

Deskriptive Statistik und Wahrscheinlichkeitsrechnung

Vorlesung an der Universität des Saarlandes

PD Dr. Martin Becker

Sommersemester 2021



Organisatorisches I

- Vorlesung: voraussichtlich nur online, Inhalte jederzeit abrufbar
- Übungen: voraussichtlich nur online, Inhalte jederzeit abrufbar
- Prüfung: *voraussichtlich* 2-stündige Klausur nach Semesterende (1. Prüfungszeitraum)
Anmeldung und Informationen zum Termin im ViPa
- Hilfsmittel für Klausur
 - ▶ „Moderat“ programmierbarer Taschenrechner, auch mit Grafikfähigkeit
 - ▶ 2 *beliebig gestaltete* DIN A 4–Blätter (bzw. 4, falls nur einseitig)
 - ▶ Benötigte Tabellen werden gestellt, aber **keine weitere Formelsammlung!**
- Durchgefallen — was dann?
 - ▶ „Wiederholungskurs“ im kommenden (Winter-)Semester
 - ▶ „Nachprüfung“ (voraussichtlich) erst März/April 2022 (2. Prüfungszeitraum)
 - ▶ „Reguläre“ Vorlesung/Übungen wieder im Sommersemester 2022

Organisatorisches II

- Informationen und Materialien über Moodle sowie unter
<https://www.lehrstab-statistik.de>
bzw. spezieller
<https://www.lehrstab-statistik.de/deskrwrss2021.html>
(bei Problemen <https://www2.lehrstab-statistik.de> versuchen!)
- Kontakt: PD Dr. Martin Becker
Geb. C3 1, 2. OG, Zi. 2.17 (im Präsenzbetrieb)
e-Mail: martin.becker@mx.uni-saarland.de
- Sprechstunde (via MS Teams) nach Terminabstimmung per e-Mail
- Vorlesungsunterlagen
 - ▶ Vorlesungsfolien
 - ▶ Erklär-Videos zu den Vorlesungsfolien
 - ▶ Zusätzlich: lehrbuchartige Aufbereitung der Inhalte der ersten drei Wochen im Online-Skript

Organisatorisches III

- Übungsunterlagen
 - ▶ Übungsblätter i.d.R. zusammen mit neuen Vorlesungsunterlagen zum Download
 - ▶ Ergebnisse (*keine Musterlösungen!*) zu den meisten Aufgaben ebenfalls unmittelbar verfügbar
 - ▶ Ausführlichere Lösungen zu den Übungsaufgaben (Online-Skript + noch ausführlichere Erklärvideos) einige Tage später, *damit Sie nicht zu sehr in Versuchung geraten, sich die Lösung vor der eigenen Bearbeitung der Übungsblätter anzuschauen!*
 - ▶ Eigene Bearbeitung der Übungsblätter (**vor** Betrachten der bereitgestellten Lösungen) wichtigste Klausurvorbereitung (eine vorhandene Lösung zu verstehen etwas **ganz** anderes als eine eigene Lösung zu finden!).
- Alte Klausuren
 - ▶ Aktuelle Klausuren inklusive der meisten Ergebnisse unter „Klausuren“ auf Homepage des Lehrstabs verfügbar
 - ▶ Prüfungsrelevant sind (natürlich) alle in Vorlesung und Übungsprogramm behandelten Inhalte, nicht nur die Inhalte der Altklausuren!

Was ist eigentlich „Statistik“?

- Der Begriff „Statistik“ hat verschiedene Bedeutungen, insbesondere:
 - ▶ Oberbegriff für die Gesamtheit der Methoden, die für die Erhebung und Verarbeitung empirischer Informationen relevant sind (→ statistische Methodenlehre)
 - ▶ (Konkrete) Tabellarische oder grafische Darstellung von Daten
 - ▶ (Konkrete) Abbildungsvorschrift, die in Daten enthaltene Informationen auf eine „Kennzahl“ (→ Teststatistik) verdichtet
- Grundlegende Teilgebiete der Statistik:
 - ▶ Deskriptive Statistik (auch: beschreibende Statistik, explorative Statistik)
 - ▶ Schließende Statistik (auch: inferenzielle Statistik, induktive Statistik)
- Typischer Einsatz von Statistik:

Verarbeitung — insbesondere Aggregation — von (eventuell noch zu erhebenden) Daten mit dem Ziel, (informelle) Erkenntnisgewinne zu erhalten bzw. (formal) Schlüsse zu ziehen.

↪ Bestimmte Informationen „ausblenden“, um neue Informationen zu erkennen

Vorurteile gegenüber Statistik

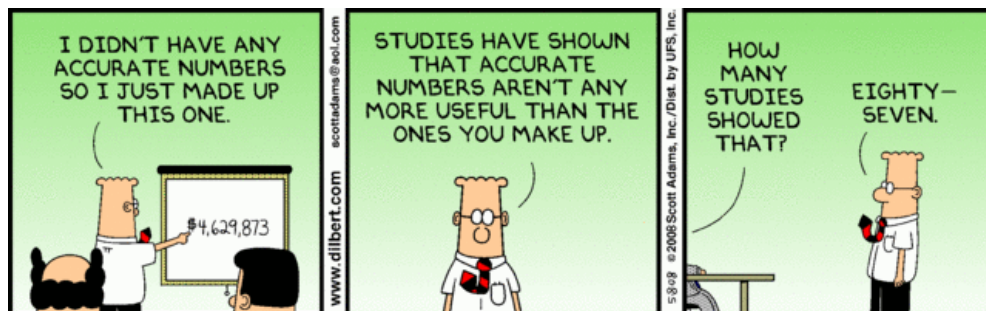
- Einige Zitate oder „Volksweisheiten“:
 - ▶ „Statistik ist pure Mathematik, und in Mathe war ich immer schlecht...“
 - ▶ „Mit Statistik kann man alles beweisen!“
 - ▶ „Ich glaube nur der Statistik, die ich selbst gefälscht habe.“
(häufig Winston Churchill zugeschrieben, aber eher Churchill von Goebbels' Propagandaministerium „in den Mund gelegt“)
 - ▶ „There are three kinds of lies: lies, damned lies, and statistics.“
(häufig Benjamin Disraeli zugeschrieben)

↪ negative Vorurteile gegenüber der Disziplin „Statistik“
- Tatsächlich aber
 - ▶ benötigt man für viele statistische Methoden nur die vier Grundrechenarten.
 - ▶ ist „gesunder Menschenverstand“ viel wichtiger als mathematisches Know-How.
 - ▶ sind nicht die statistischen Methoden an sich schlecht oder gar falsch, sondern die korrekte Auswahl und Anwendung der Methoden zu hinterfragen.
 - ▶ werden viele (korrekte) Ergebnisse statistischer Untersuchungen lediglich falsch interpretiert.

Kann man mit Statistik lügen? I

Und falls ja, wie (schützt man sich dagegen)?

- Natürlich kann man mit Statistik „lügen“ bzw. täuschen!
- „Anleitung“ von Prof. Dr. Walter Krämer (TU Dortmund):
So lügt man mit Statistik, Campus, 2015
- Offensichtliche Möglichkeit: Daten (vorsätzlich) manipulieren/fälschen:



Kann man mit Statistik lügen? II

Und falls ja, wie (schützt man sich dagegen)?

- Weitere Möglichkeiten zur Täuschung
 - ▶ Irreführende Grafiken
 - ▶ (Bewusstes) Weglassen relevanter Information
 - ▶ (Bewusste) Auswahl ungeeigneter statistischer Methoden
- Häufiges Problem (vor allem in den Medien):
Suggestion von Sicherheit durch hohe Genauigkeit angegebener Werte
↪ zusätzlich: Ablenkung vom „Adäquationsproblem“
(misst der angegebene Wert überhaupt das „Richtige“?)
- Schutz vor Täuschung:
 - ▶ Mitdenken!
 - ▶ „Gesunden Menschenverstand“ einschalten!
 - ▶ Gute Grundkenntnisse in Statistik!

Beispiel (Adäquationsproblem) I

vgl. Walter Krämer: So lügt man mit Statistik, Piper, München, 2009

- Frage: Was ist *im Durchschnitt* sicherer, Reisen mit Bahn oder Flugzeug?
- Statistik 1:

Bahn	9 Verkehrstote pro 10 Milliarden Passagierkilometer
Flugzeug	3 Verkehrstote pro 10 Milliarden Passagierkilometer

~> Fliegen sicherer als Bahnfahren!
- Statistik 2:

Bahn	7 Verkehrstote pro 100 Millionen Passagierstunden
Flugzeug	24 Verkehrstote pro 100 Millionen Passagierstunden

~> Bahnfahren sicherer als Fliegen!
- Widerspruch? Fehler?

Beispiel (Adäquationsproblem) II

vgl. Walter Krämer: So lügt man mit Statistik, Piper, München, 2009

- Nein, Unterschied erklärt sich durch höhere Durchschnittsgeschwindigkeit in Flugzeugen (Annahme: ca. 800 km/h vs. ca. 80 km/h)
- Wie wird „Sicherheit“ gemessen? Welcher „Durchschnitt“ ist geeigneter?

~> Interpretation abhängig von der Fragestellung! Hier:

 - ▶ Steht man vor der Wahl, eine gegebene Strecke per Bahn oder Flugzeug zurückzulegen, so ist Fliegen sicherer.
 - ▶ Vor einem vierstündigen Flug ist dennoch eine größere „Todesangst“ angemessen als vor einer vierstündigen Bahnfahrt.

Beispiel („Schlechte“ Statistik) I

- Studie/Pressemitteilung des ACE Auto Club Europa *anlässlich des Frauentags am 8. März 2010*: „Autofahrerinnen im Osten am besten“ (siehe https://www.ace.de/fileadmin/user_uploads/Der_Club/Dokumente/Verkehrspolitik/Handout-Booklet-ACE-Studien.pdf, S. 88–90)
- Untersuchungsgegenstand:
 - ▶ Regionale Unterschiede bei Unfallhäufigkeit mit Frauen als Hauptverursacher
 - ▶ Vergleich Unfallhäufigkeit mit Frau bzw. Mann als Hauptverursacher
- Wesentliche Datengrundlage ist eine Publikation des Statistischen Bundesamts (Destatis): „Unfälle im Straßenverkehr nach Geschlecht 2008“

Beispiel („Schlechte“ Statistik) II

- Beginn der Pressemitteilung des ACE:
„Von wegen schwaches Geschlecht: Hinterm Steuer sind Frauen besonders stark.“

Weiter heißt es:

“Auch die durch Autofahrerinnen verursachten Unfälle mit Personenschaden liegen wesentlich hinter den von Männern verursachten gleichartigen Karambolagen zurück.“

und in einer Zwischenüberschrift

„Schlechtere Autofahrerinnen sind immer noch besser als Männer“

Beispiel („Schlechte“ Statistik) III

- „Statistische“ Argumentation: Laut Destatis-Quelle sind (**angeblich!**)
 - ▶ mehr als 2/3 aller Unfälle mit Personenschaden 2008 (genauer: 217 843 von etwas über 320 000 Unfällen) durch PKW-fahrende Männer verursacht worden,
 - ▶ nur 37% aller Unfälle mit Personenschaden 2008 durch PKW-fahrende Frauen verursacht worden.
- Erste Auffälligkeit: $66.6\% + 37\% = 103.6\%$ (???)
- Lösung: **Ablesefehler** (217 843 aller 320 614 Unfälle mit Personenschaden (67.9%) wurden mit **PKW-Fahrer** (geschlechtsunabhängig) als Hauptverursacher registriert)

Beispiel („Schlechte“ Statistik) IV

- Korrekte Werte:
 - ▶ Bei 210 905 der 217 843 Hauptunfallverursacher als PKW-Fahrzeugführer wurde Geschlecht registriert.
 - ▶ 132 757 waren männlich (62.95%), 78 148 weiblich (37.05%)
- **Also:** immer noch deutlich mehr Unfälle mit PKW-fahrenden Männern als Hauptverursacher im Vergleich zu PKW-Fahrerinnen.
- **Aber:** Absolute Anzahl von Unfällen geeignetes Kriterium für Fahrsicherheit?

Beispiel („Schlechte“ Statistik) V

- Modellrechnung des DIW aus dem Jahr 2004 schätzt
 - ▶ Anzahl Männer mit PKW-Führerschein auf 28.556 Millionen,
 - ▶ Anzahl Frauen mit PKW-Führerschein auf 24.573 Millionen.
- Weitere ältere Studie (von 2002) schätzt
 - ▶ durchschnittliche Fahrleistung von Männern mit PKW-Führerschein auf 30 km/Tag,
 - ▶ durchschnittliche Fahrleistung von Frauen mit PKW-Führerschein auf 12 km/Tag.
- Damit stehen also
 - ▶ bei Männern 132 757 verursachte Unfälle geschätzten $30 \cdot 365 \cdot 28.556 = 312688.2$ Millionen gefahrenen Kilometern,
 - ▶ bei Frauen 78 148 verursachte Unfälle geschätzten $12 \cdot 365 \cdot 24.573 = 107629.74$ Millionen gefahrenen Kilometern gegenüber.

Beispiel („Schlechte“ Statistik) VI

- Dies führt im Durchschnitt
 - ▶ bei Männern zu 0.425 verursachten Unfällen mit Personenschaden pro eine Million gefahrenen Kilometern,
 - ▶ bei Frauen zu 0.726 verursachten Unfällen mit Personenschaden pro eine Million gefahrenen Kilometern.
- Pro gefahrenem Kilometer verursachen (schätzungsweise) weibliche PKW-Fahrer also durchschnittlich ca. **71% mehr** Unfälle als männliche!
- Anstatt dies zu konkretisieren, räumt die Studie lediglich weit am Ende ein entsprechendes Ungleichgewicht bei der jährlichen Fahrleistung ein.

Beispiel („Schlechte“ Statistik) VII

- Welt Online (siehe <http://www.welt.de/vermishtes/article6674754/Frauen-sind-bessere-Autofahrer-als-Maenner.html>) beruft sich auf die ACE-Studie in einem Artikel mit der Überschrift

„Frauen sind bessere Autofahrer als Männer“

und der prägnanten Bildunterschrift

„Männer glauben bloß, sie seien die besseren Autofahrer. Eine Unfall-Statistik beweist das Gegenteil.“

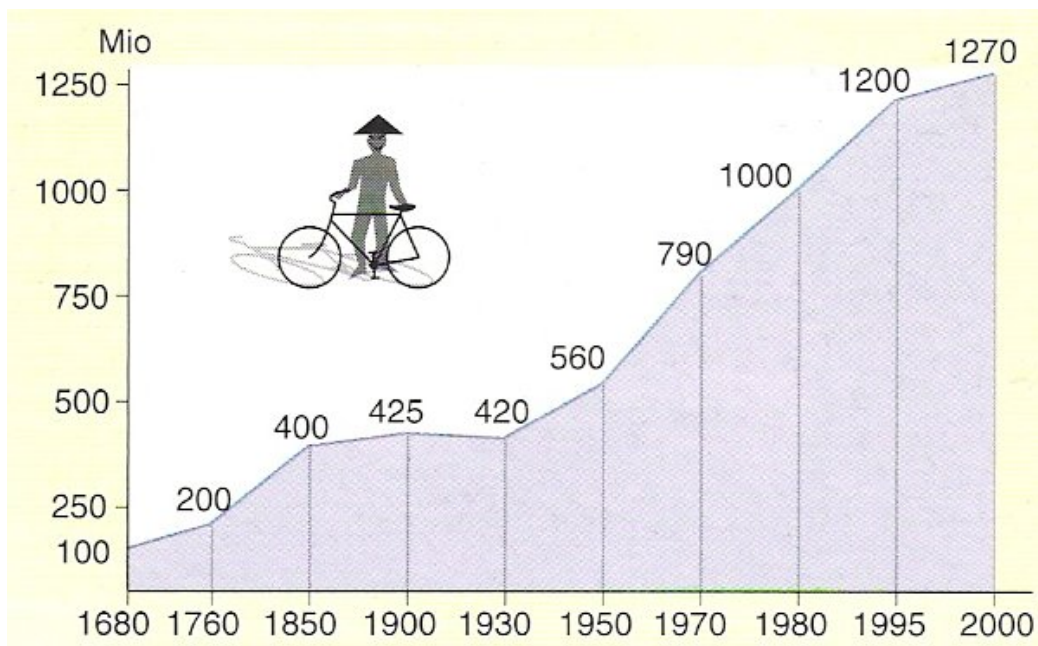
Erst am Ende wird einschränkend erwähnt:

„Fairerweise muss man erwähnen, dass Männer täglich deutlich mehr Kilometer zurücklegen. Und: Während 93 Prozent von ihnen einen Führerschein besitzen, sind es bei den Frauen lediglich 82 Prozent.“

Beispiel (Irreführende Grafik) I

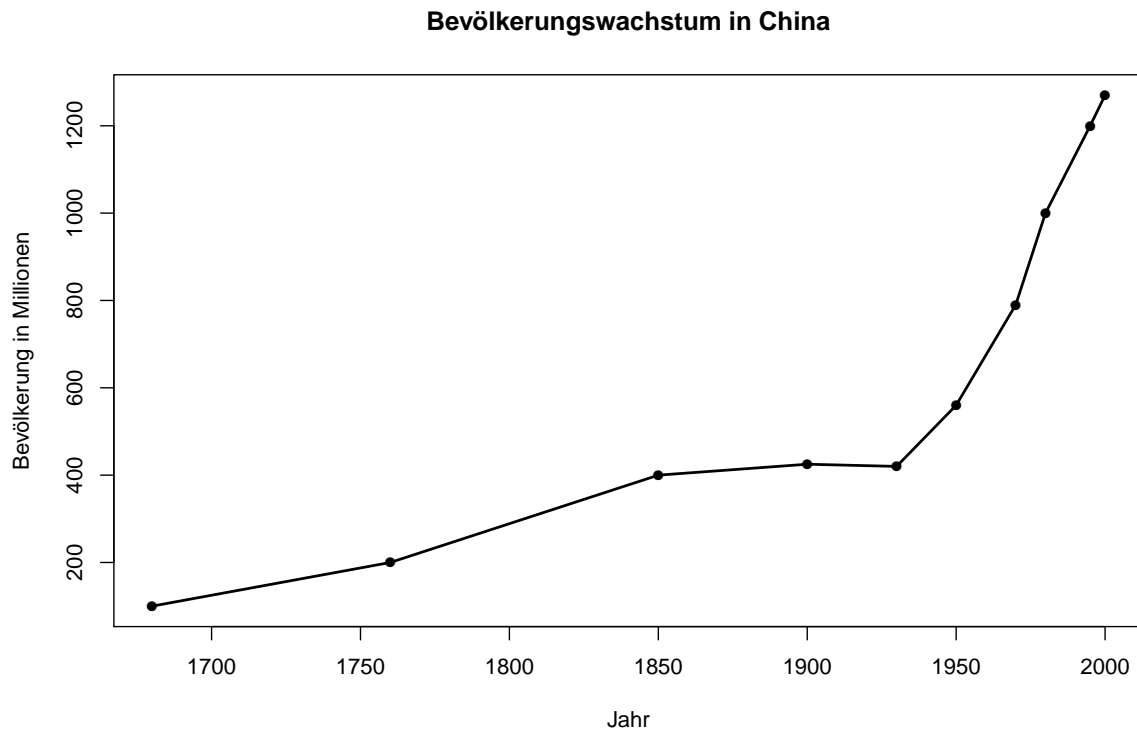
vgl. <http://www.klein-singen.de/statistik/h/Wissenschaft/Bevoelkerungswachstum.html>

Bevölkerungswachstum in China



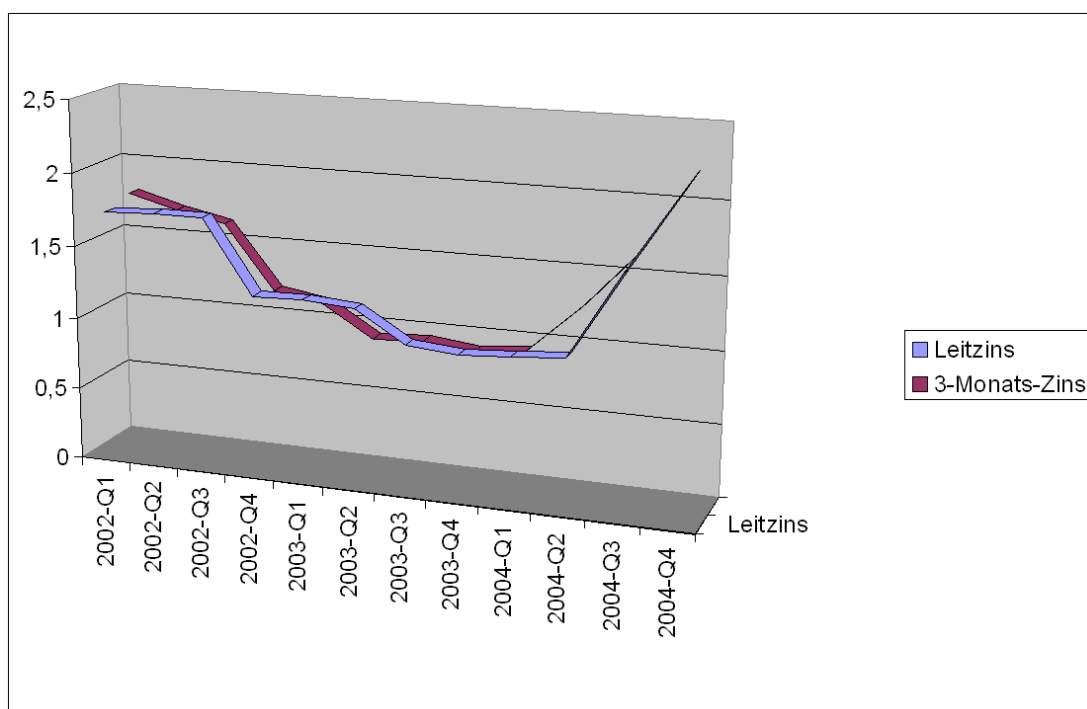
Beispiel (Irreführende Grafik) II

identischer Datensatz, angemessene Skala



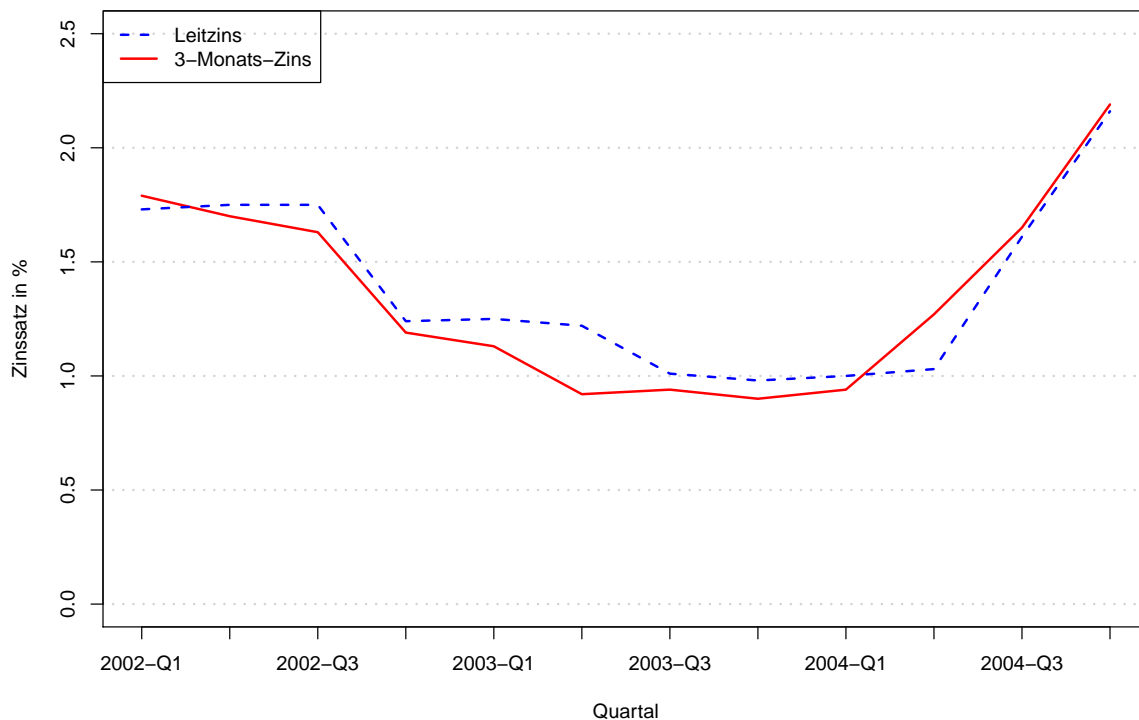
Beispiel (Chartjunk)

Microsoft Excel mit Standardeinstellung für 3D-Liniendiagramme



Beispiel (Grafik ohne Chartjunk)

Statistik-Software R, identischer Datensatz



Kann Statistik auch nützlich sein?

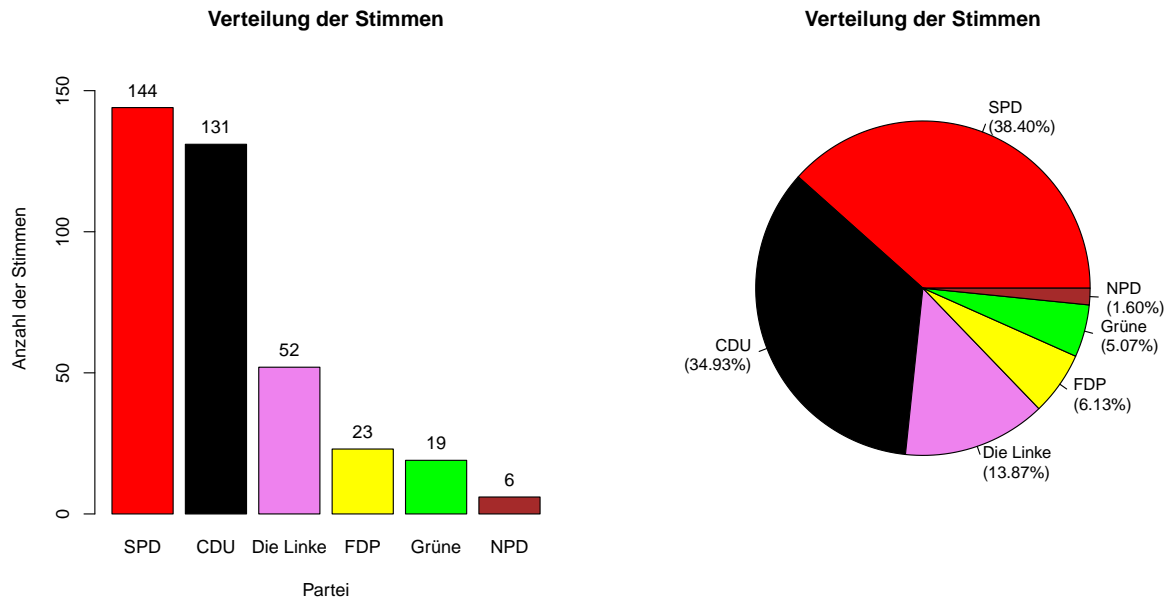
Welche Partei erhält wie viele Stimmen im Wahlbezirk 1.206 der Gemeinde Losheim am See bei den Erststimmen zur Bundestagswahl 2009? Stimmen:

Die Linke, SPD, CDU, Die Linke, SPD, SPD, Die Linke, CDU, FDP, Grüne, Die Linke, SPD, Die Linke, CDU, SPD, CDU, CDU, SPD, SPD, FDP, CDU, FDP, Die Linke, Die Linke, Grüne, CDU, CDU, CDU, CDU, Die Linke, CDU, CDU, CDU, SPD, CDU, SPD, SPD, CDU, FDP, FDP, SPD, CDU, CDU, CDU, CDU, SPD, SPD, SPD, CDU, NPD, SPD, Die Linke, CDU, CDU, FDP, Grüne, SPD, FDP, CDU, CDU, CDU, SPD, SPD, SPD, CDU, Die Linke, CDU, Die Linke, SPD, FDP, CDU, SPD, CDU, CDU, CDU, SPD, Die Linke, CDU, Die Linke, NPD, SPD, Grüne, FDP, SPD, FDP, SPD, CDU, SPD, CDU, SPD, SPD, SPD, SPD, SPD, CDU, CDU, Die Linke, CDU, CDU, SPD, CDU, CDU, Die Linke, CDU, SPD, SPD, SPD, SPD, SPD, SPD, Die Linke, Die Linke, Die Linke, CDU, Die Linke, CDU, Grüne, CDU, CDU, SPD, CDU, SPD, CDU, CDU, SPD, SPD, CDU, FDP, CDU, SPD, SPD, SPD, CDU, CDU, Die Linke, CDU, CDU, CDU, CDU, SPD, FDP, SPD, SPD, Die Linke, SPD, Grüne, SPD, Grüne, FDP, SPD, CDU, Die Linke, FDP, SPD, CDU, SPD, SPD, SPD, SPD, Die Linke, SPD, SPD, CDU, SPD, CDU, Die Linke, SPD, CDU, CDU, CDU, SPD, SPD, SPD, Die Linke, FDP, Grüne, CDU, SPD, CDU, SPD, SPD, Die Linke, SPD, CDU, CDU, CDU, SPD, SPD, SPD, Die Linke, SPD, SPD, SPD, SPD, Die Linke, CDU, CDU, Die Linke, CDU, CDU, SPD, SPD, CDU, CDU, SPD, SPD, CDU, CDU, NPD, SPD, SPD, CDU, SPD, SPD, Grüne, CDU, SPD, SPD, Die Linke, FDP, Die Linke, CDU, SPD, Grüne, SPD, CDU, SPD, Die Linke, Die Linke, SPD, CDU, Die Linke, SPD, SPD, SPD, Die Linke, Die Linke, SPD, SPD, FDP, CDU, CDU, SPD, SPD, CDU, SPD, CDU, SPD, SPD, CDU, SPD, CDU, CDU, SPD, Grüne, SPD, SPD, SPD, CDU, CDU, SPD, SPD, SPD, FDP, Die Linke, CDU, FDP, CDU, Die Linke, SPD, CDU, CDU, CDU, CDU, Grüne, CDU, CDU, CDU, SPD, Die Linke, SPD, Die Linke, NPD, CDU, Grüne, Die Linke, CDU, CDU, Die Linke, Die Linke, SPD, SPD, CDU, Grüne, SPD, Die Linke, SPD, SPD, SPD, CDU, Die Linke, SPD, SPD, NPD, SPD, CDU, SPD, SPD, SPD, Grüne, CDU, SPD, SPD, FDP, Grüne, SPD, Die Linke, CDU, SPD, SPD, CDU, SPD, Die Linke, Die Linke, CDU, FDP, CDU, SPD, Die Linke, SPD, CDU, CDU, SPD, SPD, SPD, CDU, CDU, Grüne, CDU, CDU, CDU, FDP, Die Linke, SPD, CDU, Die Linke, CDU, SPD, CDU, FDP, SPD, SPD, CDU, SPD, CDU, CDU, CDU, NPD, CDU, Grüne, SPD, SPD, CDU, Grüne, CDU, SPD, CDU, SPD

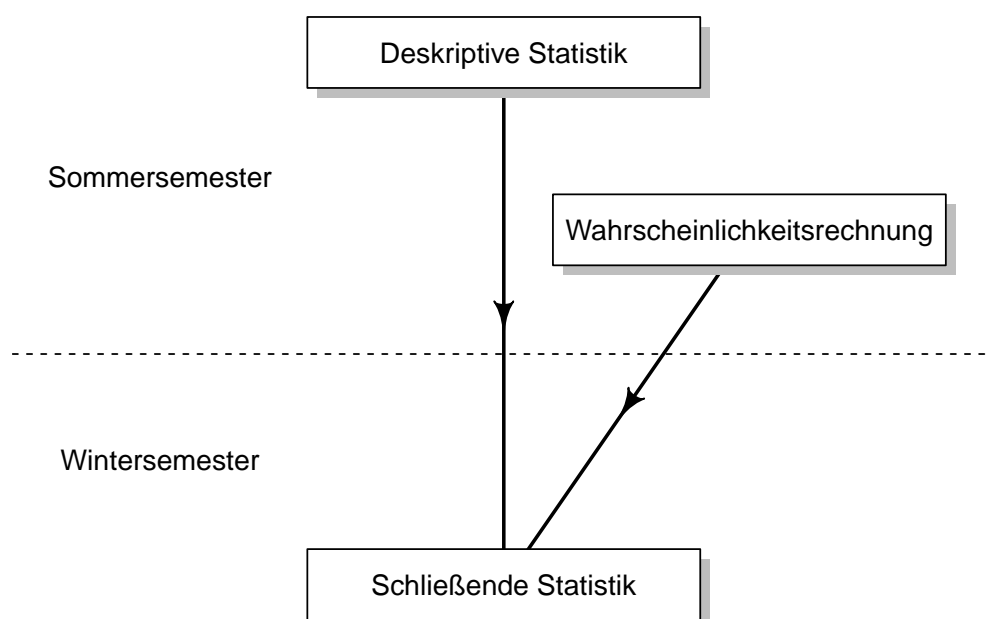
- Mit etwas (deskriptiver) Statistik in tabellarischer Form:

	SPD	CDU	Die Linke	FDP	Grüne	NPD	Summe
Anzahl der Stimmen	144	131	52	23	19	6	375
Stimmenanteil in %	38.40	34.93	13.87	6.13	5.07	1.60	100.00

- Grafisch aufbereitete Varianten:



Organisation der Statistik-Veranstaltungen



Teil I

Deskriptive Statistik

Datenerhebung I

- Beginn jeder (deskriptiven) statistischen Untersuchung: Datenerhebung
- Zu einer **Menge von Merkmalsträgern (statistische Masse)**, eventuell Teil einer größeren **Grundgesamtheit**, werden ein oder mehrere **Merkmale** erhoben
- Unterscheidung nach
 - ▶ Primärerhebung ↔ Sekundärerhebung:
Neue Erhebung oder Nutzung von vorhandenem Datenmaterial
 - ▶ Vollerhebung ↔ Teilerhebung:
Erhebung der Merkmale für ganze Grundgesamtheit oder Teilgesamtheit

Datenerhebung II

- Bei Primärerhebung: Untersuchungsziel bestimmt
 - ▶ Auswahl bzw. Abgrenzung der statistischen Masse
 - ▶ Auswahl der zu erhebenden Merkmale
 - ▶ Art der Erhebung, z.B. Befragung (Post, Telefon, Internet, persönlich), Beobachtung, Experiment
- Sorgfalt bei Datenerhebung enorm wichtig:
Fehler bei Datenerhebung sind später nicht mehr zu korrigieren!
- Ausführliche Diskussion hier aus Zeitgründen nicht möglich

Vorsicht vor „falschen Schlüssen“! I

- Deskriptive Statistik fasst lediglich Information über statistische Masse zusammen
- Schlüsse auf (größere) „Grundgesamtheit“ (bei Teilerhebung)
~> Schließende Statistik
- Dennoch häufig zu beobachten:
„Informelles“ Übertragen der Ergebnisse in der statistischen Masse auf größere Menge von Merkmalsträgern
~> Gefahr von falschen Schlüssen!

Vorsicht vor „falschen Schlüssen“! II

Beispiel: Bachelor-Absolventen (vgl. Krämer: So lügt man mit Statistik)

Hätte man am Ende des SS 2011 in der statistischen Masse der Absolventen des BWL-Bachelorstudiengangs in Saarbrücken die Merkmale „Studiendauer“ und „Abschlussnote“ erhoben, würde man wohl feststellen, dass alle Abschlüsse in Regelstudienzeit und im Durchschnitt mit einer guten Note erfolgt sind. Warum? Kann man dies ohne weiteres auf Absolventen anderer Semester übertragen?

- ↪ Zur Interpretationsfähigkeit von Ergebnissen statistischer Untersuchungen:
- ▶ Abgrenzung der zugrundeliegenden statistischen Masse **sehr** wichtig
 - ▶ (Möglichst) objektive Festlegung nach Kriterien zeitlicher, räumlicher und sachlicher Art

Definition 2.1 (Menge, Mächtigkeit, Tupel)

- ① Eine (endliche) **Menge** M ist die Zusammenfassung (endlich vieler) unterschiedlicher Objekte (Elemente).
- ② Zu einer endlichen Menge M bezeichnen $\#M$ oder auch $|M|$ die Anzahl der Elemente in M . $\#M$ bzw. $|M|$ heißen auch **Mächtigkeit** der Menge M .
- ③ Für eine Anzahl $n \geq 1$ von (nicht notwendigerweise verschiedenen!) Elementen x_1, x_2, \dots, x_n aus einer Menge M wird eine (nach ihrer Reihenfolge geordnete) Auflistung (x_1, x_2, \dots, x_n) bzw. x_1, x_2, \dots, x_n als **n -Tupel** aus der Menge M bezeichnet. 2-Tupel (x_1, x_2) heißen auch Paare.
- ④ Lassen sich die Elemente der Menge M (der Größe nach) ordnen, so sei (zu einer vorgegebenen Ordnung)
 - ① mit $(x_{(1)}, x_{(2)}, \dots, x_{(n)})$ bzw. $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ das der Größe nach geordnete n -Tupel der n Elemente x_1, x_2, \dots, x_n aus M bezeichnet, es gelte also $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$.
 - ② zu einer endlichen Teilmenge $A \subseteq M$ der Mächtigkeit m mit $(a_{(1)}, a_{(2)}, \dots, a_{(m)})$ bzw. $a_{(1)}, a_{(2)}, \dots, a_{(m)}$ das der Größe nach geordnete m -Tupel der Elemente a_1, a_2, \dots, a_m von A bezeichnet, es gelte also $a_{(1)} < a_{(2)} < \dots < a_{(m)}$.

Merkmalswerte, Merkmalsraum, Urliste I

- Bei der Datenerhebung werden den Merkmalsträgern zu jedem erhobenen Merkmal **Merkmalswerte** oder **Beobachtungswerte** zugeordnet.
- Man nimmt an, dass man (im Prinzip auch vor der Erhebung) eine Menge M angeben kann, die alle vorstellbaren Merkmalswerte eines Merkmals enthält.
- Das n -Tupel (x_1, \dots, x_n) der Merkmalswerte x_1, \dots, x_n (aus der Menge M) zu einem bei den n Merkmalsträgern erhobenen Merkmal X bezeichnet man als **Urliste**.
- Die Menge A der (verschiedenen) in der Urliste (tatsächlich) auftretenden Merkmalswerte, in Zeichen

$$A := \{a \in M \mid \exists i \in \{1, \dots, n\} \text{ mit } x_i = a\} ,$$

heißt **Merkmalsraum**, ihre Elemente **Merkmalsausprägungen**.

Merkmalswerte, Merkmalsraum, Urliste II

Beispiel Wahlergebnis

- ▶ Urliste (siehe Folie 22) aus gewählten Parteien der 375 abgegebenen gültigen Stimmen:
 $x_1 = \text{“Die Linke”}, x_2 = \text{“SPD”}, x_3 = \text{“CDU”}, x_4 = \text{“Die Linke”}, x_5 = \text{“SPD”},$
 $x_6 = \text{“SPD”}, x_7 = \text{“Die Linke”}, x_8 = \text{“CDU”}, x_9 = \text{“FDP”}, x_{10} = \text{“Grüne”}, x_{11} =$
 $\text{“Die Linke”}, x_{12} = \text{“SPD”}, x_{13} = \text{“Die Linke”}, x_{14} = \text{“CDU”}, x_{15} = \text{“SPD”}, x_{16} =$
 $\text{“CDU”}, x_{17} = \text{“CDU”}, x_{18} = \text{“SPD”}, x_{19} = \text{“SPD”}, x_{20} = \text{“FDP”}, \dots$
- ▶ Merkmalsraum: $A = \{\text{SPD, CDU, Die Linke, FDP, Grüne, NPD}\}$

Merkmalstypen I

Definition 2.2 (Merkmalstypen)

- ① Ein Merkmal heißt
 - ▶ **nominalskaliert**, wenn seine Ausprägungen lediglich unterschieden werden sollen,
 - ▶ **ordinalskaliert** oder **rangskaliert**, wenn (darüberhinaus) eine (Rang-)Ordnung auf den Ausprägungen vorgegeben ist,
 - ▶ **kardinalskaliert** oder **metrisch skaliert**, wenn (darüberhinaus) ein „Abstand“ auf der Menge der Ausprägungen vorgegeben ist, also wenn das Ausmaß der Unterschiede zwischen verschiedenen Ausprägungen gemessen werden kann.
- ② Ein Merkmal heißt **quantitativ**, wenn es kardinalskaliert ist, **qualitativ** sonst.
- ③ Ein Merkmal heißt
 - ▶ **diskret**, wenn es qualitativ ist oder wenn es quantitativ ist und die Menge der möglichen Ausprägungen endlich oder abzählbar unendlich ist,
 - ▶ **stetig**, wenn es quantitativ ist und für je zwei mögliche Merkmalsausprägungen auch alle Zwischenwerte angenommen werden können.

Merkmalstypen II

- Welche der in Definition 2.2 erwähnten Eigenschaften für ein Merkmal zutreffend sind, hängt von der jeweiligen Anwendungssituation ab.
- Insbesondere ist die Abgrenzung zwischen stetigen und diskreten Merkmalen oft schwierig (allerdings meist auch nicht besonders wichtig).
- Damit ein Merkmal (mindestens) ordinalskaliert ist, muss die verwendete Ordnung — insbesondere bei Mehrdeutigkeit — eindeutig festgelegt sein.
- Häufig findet man zusätzlich zu den in 2.2 erläuterten Skalierungen auch die Begriffe **Intervallskala**, **Verhältnisskala** und **Absolutskala**. Diese stellen eine feinere Unterteilung der Kardinalskala dar.
- *Unabhängig vom Skalierungsniveau* heißt ein Merkmal **numerisch**, wenn seine Merkmalsausprägungen Zahlenwerte sind.

Merkmalstypen III

Beispiel (Merkmalstypen)

- ▶ nominalskalierte Merkmale: Geschlecht (Ausprägungen: „männlich“, „weiblich“, „divers“), Parteien (siehe Wahlergebnis-Beispiel)
- ▶ ordinalskalierte Merkmale: Platzierungen, Zufriedenheit („sehr zufrieden“, „eher zufrieden“, „weniger zufrieden“, „unzufrieden“)
- ▶ kardinalskalierte Merkmale: Anzahl Kinder, Anzahl Zimmer in Wohnung, Preise, Gewichte, Streckenlängen, Zeiten
 - ★ davon diskret: Anzahl Kinder, Anzahl Zimmer in Wohnung,
 - ★ davon (eher) stetig: Preise, Gewichte, Streckenlängen, Zeiten

Umwandlung von Merkmalstypen I

- Umwandlung qualitativer in quantitative Merkmale durch **Quantifizierung**:
 - ▶ Ersetzen des qualitativen Merkmals „Berufserfahrung“ mit den Ausprägungen „Praktikant“, „Lehrling“, „Geselle“, „Meister“ durch quantitatives Merkmal, dessen Ausprägungen den (mindestens) erforderlichen Jahren an Berufspraxis entsprechen, die zum Erreichen des Erfahrungsgrades erforderlich sind.
 - ▶ Ersetzen des qualitativen Merkmals Schulnote mit den Ausprägungen „sehr gut“, „gut“, „befriedigend“, „ausreichend“, „mangelhaft“, „ungenügend“ (eventuell feiner abgestuft durch Zusätze „+“ und „-“) durch quantitatives Merkmal, z.B. mit den Ausprägungen 15, 14, ..., 00 oder den Ausprägungen 1.0, 1.3, 1.7, 2.0, 2.3, ..., 4.7, 5.0, 6.0.
 - ▶ **Vorsicht:** Umwandlung nur vernünftig, wenn Abstände tatsächlich (sinnvoll) interpretiert werden können!

Inhaltsverzeichnis

(Ausschnitt)

3 Eindimensionale Daten

- Häufigkeitsverteilungen unklassierter Daten
- Häufigkeitsverteilungen klassierter Daten
- Lagemaße
- Streuungsmaße
- Box-Plot
- Symmetrie- und Wölbungsmaße

Häufigkeitsverteilungen I

- Geeignetes Mittel zur Verdichtung der Information aus Urlisten vor allem bei diskreten Merkmalen mit „wenigen“ Ausprägungen: **Häufigkeitsverteilungen**
- Zur Erstellung einer Häufigkeitsverteilung: Zählen, wie oft jede Merkmalsausprägung a aus dem Merkmalsraum $A = \{a_1, \dots, a_m\}$ in der Urliste (x_1, \dots, x_n) vorkommt.

- ▶ Die **absoluten Häufigkeiten** $h(a)$ geben für die Merkmalsausprägung $a \in A$ die (absolute) Anzahl der Einträge der Urliste mit der Ausprägung a an, in Zeichen

$$h(a) := \#\{i \in \{1, \dots, n\} \mid x_i = a\} .$$

- ▶ Die **relativen Häufigkeiten** $r(a)$ geben für die Merkmalsausprägung $a \in A$ den (relativen) Anteil der Einträge der Urliste mit der Ausprägung a an der gesamten Urliste an, in Zeichen

$$r(a) := \frac{h(a)}{n} = \frac{\#\{i \in \{1, \dots, n\} \mid x_i = a\}}{n} .$$

Häufigkeitsverteilungen II

- Die absoluten Häufigkeiten sind natürliche Zahlen und summieren sich zu n auf (i.Z. $\sum_{j=1}^m h(a_j) = n$).
- Die relativen Häufigkeiten sind Zahlen zwischen 0 und 1 (bzw. zwischen 0% und 100%) und summieren sich zu 1 (bzw. 100%) auf (i.Z. $\sum_{j=1}^m r(a_j) = 1$).
- Ist die Anordnung (Reihenfolge) der Urliste unwichtig, geht durch Übergang zur Häufigkeitsverteilung keine relevante Information verloren.
- Häufigkeitsverteilungen werden in der Regel in tabellarischer Form angegeben, am Beispiel des Wahlergebnisses:

	SPD	CDU	Die Linke	FDP	Grüne	NPD	Summe
a_j	a_1	a_2	a_3	a_4	a_5	a_6	Σ
$h(a_j)$	144	131	52	23	19	6	375
$r(a_j)$	0.3840	0.3493	0.1387	0.0613	0.0507	0.0160	1.0000

Häufigkeitsverteilungen III

- Grafische Darstellung (insbesondere bei nominalskalierten Merkmalen) durch **Balkendiagramme** (auch: Säulendiagramme) oder **Kuchendiagramme** (siehe Folie 23).
- Balkendiagramme meist geeigneter als Kuchendiagramme (außer, wenn die anteilige Verteilung der Merkmalsausprägungen im Vordergrund steht)
- Oft mehrere Anordnungen der Spalten/Balken/Kreissegmente bei nominalskalierten Merkmalen plausibel, absteigende Sortierung nach Häufigkeiten $h(a_j)$ meist sinnvoll.
- Bei ordinalskalierten Merkmalen zweckmäßig: Sortierung der Merkmalsausprägungen nach vorgegebener Ordnung, also

$$a_1 = a_{(1)}, a_2 = a_{(2)}, \dots, a_m = a_{(m)}$$

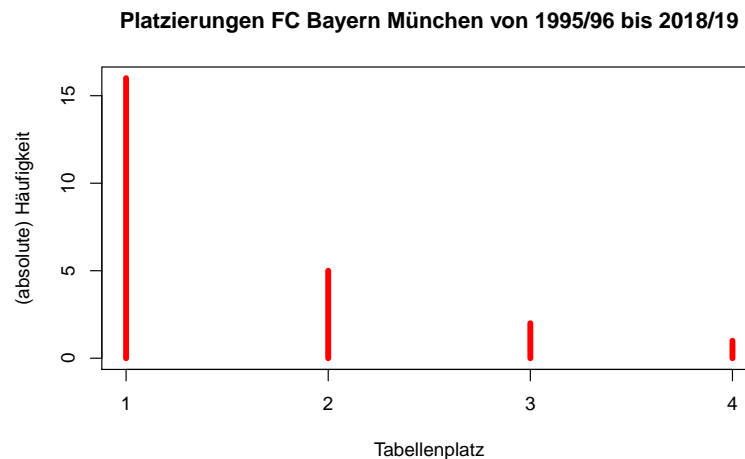
- Alternative grafische Darstellung bei (mindestens) ordinalskalierten Merkmalen mit numerischen Ausprägungen: **Stabdiagramm**

Häufigkeitsverteilungen IV

- Stabdiagramm zur Urliste

2, 1, 2, 1, 1, 1, 3, 1, 2, 1, 1, 4, 1, 2, 1, 3, 2, 1, 1, 1, 1, 1, 1

der finalen Tabellenplätze des FC Bayern München in der (ersten) Fußball-Bundesliga (Saison 1995/96 bis 2018/2019):



Empirische Verteilungsfunktion

- Bei (mindestens ordinalskalierten) numerischen Merkmalen interessante Fragestellungen:
 - ▶ Wie viele Merkmalswerte sind kleiner/größer als ein vorgegebener Wert?
 - ▶ Wie viele Merkmalswerte liegen in einem vorgegebenem Bereich (Intervall)?
- Hierzu nützlich: **(relative) kumulierte Häufigkeitsverteilung**, auch bezeichnet als **empirische Verteilungsfunktion**
- Die empirische Verteilungsfunktion $F(x)$ ordnet einer Zahl x den Anteil der Merkmalswerte x_1, \dots, x_n zu, die kleiner oder gleich x sind, also

$$F(x) := \frac{\#\{i \in \{1, \dots, n\} \mid x_i \leq x\}}{n}.$$

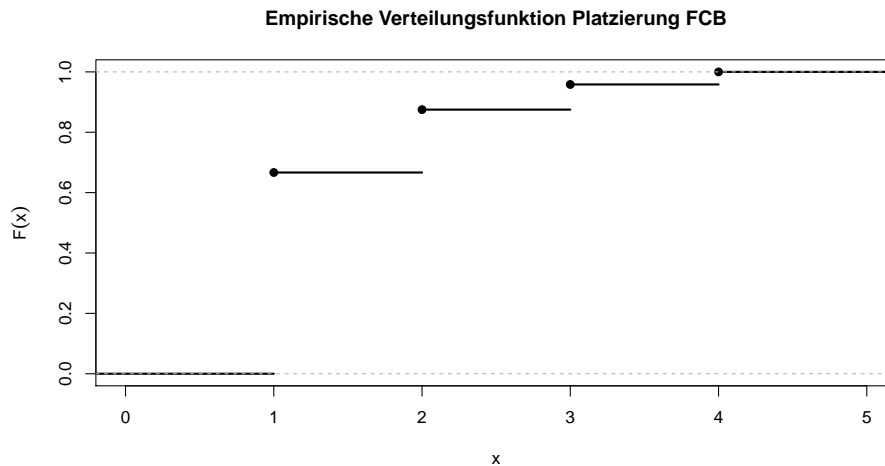
- Ein Vergleich mit den Definitionen von $h(a)$ und $r(a)$ offenbart (!), dass $F(x)$ auch mit Hilfe von $h(a)$ bzw. $r(a)$ berechnet werden kann; gibt es m Merkmalsausprägungen, so gilt:

$$F(x) = \frac{1}{n} \sum_{\substack{a_j \leq x \\ 1 \leq j \leq m}} h(a_j) = \sum_{\substack{a_j \leq x \\ 1 \leq j \leq m}} r(a_j)$$

- Beispiel: Empirische Verteilungsfunktion für FC Bayern-Platzierungen

$$F(x) = \begin{cases} 0 & \text{für } x < 1 \\ \frac{16}{24} & \text{für } 1 \leq x < 2 \\ \frac{21}{24} & \text{für } 2 \leq x < 3 \\ \frac{23}{24} & \text{für } 3 \leq x < 4 \\ 1 & \text{für } x \geq 4 \end{cases} \approx \begin{cases} 0.000 & \text{für } x < 1 \\ 0.667 & \text{für } 1 \leq x < 2 \\ 0.875 & \text{für } 2 \leq x < 3 \\ 0.958 & \text{für } 3 \leq x < 4 \\ 1.000 & \text{für } x \geq 4 \end{cases}$$

- Grafische Darstellung der empirischen Verteilungsfunktion:



Relative Häufigkeiten von Intervallen I

(bei numerischen Merkmalen)

- Relative Häufigkeit $r(a)$ ordnet Ausprägungen $a \in A$ zugehörigen Anteil von a an den Merkmalswerten zu.
- $r(\cdot)$ kann auch für $x \in \mathbb{R}$ mit $x \notin A$ ausgewertet werden ($\rightsquigarrow r(x) = 0$).
- „Erweiterung“ von $r(\cdot)$ auch auf Intervalle möglich:
- $F(b)$ gibt für $b \in \mathbb{R}$ bereits Intervallhäufigkeit

$$F(b) = r((-\infty, b]) = r(\{x \in \mathbb{R} \mid x \leq b\})$$

an.

Relative Häufigkeiten von Intervallen II

(bei numerischen Merkmalen)

- Relative Häufigkeit des offenen Intervalls $(-\infty, b)$ als Differenz

$$r((-\infty, b)) = r((-\infty, b]) - r(b) = F(b) - r(b)$$

- Analog: relative Häufigkeiten weiterer Intervalle:

- ▶ $r((a, \infty)) = 1 - F(a)$
- ▶ $r([a, \infty)) = 1 - (F(a) - r(a)) = 1 - F(a) + r(a)$
- ▶ $r([a, b]) = F(b) - (F(a) - r(a)) = F(b) - F(a) + r(a)$
- ▶ $r((a, b]) = F(b) - F(a)$
- ▶ $r([a, b)) = (F(b) - r(b)) - (F(a) - r(a)) = F(b) - r(b) - F(a) + r(a)$
- ▶ $r((a, b)) = (F(b) - r(b)) - F(a) = F(b) - r(b) - F(a)$