

# Deskriptive Statistik und Wahrscheinlichkeitsrechnung

Vorlesung an der Universität des Saarlandes

Dr. Martin Becker

Sommersemester 2016



## Organisatorisches I

- Vorlesung: Freitag, 12-14 Uhr, Gebäude B4 1, Audimax (HS 0.01)
- Übungen: nach gesonderter Ankündigung, Beginn: ab Montag, 25.04.
- Prüfung: 2-stündige Klausur nach Semesterende (1. Prüfungszeitraum)  
**Anmeldung im ViPa nur vom 12.05. (8 Uhr) – 30.05. (15 Uhr)!**  
(Abmeldung im ViPa bis 14.07., 12 Uhr)
- 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 2017 (2. Prüfungszeitraum)
  - ▶ „Reguläre“ Vorlesung/Übungen wieder im Sommersemester 2017

# Organisatorisches II

- Informationen und Materialien unter  
<http://www.lehrstab-statistik.de>  
bzw. spezieller  
<http://www.lehrstab-statistik.de/deskrwrss2016.html>
- Kontakt: Dr. Martin Becker  
Geb. C3 1, 2. OG, Zi. 2.17  
e-Mail: [martin.becker@mx.uni-saarland.de](mailto:martin.becker@mx.uni-saarland.de)
- Sprechstunde nach Vereinbarung (Terminabstimmung per e-Mail)
- Vorlesungsunterlagen
  - ▶ Vorlesungsfolien (kein Kompaktskript)
  - ▶ Download (inklusive Drucker-freundlicher 2-auf-1 bzw. 4-auf-1 Versionen)  
spätestens Donnerstags vor Vorlesung möglich

# Organisatorisches III

- Übungsunterlagen
  - ▶ Wöchentliche Übungsblätter
  - ▶ Download i.d.R. kurz nach Ende der Vorlesung Freitag nachmittags möglich
  - ▶ Ebenfalls online: Ergebnisse (*keine Musterlösungen!*) zu einigen Aufgaben
  - ▶ Besprechung der Übungsblätter mit ausführlicheren Lösungsvorschlägen in den Übungsgruppen der folgenden Woche.
  - ▶ **Übungsaufgaben sollten unbedingt vorher selbst bearbeitet werden!**
  - ▶ Eventuell: Freiwillige Bearbeitung und Abgabe von (höchstens zwei) Zusatzübungsblättern, die nach „Klausurmaßstäben“ korrigiert zurückgegeben werden.
- Alte Klausuren
  - ▶ Alte Klausuren ab Sommersemester 2010 relevant (insbesondere Aufbau)
  - ▶ Aktuelle Klausuren inklusive der meisten Ergebnisse unter „Klausuren“ auf Homepage des Lehrstabs verfügbar

# 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

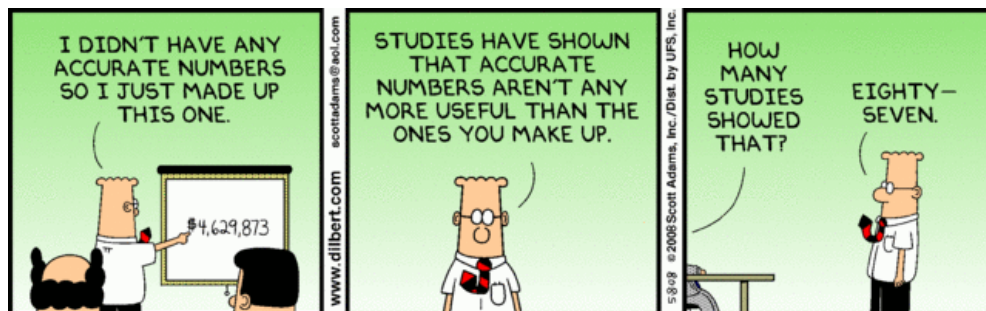
- Einige Zitate oder „Volksweisheiten“:
  - ▶ „Statistik ist pure Mathematik, und in Mathe war ich immer schlecht...“
  - ▶ „Mit Statistik kann man alles beweisen!“
  - ▶ „Ich glaube nur der Statistik, die ich selbst gefälscht habe.“  
(häufig Winston Churchill zugeschrieben, aber eher Churchill von 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, Piper, München, 2009*
- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- Studie/Pressemitteilung des ACE Auto Club Europa *anlässlich des Frauentags am 8. März 2010*: „Autofahrerinnen im Osten am besten“
- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- „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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- 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

<http://www.ace-online.de/nc/der-club/news/autofahrerinnen-im-osten-am-besten.html>

- Welt Online (siehe <http://www.welt.de/vermischtes/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 (Zugriff 04/2011, mittlerweile entfernt)

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

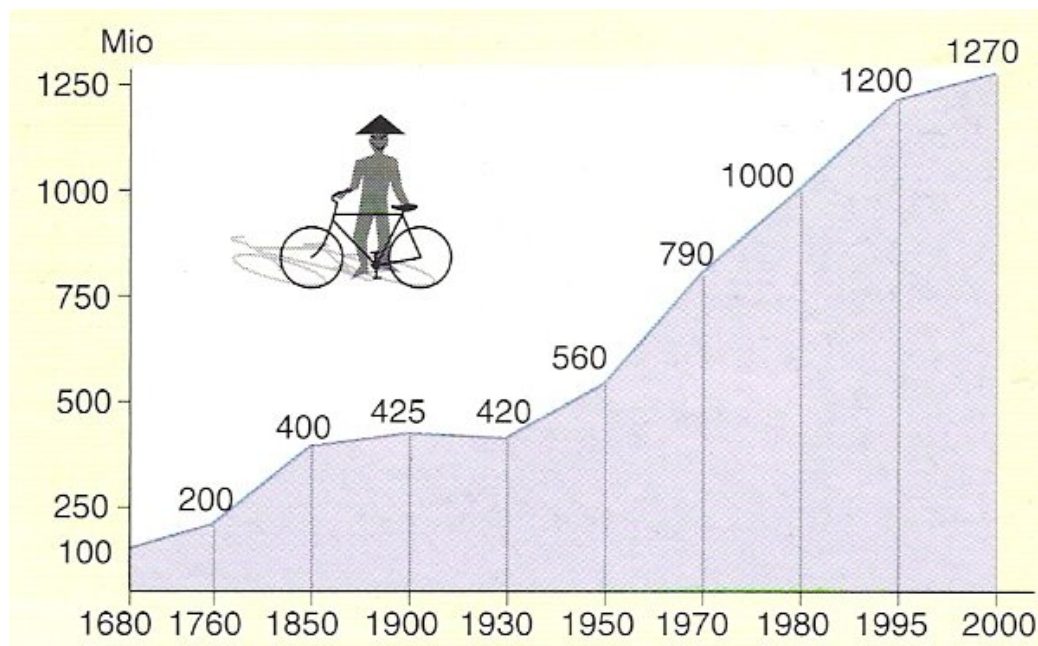
Erst am Ende wird einschränkend erwähnt:

„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

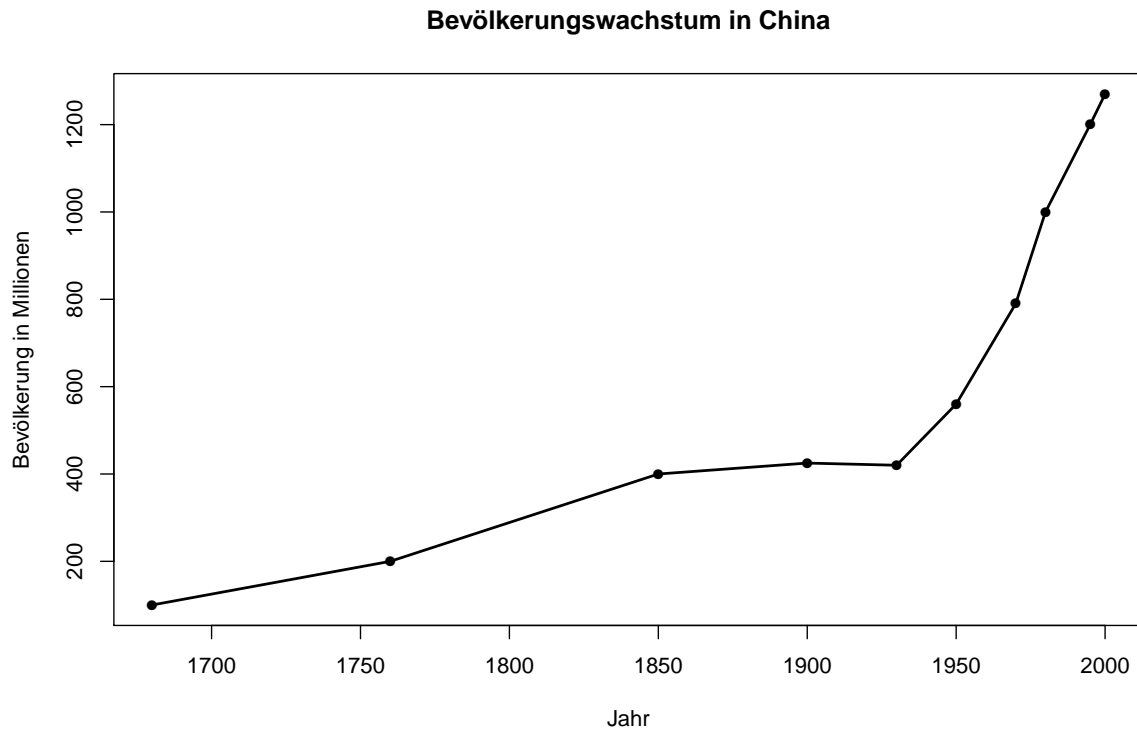
vgl. <http://www.klein-singen.de/statistik/h/Wissenschaft/Bevoelkerungswachstum.html>

### Bevölkerungswachstum in China



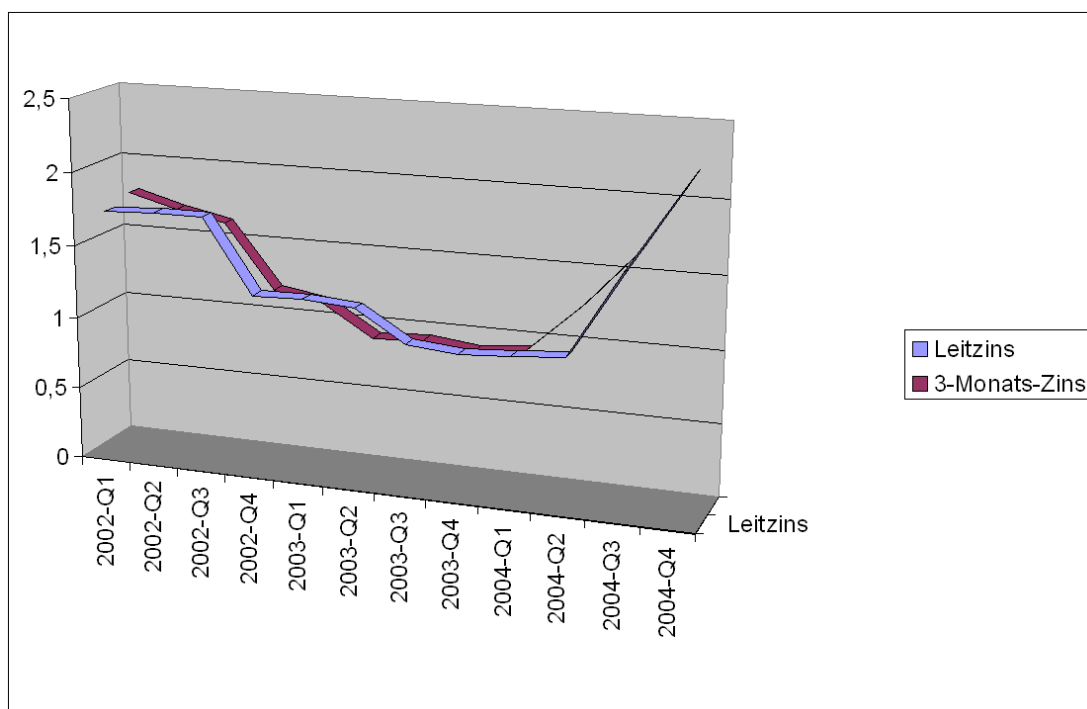
## Beispiel (Irreführende Grafik) II

identischer Datensatz, angemessene Skala



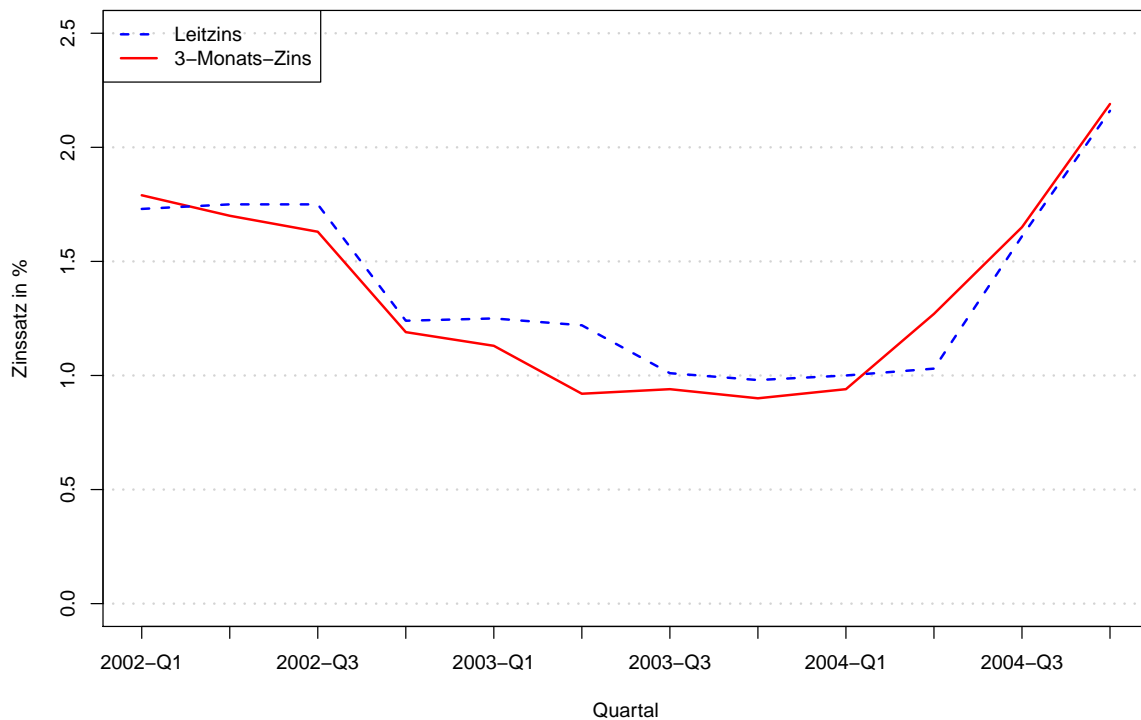
## Beispiel (Chartjunk)

Microsoft Excel mit Standardeinstellung für 3D-Liniendiagramme



# Beispiel (Grafik ohne Chartjunk)

Statistik-Software R, identischer Datensatz



## Kann Statistik auch nützlich sein?

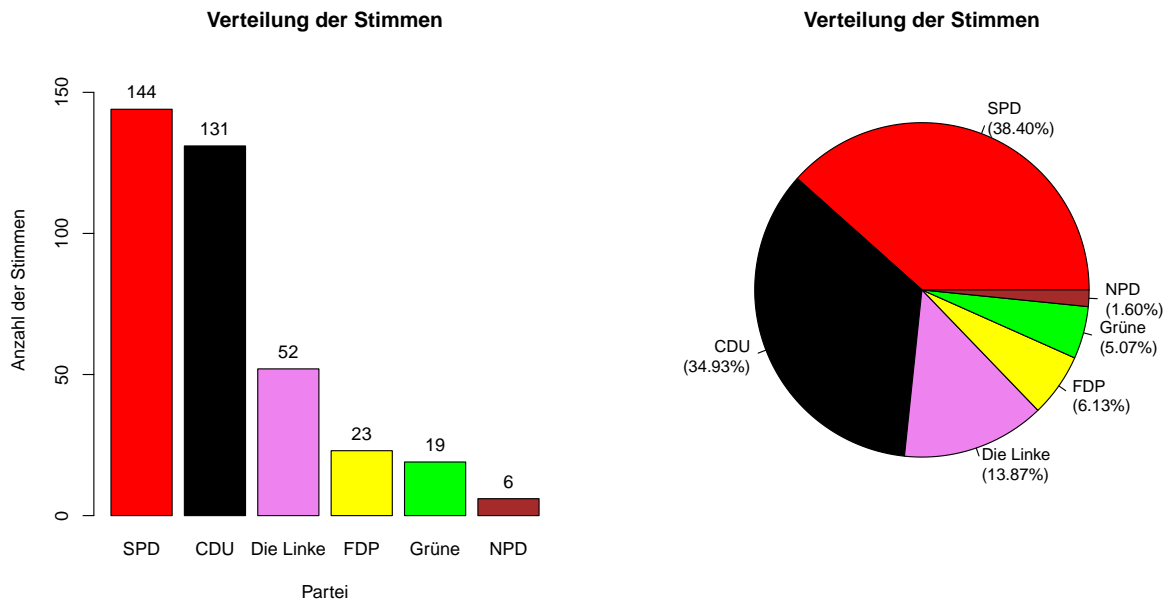
Welche Partei erhält wie viele Stimmen im Wahlbezirk 1.206 der Gemeinde Losheim am See bei den Erststimmen zur Bundestagswahl 2009? Stimmen:

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

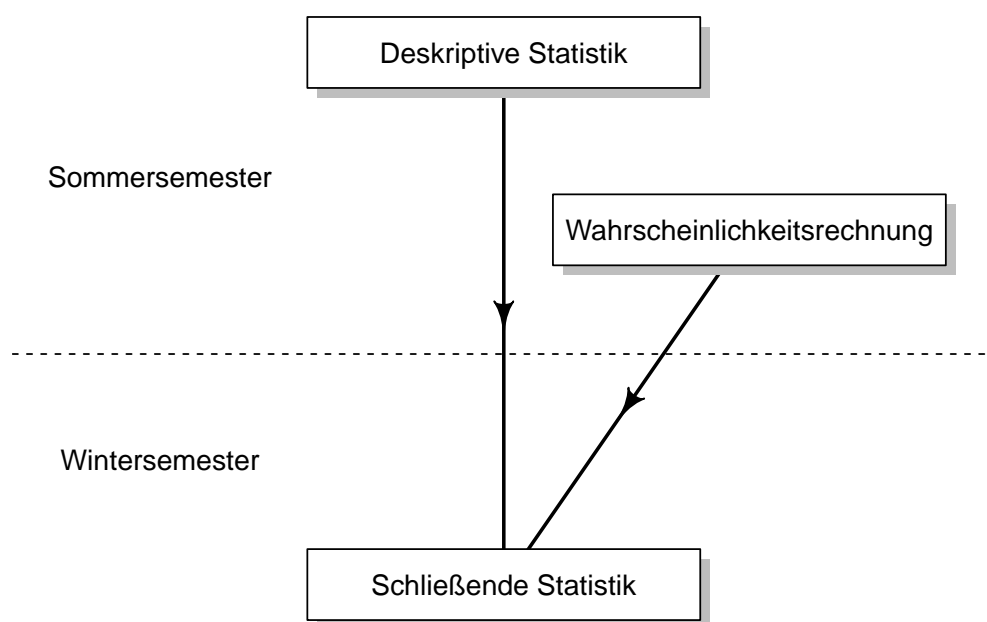
- Mit etwas (deskriptiver) Statistik in tabellarischer Form:

	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:

$x_1 = \text{"CDU"}, x_2 = \text{"SPD"}, x_3 = \text{"SPD"}, x_4 = \text{"Die Linke"}, x_5 = \text{"CDU"},$   
 $x_6 = \text{"Die Linke"}, x_7 = \text{"Die Linke"}, x_8 = \text{"SPD"}, x_9 = \text{"SPD"}, x_{10} = \text{"CDU"},$   
 $x_{11} = \text{"CDU"}, x_{12} = \text{"CDU"}, x_{13} = \text{"SPD"}, x_{14} = \text{"Grüne"}, x_{15} = \text{"FDP"},$   
 $x_{16} = \text{"SPD"}, x_{17} = \text{"SPD"}, x_{18} = \text{"NPD"}, x_{19} = \text{"SPD"}, x_{20} = \text{"FDP"}, \dots$

- ▶ 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“), 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

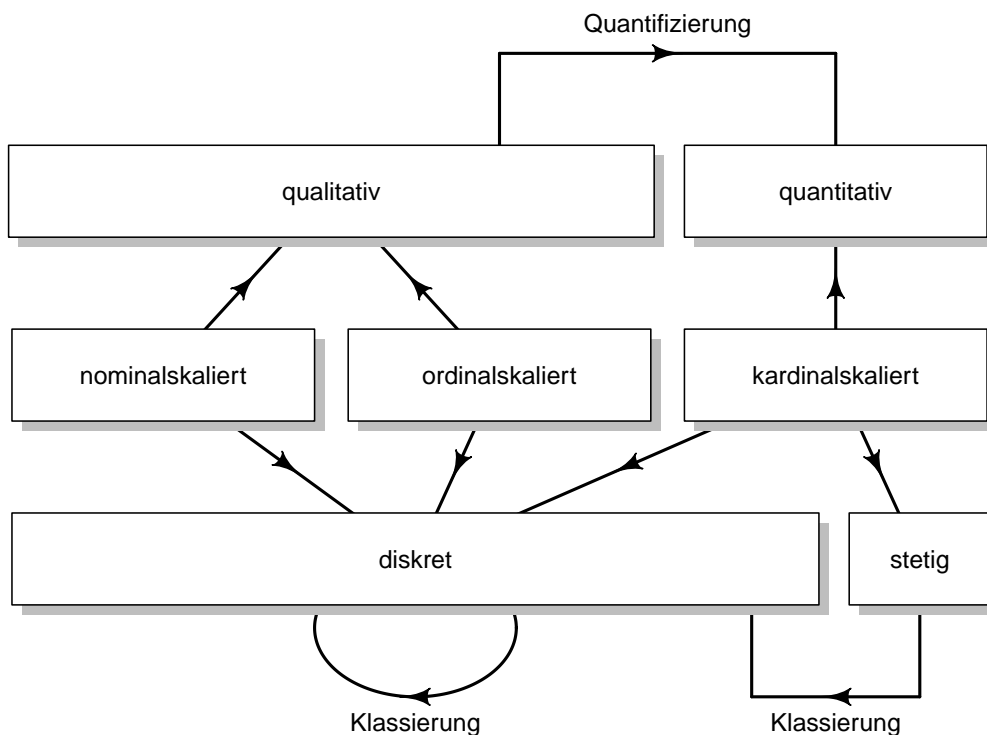
# Umwandlung von Merkmalstypen I

- 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 entspricht, 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 sinnvoll, wenn Abstände tatsächlich (sinnvoll) interpretiert werden können!

## Umwandlung von Merkmalstypen II

- Umwandlung stetiger in diskrete Merkmale durch **Klassierung** oder **Gruppierung**, d.h. Zusammenfassen ganzer Intervalle zu einzelnen Ausprägungen, z.B. Gewichtsklassen beim Boxsport.
  - ▶ Klassierung ermöglicht auch Umwandlung diskreter Merkmale in (erneut) diskrete Merkmale mit unterschiedlichem Merkmalsraum, z.B. Unternehmensgrößen kleiner und mittlerer Unternehmen nach Anzahl der Beschäftigten mit Ausprägungen „1-9“, „10-19“, „20-49“, „50-249“.
  - ▶ Klassierung erfolgt regelmäßig (aber nicht immer) bereits vor der Datenerhebung.

## Übersichtsdarstellung Merkmalstypen



# 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 **absolute 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
$a_j$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$\Sigma$
$h(a_j)$	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

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

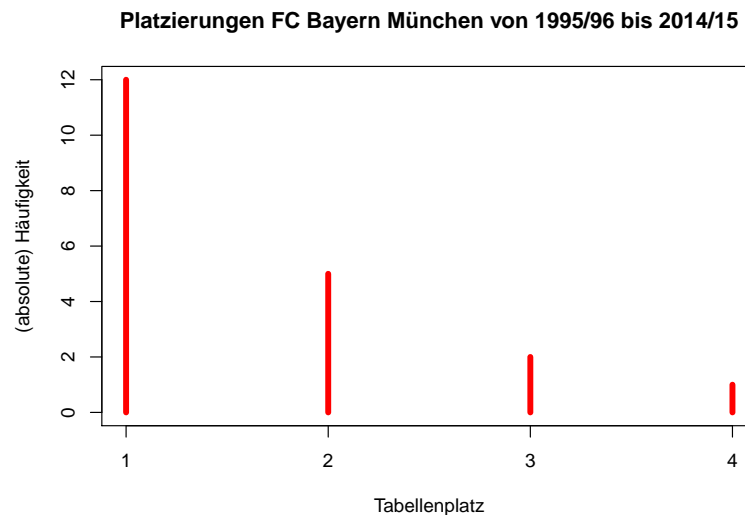
- Alternative grafische Darstellung bei (mindestens) ordinalskalierten Merkmalen mit numerischen Ausprägungen: **Stabdiagramm**

# 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

der finalen Tabellenplätze des FC Bayern München in der (ersten) Fußball-Bundesliga (Saison 1995/96 bis 2014/2015):



## 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

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

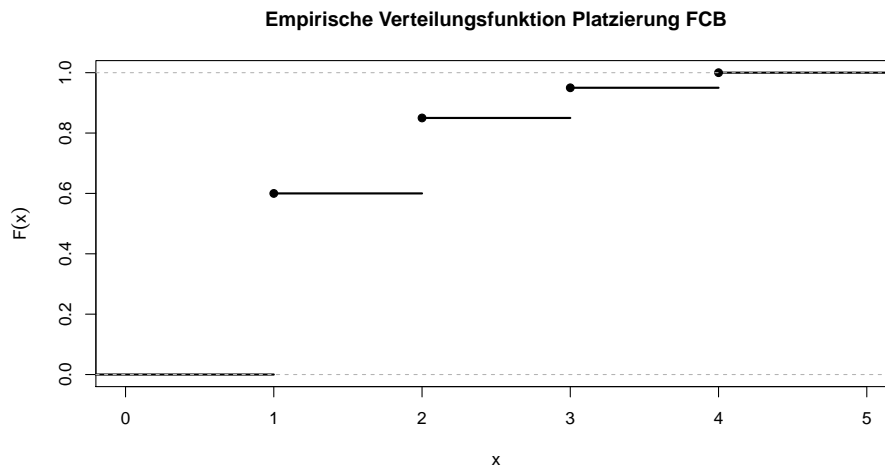
- Ein Vergleich mit den Definitionen von  $h(a)$  und  $r(a)$  offenbart (!), dass  $F(x)$  auch mit Hilfe von  $h(a)$  bzw.  $r(a)$  berechnet werden kann; gibt es  $m$  Merkmalsausprägungen, so gilt:

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

- Beispiel: Empirische Verteilungsfunktion für FC Bayern-Platzierungen

$$F(x) = \begin{cases} 0 & \text{für } x < 1 \\ \frac{12}{20} & \text{für } 1 \leq x < 2 \\ \frac{17}{20} & \text{für } 2 \leq x < 3 \\ \frac{19}{20} & \text{für } 3 \leq x < 4 \\ 1 & \text{für } x \geq 4 \end{cases} = \begin{cases} 0.00 & \text{für } x < 1 \\ 0.60 & \text{für } 1 \leq x < 2 \\ 0.85 & \text{für } 2 \leq x < 3 \\ 0.95 & \text{für } 3 \leq x < 4 \\ 1.00 & \text{für } x \geq 4 \end{cases}$$

- 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

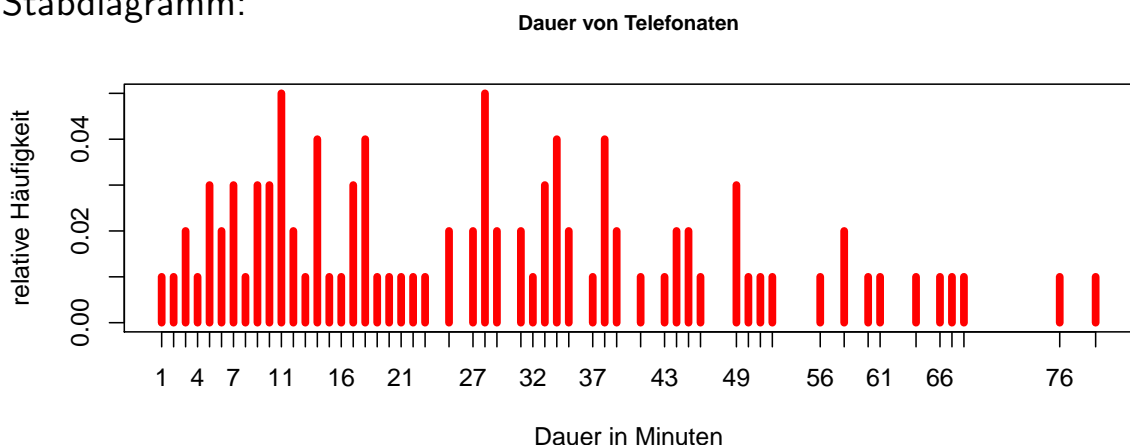
$$r((-\infty, b)) = r((-\infty, b]) - r(b) = F(b) - r(b)$$

- Analog: relative Häufigkeiten weiterer Intervalle:

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

# Häufigkeitsverteilungen klassierter Daten I

- Bisherige Analysemethoden schlecht geeignet für stetige Merkmale bzw. diskrete Merkmale mit „vielen“ Ausprägungen
- (Fiktives) Beispiel: Dauer von 100 Telefonaten (in Minuten)
  - ▶ Urliste: 44, 35, 22, 5, 50, 5, 3, 17, 19, 67, 49, 52, 16, 34, 11, 27, 14, 1, 35, 11, 3, 49, 18, 58, 43, 34, 79, 34, 7, 38, 28, 21, 27, 51, 9, 17, 10, 60, 14, 32, 9, 18, 11, 23, 25, 10, 76, 28, 13, 15, 28, 7, 31, 45, 66, 61, 39, 25, 17, 33, 4, 41, 29, 38, 18, 44, 28, 12, 64, 6, 38, 8, 37, 38, 28, 5, 7, 34, 11, 2, 31, 14, 33, 39, 12, 49, 14, 58, 45, 56, 46, 68, 18, 6, 11, 10, 29, 33, 9, 20
  - ▶ Stabdiagramm:





## Häufigkeitsverteilungen klassierter Daten II

- **Problem:** viele Merkmalswerte treten nur einmalig (oder „selten“) auf  
 $\rightsquigarrow$  Aussagekraft von Häufigkeitstabellen und Stabdiagrammen gering
- **Lösung:** Zusammenfassen mehrerer Merkmalsausprägungen in Klassen
- Zu dieser **Klassierung** erforderlich: **Vorgabe** der Grenzen  $k_0, k_1, \dots, k_l$  von  $l$  (rechtsseitig abgeschlossenen) Intervallen

$$K_1 := (k_0, k_1], K_2 := (k_1, k_2], \dots, K_l := (k_{l-1}, k_l],$$

die alle  $n$  Merkmalswerte überdecken

(also mit  $k_0 < x_i \leq k_l$  für alle  $i \in \{1, \dots, n\}$ )

## Häufigkeitsverteilungen klassierter Daten III

- Wichtige Kennzahlen der Klassierung (bzw. der klassierten Daten):

**Klassenbreiten**  $b_j := k_j - k_{j-1}$

**Klassenmitten**  $m_j := \frac{k_{j-1} + k_j}{2}$

**absolute Häufigkeiten**  $h_j := \#\{i \in \{1, \dots, n\} \mid k_{j-1} < x_i \leq k_j\}$

**relative Häufigkeiten**  $r_j := \frac{h_j}{n}$

**Häufigkeitsdichten**  $f_j := \frac{r_j}{b_j}$

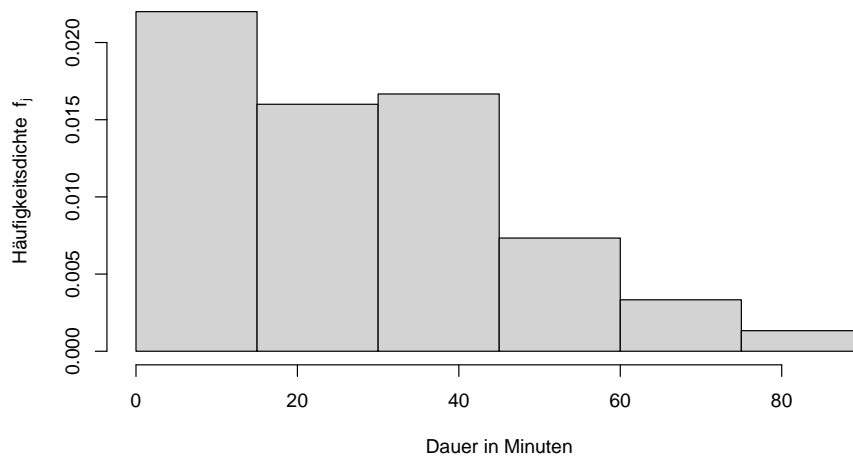
(jeweils für  $j \in \{1, \dots, l\}$ ).

- Übliche grafische Darstellung von klassierten Daten: **Histogramm**
- Hierzu: Zeichnen der Rechtecke mit Höhen  $f_j$  über den Intervallen  $K_j$  (also der Rechtecke mit den Eckpunkten  $(k_{j-1}, 0)$  und  $(k_j, f_j)$ )

- Am Beispiel der Gesprächsdauern bei 6 Klassen zu je 15 Minuten Breite:

Nr. $j$	Klasse $K_j = (k_{j-1}, k_j]$	Klassenbreite $b_j$	Klassenmitte $m_j$	absolute Häufigkeit $h_j$	relative Häufigkeit $r_j = \frac{h_j}{n}$	Häufigkeitsdichte $f_j = \frac{r_j}{b_j}$	Verteilungsfunktion $F(k_j)$
1	(0, 15]	15	7.5	33	0.33	0.022	0.33
2	(15, 30]	15	22.5	24	0.24	0.016	0.57
3	(30, 45]	15	37.5	25	0.25	0.01 $\bar{6}$	0.82
4	(45, 60]	15	52.5	11	0.11	0.007 $\bar{3}$	0.93
5	(60, 75]	15	67.5	5	0.05	0.00 $\bar{3}$	0.98
6	(75, 90]	15	82.5	2	0.02	0.001 $\bar{3}$	1.00

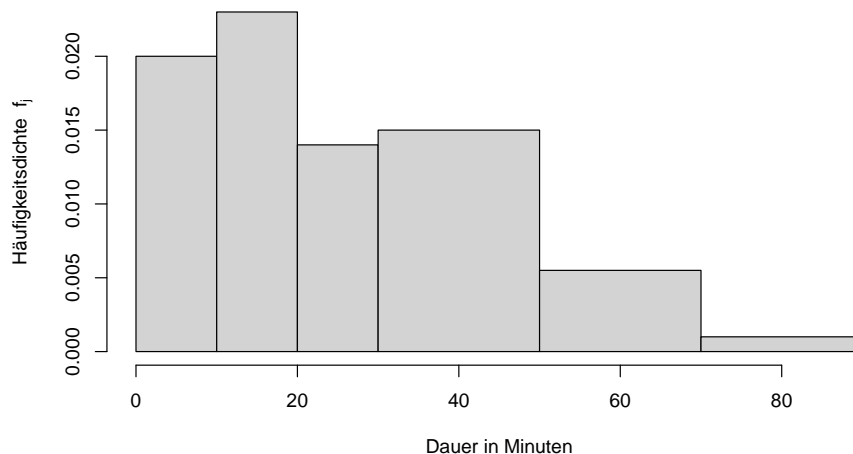
Histogramm der Gesprächsdauern



- Alternativ mit 6 Klassen bei 2 verschiedenen Breiten:

Nr. $j$	Klasse $K_j = (k_{j-1}, k_j]$	Klassenbreite $b_j$	Klassenmitte $m_j$	absolute Häufigkeit $h_j$	relative Häufigkeit $r_j = \frac{h_j}{n}$	Häufigkeitsdichte $f_j = \frac{r_j}{b_j}$	Verteilungsfunktion $F(k_j)$
1	(0, 10]	10	5	20	0.20	0.0200	0.20
2	(10, 20]	10	15	23	0.23	0.0230	0.43
3	(20, 30]	10	25	14	0.14	0.0140	0.57
4	(30, 50]	20	40	30	0.30	0.0150	0.87
5	(50, 70]	20	60	11	0.11	0.0055	0.98
6	(70, 90]	20	80	2	0.02	0.0010	1.00

Histogramm der Gesprächsdauern



## Bemerkungen I

- Der **Flächeninhalt** der einzelnen Rechtecke eines Histogramms entspricht der relativen Häufigkeit der zugehörigen Klasse
  - ↪ Die Summe aller Flächeninhalte beträgt 1
  - ↪ Die Höhe der Rechtecke ist nur dann proportional zu der relativen Häufigkeit der Klassen, falls alle Klassen die gleiche Breite besitzen!
- Die Klassierung ist abhängig von der Wahl der Klassengrenzen, unterschiedliche Klassengrenzen können einen Datensatz auch sehr unterschiedlich erscheinen lassen ↪ Potenzial zur Manipulation
- Es existieren verschiedene Algorithmen zur automatischen Wahl von Klassenanzahl und -grenzen (z.B. nach Scott, Sturges, Freedman-Diaconis)

## Bemerkungen II

- Durch Klassierung geht Information verloren!
  - ▶ Spezielle Verfahren für klassierte Daten vorhanden
  - ▶ Verfahren approximieren ursprüngliche Daten in der Regel durch die Annahme gleichmäßiger Verteilung innerhalb der einzelnen Klassen
  - ▶ (Approximative) Verteilungsfunktion (ebenfalls mit  $F(x)$  bezeichnet) zu klassierten Daten entsteht so durch lineare Interpolation der an den Klassengrenzen  $k_j$  bekannten (und auch nach erfolgter Klassierung noch exakten!) Werte der empirischen Verteilungsfunktion  $F(k_j)$
  - ▶ Näherungsweise Berechnung von Intervallhäufigkeiten dann gemäß Folie 46 f. mit der approximativen empirischen Verteilungsfunktion  $F(x)$ .

# (Approx.) Verteilungsfunktion bei klassierten Daten

## Approximative Verteilungsfunktion bei klassierten Daten

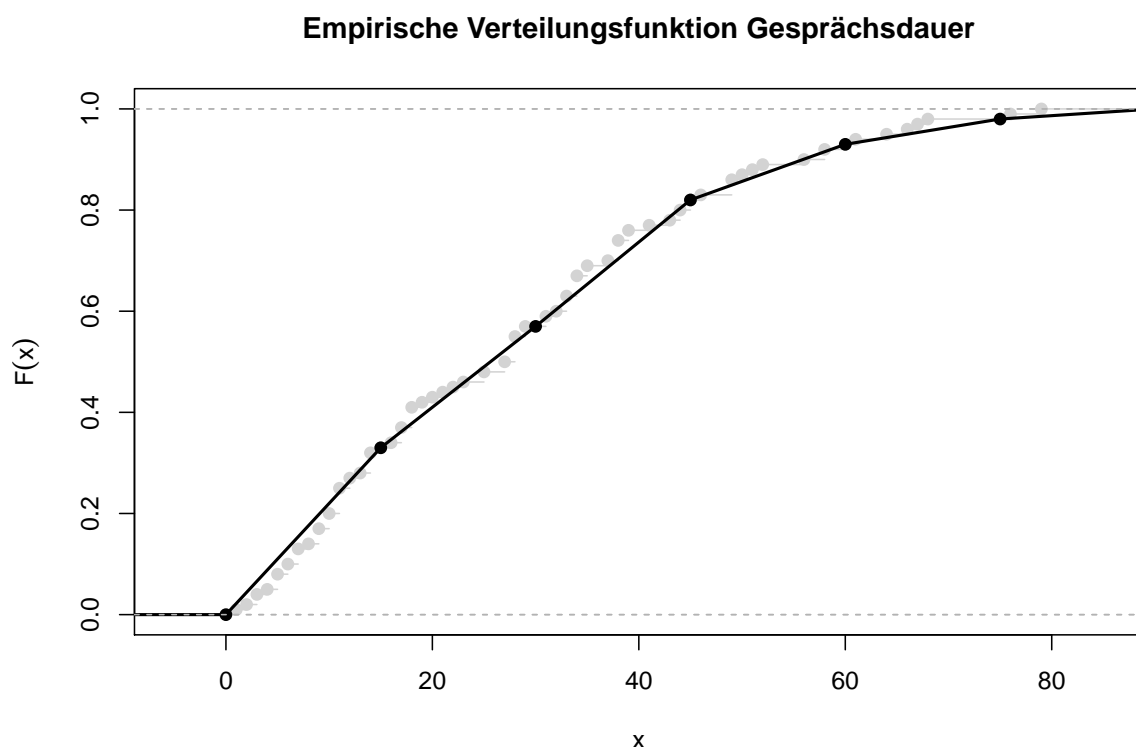
$$F(x) = \begin{cases} 0 & \text{für } x \leq k_0 \\ F(k_{j-1}) + f_j \cdot (x - k_{j-1}) & \text{für } k_{j-1} < x \leq k_j, j \in \{1, \dots, l\} \\ 1 & \text{für } x > k_l \end{cases}$$

- Am Beispiel der Gesprächsdauern (Klassierung aus Folie 52)

$$F(x) = \begin{cases} 0 & \text{für } x \leq 0 \\ 0.0200 \cdot (x - 0) & \text{für } 0 < x \leq 10 \\ 0.20 + 0.0230 \cdot (x - 10) & \text{für } 10 < x \leq 20 \\ 0.43 + 0.0140 \cdot (x - 20) & \text{für } 20 < x \leq 30 \\ 0.57 + 0.0150 \cdot (x - 30) & \text{für } 30 < x \leq 50 \\ 0.87 + 0.0055 \cdot (x - 50) & \text{für } 50 < x \leq 70 \\ 0.98 + 0.0010 \cdot (x - 70) & \text{für } 70 < x \leq 90 \\ 1 & \text{für } x > 90 \end{cases}$$

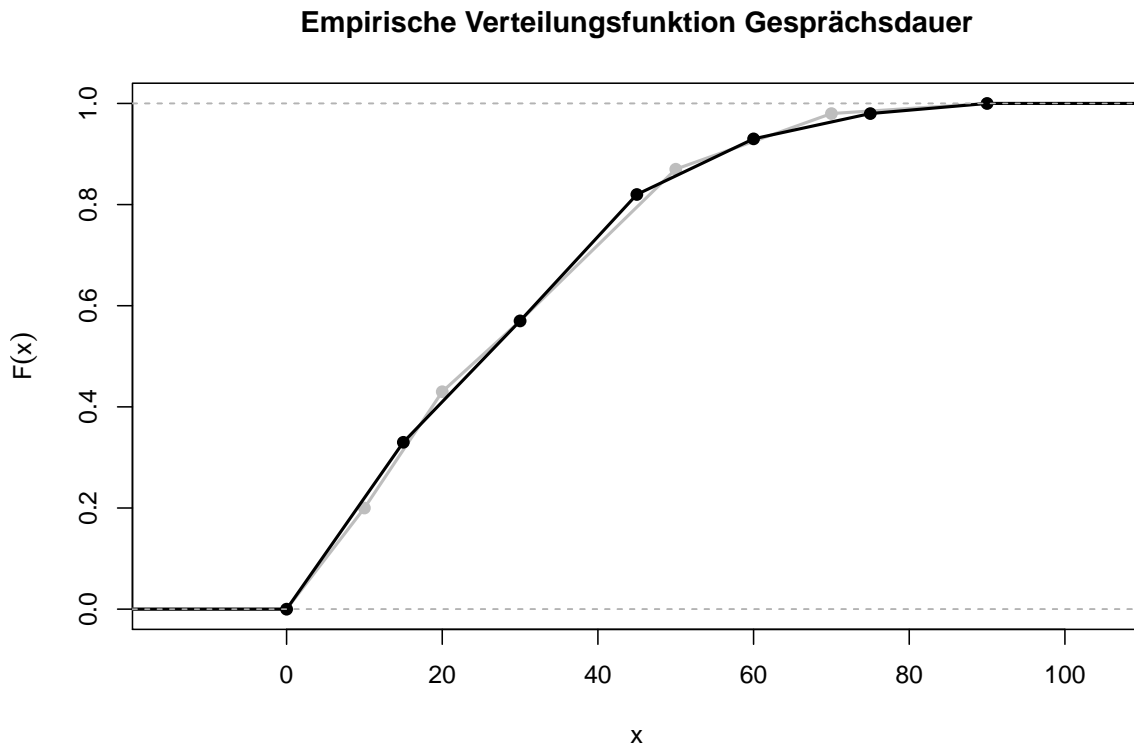
# Grafik: Verteilungsfunktion bei klassierten Daten

(Empirische Verteilungsfunktion der unklassierten Daten in hellgrau)



# Grafik: Verteilungsfunktion bei verschiedenen Klassierungen

(Klassierung aus Folie 51 in schwarz, Klassierung aus Folie 52 in grau)



## Lagemaße

- Aggregation von Merkmalswerten zu Häufigkeitsverteilungen (auch nach erfolgter Klassierung) nicht immer ausreichend.
- Häufig gewünscht: einzelner Wert, der die Verteilung der Merkmalswerte geeignet charakterisiert  $\rightsquigarrow$  „Mittelwert“
- **Aber:**
  - ▶ Gibt es immer einen „Mittelwert“?  
Was ist der Mittelwert der Merkmalswerte *rot, gelb, gelb, blau*?  
 $\rightsquigarrow$  allgemeinerer Begriff: „Lagemaß“
  - ▶ Gibt es verschiedene „Mittelwerte“?  
Falls ja, welcher der Mittelwerte ist (am Besten) geeignet?

## Lagemaße für nominalskalierte Merkmale

- Verschiedene Merkmalsausprägungen können lediglich unterschieden werden
- „Typische“ Merkmalswerte sind also solche, die häufig vorkommen
- Geeignetes Lagemaß: häufigster Wert (es kann mehrere geben!)

### Definition 3.1 (Modus, Modalwert)

Sei  $X$  ein (mindestens) nominalskaliertes Merkmal mit Merkmalsraum  $A = \{a_1, \dots, a_m\}$  und relativer Häufigkeitsverteilung  $r$ .  
Dann heißt jedes Element  $a_{\text{mod}} \in A$  mit

$$r(a_{\text{mod}}) \geq r(a_j) \text{ für alle } j \in \{1, \dots, m\}$$

**Modus** oder **Modalwert** von  $X$ .

- Beispiele:
  - ▶ Modus der Urliste *rot, gelb, gelb, blau*:  
 $a_{\text{mod}} = \text{gelb}$
  - ▶ Modalwerte der Urliste *1, 5, 3, 3, 4, 2, 6, 7, 6, 8*:  
 $a_{\text{mod},1} = 3$  und  $a_{\text{mod},2} = 6$

## Lagemaße für ordinalskalierte Merkmale I

- Durch die vorgegebene Anordnung auf der Menge der möglichen Ausprägungen  $M$  lässt sich der Begriff „mittlerer Wert“ mit Inhalt füllen.
- In der geordneten Folge von Merkmalswerten

$$X_{(1)}, X_{(2)}, \dots, X_{(n-1)}, X_{(n)}$$

bietet sich als Lagemaß also ein Wert „in der Mitte“ der Folge an.

- Ist  $n$  gerade, gibt es keine eindeutige Mitte der Folge, und eine zusätzliche Regelung ist erforderlich.

## Lagemaße für ordinalskalierte Merkmale II

### Definition 3.2 (Median)

Sei  $X$  ein (mindestens) ordinalskaliertes Merkmal auf der Menge der vorstellbaren Merkmalsausprägungen  $M$  und  $x_{(1)}, x_{(2)}, \dots, x_{(n-1)}, x_{(n)}$  die gemäß der vorgegebenen Ordnung sortierte Urliste zum Merkmal  $X$ .

- Ist  $n$  ungerade, so heißt  $x_{(\frac{n+1}{2})}$  der **Median** von  $X$ , in Zeichen  $x_{\text{med}} = x_{(\frac{n+1}{2})}$ .
- Ist  $n$  gerade, so heißen alle (möglicherweise viele verschiedene) Elemente von  $M$  *zwischen* (bezogen auf die auf  $M$  gegebene Ordnung)  $x_{(\frac{n}{2})}$  und  $x_{(\frac{n}{2}+1)}$  (einschließlich dieser beiden Merkmalswerte) **Mediane** von  $X$ .
- Bei stetigen Merkmalen kann für die Definition des Medians auch für gerades  $n$  Eindeutigkeit erreicht werden, indem spezieller der Mittelwert

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

der beiden „mittleren“ Merkmalswerte als Median festgelegt wird.

## Lagemaße für ordinalskalierte Merkmale III

- Beispiele:
  - ▶ Ist  $M = \{\text{sehr gut, gut, befriedigend, ausreichend, mangelhaft, ungenügend}\}$  als Menge der möglichen Ausprägungen eines ordinalskalierten Merkmals  $X$  mit der üblichen Ordnung von Schulnoten von „sehr gut“ bis „ungenügend“ versehen, so ist die sortierte Folge von Merkmalswerten zur Urliste  
gut, ausreichend, sehr gut, mangelhaft, mangelhaft, gut  
durch  
sehr gut, gut, gut, ausreichend, mangelhaft, mangelhaft  
gegeben und sowohl „gut“ als auch „befriedigend“ und „ausreichend“ sind Mediane von  $X$ .
  - ▶ Der oben beschriebenen Konvention für stetige Merkmale folgend ist der Median des stetigen Merkmals zur Urliste  
1.85, 6.05, 7.97, 11.16, 17.19, 18.87, 19.82, 26.95, 27.25, 28.34  
von 10 Merkmalsträgern durch  $x_{\text{med}} = \frac{1}{2} \cdot (17.19 + 18.87) = 18.03$  gegeben.

## Lagemaße für kardinalskalierte Merkmale

- Bei kardinalskalierten Merkmalen ist oft eine „klassische“ Mittelung der Merkmalswerte als Lagemaß sinnvoll, man erhält so aus der Urliste  $x_1, \dots, x_n$  das „**arithmetische Mittel**“  $\bar{x} := \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$ .

- *Beispiel:*

Die Haushalts-Nettoeinkommen (in €) von 6 Haushalten eines Mehrparteien-Wohnhauses sind:

Haushalt	1	2	3	4	5	6
Nettoeinkommen	1000	400	1500	2900	1800	2600

Frage: Wie groß ist das durchschnittliche Nettoeinkommen?

Antwort:  $\frac{1}{6} \cdot (1000 + 400 + 1500 + 2900 + 1800 + 2600) = 1700$

- Bei klassierten Daten wird der Mittelwert als gewichtetes arithmetisches Mittel der  $l$  Klassenmitten näherungsweise berechnet:

$$\bar{x} := \frac{1}{n} \sum_{j=1}^l h_j \cdot m_j = \sum_{j=1}^l r_j \cdot m_j .$$

- Arithmetisches Mittel für viele (nicht alle!) Anwendungen adäquates „Mittel“

- *Beispiel:*

Ein Wachstumssparvertrag legt folgende Zinssätze fest:

Jahr	1	2	3	4	5
Zinssatz	1.5%	1.75%	2.0%	2.5%	3.5%

Wie groß ist der Zinssatz *im Durchschnitt*?

- ▶ Aus Zinsrechnung bekannt: Kapital  $K$  inkl. Zinsen nach 5 Jahren bei Startkapital  $S$  beträgt

$$K = S \cdot (1 + 0.015) \cdot (1 + 0.0175) \cdot (1 + 0.02) \cdot (1 + 0.025) \cdot (1 + 0.035)$$

- ▶ Gesucht ist (für 5 Jahre gleichbleibender) Zinssatz  $R$ , der gleiches Endkapital  $K$  produziert, also  $R$  mit der Eigenschaft

$$K \stackrel{!}{=} S \cdot (1 + R) \cdot (1 + R) \cdot (1 + R) \cdot (1 + R) \cdot (1 + R)$$

- ▶ Ergebnis:

$$R = \sqrt[5]{(1 + 0.015) \cdot (1 + 0.0175) \cdot (1 + 0.02) \cdot (1 + 0.025) \cdot (1 + 0.035)} - 1$$

$\rightsquigarrow R = 2.2476\%$ .

- Der in diesem Beispiel für die Bruttorenditen  $(1 + \text{Zinssatz})$  sinnvolle Mittelwert heißt „**geometrisches Mittel**“.



- *Beispiel:*

Auf einer Autofahrt von insgesamt 30 [km] werden  $s_1 = 10$  [km] mit einer Geschwindigkeit von  $v_1 = 30$  [km/h],  $s_2 = 10$  [km] mit einer Geschwindigkeit von  $v_2 = 60$  [km/h] und  $s_3 = 10$  [km] mit einer Geschwindigkeit von  $v_3 = 120$  [km/h] zurückgelegt.

Wie hoch ist die durchschnittliche Geschwindigkeit?

- ▶ Durchschnittliche Geschwindigkeit: Quotient aus Gesamtstrecke und Gesamtzeit
- ▶ Gesamtstrecke:  $s_1 + s_2 + s_3 = 10$  [km] +  $10$  [km] +  $10$  [km] =  $30$  [km]
- ▶ Zeit für Streckenabschnitt: Quotient aus Streckenlänge und Geschwindigkeit
- ▶ Einzelzeiten also:

$$\frac{s_1}{v_1} = \frac{10 \text{ [km]}}{30 \text{ [km/h]}}, \quad \frac{s_2}{v_2} = \frac{10 \text{ [km]}}{60 \text{ [km/h]}} \quad \text{und} \quad \frac{s_3}{v_3} = \frac{10 \text{ [km]}}{120 \text{ [km/h]}}$$

↪ Durchschnittsgeschwindigkeit

$$\frac{s_1 + s_2 + s_3}{\frac{s_1}{v_1} + \frac{s_2}{v_2} + \frac{s_3}{v_3}} = \frac{30 \text{ [km]}}{\frac{10}{30} \text{ [h]} + \frac{10}{60} \text{ [h]} + \frac{10}{120} \text{ [h]}} = \frac{30}{\frac{7}{12}} \text{ [km/h]} = 51.429 \text{ [km/h]}$$

- Der in diesem Beispiel für die Geschwindigkeiten sinnvolle Mittelwert heißt „**harmonisches Mittel**“.

## Zusammenfassung: Mittelwerte I

### Definition 3.3 (Mittelwerte)

Seien  $x_1, x_2, \dots, x_n$  die Merkmalswerte zu einem kardinalskalierten Merkmal  $X$ . Dann heißt

- 1  $\bar{x} := \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$  das arithmetische Mittel,
- 2  $\bar{x}^{(g)} := \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} = \sqrt[n]{\prod_{i=1}^n x_i} = \left(\prod_{i=1}^n x_i\right)^{\frac{1}{n}}$  das geometrische Mittel,
- 3  $\bar{x}^{(h)} := \frac{1}{\frac{1}{n}\left(\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}\right)} = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}}$  das harmonische Mittel

von  $x_1, \dots, x_n$ .

## Zusammenfassung: Mittelwerte II

### Bemerkung 3.4

Liegt die absolute (bzw. relative) Häufigkeitsverteilung  $h$  (bzw.  $r$ ) eines kardinalskalierten Merkmals  $X$  mit Merkmalsraum  $A = \{a_1, \dots, a_m\}$  vor, so gilt

$$1 \quad \bar{x} = \frac{1}{n} \sum_{j=1}^m h(a_j) \cdot a_j = \sum_{j=1}^m r(a_j) \cdot a_j$$

$$2 \quad \bar{x}^{(g)} = \sqrt[n]{\prod_{j=1}^m a_j^{h(a_j)}} = \prod_{j=1}^m a_j^{r(a_j)}$$

$$3 \quad \bar{x}^{(h)} = \frac{1}{\frac{1}{n} \sum_{j=1}^m \frac{h(a_j)}{a_j}} = \frac{n}{\sum_{j=1}^m \frac{h(a_j)}{a_j}} = \frac{1}{\sum_{j=1}^m \frac{r(a_j)}{a_j}}$$

- Die in Bemerkung 3.4 berechneten Mittelwerte können als sogenannte *gewichtete Mittelwerte* der aufgetreten Merkmalswerte  $a_1, \dots, a_m$  aufgefasst werden, wobei die Gewichte durch die absoluten Häufigkeiten  $h(a_1), \dots, h(a_m)$  (bzw. durch die relativen Häufigkeiten  $r(a_1), \dots, r(a_m)$ ) der aufgetretenen Merkmalswerte gegeben sind.

## Weitere Beispiele I

- Pauschale Aussagen, wann welcher Mittelwert geeignet ist, nicht möglich!
- Beispiel Zinssätze:*  
Aufgrund Begrenzungen der europäischen Einlagensicherung möchte ein Anleger Kapital von 100 000 € gleichmäßig auf 5 Banken verteilen, die für die vorgegebene Anlagedauer folgende Zinsen anbieten:

Bank	1	2	3	4	5
Zinssatz	2.5%	2.25%	2.4%	2.6%	2.55%

Frage: Wie groß ist der durchschnittliche Zinssatz?

Antwort:  $\frac{1}{5} \cdot (2.5\% + 2.25\% + 2.4\% + 2.6\% + 2.55\%) = 2.46\%$

## Weitere Beispiele II

- *Beispiel Geschwindigkeiten:*

Auf einer Autofahrt von insgesamt 30 [Min.] Fahrzeit werden  $t_1 = 10$  [Min.] mit einer Geschwindigkeit von  $v_1 = 30$  [km/h],  $t_2 = 10$  [Min.] mit  $v_2 = 60$  [km/h] und  $t_3 = 10$  [Min.] mit  $v_3 = 120$  [km/h] zurückgelegt. Wie hoch ist die durchschnittliche Geschwindigkeit?

- ▶ Durchschnittliche Geschwindigkeit: Quotient aus Gesamtstrecke und -zeit
- ▶ Gesamtzeit:  $t = t_1 + t_2 + t_3 = 10$  [Min.] +  $10$  [Min.] +  $10$  [Min.] =  $30$  [Min.]
- ▶ Länge der Streckenabschnitte: Produkt aus Geschwindigkeit und Fahrzeit
- ↪ Durchschnittsgeschwindigkeit

$$\frac{v_1 \cdot t_1 + v_2 \cdot t_2 + v_3 \cdot t_3}{t} = \frac{1}{3} \cdot 30 \text{ [km/h]} + \frac{1}{3} \cdot 60 \text{ [km/h]} + \frac{1}{3} \cdot 120 \text{ [km/h]} = 70 \text{ [km/h]}$$

## Bemerkungen I

- Insbesondere bei diskreten Merkmalen wie z.B. einer Anzahl muss der erhaltene (arithmetische, geometrische, harmonische) Mittelwert weder zum Merkmalsraum  $A$  noch zur Menge der vorstellbaren Merkmalsausprägungen  $M$  gehören (z.B. „im Durchschnitt 2.2 Kinder pro Haushalt“).
- Auch der/die Median(e) gehören (insbesondere bei numerischen Merkmalen) häufiger nicht zur Menge  $A$  der Merkmalsausprägungen; lediglich der/die Modalwert(e) kommen stets auch in der Liste der Merkmalswerte vor!
- **Vorsicht** vor falschen Rückschlüssen vom Mittelwert auf die Häufigkeitsverteilung!

## Bemerkungen II

### Mobilfunknutzung Europa in 2006

In einem aktuell nicht mehr verfügbaren Online-Artikel auf der Homepage <http://www.digital-world.de> wurde aus der Tatsache, dass die Anzahl der Mobiltelefone in Europa größer ist als die Anzahl der Europäer, also das arithmetische Mittel des Merkmals *Anzahl Mobiltelefone pro Person* in Europa größer als 1 ist, die folgende Aussage in der Überschrift abgeleitet:

**Statistik: Jeder Europäer telefoniert mobil**

Zusammenfassend heißt es außerdem:

**Laut einer aktuellen Studie telefoniert jeder Europäer mittlerweile mit mindestens einem Mobiltelefon.**

Wie sind diese Aussagen zu beurteilen? Welcher Fehlschluss ist gezogen worden?

## Optimalitätseigenschaften

einiger Lagemaße bei kardinalskalierten Daten

- Für kardinalskalierte Merkmale besitzen Mediane und arithmetische Mittelwerte spezielle (Optimalitäts-)Eigenschaften.
- Für jeden Median  $x_{\text{med}}$  eines Merkmals  $X$  mit den  $n$  Merkmalswerten  $x_1, \dots, x_n$  gilt:

$$\sum_{i=1}^n |x_i - x_{\text{med}}| \leq \sum_{i=1}^n |x_i - t| \quad \text{für alle } t \in \mathbb{R}$$

- Für das arithmetische Mittel  $\bar{x}$  eines Merkmals  $X$  mit den  $n$  Merkmalswerten  $x_1, \dots, x_n$  gilt:

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

$$\textcircled{2} \quad \sum_{i=1}^n (x_i - \bar{x})^2 \leq \sum_{i=1}^n (x_i - t)^2 \quad \text{für alle } t \in \mathbb{R}$$