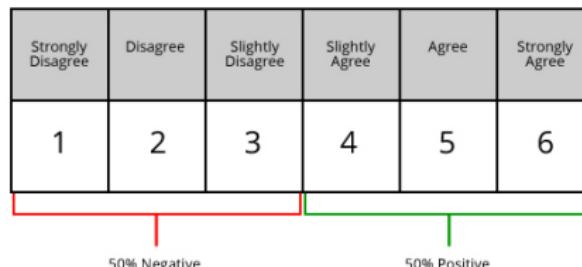
continuous

Literature

single items can be binomial distributed (! however the real distribution is never known)

Imagine that participants should answer the following two items on a 6 point Likert Scale:

- Should something be done about climate change for future generations?
- Would you be willing to accept personal restrictions for the fight against climate change?



→ What would you expect?

Take-Home-Message: surveys

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- individual items of a Likert scale result from a binomially distributed process ($Y_1 \sim Bin(n, p)$, ...)
 - we want to achieve a high variance in the items (all response options are used)
 - variance or covariance (if multiple items) is the most important building block for the statistical models
 - the answering of items should be independent from other people or the ordering of the items
- our statistical models should reflect the ongoing in the system of interest properly / sufficiently and should be theoretically derived

Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful

George Box (1987)

Descriptive statistics: Motivation

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- to be able to analyse data, one needs an impression of the available information
- data sets are so large that direct viewing of the data is not useful
- goals include:
 - description and distribution of individual variables
 - detection of data errors (implausible values, outliers,...)
 - investigation of correlations between variables

Descriptive statistics: Fundamentals

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

In general, statistical procedures can be divided into how many variables are considered at once:

- univariate: Representation of values one variable (also: denotes dependence on only one variable, one-dimensional)
- bivariate: Representation of the values of two variables
- multivariate: Representation of the values of several variables

→ the statistical procedures and the assumed distributions behind them (see from slide ??) become increasingly complex

Descriptive statistics: univariate

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

The aim is the representation, description of data consisting of observations of a single characteristic (univariate, one-dimensional data)
→ Comparison of several univariate results enables impressions about possible dependencies that can be further analysed using multivariate procedures

univariate representations / statistical methods:

- frequency distributions
 - frequency table
 - bar chart / histogram
- central measures (from slide 48)
- symmetry of distributions (from slide 82)

Frequency table I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Initial situation: Observed characteristics x_1, \dots, x_n of a feature X
→ Raw data is summarised according to different occurring counts / characteristics a_1, \dots, a_k , where $k \leq n$ (for categorical characteristics k is equal to the number of categories)

Representation of a part of the original list for the variable IQ.cat (grouped intelligence quotient in 3 groups):

1 below average	2 below average	3 - 18 ...	19 medium	20 above average
--------------------	--------------------	---------------	--------------	---------------------

1	2	3 - 18	19	20
1	1	...	2	3

→ here is an ordinal scaled variable, thus $3 > 2 > 1$ (this would not be the case with a nominal scaled variable)

Frequency table II

[Motivation](#)[Content](#)

Fundamentals

[Data collection](#)[Level of measurement](#)[Surveys](#)

Descriptive statistics

[univariat](#)

Graphs

[Measures](#)[Symmetry](#)[multivariate](#)

Graphs

[Measures](#)

Theory

Probability

[Randomness](#)

Distributions

[discrete](#)[continuous](#)

Literature

Frequency table for the variable IQ.cat:

Characteristic	class	absolute frequency h	relative frequency f
below average	$85 < \dots \leq 95$	11	$11/20 = .55$
medium	$96 < \dots \leq 113$	5	$5/20 = .25$
above average	$114 < \dots \leq 127$	4	$4/20 = .20$

- absolute frequency h : $h(a_j) = h_j$, number of x_i from x_1, \dots, x_n with $x_i = a_j$.
- relative frequency f : $f(a_j) = f_j = h_j/n$

Bar chart

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Row chart: a bar perpendicular to the x-axis with absolute heights h_1, \dots, h_k or relative heights f_1, \dots, f_k is plotted over different occurring characteristics a_1, \dots, a_k .
- Bar chart: like bar chart but with x-axis placed vertically instead of horizontally

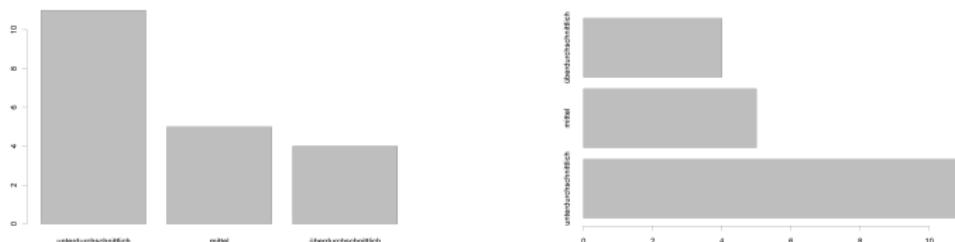


Abbildung: left bar chart, right row chart

Histogram I: construction

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

The construction of a histogram is done under the principle of area fidelity, so that the areas are proportional to the absolute h_j or relative frequencies f_j :

- ① Data are grouped into classes c_0, \dots, c_k , where these have adjacent intervals: $[c_0, c_1), \dots, [c_{k-1}, c_k)$
 - Rectangles in the histogram have the width $d_j = c_j - c_{j-1}$
- ② Height of the rectangles in the histogram are defined by *frequency/classwidth*
 - absolute frequencies: h_j/dj
 - relative frequencies: f_j/dj
- ③ Area over the intervals results from $\text{Area} = \text{Width} \times \text{Height}$
→ if class widths d_j are equal, the absolute h_j or relative frequencies f_j are simply deducted as height

Histogram II: example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

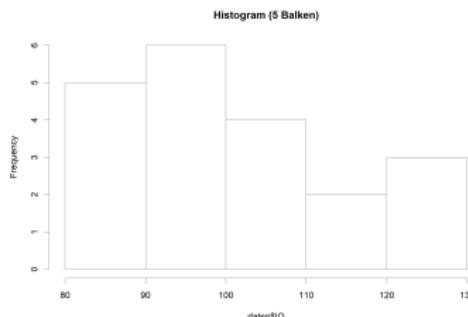
continuous

Literature

Representation of a histogram for the variable IQ:

$(c_{j-1}, c_j]$	h_j	b_j	h_j / b_j
$80 \leq x \leq 90$	5	10	5
$90 \leq x \leq 100$	6	10	6
$100 \leq x \leq 110$	4	10	4
$110 \leq x \leq 120$	2	10	2
$120 \leq x \leq 130$	3	10	3

→ all class widths d_j are the same size



Histogram III: number of classes

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Formation of classes always represents a certain loss of information (*data aggregation*) → Number of classes has a decisive influence on the visual impression of a distribution

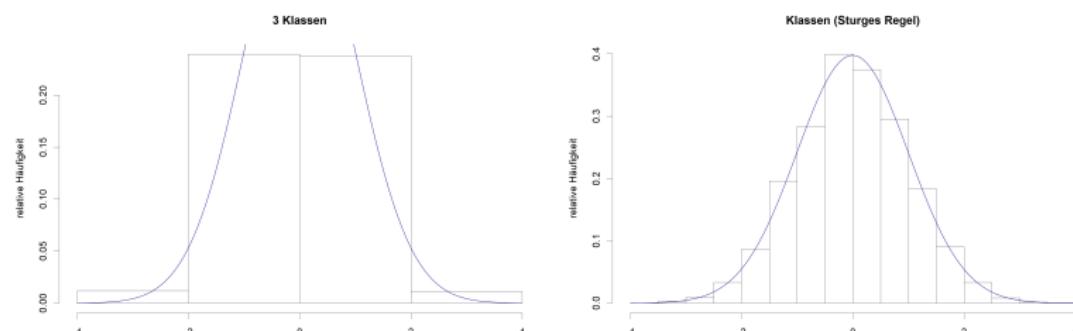


Abbildung: left 3 classes, right Sturges rule*

→ blue line represents the true distribution, this is a standard normal distribution $X \sim \mathcal{N}(0, 1)$ (see later from slide 153)

Histogram IV: number of classes

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

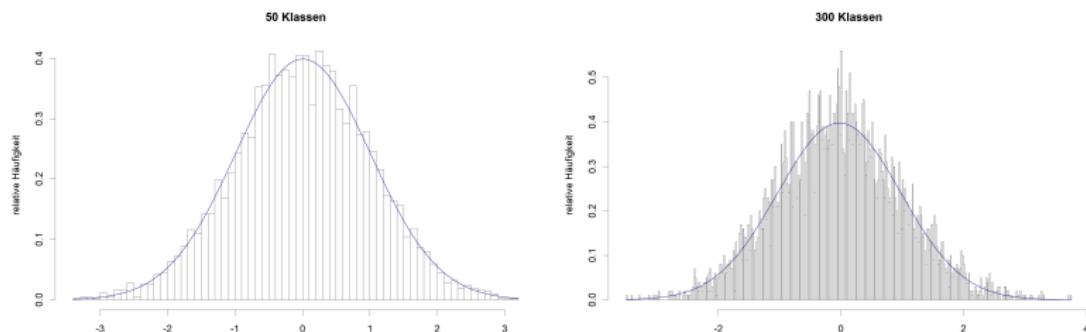


Abbildung: left 50 classes, right 300 classes

Histogram V: number of classes

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Number of classes significantly influences the shape of the histogram

- too low number give only a smooth estimate → low variability
- too high number give a good approximation of the data, but a very rough estimate → large variability

this is a classic **bias-variance trade off**, i.e.

- J small: large bias but low variability
- J large: small bias but large variability

Number of classes is often estimated by the statistics programme using rules of thumb, e.g. using the *Sturges rule: $\lfloor(1 + \log_2 n)$

Cumulative frequency distribution: construction

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

In contrast to the histogram, allows the most exact representation of the measured values (! no loss of information) with metric variables

- absolute cumulative frequency distribution: determine the number of observations (characteristic values) that are less than or equal to x : $H(x) = \text{Anzahl der Werte } x_i \text{ mit } x_i \leq x$
 - these absolute frequencies can be written as $H(x) = h(a_1) + \dots + h(a_j)$ (sum of the steps)
- relative cumulative frequency distribution (also empirical distribution function): division of the absolute frequencies by the number of observations: $F(x) = H(x)/n$

→ these are monotonically increasing staircase functions which jump upwards at the values a_1, \dots, a_k .

Cumulative frequency distribution: example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

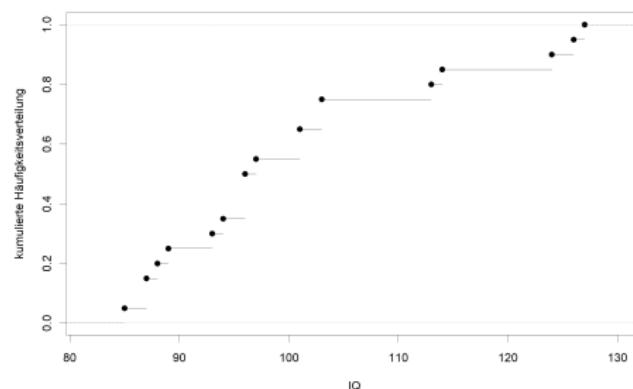
Randomness

Distributions

discrete

continuous

Literature



Representation of the variable IQ in the data set of Trame et al. 1998
(see slide 10)

Take Home Messages: Choice of presentation

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- discrete variables with few expressions: pie chart and bar chart
- metric variables, but also discrete variables with many values: Histograms (! certain loss of information)
- metric variables with many expressions: cumulative frequencies (! no loss of information, but the shape of the distribution is not visible - see from slide 82)

→ Linking level of measurement with forms of representation

R-commands section univariate statistics using graphs

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Create a frequency table: `table(x)`
- Bar chart: `barplot(table(x), horiz = TRUE)`, Säulendiagramm: `barplot(table(x), horiz = FALSE)`; here, the argument `horiz` only specifies whether the bars are horizontal or vertical
- Histogram with absolute frequencies: `hist(x)`, Histogram with relative frequencies: `hist(x, freq = FALSE)`
- Cumulative frequency distribution: `plot(x = ecdf(x))`

Measures: Basics I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Measures (*also parameters*) summarise the essential properties of a frequency distribution, these are divided into:
 - Location parameters / measures of central tendency: represent central tendency, centre of gravity of data.
 - dispersion parameter / dispersion: represent distributions / dispersions of data around the centre of gravity
 - shape parameter / skewness: represent the shape of the distribution of data (statements about symmetry, see from slide 82)
- Note: some of the measures correspond to moments of random variables (see from slide 143)

Measures: Basics II

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Which measure is appropriate depends on the scale level (and context, as well as the question):

Beschreibung der	Statistik	Skalenniveau
zentralen Tendenz	Mittelwert	mind. Intervallskala
	Median	mind. Ordinalskala
	Modus	mind. Nominalskala
Position in der Verteilung	Quantile (Quartile, Dezile)	mind. Ordinalskala
Streuung	Varianz	mind. Intervallskala
	Standardabweichung	mind. Intervallskala
	Interquartilabstand	mind. Intervallskala
	Spannweite	mind. Intervallskala
Schiefe	Schiefekoeffizient	mind. Intervallskala
	Quartilskoeffizient	mind. Ordinalskala
Wölbung	Kurtosiskoeffizient	mind. Intervallskala

Measures: Basics II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Which measure is appropriate depends on the scale level (and context, as well as the question):

Scale type	Empirical properties	Permissible statistics	Examples
Nominal	Equivalence	Mode, chi square	Eye colour, place of birth, etc...
Ordinal	Equivalence, order (greater or less)	Median, percentile	Surface hardness, military rank, etc...
Interval	Equality, order, distance (addition or subtraction)	Mean, standard deviation, correlation, regression, analysis of variance	Temperature in °C, serial numbers, etc...
Ratio	Equality, order, distance, ratio (multiplication or division)	All statistics permitted for interval scales plus the following: geometric mean, harmonic mean, coefficient of variation, logarithms	Temperature in K, weight, age, number of children, etc...

Measures of central tendency: Mean - formula

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: metric variables

Formula:

$$\bar{x} = \frac{1}{n}(x_1 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$$

Formula for stratification for r group averages:

$$\bar{x} = \frac{1}{n}(n_1\bar{x}_1 + \dots + n_r\bar{x}_r) = \frac{1}{n} \sum_{j=1}^r n_j x_j$$

briefly sum sign

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

The scary-looking Sum symbol

$$\sum_{i=1}^5 i^2$$

The last value of i → 5

The thing to add up → i^2

The first value of i → $i=1$

To be read as: "The sum from 1 to 5 of i squared"

- Sum sign consists of the capital greek letter Σ followed by a function
- running index (also count variable), this appears again in the function
- the running index starts from its start value and ends at the given end value (*here 5 steps*)

Measures of central tendency: Mean - properties I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

the mean value has a so-called **center of gravity property**, sum of deviations between x_i and \bar{x} disappears:

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

→ no statement about the dispersion of the data possible, therefore choice of squared differences - the smaller this value, the better the mean value represents the data:

$$\bar{x} = \min_{c \in \mathbb{R}} \left(\sum_{i=1}^n (x_i - c)^2 \right) \quad (1)$$

Note: Term states that no value c from the real numbers \mathbb{R} other than the mean \bar{x} leads to a smaller sum of the squares of the deviations

Measures of central tendency: Mean - properties II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

How representative are the group means (green) for the available data (example: test performance of 5 schools)?

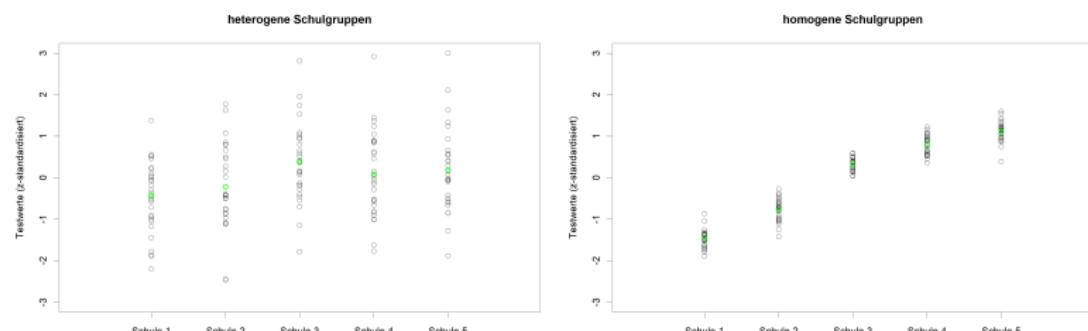


Abbildung: Heterogeneous school groups on the left, homogeneous school groups on the right (group mean representative)

Measures of central tendency: Median

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: at least ordinal scaled variables

→ the median is **robust** / **resistant** to extreme values (also outliers)³

„calculation“: Median is placed in the middle of the data so that at least 50% of the values are less than / equal to and at the same time at least 50% of the values are greater than / equal to the median.

$$x_{med} \begin{cases} = x_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \in [x_{(\frac{n}{2})}, x_{(\frac{n+1}{2})}] & \text{if } n \text{ even} \end{cases}$$

→ for even n, the following rule is often used:

$$x_{med} = \frac{1}{2}(x_{(\frac{n}{2})} + x_{(\frac{n+1}{2})})$$

³extreme values: measured value that does not meet expectations (! subjective), often values (arbitrarily) are called outliers that lie outside 1.5 times the quartile distance (see from slide 87, as well as from slide ??)

Measures of central tendency: Median - properties I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

the median is exactly in the centre of the data

$$x_{med} = \min_{c \in \mathbb{R}} \left(\sum_{i=1}^n |x_i - c| \right) \quad (2)$$

Note: Term states that no value c from the real numbers \mathbb{R} other than the median leads to a smaller sum of the deviation

Measures of central tendency: Median - properties II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

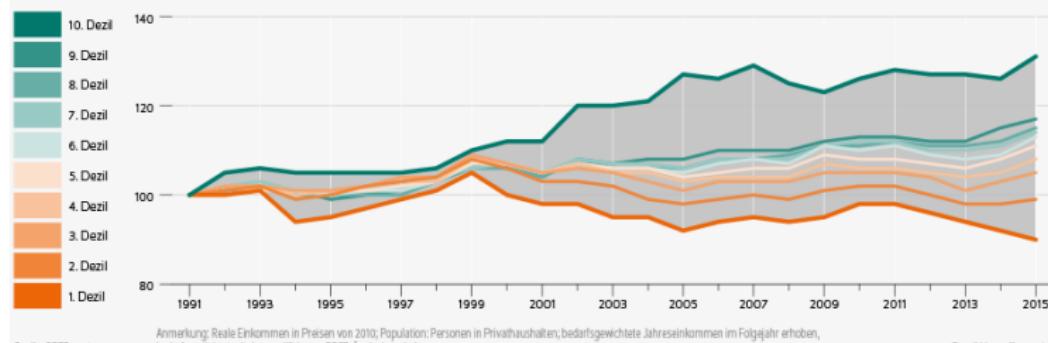
continuous

Literature

Is the mean income in Germany determined by the mean or median?

Die Einkommen der Gruppen am unteren Ende der Verteilung sind seit 1991 gesunken

Entwicklung des durchschnittlichen verfügbaren Haushaltseinkommens nach Dezilen, in Prozent (1991 = 100)



German Institute for Economic Research: Real Income Germany between 1991 and 2015

vs.

Wikipedia Income distribution in Germany
current report: WSI-Verteilungsbericht 2019

Measures of central tendency: Mode

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: at least nominally scaled variables, most important positional measure for categorical variables
the mode, like the median, is **robust / resistant** to extreme values

„Calculation“: x_{mod} is the expression with the highest frequency, it is only unique if the frequency distribution has only one maximum

Measures of central tendency: Summary

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

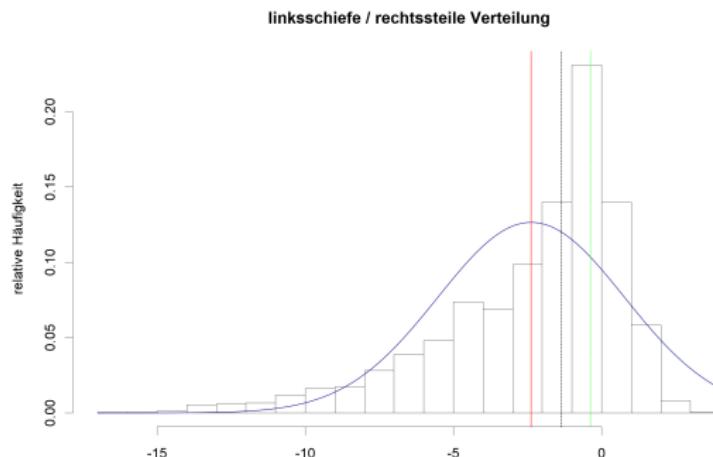
Distributions

discrete

continuous

Literature

- mean indicates centre of gravity of the data, cf. equation (1)
- Median indicates the midpoint (50%) of the data, cf. equation (2)
- Mode indicates the most frequently occurring value



vertical lines from left to right: mean, median, mode

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs
Measures

Symmetry

multivariate

Graphs
Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- mean

- metric features
- also useful for binary 0 / 1 coded characteristics
- partly used for ordinal scaled variables (like grades) (! interpret with care)

- median

- ordinal scaled variables

- mode

- Applicable to any scale level, but rarely interpretable for metric variables

Measures of central tendency: allowed transformations

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- mean allows positive affine transformations (= interval scale, see from slide 20):

$$f(x) = a + b * x \rightarrow \bar{x} = a + b * \bar{x} \quad \text{mit } b > 0$$

- Median allows strictly monotonic transformations (= ordinal scale):

$$x < y \rightarrow f(x) < f(y)$$

- mode allows injective transformations (= nominal scale):
 a, b, x, y values of a nominal scale where $a = b$ and $x \neq y$ and $f(a) = f(b)$ and $f(x) \neq f(y)$

Position in the distribution: quantiles

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: at least ordinal scaled variables

Measure is a generalisation of the median (median divides the data into the 50% quantile)

„Calculation“: at least $p * 100\%$ of the data are less than or equal to and $(1 - p) * 100\%$ are greater than or equal to the p -quantile x_p , where $0 \leq p \leq 1$

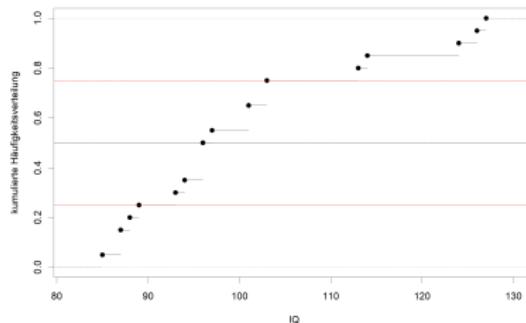
$$x_p \begin{cases} = x_{(np+1)} & , \text{ if } n*p \text{ is not integer} \\ \in [x_{(np)}, x_{(np+1)}] & , \text{ if } n*p \text{ is integer} \end{cases}$$

Position in the distribution: quantiles - example

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

quartiles (25% p -quantiles) of the variable IQ in the data set of Trame et al. 1998 (see slide 10):

$$\begin{array}{l} p*100 = \%: \quad 0\% \quad 25\% \quad 50\% \quad 75\% \quad 100\% \\ x: \quad \quad \quad 85.0 \quad 92.0 \quad 96.5 \quad 105.5 \quad 127.0 \end{array}$$



Position in the distribution: quantiles - types

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **Quartile:** $x_{.00}, x_{.25}, x_{.50}, x_{.75}, x_{1.00}$
 - **Decile:** $x_{.10}, x_{.20}, \dots, x_{.90}$
 - **Percentile:** $x_{.01}, x_{.02}, \dots, x_{.99}$
 - $x_{.25}$: lower quantile
 - $x_{.50} = y_{med}$: Median
 - $x_{.75}$: upper quantile
 - as well as *Inter Quartile Range*, IQR $d_Q = x_{.75} - x_{.25}$
 - in this range are the 50% of the data; this is a measure of dispersion; answers the question: *How far apart are the data?*
 - as well as **span** $d_Q = x_{1.00} - x_{0.00}$ (maximum - minimum)
- relevant for the 5-point summary (see slide 87)

Dispersion: Variance - Formula

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Definition: Measure of the dispersion of the probability density (values) around its centre of gravity;
Prerequisite: metric variables

Formulas, compare equation (1):
empirical variance:

$$\tilde{s}^2 = \frac{1}{n} ((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Sampling variance:

$$s^2 = \frac{1}{n-1} ((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

→ mostly the sample variance is calculated, since the empirical variance underestimates the population variance (practical reason); as well as only $n - 1$ variances vary freely and thus one averages by the number of degrees of freedom (theoretical reason)

Dispersion: Variance - properties I

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

The idea of defining variance

- ❶ Deviations between $x_i - \bar{x}$ measure how much the data scatter around mean \bar{x} .
- ❷ by squaring, all deviations receive a positive sign: $(x_i - \bar{x})^2$
 - sensitive to outliers: large x -values enter the formula strongly
 - ! the variance s^2 no longer has the same unit of measurement due to squaring, e.g. body height measured in cm , so the unit of variance is cm^2
- ❸ variance is the mean of the squared deviations and is small if the values scatter closely around the mean value

Dispersion: Variance - properties II: calculation

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

1 $x_i - \bar{x}$

2 $(x_i - \bar{x})^2$

3 Interpretation, *Is there an outlier?*

Example for the variable IQ: $\bar{x} = 101$, $s^2 = 174.53$, $s = 13.21$

IQ	$IQ_i - \bar{IQ}$	$(IQ_i - \bar{IQ})^2$
96	-5	25
89	-12	144
87	-14	196
...
113	12	144
124	23	529
$\sum_{i=1}^n$		3316

$$s^2 = \frac{1}{20-1} * 3316 = 174.53, \sqrt{s^2} = 13.21$$

Dispersion: Standard deviation

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Objective: Variance of a variable is always given in square units,
standard deviation has the same unit as the variable itself
Prerequisite: metric variables

Formula: Root of the sample variance

$$s = \sqrt{s^2}$$

→ the units of the standard deviation s now correspond to the original unit of measurement

Dispersion: variance, standard deviation - properties

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- sensitive to outliers / extreme values: due to squaring, strong deviations have a high weight.
 - e.g. *last value of variable IQ is replaced by 220, change variance: $s^2 = 174.53 \rightarrow s^2 = 867.75$ (also mean: $\bar{x} = 101 \rightarrow \bar{x} = 105.8$)*
- variance can be linearly transformed: $f y_i = a + b * x_i$, so variance is unchanged for a and for b : $s_y^2 = b^2 s_x^2$
- by means of the standard deviation and the arithmetic mean under the assumption that the characteristic is normally distributed (see from slide 153), classic confidence intervals can be formed (see from slide ??)

Skewness: basics

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

! link to symmetry of one-dimensional distributions (see from slide 82)
the dimensions characterise the following:

- skewness: statistical ratio describing the type and strength of the asymmetry of a probability distribution.
 - $g_m < 0$: distribution skewed to the right (right skew, left skew, negative skew).
 - $g_m > 0$: distribution skewed to the left (left-skewed, right-skewed, positive skewness)
- kurtosis: measure of the steepness of a (single-peaked) probability function (also density function, frequency distribution).
 - $\gamma > 0$: steep-peaked (also leptokurtic), more peaked distribution compared to normal distribution
 - $\gamma < 0$: flat-peaked (also platykurtic), flattened distribution compared to normal distribution

Skewness: coefficient I

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Prerequisite: metric variables

Skewness of a random variable X is the central moment⁴ 3rd order, normalised to the standard deviation s formula:

$$g_m = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3 = \frac{m_3}{\tilde{s}^3}$$

→ Note second term: this is nothing more than a z-standardisation raised to the power of 3 (see slide 153)

⁴central moments describe the distribution of the probability mass around the expected value $\mu = \mathbb{E}(X)$ (mean) of the random variable X , which are defined as $\mu_k := \mathbb{E}((X - \mu)^k)$

briefly importance to understand fractions

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature



$$\frac{3}{4}$$

Numerator

How many equal parts do you have?

Denominator

How many equal parts is the whole divided into?

→ by subtracting the mean in the numerator and dividing the denominator (often standard deviation), a variety of measures are **normalised** (also standardised) in statistics to make their range of values independent of the metric of the variable (e.g. body weight measured in g vs. kg), examples are:

- Coefficient of variation: $v = \frac{s}{\bar{x}}, \bar{x} > 0$ - Fahrmeir et al. (2016), 68
 - normalised to: $0 \leq v \leq 1$
- Correlation coefficient
 - normalised to: $-1 \leq r_{XY} \leq 1$
- ...

Skewness: coefficient II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

the skewness coefficient g_m is constructed as follows:

- in the numerator, the deviations from the mean are calculated and raised to the power of 3
- in the denominator the standard deviation is raised to the power of 3
 - denominator is always positive in contrast to the numerator (*e.g.* $-2^3 = -8$)
 - symmetric distribution: sum of negative and positive deviation terms in the numerator is approximately equal
 - left skewed distribution: sum of the deviation terms in the numerator is larger for the observations below the mean; vice versa for right skewed distributions
 - the further the skewness coefficient g_m deviates from 0, the more left-skewed or right-skewed is the distribution

Skewness: coefficient III: example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

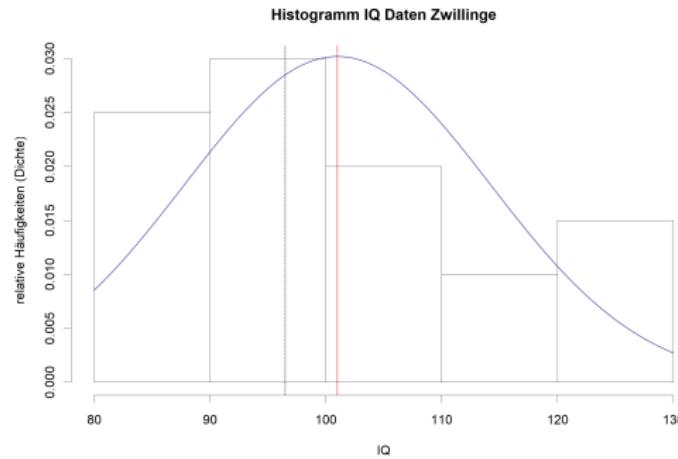
Distributions

discrete

continuous

Literature

Example of the variable IQ - graphical representation:
positive skewness: Schiefekoeffizient = .81 (also right-skewed,
right-tailed, left-leaning curve)



vertical lines from left to right: median, **mean**

Skewness: quartile coefficient I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: at least ordinal scaled variables

→ Advantage: in contrast to the skewness coefficient, the quartile coefficient is not influenced by outliers

formula:

$$g_p = \frac{(x_{1-p} - x_{med}) - (x_{med} - x_p)}{x_{1-p} - x_p}$$

often $p = .25$ is set to obtain the classical quartile coefficient:

$$g_p = \frac{(x_{.75} - x_{.50}) - (x_{.50} - x_{.25})}{x_{.75} - x_{.25}}$$

Skewness: quartile coefficient II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

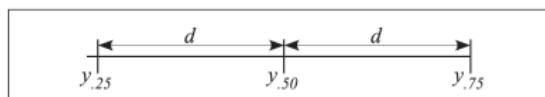
discrete

continuous

Literature

the quartile coefficient g_p is constructed as follows:

- this calculates the difference between the upper quantile and the median, as well as the difference between the median and lower quantile, example symmetrical distribution:



- $g_{.25} < 0$: Distribution skewed to the right (negative skewness).
- $g_{.25} > 0$ distribution tilted to the left (positive skewness)

Skewness: Example of skewness and quartile coefficient

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Representation of two boxplots for the variable IQ (with plotted values):

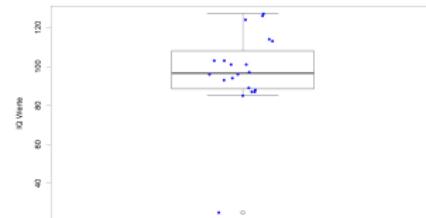
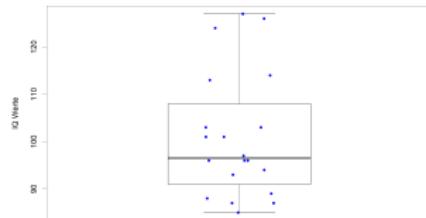


Abbildung: links Variable IQ, rechts Variable IQ mit Datenfehler

Result:

- $g_m = .81$: right-skewed distribution (positive skewness)
respectively $g_m = -1.67$ left-skewed distribution (negative skewness)
- $g_{.25} = .33$: right-skewed distribution (positive skewness)
respectively $g_{.25} = .07$: symmetrical distribution

Kurtosis (curvature): coefficient I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: metric variables

kurtosis of a random variable X is the 4th order central moment normalised to the standard deviation s formula:

$$g_m = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}^4} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^4 = \frac{m_4}{\tilde{s}^4} := \underbrace{\frac{m_4}{\tilde{s}^4}}_{\text{comparison to}} - 3 \quad N()$$

→ Note second term: this is nothing more than a z-standardisation raised to the power of 4 (see from slide 153).

→ to better assess kurtosis, it is compared with the kurtosis of a normal distribution (last term) - e.g. *implemented like this in SPSS.*

Kurtosis : coefficient II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- $\gamma > 0$: steep peaked (also leptokurtic), distribution more peaked compared to the normal distribution.
- $\gamma < 0$: flat-peaked (also platykurtic), distribution flattened compared to normal distribution
- $\gamma \approx 0$: approximately normally distributed (also mesokurtic)

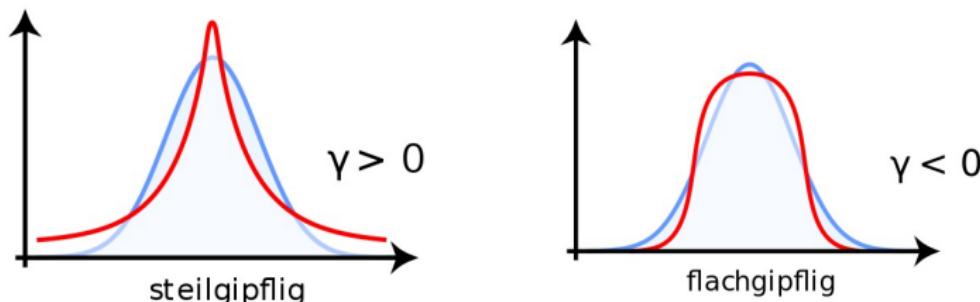


Abbildung: left steep peaked, right flat-peaked

Take Home Messages: measures

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- in dependence on the scale level, different measures can be calculated for individual variables, which are subdivided into:
 - Measures of central tendency (centre of gravity of the data)
 - measures of dispersion (variation in the data from the mean)
 - Measures of skewness / kurtosis (statements about symmetry of a distribution compared to the normal distribution)
- often underestimated is the possibility to calculate quantiles of data (*example income*)

→ linking level of measurement with measures

R-commands section measures

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- computation mean: `mean(x)`
- median: `median(x)`
- mode: `getmode(x)` → is not implemented in **R**, so it is necessary to write own function
- quantiles: `quantile(x, probs = ...)`
 - using the argument `probs` you can determine which quantiles are calculated (quartiles by default)
- variance: `var(x)`
- standard deviation: `sd(x)`
- quartile coefficient: `getquartilecoef(x)` → not implemented in **R**
- skewness: `moments::skewness(x)`
- kurtosis: `moments::kurtosis(x)`

Symmetry: Examples

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- often one encounters certain (partly undesirable) forms of distribution in reality / in surveys:
 - income distributions are almost always **right-skewed** (high incomes run out to the right).
 - questions that are answered in a socially desirable way are often **right-skewed** (question is strongly rejected) or **left-skewed** (question is strongly favoured)
 - standardised tests are often said to have a **symmetrical** distribution, but can be right-skewed if questions are too easy for example
 - ...

→ in addition to the shape of a distribution, extreme values or outliers can also be detected

Symmetry: basics

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Symmetry or skewness can be assessed via so-called (Fechner's) position rules:

- Symmetrical distribution: $\bar{x} \approx x_{med} \approx x_{mod}$
 - only approximate, since empirical distributions are rarely exactly symmetric
- Left-skewed distribution: $\bar{x} > x_{med} > x_{mod}$
- Right-skewed distribution: $\bar{x} < x_{med} < x_{mod}$

→ mainly relevant for **unimodal** distributions (see from slide 88)

Symmetry: negative skew (left-skewed)

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

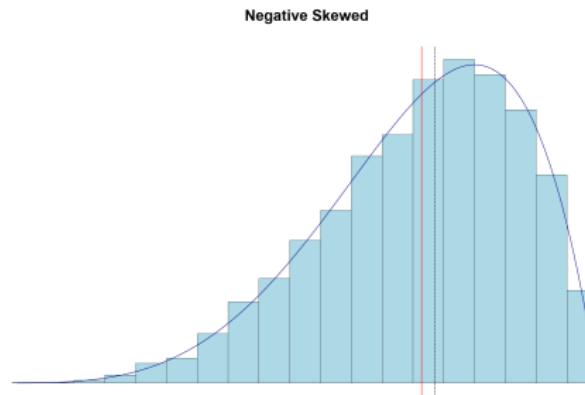
Distributions

discrete

continuous

Literature

negative skew: skewness coefficient < 0
→ falls more flat on the left side than on the right side



vertical lines from left to right: **mean**, median

Symmetry: positive skew (right-skewed)

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

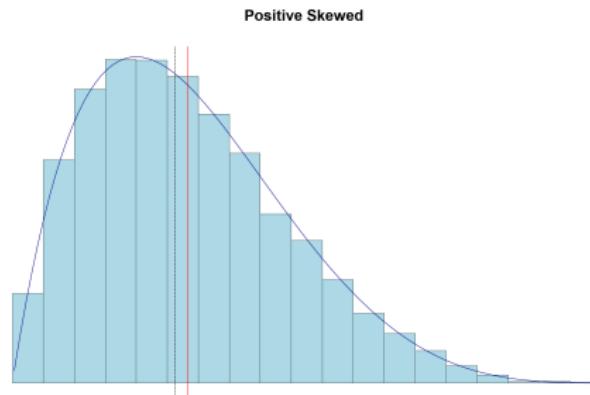
Distributions

discrete

continuous

Literature

positive skew: skewness coefficient > 0
→ falls more flat on the right side than on the left side



vertical lines from left to right: median, median

Symmetry: symmetrical distribution

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

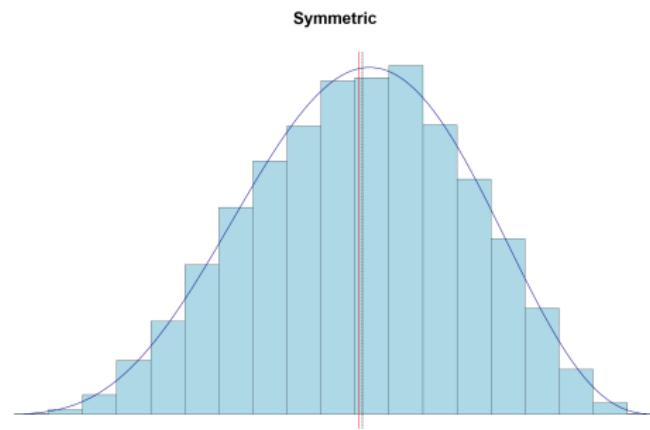
discrete

continuous

Literature

positive skew: skewness coefficient ≈ 0

\rightarrow falls evenly on both sides



vertical lines: median \approx **median**

Boxplot: 5-point summary

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

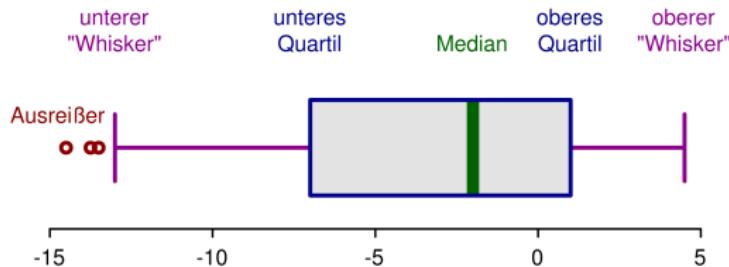
Distributions

discrete

continuous

Literature

a boxplot is derived from quantiles and visually represents the following measures: $x_{\min}, x_{0.25}, x_{\text{med}}, x_{0.75}, x_{\max}$
outlier defined as $[x_{0.25} - 1.5 * \text{IQR}, x_{0.75} + 1.5 * \text{IQR}]$



Interpretation:

Median within the box slightly shifted to the right → right-sided distribution

several outliers → no normal distribution

[Wikipedia Box-Plot](#)

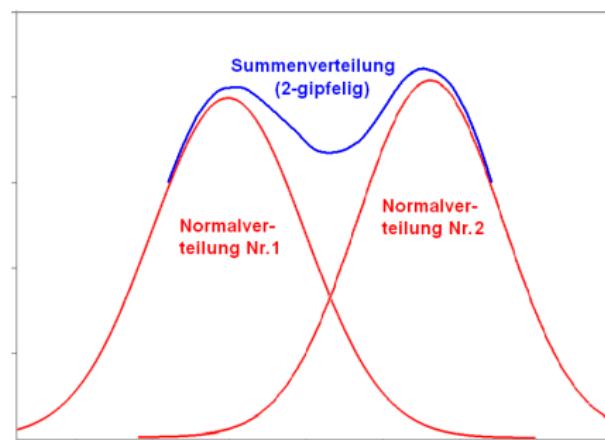
Symmetry: peak of a distribution

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

a distribution may have the following peaks:

- unimodal (single-peaked): distribution with one peak / maximum
- bimodal (two-peaked): has several maxima; as well as multimodal (multi-peaked) distributions

→ data come from different subpopulations, but are mixed in the representation; subpopulations can be separated by the minimum between the two maxima



Take Home Messages: symmetry

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

unimodal distributions can have the following possible forms:

- Symmetric distribution
- Left-skewed distribution
- Right-skewed distribution

multimodal distributions, on the other hand, are an indication that the observed phenomenon comes from several subpopulations that can (possibly) be separated into unimodal distributions

→ it is important to theoretically explain properties of symmetry of distributions by possible causes (such as response tendencies, intrinsically non-symmetrical phenomenon, ...)

R-commands section symmetry

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Calculation of the skewness: `moments::skewness(x)`
- Boxplot: `boxplot(x)`; several boxplots next to each other: `boxplot(x ~ y)`, where y is a categorical (discrete) variable
- 5-point summary (several descriptive measures): `summary(x)`

Multivariate descriptive statistics: Basics I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Looking at two or more variables at the same time can answer questions about **relationships**:

- Do the values of one variable tend to go hand in hand with certain values of another variable
 - Example: *People who are taller have higher body weight*
- Measures of correlation are also to be interpreted as measures of difference, example:
 - *Average weight of taller people is higher than average weight of people with medium or short height*

Multivariate descriptive statistics: Basics II

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

The aim is the representation, description of data consisting of observations of two or more characteristics (variables) (bivariate, two-dimensional data or multivariate, multidimensional data).

→ also allows comparison of several univariate results to gain impressions about possible dependencies, which can be further analysed with multivariate procedures

bivariate and multivariate representations:

- contingency table
- two-dimensional bar chart
- two-dimensional histogram
- scatter plot
 - pairwise scatterplots (scatterplot matrix)

Contingency table I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: discrete variables with relatively few characteristics (also possible to categorise metric characteristics)

Contingency tables use only the nominal scale level, thus the lowest information content

Example: Are there differences between the sexes on the categorised IQ variable?

Variable	IQ.cat	male (1)	female (2)
	below average	6	5
	medium	2	3
	above average	2	2

Contingency table II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

		Outcome: lung cancer		
		Yes	No	Total
Exposure: smoked	Yes	75	25	100
	No	10	90	100
	Total	85	115	N

measure of risk; Basic idea frequency table from slide 36

Contingency table III

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Contingency tables represent the marginal distributions of two variables X, Y , where the

- marginal frequency of feature X is given by: $h_{i\cdot} = h_{i1} + \dots + h_{im}$, where $i = 1, \dots, k$ is the marginal frequency of feature Y given by: $h_{\cdot j} = h_{1j} + \dots + h_{kj}$, where $j = 1, \dots, m$

→ Dot notation: $h_{i\cdot} = \sum_j h_{ij}$ (sum over all columns)

Example:

Variable IQ.cat	male (1)	female (2)	$h_{i\cdot}$
below average	6	5	11
medium	2	3	5
above average	2	2	4
$h_{\cdot j}$	10	10	20

Contingency table III - conditional frequency

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Groups are (almost always) of different sizes, so it requires conditional frequencies:

conditional frequency distribution of Y under the condition $X = a_i$:

$$f_Y(b_1|a_i) = \frac{h_{i1}}{h_{i\cdot}}, \dots, f_Y(b_m|a_i) = \frac{h_{im}}{h_{i\cdot}}$$

conditional frequency distribution of X under the condition $Y = b_j$:

$$f_X(a_1|b_j) = \frac{h_{1j}}{h_{\cdot j}}, \dots, f_X(a_k|b_j) = \frac{h_{kj}}{h_{\cdot j}}$$

Example of fixed level of variable IQ.cat ($X = a_i$):

Variable IQ.cat	male (1)	female (2)	
below average	.55	.45	1
medium	.4	.6	1
above average	.5	.5	1

Three-dimensional bar chart: example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

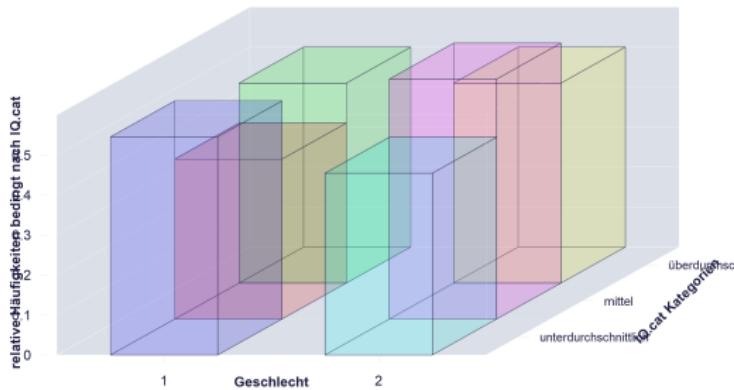
Distributions

discrete

continuous

Literature

Example of the calculated conditional frequencies from the previous slide:



Three-dimensional histogram I: construction

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Prerequisite: discrete variables with relatively few expressions

Analogous to the construction of a one-dimensional histogram:

- ① data are grouped into rectangle classes c_0, \dots, c_k , as well as e_0, \dots, e_j , where these have adjacent intervals:
 $[c_0, c_1] \times [e_0, e_1], \dots, [c_{i-1}, c_i] \times [e_{j-1}, e_j]$
 - Blocks in the histogram have the base edge $[c_{i-1}, c_i]$ in the x coordinate and $[e_{j-1}, e_j]$ in the y coordinate, respectively.
- ② absolute height of the blocks: $\frac{h_{ij}}{(c_i - c_{i-1})(e_j - e_{j-1})}$
 - Note: Term divides the frequencies within the intervals (h_{ij}) by the base area.
 - is done under the principle of area fidelity (volume is proportional to height and thus to absolute frequencies)

→ same disadvantages as with a one-dimensional histogram (see from slide 39)

Three-dimensional histogram II: example

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

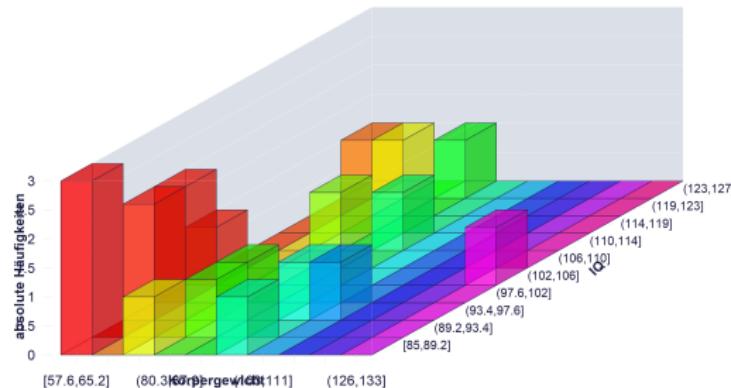
Distributions

discrete

continuous

Literature

Representation of the variables IQ x body weight:



→ interrelationships would be shown by „mountains and valleys“ in the data (data centroids)

Scatter plot I: construction, example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

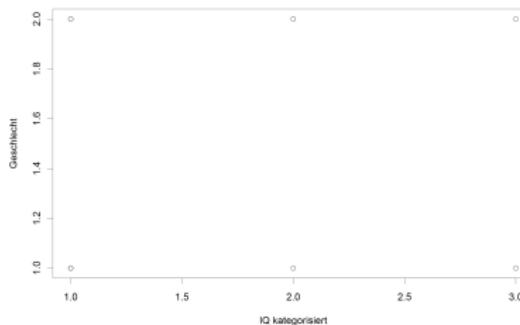
discrete

continuous

Literature

Prerequisite: metric variables with relatively few identical values.
pairs of measurements $(x_1, y_1), \dots, (x_n, y_n)$ are represented in an x-y-coordinate system (same index n : need equal number of observations for both variables)

Visualisation scatter plot with multiple identical values:



Scatter plot II: example

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

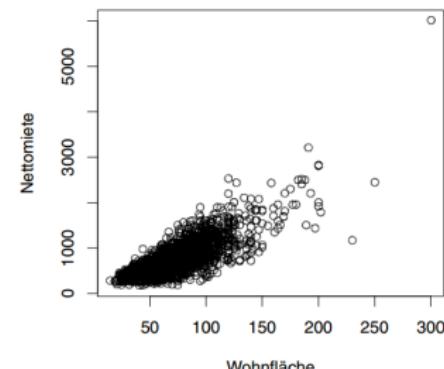
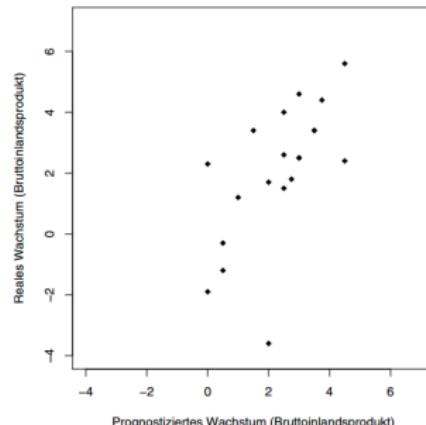


Abbildung: left projected and real growth rate for the Council of Economic Experts' (Sachverständigenrat) forecast in 1975-94, right living space times net rent in Munich

pairwise scatterplots (scatterplot matrix) I: notes

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: pairwise metric variables with relatively few identical values

- scatter plots of two variables are formed in pairs
- on the diagonal is the variable name from the data set
- ! Relationships are shown without taking the other variables into account.

pairwise scatterplots (scatterplot matrix) II: example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

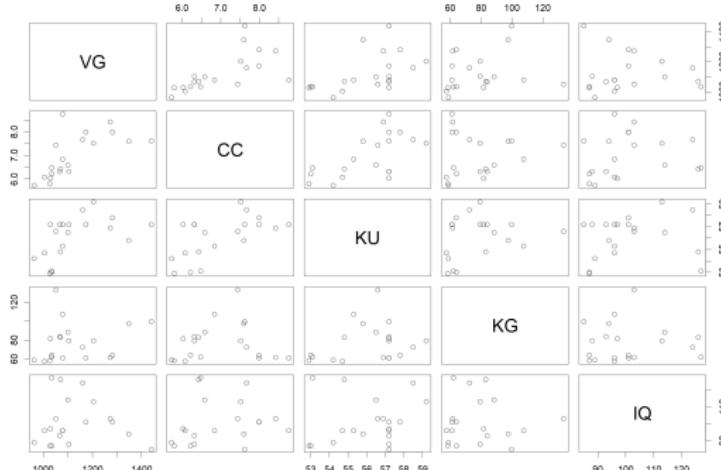
Distributions

discrete

continuous

Literature

Plot scatterplot matrix for the variables VG (volume of the forebrain in cm^3), CC (area of the corpus callosum in cm^2), KU (head circumference in cm), KG (body weight in kg), IQ



pairwise scatterplots II: correlation matrix

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Correlation matrix (empirical correlation coefficient according to Bravais-Pearson) for the variables VG (volume of the forebrain in cm³), CC (area of the corpus callosum in cm²), KU (head circumference in cm), KG (body weight in kg), IQ:

	VG	CC	KU	KG	IQ
VG	1.00	0.66	0.51	0.21	-0.06
CC	0.66	1.00	0.61	0.08	0.16
KU	0.51	0.61	1.00	0.24	0.14
KG	0.21	0.08	0.24	1.00	0.00
IQ	-0.06	0.16	0.14	0.00	1.00

→ VG, CC, KU are highly correlated with each other, but no (low) correlation with KG and with IQ

Multivariate descriptive statistics: measures

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

if individual variables show a (linear) correlation, this can be quantified using statistical measures

bivariate (and multivariate) statistical measures:

- for two variables with the same scale level:
 - Empirical correlation coefficient according to Bravais-Pearson
 - Spearman's correlation coefficient
 - Alternative rank correlation measures
- *not mentioned*: there are statistical procedures if two variables have a different scale level:
 - η^2 (Eta-squared)
 - r_{pbis} (point biserial correlation)

Empirical correlation coefficient I - formula

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: two metric variables

Formula Empirical correlation coefficient according to Bravais-Pearson:

$$r = r_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\tilde{s}_{XY}}{\tilde{s}_X \tilde{s}_Y}$$

→ the denominator is used for normalisation⁵, range of values between perfectly positive correlation ($r_{XY} = 1$), no correlation ($r_{XY} = 0$) and perfectly negative correlation ($r_{XY} = -1$)

⁵the common variation of the values can be at most as large as the variation of the individual values

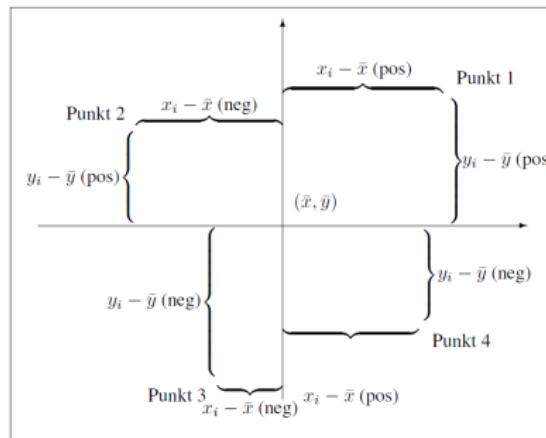
Empirical correlation coefficient II - properties

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

equidistant **linear** connection

Empirical covariance:

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



briefly properties covariance

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

the following properties of the covariance Cov also apply to the correlation:

- symmetry: $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- sensitive to outliers
- independence: if two random variables X and Y are independent, then the covariance is $\text{Cov}(X, Y) = 0$.
- constant parameter: if the data consists of constant values a in addition to a random variable X , the covariance $\text{Cov}(X, a) = 0$.

Relationship to variance:

- $\text{Cov}(X, X) = \text{Var}(X)$
- with the covariance, in contrast to the variance, the unit of measurement remain
 - , however, the correlation is dimensionless (regardless of the unit of measurement, it takes the values $[-1 \leq r_{SP} \leq 1]$)

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

equidistant **linear** connection

- clouds of points in the 1st, 3rd quadrant $r_{XY} > 0$, same-sense correlation
 - extreme case $r_{XY} = 1$: all points lie on a straight line with positive slope.
- clouds of points in 2nd, 4th quadrant $r_{XY} < 0$, opposite relation
- point clouds in all quadrants $r_{XY} \approx 0$ similarly strongly distributed

Empirical correlation coefficient IV - examples

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

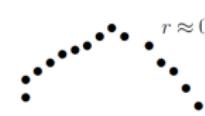
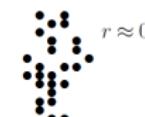
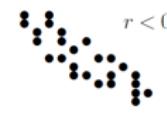
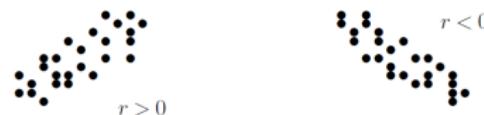
Randomness

Distributions

discrete

continuous

Literature



→ Empirical correlation coefficient only maps **linear** correlations

Empirischer Korrelationskoeffizient V - Interpretation Stärke

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Cohen (1992) suggests the following classification to determine the strength of correlations:

- $r_{XY} = .1$: corresponds to a weak effect
- $r_{XY} = .3$: corresponds to a medium effect
- $r_{XY} = .5$: corresponds to a strong effect

→ for a linear regression, the coefficient of determination is mathematically equal to the squared empirical correlation coefficient

$$R^2 = r_{XY}^2$$

$r_{XY}^2 * 100$ indicates how much variance can be explained by the included variables

Cohen (1992)

Spearman's correlation coefficient I - formula

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Prerequisite: two ordinal scaled variables

Spearman's correlation coefficient formula:

$$r = r_{SP} = \frac{\sum_{i=1}^n (rg(x_i) - rg(\bar{x}))(rg(y_i) - rg(\bar{y}))}{\sqrt{\sum_{i=1}^n (rg(x_i) - rg(\bar{x}))^2 \sum_{i=1}^n (rg(y_i) - rg(\bar{y}))^2}}$$

→ the denominator is used to normalise the range of values:

$$[-1 \leq r_{SP} \leq 1]$$

Example ranks:

x_i	2.17	8	1.09	2.16
$rg(x_i)$	3	4	1	2

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

gleichsinniger **monotoner** Zusammenhang

- $r_{SP} > 0$, if in tendency x large \rightarrow also y large and vice versa, monotonic relationship in the same direction
- $r_{SP} < 0$, if in tendency x large \rightarrow y small and vice versa, opposite monotonic relationship
- $r_{SP} \approx 0$, if no tendency is visible, no monotonic relationship

Spearman's correlation coefficient III - examples

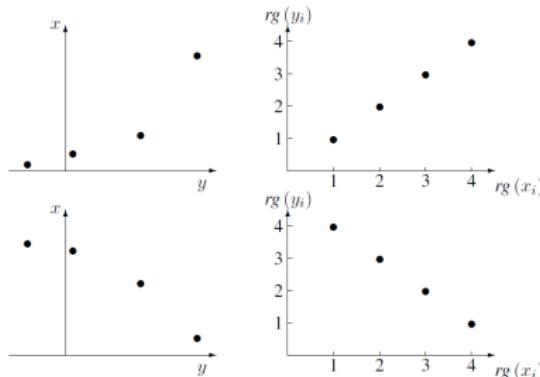
[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Abbildung: extreme cases for Spearman's correlation coefficients: top $r_{SP} = 1$, bottom $r_{SP} = -1$

→ Empirical correlation coefficient only maps **monotone** correlations , e.g. for upper two graphs always holds if $x_i < x_j$, then $y_i < y_j$ for any $i \neq j$

Alternative rank correlation measures: basics I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

in contrast to Spearman's correlation coefficients, only the discrepancy and not the difference between the ranks is considered

- compared to Spearman's correlation coefficients, the values of **concordance measures** are consistently smaller
- to be applied especially when the data are not normally distributed, the scales have unequal answering options, for very small sample sizes

Alternative rank correlation measures: basics II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Concordance measures build on the comparison of observation tuples:

- concordant (C): if $x_i < x_j$ then $y_i < y_j$, values behave in the same way, these are summed up in N_c .
- discordant (D): if $x_i < x_j$ then $y_i > y_j$, values behave in opposite senses, these are summed up in N_d
- bound pairs (T, ties) exist in 3 different forms:
 - T_{XY} if $x_i = x_j$ then $y_i = y_j$
 - T_X if $x_i = x_j$ then $y_i \neq y_j$
 - T_Y if $x_i \neq x_j$ then $y_i = y_j$

Concordance measures I - formula

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Prerequisite: two ordinal scaled variables

Kendall's τ_α :

$$\tau_\alpha = \frac{N_c - N_d}{n(n-1)/2}$$

Goodman and Kruskal γ :

$$\gamma = \frac{N_c - N_d}{N_c + N_d}$$

→ the denominator is used to normalise the range of values:
 $[-1 \leq r_{SP} \leq 1]$

Concordance measures II - example

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- ➊ in the first step, the first person (statistical object) is compared with all others that are below it.
- ➋ in the second step the second person is compared with all others that are below it
- ➌ ...

Nr.	Vergleiche								
	1	2	3	4	5	6	7	8	9
1									
2	D								
3	C	C							
4	C	C	C						
5	C	C	C	T _{XY}					
6	C	C	C	C	C				
7	C	C	C	C	C	T _Y			
8	C	C	C	C	C	C	D		
9	C	C	C	C	C	C	C	C	
10	C	C	C	C	C	C	C	C	C

Daten	
X	Y
x ₁ 2	y ₁ 1
x ₂ 1	y ₂ 2
x ₃ 3	y ₃ 3
x ₄ 4,5	y ₄ 4,5
x ₅ 4,5	y ₅ 4,5
x ₆ 6	y ₆ 6,5
x ₇ 8	y ₇ 6,5
x ₈ 7	y ₈ 8
x ₉ 9	y ₉ 9
x ₁₀ 10	y ₁₀ 10

	Summe								
	D	1					1		2
D	1						1		2
C	8	8	7	5	5	3	2	2	1
T _X									0
T _Y						1			1
T _{XY}				1					1

Take Home Messages: multivariate descriptive statistics

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

a large number of questions can be described with the following methods:

Graphs:

- grouped or discrete variables: Contingency tables, three-dimensional bar chart.
- quantitative variables: Scatterplot, three-dimensional histogram, scatterplot matrix (pairwise scatterplots).

measures:

- intervallscaled variables: Empirical correlation coefficient according to Bravais-Pearson.
- ordinal scaled variables: Spearman's correlation coefficient; deviation normal distribution: concordance measures such as Kendall's τ_α

R-commands section multivariate descriptive statistics I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Graphs:

- contingency table: `table(x, y)`
- contingency table conditional frequencies:
`prop.table(x, margin =)`, depending on whether
`margin = 1` or `margin = 2` is conditioned on variable `x` or `y`.
- two-dimensional bar chart: `epade::bar3d.ade(x)`
- two-dimensional histogram:
`epade::bar3d.ade(x = gplots::hist2d(x, y)$counts)`
- scatter plot: `plot(x, y)`
- scatter plot matrix (pairwise scatterplots): `pairs(x)`

R-commands section multivariate descriptive statistics II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

statistical correlation measures:

- Empirical correlation coefficient according to Bravais-Pearson:
`cor(x, y, method = "pearson")`.
- Spearman's correlation coefficient:
`cor(x, y, method = "spearman")`
- Alternative rank correlation measures:
`cor(x, y, method = "kendall")`

further correlation measures such as r_{pbis} (point biserial correlation) can be found in chapter 6 in [Dormann \(2013\)](#)

Correlation measures - no causality

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful

George Box (1987)

Correlation measures (generally statistical methods) never (!) indicate causality unless a correlation meets very strong assumptions

→ Important contributions on the topic can be found, for example, in Shadish, Cook und Campbell (2002) or in contributions to complex thinking Elbow (1993); Tsoukas und Hatch (2001)

Causality assumptions

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Necessary assumptions to attribute causality to a relationship between two variables X and Y :

- if the effect Y is assumed to be caused by the cause X , the cause must precede the effect ($X \xrightarrow{\text{cause}} Y$).
- there is a correlation or observable covariation between the occurrence of X and Y .
- the correlation between X and Y is **isolated**. Thus, there are no other possible causes that cause the correlation of X and Y . The relationship between the two variables is also evident when other variables that may cause Y are statistically controlled for.
- the estimation procedure chosen for the estimation of the relationship between the two variables is correct. The distribution of the observed data is consistent with the distributional assumptions of the estimation procedure.
- The direction of the causal relationship is correctly specified.

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Relationship structures between statements on integration and on political culture at an overall level (*correlation matrix*):

Tab. 3 Einstellungen gegenüber Integration und Haltung zur Demokratie in Europa. (Quelle: Eigene Berechnungen; European Social Surveys 2012–2016)

	Soziales Ver-trauen	Wichtigkeit Leben in einer Demokratie ^a	Wie demokra-tisch ist das eigene Land ^a	Zufriedenheit wie Demokratie im Land arbeitet
Erlaubnis zur Immigrati-on von Mitgliedern der eigenen „Ethnie“/ Mehrheitsgesellschaft	+0,19	+0,16	+0,14	+0,16
Erlaubnis zur Immigrati-on von Mitgliedern anderer „Ethnien“ Minderheiten	+0,24	+0,14	+0,14	+0,19
Erlaubnis zur Immigrati-on von Muslimen ^b	+0,23	–	–	+0,24
Einschätzung von Immigrati-on als gut	+0,26	+0,17	+0,18	+0,27
Immigration macht Land weniger lebens-wert	-0,27	-0,18	-0,18	-0,28
Immigration unter-wandert Kultur des Landes	-0,27	-0,20	-0,24	-0,28

Pearsons Produkt-Moment-Korrelationen über alle Länder der Datensätze (kontrolliert auf Länder)

Alle Skalen 11-Punkti-Antwortvorgaben; $n=32000-37000$

^a2012

^b2014

→ not test for measurement invariance (correlations have the same patterns between countries)

stay critical! own calculations

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Relationship structures (*correlation matrix*) between statements on integration and on political culture separately for the countries Hungary and Germany:

	St H.	Swd H.	St G.	Swd G.
Migration eigener Ethnien	0.13	0.06	0.27	0.28
Migration fremder Ethnien	0.15	-0.03	0.27	0.31
Besser Ort zu leben	0.27	0.14	0.34	0.42
Bereicherung Kultur	0.23	0.06	0.30	0.37

Tabelle: St = Social trust, Swd = Satisfaction with democracy, H. = Hungary, G. = Germany

→ Relationship is **not stable** between countries

Standards of evidence

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

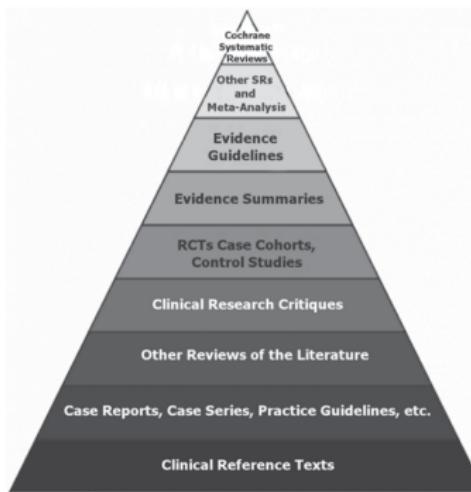
Distributions

discrete

continuous

Literature

Hierarchy of evidence in the context of „evidence-based practice“ in medicine::



→ from the 6th level of the pyramid onwards, studies often have a much lower relevance

Probability?

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Fragen:

- How likely is it that the sun will rise tomorrow?
- What is the probability that it will rain tomorrow?
- How likely is it to roll the number 6 on a fair die?
- Bet on coin tosses: If tails comes up you win 2€, if heads comes up you lose 1€. How much money did you win after 100 rolls?

→ Probabilities are used to quantify the degree of uncertainty of the occurrence of events

Zucchini et al. (2009), 73ff.; Fahrmeir et al. (2016), 165ff.

Probability concepts I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

There are three approaches in everyday life and in science for determining **probabilities**:

- subjective probability: assignment of probabilities based on subjective beliefs
- logical probability (also classical, Laplace probability): logically seen / on the basis of plausible considerations there is nothing to be said against the assignment of probabilities
- frequentist (objective) probability: relative frequency with which a certain **event** occurs in an independent repetition of a **random experiment**

Probability concepts II

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **random experiment:** thought or actual experiment whose outcome (event) cannot be predicted with certainty
- **(outcome) set:** Set of all possible outcomes of a random experiment, denoted by symbol Ω (*Omega*)
 - normal coin: $\Omega = \{\text{head}, \text{number}\}$
 - normal die: $\Omega = \{1, 2, 3, 4, 5, 6\}$
 - flight from Hamburg to Rome: $\Omega = \{0, 1\}$ (0 do not land in Rome, 1 land in Rome)
 - Burning time of a light bulb which is not countable⁶ has many outputs: $\Omega = x | x \geq 0$

⁶only possible probabilities with the following formula „ $P_{x_1} \leq \text{result} \leq x_2$ “ for any subintervals of the result set, where x_1, x_2 are real numbers

Probability concepts III

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive

statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- **subsets** of the result set: are denoted with capital Latin letters
 - normal dice: $A = \{\text{Augenzahl ist gerade}\} = \{2, 4, 6\}$
 - burning time of a light bulb:
 $B = \{\text{Brenndauer ist länger als 500 Stunden}\} = \{x | x > 500\}$
- an **empty set** \emptyset is an impossible event
 - normal dice: $\Omega = \{1, 2, 3, 4, 5, 6\}$ has event
 $C = \{\text{number of dice is a 7}\} = \{7\} \rightarrow C \notin \Omega$
- the complete result set Ω is always a **safe event**
 - normal dice: $\Omega = \{1, 2, 3, 4, 5, 6\}$ hat Ereignis
 $D = \{1, 2, 3, 4, 5, 6\} \rightarrow D \in \Omega$, also: $D \subseteq \Omega$

Definition of probability I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

probability function is a function P which assigns a real number to all events from Ω and which must satisfy the free following axioms of Kolmogoroff:

$$\text{A1 } 0 \leq P(A) \leq 1$$

$$\text{A2 } P(\Omega) = 1$$

$$\text{A3 wenn } A \cap B = \emptyset \text{ (disjunkt), gilt } P(A \cup B) = P(A) + P(B)$$

→ what conditions must the function P fulfil in order to represent a probability (also applies to subjective probabilities)

Definition of probability II - example

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Are these functions probabilities?

Ereignis	\emptyset	{Kopf}	{Zahl}	Ω
Wahrscheinlichkeit	$P(\emptyset)$	$P(\{\text{Kopf}\})$	$P(\{\text{Zahl}\})$	$P(\Omega)$
Fall 1	0.0	0.5	0.5	1.0
Fall 2	0.0	0.6	0.4	1.0
Fall 3	0.0	0.0	1.0	1.0
Fall 4	0.0	0.5	0.6	1.0
Fall 5	0.1	0.4	0.5	1.0

Verletzung folgender Axiome:

- case 4 violates A2: assuming that $A \cap B = \emptyset$, $P(\text{head}) + P(\text{number}) = 1.0 \neq 1.1$ holds
- case 5 violates A2: an impossible event never has a probability and it should hold $P(\Omega) = P(\Omega) + P(\emptyset)$ since $P(\emptyset) = 0$.

Frequentist (objective) probability

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

10,000 times the toss was simulated with a coin, this has the result set $\Omega = \{\text{head}, \text{tail}\}$, here we are only interested in the event $A = \{\text{coin shows head}\} = \{\text{head}\}$.

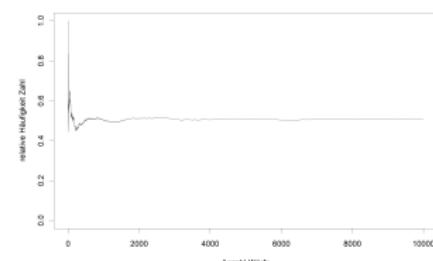
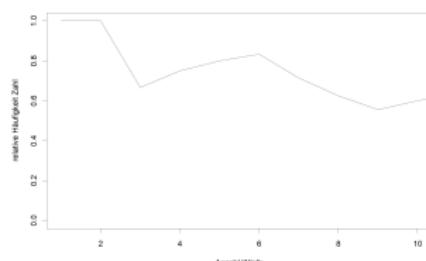


Abbildung: left relative frequency number to throw head at 10 throws, right at 10,000 throws

→ relative frequency with which the event occurs in 10,000 independent repetitions of a random experiment is as expected at $P(A) = 0.5$

Frequentist (objective) probability: law of large numbers

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

law of large numbers: the final value of the relative frequency of an event A can be interpreted as the probability of A (relative frequency interpretation of probability)

$$f_n(A) \xrightarrow{n \rightarrow \infty} P(A)$$

Problem: real random experiments are not symmetrical, i.e. all elementary events do not always have the same probability (one cannot repeat experiments under always the same conditions) → no Laplace probability

Random variable I: definition

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

random variable X is not a single value but has a whole range of possible values

→ behaviour of a random variable X can only be described by probabilities; justification:

- in reality, most situations are **stochastic**, i.e. only probabilities and no certainties can be stated, since correlations are subject to random fluctuations
 - example: *flight duration of the flight connection between Hamburg and Rome*
- if an event is **deterministic**, for each value of variable X, an exact value of variable Y can be determined.
 - example: *period of oscillation of a pendulum in a vacuum*

Random variable II: definition

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

if the desired event (*flight duration of flight connection A between Hamburg and Rome*) occurred, then a **realisation** of the random variable X took place

- before the aircraft has landed, the actual flight duration of flight connection A is unknown.
- after the aircraft has landed safely, the flight duration is a simple number

→ Probabilities allow the degree of uncertainty to be predicted (more or less)

Random variable III: conditional probability - basic idea

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

A variety of factors influence the likelihood of landing safely in Rome:

- Random experiment: Flight to Rome.
 - Event of interest: I arrive in Rome in one piece.
 - Initial assessment: My chance of arriving safely is good (high probability).
 - Additional information:
 - One engine appears to be defective.
 - The pilot staggers through the cabin shouting pop songs
- new assessment: my chances have become worse (lower conditional probability)

conditional probability: $P(A|B)$, where A is event „land in Rome“ and B is additional information

Random variable: example light bulbs I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

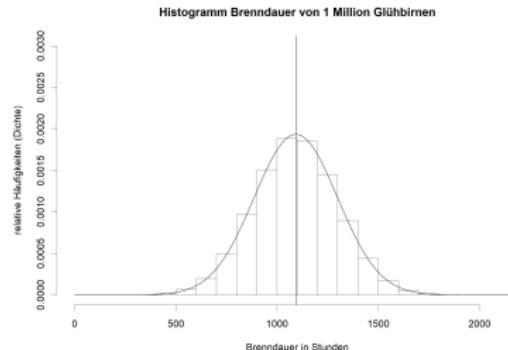
Distributions

discrete

continuous

Literature

You are the operator of a small ghost train at the funfair / kermis and you have to buy 30 light bulbs. The light bulbs of a certain type have the following burning time:



→ Sadly, nothing compared to the centennial light at 4550 East Avenue, Livermore in California, which has been burning since 1901 (see [Wikipedia Centennial Light](#))

Random variable: example light bulbs II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

You have now been to two shops and your 30 randomly selected light bulbs have the following burning time:

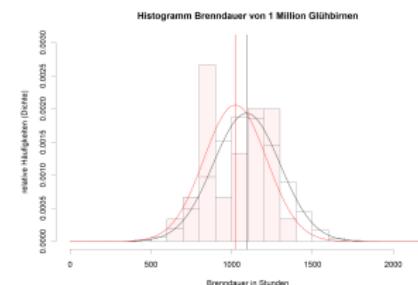
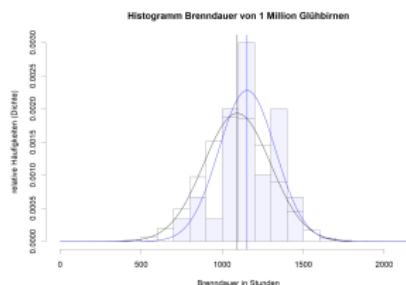


Abbildung: on the left 30 light bulbs in blue, on the right 30 light bulbs in red in comparison to the theoretical purchase of 1 million light bulbs

Where would you buy light bulbs in the future?

Random variable: example light bulbs III

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

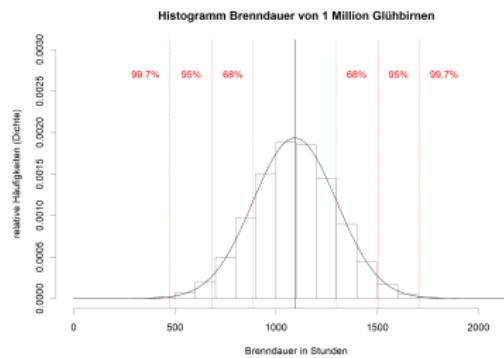
Distributions

discrete

continuous

Literature

Probabilities can be calculated via the density function*:



this corresponds to the 68-95-99.7 rule of a **normal distribution**
 $(\mu \hat{=} \bar{x}, \sigma \hat{=} s:)$

- $P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 0.6828$
- $P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 0.9545$
- $P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 0.9973$

Random variable: example light bulbs VI

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Calculation of the 68-rule (further explanation see slide 153):
Normal distribution:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} \quad [-\infty < x < \infty]$$

Integral to be calculated for the 68 rule:

$$P(886.804 < X < 1298.961) = \int_{886.804}^{1298.961} \frac{1}{206.1\sqrt{2\pi}} e^{\frac{-(x-1092.9)^2}{2206.1^2}} dx = 0.6828$$

basic step-by-step calculation:

- ① calculate root function
 - ② determine individual integrals and subtract them from each other
- Calculation for Normal Distribution [youtube: Integrating Normal Density Function](#)

Take Home Messages: Probabilities

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- many processes in reality are stochastic, thus they are (more or less) random processes.
- random processes are described by probabilities
- probabilities serve to quantify the degree of uncertainty of the occurrence of events and thus prevent surprises
- conditional probabilities describe the change of probabilities when prior knowledge is known

Distributions: describe random processes I

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

the aim is the description of results (events) of random processes by numbers

- numbers are here the characteristic value of a variable
 - Variable X whose characteristic x is the result of a random process is called random variable X , $x \in \mathbb{R}$ is the realization of X
 - behaviour of X can only be described by probabilities

example of discrete random variables Y, Z :

$Y = \text{Sum of two dice}$

$$Z = \begin{cases} 1 & \text{if sum of two dice} = 7 \\ 0 & \text{else} \end{cases}$$

→ in the example, the underlying random process is not of primary interest or cannot always be described, otherwise the value range of the variable is directly considered as the result space

Distributions: describe random processes II

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Random variable is defined by assigning numbers to an event of a random process (realisation), for which **distributions** are necessary:

- it is not desirable or feasible to assign a probability to every possible event
 - Random processes can only be described by functional relations, types of representation:
 - for discrete random variables → probability function
 - for continuous random variables → density function
- Distributions are determined by functions and their parameters

A random variable $f(X)$ is a function that assigns a real number to each possible outcome of a random experiment / process.

Example discrete random variables

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Random variables can only take on finitely many values or countably infinitely many values

→ every possible event can be assigned a natural number

Example: rolling 6 or 30 fair dices several times at once and we want to know $A = \{\text{result is 1}\} = \{1\} \rightarrow P(\{1\}) = 1/6$ (binomial distributed):

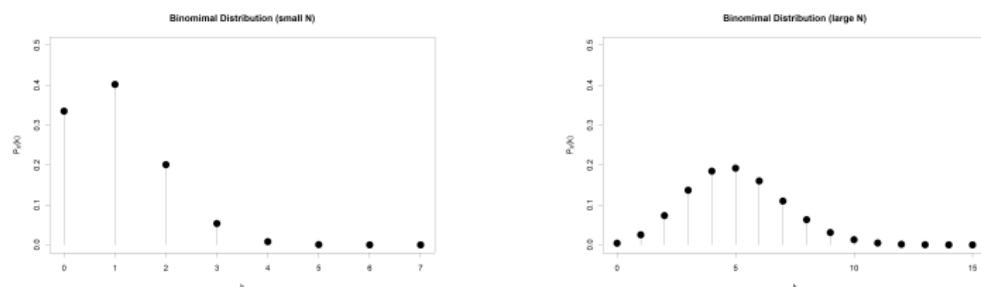


Abbildung: left rolling 6 dices, right 30 dices

Is rolling a dice a Gauss distribution?

Example continuous random variables

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

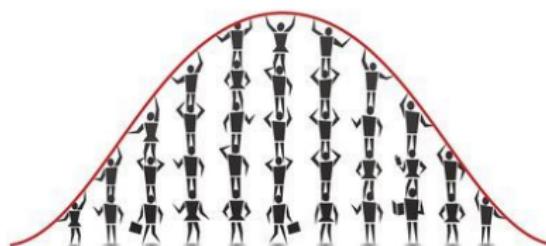
discrete

continuous

Literature

Random variables can take on an infinite number of values
→ no probability can be assigned to an possible single event, e.g.
Probability that a machine will fail after 7.483428383... days is zero (real numbers)

Example body size of a specific population (normal distribution):



Which distributions exist?

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Probability distributions that are significant in theory or in practical application have specific names, (non-exhaustive) list of distributions:

- [List of probability distributions](#)
- [List of convolutions of probability distributions](#)
- [Flowchart of relationships between probability distributions \(last page\)](#)

Distribution of discrete random variables

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

Random variable X is called discrete if it can only take on finitely or countably infinitely many values x_1, x_2, \dots, x_k .

determined by probabilities:

$$P(X = x_i) = p_i, \quad i = 1, 2, \dots, k$$

Probability function:

$$f(x) = \begin{cases} P(X = x_i) = p_i & x = x_i \in \{x_1, x_2, \dots, x_k, \dots\} \\ 0 & \text{else} \end{cases}$$

Example: 2 dice were thrown and the event of interest was „The sum of the eyes of both dice is at most equal to 4“:

$$P(X \leq 4) = \{(1, 1), (1, 2), (1, 3), (2, 1), (2, 2), (3, 1)\}$$

$$= P(1, 1) + \dots + P(3, 1) = \frac{6}{36} = \frac{1}{6}$$

Binomial distribution I

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

$$P(X = x) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}, \quad x = 0, 1, \dots, n$$

- parameter π constant probability of success of each run n .
- number n : independent repetitions of a Bernoulli experiment

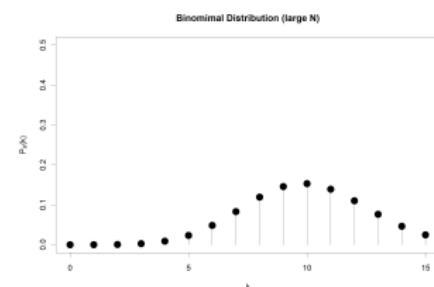
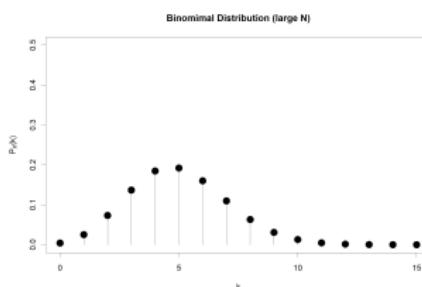


Abbildung: left $\pi = 1/6$, right $\pi = 1/3$ with 30 dice each

Binomial distribution II

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Discrete distributions have formally very similar properties to empirical distributions:

Binomial distribution: Expected value:

$$E(X) = E(X_1) + \dots + E(X_n) = n\pi$$

Variance:

$$\text{Var}(X) = \text{Var}(X_1) + \dots + \text{Var}(X_n) = n\pi(1 - \pi)$$

→ expected value $E(X)$ is analogous to the arithmetic mean \bar{x}

Interpretation: if a random experiment is repeated very often, the arithmetic mean approaches the expected value

Distribution of continuous random variables

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

Random variable X is called continuous if it can take on an infinite number of values x_1, x_2, \dots, x_k (is continuous)

Probabilities are determined by the area $[a, b]$ and the overlying function $f(x)$:

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

Function $f(x)$ is called (probability) density of X

→ for continuous variables, the x -values in $[a, b]$ are no longer countable and a single value x is $P(X = x) = 0$, since the limits $a = b$ would be

briefly density

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

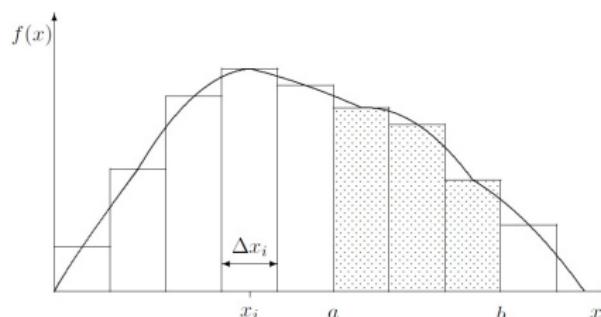
discrete

continuous

Literature

following is an approximating probability histogram, where the rectangular areas over Δx_i can be calculated using $P(X_d = x_i) = f(x_i)\Delta x_i$.

→ If the class widths Δx_i approach zero, X_d becomes a continuous random variable



Normal distribution I

Motivation

Content

Fundamentals

Data collection

Level of measurement

Surveys

Descriptive statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

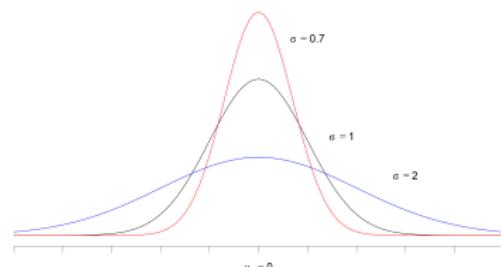
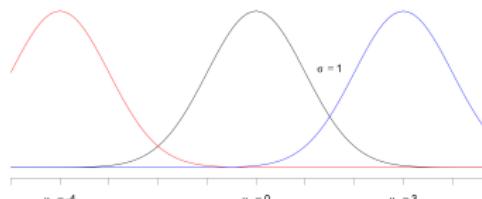
continuous

Literature

Normal distribution (also Gaussian bell curve):

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad [-\infty < x < \infty]$$

- parameter μ (left figure): maximum of the normal distribution.
- parameter σ (right figure): dispersion of the normal distribution



Normal distribution II

[Motivation](#)[Content](#)[Fundamentals](#)[Data collection](#)[Level of measurement](#)[Surveys](#)[Descriptive statistics](#)[univariat](#)[Graphs](#)[Measures](#)[Symmetry](#)[multivariate](#)[Graphs](#)[Measures](#)[Theory](#)[Probability](#)[Randomness](#)[Distributions](#)[discrete](#)[continuous](#)[Literature](#)

a normally distributed random variable X can be standardised as follows:

$$Z = \frac{X - \mu}{\sigma}$$

Standard normal distribution $X \sim N(\mu = 0, \sigma = 1)$:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} \quad [-\infty < x < \infty]$$

Standardisation of normally distributed random variables has multiple advantages: Comparison of variables of different metrics, $Cov(X, Y) = r_{xy}$, in a regression coefficients directly correspond to effect sizes, ...

Take Home Message: Distributions

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariat

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

- Random variables describe the outcome of a random experiment by numbers
- the most important continuous distribution is the normal distribution
- Note: Binomial distribution can be approximated by the normal distribution under the following conditions (*limiting form*)
 - number of independent trials goes towards infinity $n \rightarrow \infty$
 - probability of success π does not go towards 0 or 1

Literature

Motivation

Content

Fundamentals

Data collection

Level of
measurement

Surveys

Descriptive
statistics

univariate

Graphs

Measures

Symmetry

multivariate

Graphs

Measures

Theory

Probability

Randomness

Distributions

discrete

continuous

Literature

-  Bromme, Rainer, Manfred Prenzel und Michael Jäger (2014). „Empirische Bildungsforschung und evidenzbasierte Bildungspolitik“. In: *Zeitschrift für Erziehungswissenschaft* 17.4, S. 3–54.
-  Cohen, Jacob (1992). „A Power Primer.“. In: *Psychological Bulletin* 112.1, S. 155–159.
-  Dormann, Carsten (2013). *Parametrische Statistik Verteilungen, maximum likelihood und GLM in R*. Springer Spektrum.
-  Elbow, Peter (1993). „The uses of binary thinking“. In: *Journal of Advanced Composition*, S. 51–78.
-  Fahrmeir, Ludwig et al. (2016). *Statistik: Der Weg zur Datenanalyse*. Springer-Verlag.
-  Hoyle, Rick (2012). *Handbook of structural equation modeling*. Guilford Press.
-  Kaldor, Mary (2018). „Cycles in world politics“. In: *International Studies Review* 20.2, S. 214–222.
-  Kelle, Udo (2008). *Die Integration qualitativer und quantitativer Methoden in der empirischen Sozialforschung*. Springer.
-  Leeuw, Edith, Joop Hox und Don Dillman (2012). *International handbook of survey methodology*. Routledge.
-  Lintorf, Katrin (2012). *Wie vorhersehbar sind Grundschulnoten? Prädiktionskraft individueller und kontextspezifischer Merkmale*. Springer-Verlag.
-  Lütkepohl, Helmut und Markus Krätzig (2004). *Applied time series econometrics*. Cambridge University Press.
-  Macho, Siegfried (2016). „Moderne Testtheorie“. Université de Fribourg.
-  Neumann, Astrid (2007). *Briefe schreiben in Klasse 9 und 11. Beurteilungskriterien, Messungen, Textstrukturen und Schülerleistungen*. Waxmann Verlag.
-  O'Connell, Michael (2019). „Is the impact of SES on educational performance overestimated? Evidence from the PISA survey“. In: *Intelligence* 75, S. 41–47.
-  Pickel, Gert und Susanne Pickel (2018). „Migration als Gefahr für die politische Kultur?“ In: *Zeitschrift für Vergleichende Politikwissenschaft* 12.1, S. 297–320.
-  Shadish, William R., Thomas D Cook und Donald Thomas Campbell (2002). *Experimental and quasi-experimental designs for generalized causal inference*. Wadsworth Cengage Learning.
-  Tsoukas, Haridimos und Mary Hatch (2001). „Complex thinking, complex practice: The case for a narrative approach to organizational complexity“. In: *Human relations* 54.8, S. 979–1013.
-  Wolf, Sabrina (2008). „Der Methodenstreit quantitativer und qualitativer Sozialforschung unter besonderer Berücksichtigung der grundlegenden Unterschiede beider Forschungstraditionen“. Bachelor Thesis. Universität Augsburg.
-  Zucchini, Walter et al. (2009). *Statistik für Bachelor- und Masterstudenten: eine Einführung für Wirtschafts- und Sozialwissenschaftler*. Springer-Verlag.