

Statistische Verfahren in der Geographie

Skript für den Theorieteil

Till Straube
straube@geo.uni-frankfurt.de
Institut für Humangeographie
Goethe-Universität Frankfurt

Sommersemester 2022

Inhaltsverzeichnis

| | |
|--|-----------|
| Terminüberblick | 2 |
| Vorbesprechung | 3 |
| 1 Datenerhebung und Häufigkeiten | 5 |
| 1.1 Statistische Praxis | 5 |
| 1.2 Grundlagen der Datenerhebung | 8 |
| 1.3 Häufigkeitsverteilungen | 10 |
| Tipps zur Vertiefung | 14 |
| Übungsaufgaben | 16 |
| 2 Maßzahlen | 19 |
| 2.1 Einleitende Bemerkungen | 19 |
| 2.2 Lagemaße | 21 |
| 2.3 Streumaße | 23 |
| 2.4 Boxplot | 27 |
| Tipps zur Vertiefung | 27 |
| Lösungen der Übungsaufgaben | 29 |
| Sitzung 1 | 29 |
| Quellenverzeichnis | 32 |

Terminüberblick

Alle Sitzungen finden von 14 bis 16h c. t. statt, die Klausuren s. t.

| Datum | Sitzung | Inhalt | Ort |
|------------------|---------|--------------------------------------|------|
| 12. April 2022 | | Vorbesprechung | Zoom |
| 19. April 2022 | 1 | Datenerhebung und Häufigkeiten | Zoom |
| 26. April 2022 | 2 | Maßzahlen | HZ10 |
| 3. Mai 2022 | 3 | [z-Werte und Normalverteilung] | HZ10 |
| 10. Mai 2022 | 4 | [Schätzstatistik] | HZ10 |
| 17. Mai 2022 | 5 | [Grundlagen der Teststatistik] | HZ10 |
| 24. Mai 2022 | 6 | [Testverfahren mit zwei Stichproben] | HZ10 |
| 31. Mai 2022 | 7 | [Korrelation] | HZ10 |
| 7. Juni 2022 | | entfällt | |
| 14. Juni 2022 | 8 | [Lineare Regression] | HZ10 |
| 21. Juni 2022 | 9 | [Kreuztabellen] | HZ10 |
| 28. Juni 2022 | 10 | [χ^2 -Tests] | HZ10 |
| 5. Juli 2022 | | Klausurvorbereitung | HZ10 |
| 12. Juli 2022 | | Klausur (14h s. t.) | |
| 11. Oktober 2022 | | Nachklausur (14h s. t.) | |

Vorbesprechung

Lernziele der Veranstaltung

Sie können...

- Grundbegriffe der Statistik sinnvoll verwenden.
- die wichtigsten statistischen Kennzahlen berechnen.
- gängige Diagramme interpretieren.
- einfache statistische Schätz- und Prüfverfahren anwenden.
- passende Verfahren für verschiedene Aufgaben wählen.

Konzept der Veranstaltung

- Die gesamte Veranstaltung dient als Klausurvorbereitung
- Die selbständige Anwendung der Verfahren steht im Vordergrund

Sitzungsvorbereitung

- Materialien werden zur eigenständigen Vorbereitung bereit gestellt
- Dieses Online-Skript mit den Kerninhalten
- Darin: Videos (aus 2020) mit Beispielen und Übungen
- Darin: Verweise auf weiterführende Literatur, YouTube-Videos, etc.
- Fehler und Unklarheiten bitte per E-Mail melden!

Sitzungsablauf

- Dienstags, 14 h c. t..
 - Sitzung 1 auf Zoom (Link in OLAT)
 - Folgende Sitzungen in HZ10
- Übungsaufgaben (und Lösungen) werden online bereit gestellt
- Teilnehmer*innen bearbeiten die Aufgaben in Break-Out-Sessions bzw. Kleingruppen
- Bei Problemen fragen Sie sich erst mal gegenseitig
- Sonst bin ich ansprechbar (Melden oder Zoom-Funktion: Um Hilfe bitten)

Empfehlungen

- Lassen Sie sich auf den wöchentlichen Rhythmus ein
- Bereiten Sie die Sitzungen vor und nach
- Bilden Sie Lerngruppen
- Gleichen Sie in Lerngruppen Ihre Ziele ab

- Machen Sie sich mit Ihrem Taschenrechner vertraut

Literaturempfehlungen

- Ganz besonders:
 - Bortz und Schuster (2010) (als E-Book bei der UB erhältlich; dieselben Notationskonventionen wie in der Veranstaltung)
- Ergänzend:
 - Lange und Nipper (2018) (geographiebezogen)
 - Bahrenberg, Giese und Nipper (2010) (geographiebezogen)
 - Benninghaus (2007) (als E-Book bei der UB erhältlich)
- Bedingt:
 - Zimmermann-Janschitz (2014) (geographiebezogen; als E-Book bei der UB erhältlich)
- Englisch:
 - Burt und Barber (1996)

Taschenrechner

- Zulassungsregeln für Klausur wie für Mathe-Abi (Hessen)
- Also kein „programmierbarer“ Taschenrechner
- Erlaubt ist z. B. CASIO FX-991DE Plus
- „Wissenschaftlicher“ Taschenrechner kann von großem Vorteil sein... aber den statistischen Funktionen nicht blind vertrauen!

Sitzung 1

Datenerhebung und Häufigkeiten

Lernziele dieser Sitzung

Sie können...

- einige Grundbegriffe der Statistik definieren.
- Typen von Stichproben unterscheiden.
- Skalenniveaus von Variablen bestimmen.
- Häufigkeitsverteilungen beschreiben.

Lehrvideos (Sommersemester 2020)

- [1a\) Grundbegriffe](#)
- [1b\) Skalenniveaus](#)
- [1c\) Grundbegriffe](#)

1.1 Statistische Praxis

Was ist Statistik? Je nach Perspektive kann Statistik vieles sein: ein Teilgebiet der Mathematik, ein Untersuchungsobjekt kritischer Forschung oder ein unbeliebtes Studienfach.

Im Rahmen dieser Veranstaltung soll Statistik als eine Zusammenstellung von Praktiken in der quantitativen Forschung verstanden werden, wobei ihre Anwendung stets im Mittelpunkt steht. Eine hilfreiche Definition findet sich bei Haseloff et al. (1968):

„Allgemein kann gesagt werden: Die Statistik hat es mit Zahlen zu tun, die entweder aus Abzählvorgängen oder aus Messungen gewonnen wurden. Ihre Aufgabe ist es, ein solches Zahlenmaterial in eine optimal übersichtliche und informationsreiche Form zu bringen, aus ihnen methodische Schlußfolgerungen zu ziehen und gegebenenfalls auch die Ursachen der analysierten Zahlenverhältnisse mit sachlichen Methoden aufzudecken.“ (Haseloff et al. 1968: 27)

1.1.1 Grundbegriffe der Statistik

1.1.1.1 Untersuchungselement

Untersuchungselemente (auch Untersuchungseinheiten, Merkmalsträger, bei Personen: Proband*innen, engl. *sampling unit*) sind die individuellen Gegenstände empirischer Untersuchungen. Bei einer Hochrechnung zur Bundestagswahl ist dies z.B. eine befragte Wählerin.

1.1.1.2 Stichprobe

Eine Stichprobe (engl. *sample*) ist die Menge aller Untersuchungselemente, deren Daten direkt erhoben werden. Die Anzahl der Untersuchungselemente in der Stichprobe wird in Formeln mit n bezeichnet. Bei einer Hochrechnung z.B. bilden alle tatsächlich befragten Wähler*innen die Stichprobe.

1.1.1.3 Grundgesamtheit

Die Grundgesamtheit (auch Population, engl. *population*) ist die Menge aller potentiell untersuchbaren Elemente, über die Aussagen getroffen werden sollen. Die Stichprobe ist eine Teilmenge der Grundgesamtheit. Die Anzahl der Elemente in der Grundgesamtheit wird in Formeln mit N bezeichnet. Bei einer Hochrechnung zur Bundestagswahl sind dies z.B. alle Wähler*innen (bzw. alle Wahlberechtigten, wenn Wahlbeteiligung von Interesse ist).

1.1.1.4 Variable

Variablen (auch Merkmale, engl. *variable*) sind Informationen über die Untersuchungselemente, die in einer Untersuchung von Interesse sind. Typischerweise unterscheiden sie sich von Untersuchungselement zu Untersuchungselement, sind also variabel. Bei einer Hochrechnung ist dies die Antwort auf die Frage: „Welche Partei haben Sie gerade gewählt?“

1.1.1.5 Wert

Ein Wert (auch Merkmalsausprägung, engl. *observation*) ist die erfasste Ausprägung einer Variable bei einem Untersuchungselement. In Formeln werden Werte mit $x_1, x_2, x_3, \dots, x_n$ durchnummeriert. Bei einer Hochrechnung kann die Variable „gewählte Partei“ für ein Untersuchungselement z.B. den Wert „CDU“ annehmen.

1.1.1.6 Kennwert

Kennwerte (auch Maßzahlen, Kennzahlen, engl. *summary statistics*) sind Zahlen, die aus den beobachteten Werten errechnet werden. Sie können beispielsweise Aufschluss über Mittelwerte und Verteilung einer Variable oder den Zusammenhang mehrerer Variablen geben. Bei einer Hochrechnung sind z.B. die relativen Häufigkeiten (in Prozent) der Variable „gewählte Partei“ von besonderem Interesse.

1.1.2 Taxonomien statistischer Verfahren

Statistische Verfahren werden in mehrerlei Hinsicht unterschieden, wie im Folgenden beschrieben. Dabei schließen sich verschiedene Kategorien nicht unbedingt aus, es gibt also durchaus statistische Verfahren, die z.B. als univariat *und* deskriptiv bezeichnet werden.

1.1.2.1 Uni-, bi- und multivariate Statistik

Bei diesen Bezeichnungen ist entscheidend, wie viele Variablen bei den jeweiligen Verfahren zum Einsatz kommen. Im Allgemeinen spricht man bei einer Variable von univariater Statistik, bei zwei Variablen von bivariater Statistik und bei mehr als zwei Variablen von multivariater Statistik. (Manchmal werden allerdings auch Verfahren mit nur zwei Variablen als multivariat bezeichnet.)

In dieser Veranstaltung beschäftigen wir uns zunächst mit univariaten, dann mit bivariaten Verfahren. Verfahren mit mehr als zwei Variablen werden nicht behandelt.

1.1.2.2 Deskriptive und schließende Statistik

Unabhängig von der Anzahl der Variablen unterscheidet man auch nach der Art und Weise des Vorgehens:

1.1.2.2.1 Deskriptive Statistik Die deskriptive Statistik (auch: beschreibende Statistik) dient der Beschreibung der Verteilung von Merkmalen, indem sie z. B. Durchschnittswerte bildet, Häufigkeiten bestimmt oder etwas über die Streuung eines Merkmals aussagt. Sie kann so große Datenmengen übersichtlicher machen, indem sie diese ordnet, gruppiert oder verdichtet. Sie erleichtert es also, das Charakteristische, Wichtige zu erkennen.

1.1.2.2.2 Schließende Statistik Die schließende Statistik (auch: analytische, operative Statistik, Inferenzstatistik, Prüfstatistik) verhilft dazu, von Eigenschaften einer Stichprobe auf Eigenschaften der Grundgesamtheit verallgemeinern bzw. schließen zu können (deshalb eben auch: schließende Statistik) und diese Einschätzung überprüfen zu können.

Die schließende Statistik wird weiter unterteilt in Schätz- und Teststatistik:

1.1.2.2.2.1 Schätzende Statistik Die Schätzstatistik schätzt Kennwerte der Grundgesamtheit aus den Kennwerten einer Stichprobe.

1.1.2.2.2.2 Testende Statistik Die Teststatistik überprüft, als wie wahrscheinlich oder unwahrscheinlich gemachte Schätzungen bzw. Hypothesen gelten können.

1.1.3 Ablauf einer statistischen Untersuchung

Eine typische Anwendung statistischer Verfahren in der Forschung folgt diesem Schema:

1.1.3.1 Datenerhebung

- Eigene Erhebung z.B. durch Zählen, Messen, Befragung (primärstatistische Daten)
 - Auswahl von Untersuchungseinheiten
 - Wahl der Datenniveaus
- Rückgriff auf vorhandenes Datenmaterial (sekundärstatistische Daten)

1.1.3.2 Datenaufbereitung

- Verdichtung des gewonnenen Datenmaterials und Digitalisierung in Form einer Datenmatrix
- Verschneidung von mehreren Datensätzen

- Vereinheitlichung und Säuberung der Daten
- Überblick verschaffen durch einfache Beschreibung von Häufigkeiten und Maßzahlen (deskriptive Statistik)

1.1.3.3 Datenauswertung

- Verdichtete Beschreibung von Verteilungsmustern einer Variable (univariate deskriptive Statistik)
- Verdichtete Beschreibung der Beziehung zwischen zwei Variablen (bivariate deskriptive Statistik)
- Schluss von Stichprobe auf Grundgesamtheit (Schätzstatistik)
- Testen von Hypothesen über die Grundgesamtheit (Teststatistik)

1.2 Grundlagen der Datenerhebung

1.2.1 Typen von Stichproben

1.2.1.1 Reine Zufallsstichprobe

Bei endlichen Grundgesamtheiten können Lotterieverfahren angewendet werden. Dabei wird allen Elementen der Grundgesamtheit eine Zahl zwischen 1 und N zugeordnet. Anschließend werden Zufallszahlen ausgewählt und die entsprechenden Elemente in die Stichprobe übernommen.

1.2.1.2 Systematische Zufallsstichprobe

Die Elemente einer endlichen Grundgesamtheit werden in eine Rangordnung gebracht (Nummerierung 1 bis N). Anschließend wählt man jedes (N/n) -te Element aus. So entsteht eine Stichprobe der Größe n .

1.2.1.3 Geschichtete Zufallsstichprobe

Die Elemente einer endlichen Grundgesamtheit werden in Schichten (Klassen) zusammengefasst. Anschließend zieht man eine Zufallsstichprobe aus jeder Schicht. Geschichtete Stichproben setzen die Kenntnis einiger Parameter der Grundgesamtheit voraus. Zur Aufteilung des Stichprobenumfangs auf die einzelnen Schichten wird in der Regel die proportionale Aufteilung gewählt.

1.2.1.4 Klumpenstichprobe

Hier ist die Grundgesamtheit schon in „natürliche“ Gruppen aufgeteilt (z.B. Schulklassen) und es werden mehrere dieser Gruppen (Klumpen, engl. *cluster*) nach einem Zufallsverfahren als Stichprobe gewählt.

„Man beachte, dass ein einzelner Klumpen (...) keine Klumpenstichprobe darstellt, sondern eine Ad-hoc-Stichprobe, bei der zufällige Auswahlkriterien praktisch keine Rolle spielen. Die Bezeichnung „Klumpenstichprobe“ ist nur zu rechtfertigen, wenn mehrere zufällig ausgewählte Klumpen vollständig untersucht werden.“ (Bortz und Schuster 2010: 81)

1.2.2 Variablentypen

Tabelle 1.1: Die vier wichtigsten Skalenniveaus

| Skalenart | Beispiel | mögliche Aussagen | gültige Lagemaße |
|-----------------|--------------------|--------------------------------|-------------------------|
| Nominalskala | Postleitzahl | Gleichheit, Verschiedenheit | Modus |
| Ordinalskala | Militärischer Rang | + Größer-kleiner-Relationen | + Median |
| Intervallskala | Temperatur in °C | + Gleichheit von Differenzen | + arithmetisches Mittel |
| Verhältnisskala | Körpergröße | + Gleichheit von Verhältnissen | + geometrisches Mittel |

1.2.2.1 Qualitative Variablen

Qualitative Variablen können nicht der Größe nach, sondern nur im Hinblick auf ihre Eigenschaft/Art („Qualität“) unterschieden werden (z.B. Parteizugehörigkeit, Telefonnummer, Automarke).

Qualitative Variablen, die nur zwei mögliche Werte annehmen können, nennt man „dichotome“ Variablen (etwa Antworten auf Ja-Nein-Fragen).

1.2.2.2 Quantitative Variablen

Quantitative Variablen können der Größe nach unterschieden werden (Bsp. Geburtenzahl, Arbeitslosenrate).

Quantitative Variablen können diskret oder stetig sein:

1.2.2.2.1 Diskrete Variablen Diskrete Variablen (auch diskontinuierliche Variablen) können nur endlich viele, ganzzahlige Werte annehmen. Zwischen zwei Ausprägungen befindet sich eine abzählbare Menge anderer Ausprägungen (z.B. Anzahl eigener Kinder, Haushaltsgröße in Personen).

1.2.2.2.2 Stetige Variablen Stetige Variablen (auch: kontinuierliche Variablen) können in einem bestimmten Bereich jede beliebige Ausprägung annehmen. Der Ausdehnungsbereich kennt keine Lücken, sondern ist als ein fortlaufendes Kontinuum vorstellbar: Bei stetigen Variablen können zwischen zwei Werten oder Ausprägungen unendlich viele weitere Ausprägungen oder Werte liegen (z.B. Körpergröße, Längengrad in Dezimalform).

1.2.3 Skalenniveaus

Eine Variable lässt sich aufgrund ihrer Eigenschaften einem Skalenniveau (auch Skalentyp, Messniveau, Datenniveau, engl. *level of measurement*) zuordnen. Bestimmte Rechenoperationen und statistische Verfahren setzen bestimmte Skalenniveaus voraus. Deshalb ist es wichtig zu wissen, welchem Skalenniveau eine Variable zuzuordnen ist.

Variablen lassen sich immer auch einem niedrigeren Skalenniveau zuordnen. Dies geht allerdings mit Informationsverlust einher.

Die im Folgenden beschriebenen Skalenniveaus sind nicht deckungsgleich mit den o.g. Variablentypen. Intervall- und Verhältnisskalen können z.B. jeweils diskret oder stetig sein.

In Tabelle 1.1 sind die wichtigsten Skalenniveaus im Überblick aufgeführt. „Gültige Lagemaße“ sind dabei als Zusatzinformation aufgelistet und werden erst in der [nächsten Sitzung](#) behandelt.

1.2.3.1 Nominalskala

Die Merkmalsausprägungen einer Variable stehen je ‚für sich‘; sie lassen sich nicht sinnvoll in eine Rangordnung bringen oder gar miteinander verrechnen.

Die einzige Aussage, die sich über zwei Werte in einer Nominalskala treffen lässt, ist dass sie gleich oder nicht gleich sind.

Beispiele: Postleitzahlen, Telefonnummern, Staatsangehörigkeit, Krankheitsklassifikationen

1.2.3.2 Ordinalskala

Die Merkmalsausprägungen einer Variablen lassen sich sinnvoll in eine Rangordnung bringen, die Abstände zwischen den Merkmalsausprägungen aber lassen sich nicht sinnvoll quantifizieren.

Über zwei Werte in einer Ordinalskala lässt sich nicht nur sagen, ob sie gleich oder verschieden sind (wie in der Nominalskala), sondern darüber hinaus, welcher Wert bei Verschiedenheit größer ist.

Beispiele: Militärische Ränge, Windstärken, pauschale Häufigkeitsangaben (sehr oft ... nie), Zufriedenheitsangaben (sehr zufrieden ... unzufrieden)

1.2.3.3 Metrische Skalen (oder Kardinalskalen)

Abstände zwischen den Merkmalsausprägungen lassen sich exakt angeben.

Zusätzlich zu den Möglichkeiten der Ordinalskala können auf einer metrischen Skala Rechenoperationen auch sinnvoll auf die Differenzen zwischen den Merkmalsausprägungen angewendet werden.

Metrische Skalen werden unterteilt in Intervall- und Verhältnisskalen:

1.2.3.3.1 Intervallskala

Maßeinheit und Wahl des Nullpunktes sind willkürlich gewählt.

Beispiele: Grad Celsius, Geburtsjahr als Jahreszahl („1961“), in der Praxis häufig: subjektive Bewertung auf einer Skala von 1 bis 10.

1.2.3.3.2 Verhältnisskala (auch Ratioskala)

Es gibt einen invarianten (absoluten, natürlichen) Nullpunkt.

In einer Verhältnisskala lassen sich über alle o.a. Möglichkeiten hinaus auch Aussagen über Verhältnisse zwischen Werten treffen (z.B. „ x_1 ist doppelt so groß wie x_2 “).

Beispiele: Lebensalter in Jahren, Haushaltgröße, Körpergröße, Körpergewicht

1.3 Häufigkeitsverteilungen

1.3.1 Urliste

Die Urliste ist eine ungeordnete Liste aller erfassten Werte.

Für die statistische Erhebung „Anfangsbuchstaben der Vornamen von Teilnehmenden an einer Statistikvorlesung“ könnte die Urliste z.B. so aussehen:

T J D T E N D F F M A J V T T V A L V P J K P M F M A J N A C I T P B A P H T L
N S P C K J K L J R E Y M K H M N L A A L L M L J G P L B F L J J V M P C J M J

S A M M M P A A L L O C J L P L V F J R M A V K S B B B N C A A T J P C F L E B
 L C A K A L T V Y P F L J S T T N R J A S E L M L T A E B M N M V D P P L N L B
 A A J M L N N S H M

1.3.2 Geordnete Liste

Die geordnete Liste bringt die Werte der Urliste in eine geeignete Reihenfolge, so dass die unterschiedlichen Werte leicht gezählt werden können:

A A A A A A A A A A A A A A A A A A B B B B B B B B C C C C C C C D D D E E E
 E E F F F F F F F G H H H I J J J J J J J J J J J J J J J J K K K K K K L L L
 L L L L L L L L L L L L L L L L L M M M M M M M M M M M M M M M M N N N N
 N N N N N O P P P P P P P P P P P P R R R S S S S S S T T T T T T T T T T
 V V V V V V V Y Y

1.3.3 Häufigkeiten

Die absoluten Häufigkeiten erhält man durch einfaches Abzählen der jeweiligen Werte. Für die relativen Häufigkeiten teilt man diese Zahl durch n . Kumulierte Häufigkeiten zählen die bisherigen Summen bzw. Anteile zusammen (s. Tabelle 1.2).

Softwarehinweis

In R lässt sich mit dem Befehl `table()` eine einfache Häufigkeitstabelle aus Rohdaten erstellen.

1.3.4 Stabdiagramme

Die so ermittelten Häufigkeiten lassen sich als Stabdiagramm (auch Säulen-, Streifen-, Balkendiagramm, engl. *bar chart*) darstellen (s. Abbildung 1.1).

Softwarehinweis

In R lautet der Standardbefehl zur Erstellung eines Stabdiagramms `barplot()`.

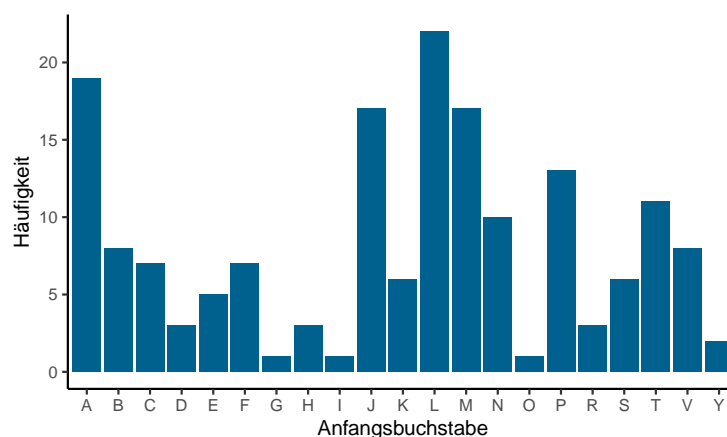


Abbildung 1.1: Stabdiagramm

Tabelle 1.2: Tabelle mit kumulierten Häufigkeiten

| Buchstabe | Absolute Häufigkeit f | f_{kum} | Relative Häufigkeit | $\%_{kum}$ |
|------------------|---|-----------------------------|----------------------------|------------------------------|
| A | 19 | 19 | 11,2% | 11,2% |
| B | 8 | 27 | 4,7% | 15,9% |
| C | 7 | 34 | 4,1% | 20% |
| D | 3 | 37 | 1,8% | 21,8% |
| E | 5 | 42 | 2,9% | 24,7% |
| F | 7 | 49 | 4,1% | 28,8% |
| G | 1 | 50 | 0,6% | 29,4% |
| H | 3 | 53 | 1,8% | 31,2% |
| I | 1 | 54 | 0,6% | 31,8% |
| J | 17 | 71 | 10% | 41,8% |
| K | 6 | 77 | 3,5% | 45,3% |
| L | 22 | 99 | 12,9% | 58,2% |
| M | 17 | 116 | 10% | 68,2% |
| N | 10 | 126 | 5,9% | 74,1% |
| O | 1 | 127 | 0,6% | 74,7% |
| P | 13 | 140 | 7,6% | 82,4% |
| R | 3 | 143 | 1,8% | 84,1% |
| S | 6 | 149 | 3,5% | 87,6% |
| T | 11 | 160 | 6,5% | 94,1% |
| V | 8 | 168 | 4,7% | 98,8% |
| Y | 2 | 170 | 1,2% | 100% |

Tabelle 1.3: Häufigkeitstabelle mit klassierten Werten

| Durchmesser | Absolute Häufigkeit f | f_{kum} | Relative Häufigkeit | $\%_{kum}$ |
|---------------------|-------------------------|-----------|---------------------|------------|
| über 8 bis 10 Zoll | 3 | 3 | 9,7% | 9,7% |
| über 10 bis 12 Zoll | 12 | 15 | 38,7% | 48,4% |
| über 12 bis 14 Zoll | 6 | 21 | 19,4% | 67,7% |
| über 14 bis 16 Zoll | 3 | 24 | 9,7% | 77,4% |
| über 16 bis 18 Zoll | 6 | 30 | 19,4% | 96,8% |
| über 18 bis 20 Zoll | 0 | 30 | 0% | 96,8% |
| über 20 bis 22 Zoll | 1 | 31 | 3,2% | 100% |

1.3.5 Quantitative Variablen

Das oben beschriebene Verfahren funktioniert gut für qualitative Variablen (und diskrete Variablen mit wenigen unterschiedlichen Werten). Für quantitative Variablen wird ein anderes Verfahren empfohlen.

Zur Veranschaulichung soll diese geordnete Liste von Messwerten des Stammdurchmessers von Schwarzkirschen (Beispieldatensatz `trees` aus [R Core Team 2018](#)) dienen:

8,3 8,6 8,8 10,5 10,7 10,8 11,0 11,0 11,1 11,2 11,3 11,4 11,4 11,7 12,0 12,9
12,9 13,3 13,7 13,8 14,0 14,2 14,5 16,0 16,3 17,3 17,5 17,9 18,0 18,0 20,6

Für solche Verteilungen müssen zuerst Klassen (engl. *bins*) gebildet werden, in denen die Werte dann zusammengefasst werden (s. Tabelle 1.3).

Für die Wahl der Klassengrenzen gibt es zwei feste Regeln:

- Alle Werte müssen abgedeckt sein.
- Die Klassen dürfen sich nicht überlappen.

Zusätzlich sollten die folgenden Konventionen nach Möglichkeit befolgt werden:

- Klassen sollten gleich große Wertebereiche abdecken.
- Alle Klassen sollten besetzt sein.
- Klassengrenzen sollten möglichst glatte Zahlen sein.
- Aus Gründen der Übersichtlichkeit sollten nicht mehr als 20 Klassen gewählt werden.
- Klassengrenzen sollten „Klumpen“ mit ähnlichen Werten nicht trennen.

Die Darstellung erfolgt in so genannten Histogrammen (engl. *histogram*). Abbildung 1.2 enthält ein Beispiel für ein Histogramm.

Softwarehinweis

In R können Histogramme mit `hist()` erstellt werden.

1.3.6 Polygone

Statt ausgefüllten Flächen wie im Histogramm lassen sich für die Häufigkeiten auch Punkte setzen, die dann mit Linien verbunden werden. So entsteht ein Häufigkeitspolygon (s. Abbildung 1.3).

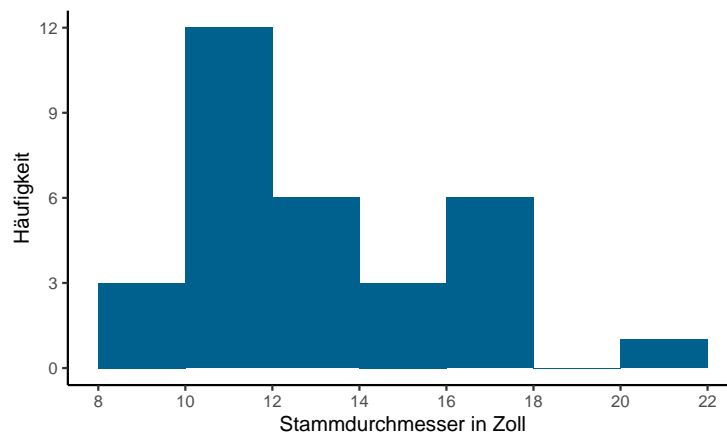


Abbildung 1.2: Histogramm

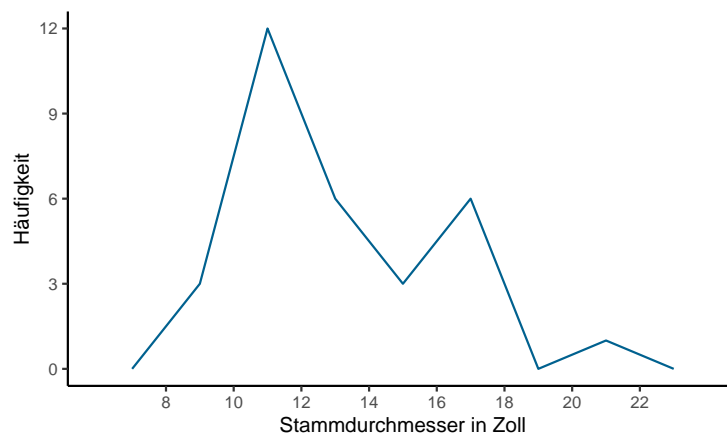


Abbildung 1.3: Polygonzug

1.3.7 Eigenschaften von Häufigkeitsverteilungen

Polygone von Häufigkeitsverteilungen (insbesondere in geglätteter Form) ergeben Annäherungen an so genannte Dichtefunktionen (engl. *density functions*). Diese lassen sich mit Attributen (uni-/bimodal, schmal-/breitgipflig, etc.) beschreiben, wie in Abbildung 1.4 veranschaulicht.

Tipps zur Vertiefung

1.3.8 Grundbegriffe

- YouTube-Kanal „Kurzes Tutorium Statistik“: [Statistische Grundbegriffe](#)
- Kapitel 1.1 in Bortz und Schuster (2010)
- Kapitel 1.1 in Benninghaus (2007)
- Kapitel 2.1 in Bahrenberg, Giese und Nipper (2010)
- *Englisch*: Kapitel 1 in Burt und Barber (1996)

1.3.9 Stichproben

- Kapitel 6.1 in Bortz und Schuster (2010)

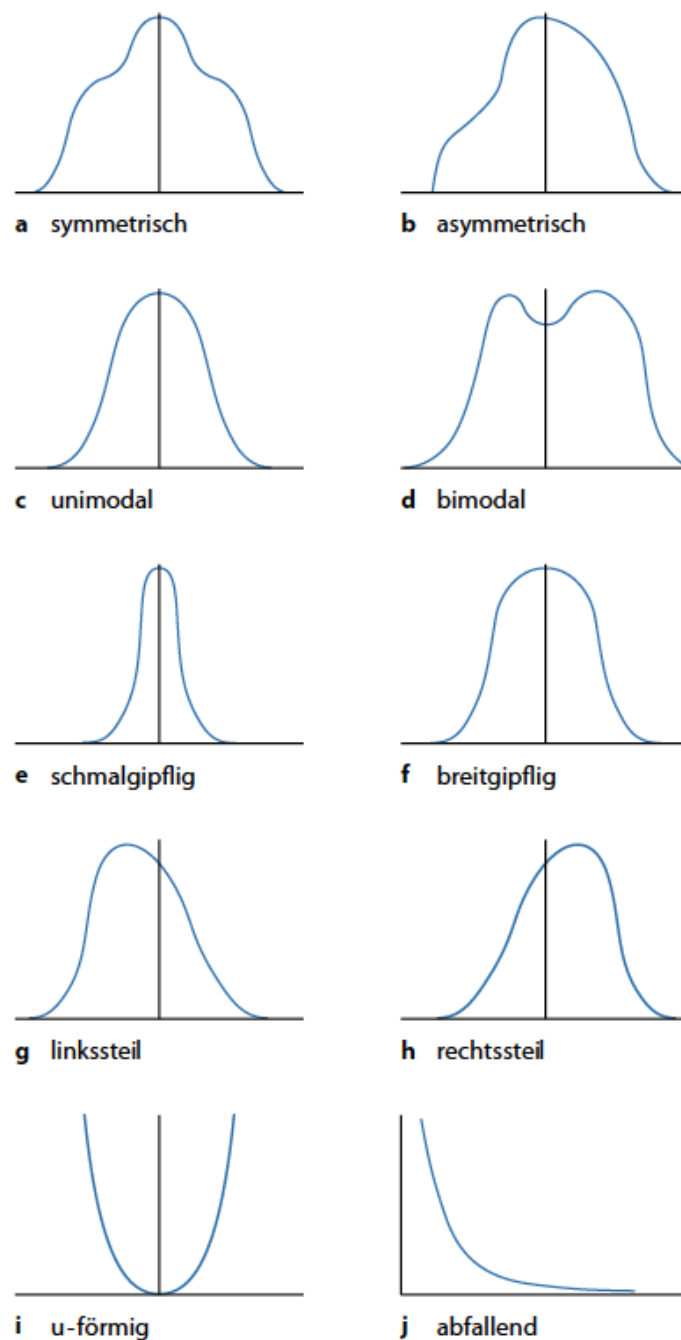


Abbildung 1.4: Merkmale von Verteilungen [aus: @bortz: 42]

- Kapitel 2.5 in Lange und Nipper (2018)
- Kapitel 2.3 in Bahrenberg, Giese und Nipper (2010)
- *Englisch*: Kapitel 1 in Burt und Barber (1996)

1.3.10 Skalenniveaus

- Kapitel 1.2 in Bortz und Schuster (2010)
- Kapitel 2.5 in Lange und Nipper (2018)

- Kapitel 2.1 in Benninghaus (2007)
- Kapitel 2.2 in Bahrenberg, Giese und Nipper (2010)
- YouTube-Kanal „Kurzes Tutorium Statistik“: [Skalenniveaus](#)
- *Englisch*: Kapitel 1.3 in Burt und Barber (1996)

1.3.11 Häufigkeiten und Diagramme

- YouTube-Kanal „Kurzes Tutorium Statistik“: [Stabdiagramme und Histogramme](#)
- Kapitel 3.1 und 3.2 in Bortz und Schuster (2010)
- Kapitel 2.5 in Lange und Nipper (2018)
- Kapitel 1.2 in Benninghaus (2007)
- Kapitel 4.1 in Bahrenberg, Giese und Nipper (2010)
- *Englisch*: Kapitel 2.1 in Burt und Barber (1996)

Übungsaufgaben

1.3.12 Aufgabe 1-1

[zur Lösung](#)

Teilen Sie in Ihrer Kleingruppe folgende Begriffe untereinander auf:

- Variable
- Kennwert
- Wert
- Grundgesamtheit
- Stichprobe
- Untersuchungselement

Gehen Sie nun für jeden Begriff wie folgt vor:

1. Erklären Sie der Reihe nach „Ihren“ Begriff den anderen Gruppenmitgliedern, gerne auch mit Beispielen.
2. Die anderen Gruppenmitglieder nehmen die Rolle von unwissenden Dritten ein und stellen bei Bedarf Nachfragen.
3. Die anderen Gruppenmitglieder geben direkt danach Feedback auf die Erklärung:
 - Was fanden Sie gut erklärt?
 - Was fanden Sie unverständlich?
 - Was hat Ihnen gefehlt?

1.3.13 Aufgabe 1-2

[zur Lösung](#)

Finden Sie als Gruppe jeweils zwei Beispiele für:

- systematische Zufallsstichproben
- geschichtete Zufallsstichproben
- Klumpenstichproben

1.3.14 Aufgabe 1-3[zur Lösung](#)

Bestimmen Sie das Skalenniveau der folgenden Variablen. Kennzeichnen Sie darüber hinaus, ob die Variable qualitativ, diskret oder stetig ist.

- a) Lebensalter in Jahren
- b) Regenmenge in mm
- c) Güteklasse
- d) Passagieraufkommen
- e) Baujahr
- f) Geschwindigkeit in km/h
- g) Sozialstatus (Unter-, Mittel und Oberschicht)
- h) Temperatur in °F
- i) Fläche eines Bundeslands in km²
- j) Temperatur in K
- k) Einwohnerzahl
- l) Pegelstand
- m) Staatsangehörigkeit
- n) Interesse an Statistik (gering bis hoch)
- o) Klausurnote
- p) Bodentyp
- q) Entfernung zum Stadtzentrum in km
- r) Körpergröße
- s) Kleidergröße (S bis XXL)
- t) Monatliches Nettoeinkommen

1.3.15 Aufgabe 1-4[zur Lösung](#)

Folgende Werte seien erfasst über die Lebensdauer von Klimaanlage in Stunden (Beispieldatensatz `aircondit7` aus [R Core Team 2018](#)):

14 23 15 139 13 39 188 22 50 3 36 46 30 5 102 5 88 22 197 72 210 97 79 44

- a) Erstellen Sie eine Häufigkeitstabelle. Welche Klassen wählen Sie und warum?
- b) Zeichnen Sie ein Histogramm.
- c) Beschreiben Sie die Verteilung.

1.3.16 Aufgabe 1-5[zur Lösung](#)

Sind die folgenden Aussagen wahr oder unwahr?

- a) Die Auswahl z. B. jedes 100. Merkmalsträgers nennt man „systematische Stichprobe“.
- b) Eine Stichprobe kann eine Grundgesamtheit niemals völlig richtig repräsentieren, es gibt immer einen Zufallsfehler.
- c) Die Größe der Stichprobe wird auch mit N bezeichnet.
- d) Klassengrenzen müssen so gewählt werden, dass alle Werte abgedeckt sind.

- e) Je stärker die Werte der Variablen streuen, desto kleiner sollte die Stichprobe sein.
- f) Variablen auf der Verhältnisskala sind immer metrisch und stetig.
- g) Verhältnisskala und Intervallskala unterscheiden sich durch den natürlichen Nullpunkt.
- h) Intervallskalierte Daten können immer auf die Nominalskala transformiert werden.
- i) Ordinalskalierte Daten können immer auf die Intervallskala transformiert werden.
- j) Eine stetige Variable ist nicht zwingend auch metrisch.
- k) Im Gegensatz zu nominalskalierten Variablen lassen sich Werte von ordinalskalierten Variablen in eine sinnvolle Reihenfolge bringen.
- l) Die relative Häufigkeit eines Werts ist nie größer als 100%.
- m) Verfahren der deskriptiven Statistik sind immer auch univariat.
- n) Klassengrenzen dürfen sich in Ausnahmefällen überlappen.
- o) x_3 ist immer kleiner als x_4 .
- p) Variablen auf der Verhältnisskala haben einen natürlichen Nullpunkt.
- q) Die absolute Häufigkeit eines Werts ist immer eine positive ganze Zahl.
- r) Wenn man die Urliste ordnet, erhält man die geordnete Liste.

Sitzung 2

Maßzahlen

Lernziele dieser Sitzung

Sie können...

- die wichtigsten Lagemaße von Stichproben bestimmen.
- die wichtigsten Streumaße von Stichproben bestimmen.
- Boxplots interpretieren.

Lehrvideos (Sommersemester 2020)

- [2a\) Lagemaße](#)
- [2b\) Streumaße](#)
- [2c\) Klassierte Verteilungen](#)
 - In diesem Video ist mir ein Fehler unterlaufen: Bei Minute 6:30 muss das arithmetische Mittel $\bar{x} \approx 4,59$ betragen. Daraus ergibt sich ein Folgefehler: Die Varianz müsste den Wert $s^2 \approx 14,56$ haben.

2.1 Einleitende Bemerkungen

Die im Folgenden besprochenen Maßzahlen (oder Kennzahlen, Parameter) verdichten (oder aggregieren) Häufigkeitsverteilungen einer Variable. Durch diese Parameter kann das Charakteristische einer Verteilung schnell erfasst und vergleichbar gemacht werden. Die Verdichtung auf Maßzahlen geht jedoch immer auch mit Informationsverlust einher.

Die Möglichkeit der Angabe statistischer Maßzahlen ist abhängig vom Skalenniveau der Daten, wie der Überblick in Tabelle [2.1](#) zeigt.

2.1.1 Beispielverteilung

Alle Berechnungen von Maßzahlen werden am folgenden Beispiel illustriert: Für die 14 Gemeinden im Landkreis Rothenberge wurde die jeweilige Anzahl an Gaststätten erhoben. Die Zählung ergab die Wertereihe in Tabelle [2.2](#).

Tabelle 2.1: Die wichtigsten Maßzahlen

| Parameter | Typ | Mindestes Skalenniveau | Formel |
|-----------------------|----------|------------------------|---|
| Modalwert | Lagemaß | nominal | Mo |
| Median | Lagemaß | ordinal | $Md = \begin{cases} \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2} \\ x_{(\frac{n+1}{2})} \end{cases}$ für |
| Arithmetisches Mittel | Lagemaß | metrisch | $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ |
| Spannweite | Streumaß | ordinal | $R = x_{(n)} - x_{(1)}$ |
| Quartilsabstand | Streumaß | ordinal | $IQR = Q_3 - Q_1$ |
| Varianz | Streumaß | metrisch | $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$ |
| Standardabweichung | Streumaß | metrisch | $s = \sqrt{s^2}$ |

Tabelle 2.2: Beispielverteilung

| x_1 | x_2 | x_3 | x_4 | x_5 | x_6 | x_7 | x_8 | x_9 | x_{10} | x_{11} | x_{12} | x_{13} | x_{14} |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|----------|
| 4 | 1 | 4 | 1 | 5 | 5 | 0 | 1 | 8 | 5 | 1 | 25 | 3 | 3 |

Tabelle 2.3: Sortierte Wertereihe

| $x_{(1)}$ | $x_{(2)}$ | $x_{(3)}$ | $x_{(4)}$ | $x_{(5)}$ | $x_{(6)}$ | $x_{(7)}$ | $x_{(8)}$ | $x_{(9)}$ | $x_{(10)}$ | $x_{(11)}$ | $x_{(12)}$ | $x_{(13)}$ | $x_{(14)}$ |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|------------|------------|------------|------------|
| 0 | 1 | 1 | 1 | 1 | 3 | 3 | 4 | 4 | 5 | 5 | 5 | 8 | 25 |

2.2 Lagemaße

Lagemaße (auch Maße der Zentraltendenz, Lokalisationsparameter, Mittelwerte, engl. *measures of central tendency*) bezeichnen alle statistischen Maßzahlen, die eine Verteilung repräsentieren, indem sie die Lage der mittleren oder häufigsten Variablenwerte angeben.

Im Falle einer unimodalen, perfekt symmetrischen Verteilung (z. B. Glockenform) haben alle drei Lageparameter den gleichen Wert. Je weiter Verteilungen von dieser Form abweichen – durch Mehrgipfligkeit oder Asymmetrie – desto unpräziser ist die Beschreibung der Verteilung durch einen einzigen Parameter.

2.2.1 Median

Der Median (engl. *median*) einer Verteilung ist der Wert, der größer als genau 50% aller Werte ist.

Da dies eine Größer-kleiner-Relation der Werte voraussetzt, kann der Median nur für ordinale und metrische Skalenniveaus angegeben werden.

Im Folgenden wird die (einfachere) Bestimmung des Medians nach Bortz und Schuster (2010) verwendet. Benninghaus (2007) beschreibt ein anderes Verfahren, welches zu anderen Ergebnissen kommen kann.

Um den Median zu bestimmen, wird zunächst eine geordnete Liste angefertigt, indem die Werte aufsteigend sortiert werden. Diese sortierten Werte werden mit $x_{(1)}, x_{(2)}, x_{(3)}, \dots, x_{(n)}$ bezeichnet (also mit Klammern). Für unsere Beispielveilung ergibt sich Tabelle 2.3.

Bei einer ungeraden Stichprobengröße n teilt der $(\frac{n+1}{2})$ -te Wert (also der Wert genau in der Mitte) die Stichprobe in zwei Hälften, weshalb gilt:

$$Md = x_{(\frac{n+1}{2})} \quad \text{falls } n \text{ ungerade.}$$

Bei geradem n entstehen zwei gleich große Hälften der Stichprobe: $x_{(1)}$ bis $x_{(\frac{n}{2})}$ einerseits, und $x_{(\frac{n}{2}+1)}$ bis $x_{(n)}$ andererseits. Der Durchschnitt zwischen $x_{(\frac{n}{2})}$ und $x_{(\frac{n}{2}+1)}$ teilt die Stichprobe in zwei Hälften. Es gilt:

$$Md = \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2} \quad \text{falls } n \text{ gerade.}$$

In unserem Beispiel ist $n = 14$ und damit gerade. Der Median errechnet also nach Formel (2.2.1) wie folgt:

Tabelle 2.4: Häufigkeiten der Beispielverteilung

| Wert x_i | Häufigkeit f_i |
|------------|------------------|
| 0 | 1 |
| 1 | 4 |
| 3 | 2 |
| 4 | 2 |
| 5 | 3 |
| 8 | 1 |
| 25 | 1 |

$$\begin{aligned}
 Md &= \frac{x_{(7)} + x_{(8)}}{2} \\
 &= \frac{3 + 4}{2} \\
 &= 3,5
 \end{aligned}$$

Softwarehinweis

In R gibt die Funktion `median()` den Median einer Verteilung aus.

2.2.2 Modalwert

Der Modalwert Mo (auch Modus, engl. *mode*) gibt den häufigsten Wert oder die häufigsten Werte einer Verteilung an.

Der Modalwert kann so auch (als einziger Mittelwert) für nominalskalierte Variablen angegeben werden.

Bei ordinalen und metrischen Skalenniveaus sind folgende Besonderheiten zu beachten:

- Wird der Modus einer Verteilung durch unmittelbar benachbarte Werte gebildet, wird er als Kombination (bei metrischen Variablen als arithmetisches Mittel) dieser Werte angegeben.
- Bei bimodalen (multimodalen) Verteilungen werden beide (alle) Modalwerte angegeben.

Hierzu müssen die Häufigkeiten der Werte bekannt sein, bzw. bestimmt werden (s. Tabelle 2.4).

Der Modalwert der Beispielverteilung beträgt 1, da der Wert 1 am häufigsten (viermal) vorkommt.

2.2.3 Arithmetisches Mittel

Das arithmetische Mittel (auch Mittelwert, Durchschnitt, engl. *mean*) ist das gebräuchlichste Lagemaß und Grundlage für viele statistische Verfahren.

Das arithmetische Mittel setzt ein metrisches Skalenniveau voraus.

Die Berechnung des arithmetischen Mittels einer Stichprobe erfolgt durch die Formel:

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

Für unsere Beispielverteilung ergibt sich durch einsetzen in Formel (2.2.3):

$$\begin{aligned}\bar{x} &= \frac{\sum_{i=1}^{14} x_i}{14} \\ &= \frac{4 + 1 + 4 + 1 + 5 + 5 + 0 + 1 + 8 + 5 + 1 + 25 + 3 + 3}{14} \\ &= \frac{63}{14} \\ &\approx 4,71\end{aligned}$$

Softwarehinweis

Der Befehl für die Ermittlung des arithmetischen Mittels in R lautet `mean()`.

2.3 Streumaße

Streumaße (auch Streuungs-, Variabilitäts-, Dispersionswerte, engl. *measures of variability*) geben Auskunft darüber, wie heterogen die Werte einer Verteilung sind, d. h. wie breit sie gestreut sind. Während Lagemaße den typischen Wert einer Verteilung ermitteln, zeigen Streumaße, wie gut (oder eigentlich: wie schlecht) dieser typische Wert die Verteilung repräsentiert.

2.3.1 Spannweite

Die Spannweite (engl. *range*) gibt Auskunft darüber, wie groß der Wertebereich ist, der von einer Verteilung abgedeckt wird. Sie wird (für metrische Skalen) als die Differenz vom größten zum kleinsten Wert (also vom letzten zum ersten Wert einer geordneten Werteliste) angegeben:

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

Für unsere Beispielstichprobe ergibt sich (mit Blick auf Tabelle 2.3):

$$\begin{aligned}R &= x_{(14)} - x_{(1)} \\ &= 25 - 0 \\ &= 25\end{aligned}$$

Softwarehinweis

In R gibt die Funktion `range()` die Werte für $x_{(1)}$ und $x_{(n)}$ aus.

| $x_{(1)}$ | $x_{(2)}$ | $x_{(3)}$ | $x_{(4)}$ | $x_{(5)}$ | $x_{(6)}$ | $x_{(7)}$ |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 0 | 1 | 1 | 1 | 1 | 3 | 3 |

| $x_{(8)}$ | $x_{(9)}$ | $x_{(10)}$ | $x_{(11)}$ | $x_{(12)}$ | $x_{(13)}$ | $x_{(14)}$ |
|-----------|-----------|------------|------------|------------|------------|------------|
| 4 | 4 | 5 | 5 | 5 | 8 | 25 |

2.3.2 Quartilsabstand

Der Quartilsabstand (auch Interquartilsabstand, engl. *interquartile range*, *IQR*) gibt die Größe des Wertebereichs der mittleren 50% einer Verteilung an.

Genau so wie der Median eine Messwertreihe in zwei gleich große Hälften „schneidet“, schneiden die Quartile die Werte in Viertel. Dabei liegt der so genannte untere Angelpunkt Q_1 genau über 25% der Werte, Q_2 ist identisch mit dem Median und der obere Angelpunkt Q_3 liegt genau über 75% der Werte.

Der Angelpunkt Q_1 wird ermittelt, indem der Median für die unteren 50% (Q_3 : die oberen 50%) der Werte bestimmt wird – also jener Werte, die theoretisch unterhalb des Medians der Gesamtverteilung liegen.

Dabei folgen wir Bortz und Schuster (2010) und nehmen im Fall eines ungeraden n den Median auf beiden Seiten hinzu.

Die Formel für den Quartilsabstand lautet:

$$IQR = Q_3 - Q_1$$

Der Quartilsabstand ist Ausreißern gegenüber stabiler als die Spannweite, da extreme hohe oder niedrige Wert nicht in die Berechnung einfließen.

In unserem Beispiel (mit $n = 14$) ist die untere Hälfte der Verteilung:

Q_1 ist der Median dieser Werte, also $x_{(4)} = 1$.

Die oberen 7 Werte lauten:

Q_3 ist also $x_{(11)} = 5$.

Für den Quartilsabstand ergibt sich durch einsetzen in Formel (2.3.2):

$$\begin{aligned} IQR &= 5 - 1 \\ &= 4 \end{aligned}$$

Softwarehinweis

In R werden die Quartile üblicherweise mit `quantile()` und der Quartilsabstand mit `IQR()` bestimmt.

Tabelle 2.5: Häufigkeitstabelle zur Berechnung der Varianz

| Werte x_i | Häufigk. f_i | $(x_i - \bar{x})$ | $(x_i - \bar{x})^2$ | $f_i \cdot (x_i - \bar{x})^2$ |
|-------------|----------------|-------------------|---------------------|-------------------------------|
| 0 | 1 | -4,71 | 22,18 | 22,18 |
| 1 | 4 | -3,71 | 13,76 | 55,04 |
| 3 | 2 | -1,71 | 2,92 | 5,84 |
| 4 | 2 | -0,71 | 0,50 | 1,00 |
| 5 | 3 | 0,29 | 0,08 | 0,24 |
| 8 | 1 | 3,29 | 10,82 | 10,82 |
| 25 | 1 | 20,29 | 411,68 | 411,68 |

Achtung: Genau wie für den Median gibt es auch für die Ermittlung der Quartile bzw. des Quartilsabstands unterschiedliche Verfahren. Die Ergebnisse dieser R-Funktionen weichen hier deshalb meist leicht vom hier besprochenen Verfahren ab!

2.3.3 Varianz

Die Varianz einer Messwertreihe (engl. *variance*) kann verstanden werden als der durchschnittliche quadrierte Abstand der Werte zum arithmetischen Mittel.

Die Formel lautet:

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

Die Quadrierung der Differenz hat dabei einen doppelten Effekt: Zum einen bekommen auch negative Differenzen ein positives Vorzeichen, so dass sich positive und negative Differenzen nicht neutralisieren. Zum anderen werden hierdurch besonders große Abweichungen zum arithmetischen Mittel stärker gewichtet als dies ohne Quadrierung der Fall wäre.

Zudem fällt auf, dass im Gegensatz zur Formel für das arithmetische Mittel im Nenner $n - 1$ steht und nicht etwa n . Dies hat mit so genannten Freiheitsgraden zu tun, die wir allerdings erst in [Sitzung 5](#) genauer kennenlernen.

Für unsere Beispielstichprobe wird die Berechnung für alle einzelnen $(x_i - \bar{x})^2$ schnell aufwendig und unübersichtlich. Deshalb berechnen wir ihre Summe hier mit Hilfe einer Häufigkeitstabelle (s. Tabelle 2.5). Dabei werden alle distinkten Werte einzeln transformiert und in der letzten Spalte mit ihrer Häufigkeit multipliziert.

Schließlich werden die Werte in Formel (2.3.3) eingesetzt:

$$\begin{aligned}s^2 &= \frac{\sum_{i=1}^{14} (x_i - \bar{x})^2}{14 - 1} \\&\approx \frac{22,18 + 55,04 + 5,84 + 1 + 0,24 + 10,82 + 411,68}{13} \\&= \frac{506,80}{13} \\&\approx 38,98\end{aligned}$$

Eine solche Tabelle lässt sich analog auch für die Berechnung von Summen größerer Messwertreihen für das arithmetische Mittel verwenden.

Zudem lässt dieses Verfahren sich auf klassierte Daten anwenden, wenn für x_i der Mittelwert der Klassen eingesetzt wird (womit allerdings Informations- und Präzisionsverlust einhergeht).

Softwarehinweis

In R lautet der Befehl für die Errechnung der Varianz `var()`.

2.3.4 Standardabweichung

Die Standardabweichung (engl. *standard deviation*) ist das gebräuchlichste Streumaß und spielt eine herausragende Rolle in den allermeisten statistischen Verfahren.

Die Standardabweichung einer Messwertreihe ist definiert als die Quadratwurzel ihrer Varianz:

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

Indem hier die Wurzel gezogen wird, wird in gewisser Weise die Quadrierung der Differenzen für die Varianz wieder „korrigiert“. Insbesondere wird die Quadrierung der Maßeinheit wieder aufgehoben – die Standardabweichung hat also die gleiche Einheit wie die Messreihe selbst.

In unserem Beispiel beträgt die Standardabweichung also:

$$s \approx \sqrt{38,98} \approx 6,24$$

Softwarehinweis

Die Standardabweichung wird in R mit der Funktion `sd()` berechnet.

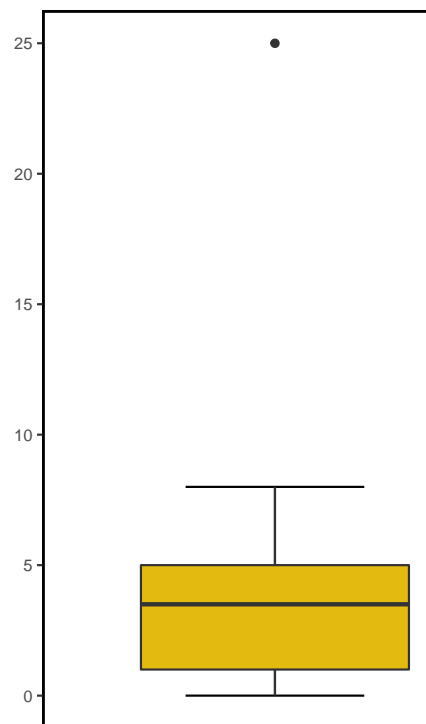


Abbildung 2.1: Boxplot der Beispielveilung

2.4 Boxplot

Der Boxplot (auch Box-and-whisker-plot) kombiniert einige der gebräuchlichsten Maßzahlen in einer übersichtlichen Grafik (s. Abbildung 2.1).

Die Höhe der „Box“ definiert sich durch den Quartilsabstand, der mittlere Strich markiert den Median und die „Whisker“ markieren den Wertebereich insgesamt – wobei Ausreißer, deren Abstand zur Box mehr als das 1,5-Fache des Quartilsabstands beträgt, üblicherweise gar nicht oder (wie hier) gesondert mit Punkten markiert werden.

Softwarehinweis

In R lässt sich ein Boxplot mit dem Befehl `boxplot()` ausgeben.

Tipps zur Vertiefung

2.4.1 Lagemaße

- Kapitel 2.1 in Bortz und Schuster (2010)
- Kapitel 3.3.2 in Lange und Nipper (2018)
- Kapitel 3.3.1 in Benninghaus (2007)
- Kapitel 4.2.1 in Bahrenberg, Giese und Nipper (2010)
- YouTube-Kanal „Kurzes Tutorium Statistik“: [Arithmetisches, harmonisches und geometrisches Mittel](#)
- YouTube-Kanal „Kurzes Tutorium Statistik“: [Boxplots, Median, Quartile](#)

- *Englisch*: Kapitel 2.2 in Burt und Barber (1996)

2.4.2 Streumaße

- Kapitel 2.2 in Bortz und Schuster (2010)
- Kapitel 3.3.3 in Lange und Nipper (2018)
- Kapitel 3.1.2 in Benninghaus (2007)
- Kapitel 4.2.2 in Bahrenberg, Giese und Nipper (2010)
- YouTube-Kanal „Kurzes Tutorium Statistik“: [Streumaße - Varianz, Standardabweichung, Variationskoeffizient und mehr!](#)
- *Englisch*: Kapitel 2.3 in Burt und Barber (1996)

2.4.3 Boxplot

- Kapitel 3.4 in Bortz und Schuster (2010)
- Kapitel 5.3.1 in Lange und Nipper (2018)
- YouTube-Kanal „Kurzes Tutorium Statistik“: [Boxplots, Median, Quartile](#)
- *Englisch*: Kapitel 16.3 in Burt und Barber (1996)

Lösungen der Übungsaufgaben

Sitzung 1

2.4.4 Lösung 1-1

[zur Aufgabenstellung](#)

– keine Musterlösung –

2.4.5 Lösung 1-2

[zur Aufgabenstellung](#)

– keine Musterlösungen –

2.4.6 Lösung 1-3

[zur Aufgabenstellung](#)

| | Variable | Skalenniveau | Variablentyp | Anmerkungen |
|----|---|---------------------|---------------------|------------------------------------|
| a) | Lebensalter in Jahren | Verhältnisskala | diskret | ganze Zahlen vorausgesetzt |
| b) | Regenmenge in mm | Verhältnisskala | stetig | |
| c) | Gütekategorie | Ordinalskala | qualitativ | |
| d) | Passagieraufkommen | Verhältnisskala | diskret | |
| e) | Baujahr | Intervallskala | diskret | |
| f) | Geschwindigkeit in km/h | Verhältnisskala | stetig | bei ganzzahligen Werten |
| g) | Sozialstatus (Unter-, Mittel und Oberschicht) | Ordinalskala | qualitativ | |
| h) | Temperatur in °F | Intervallskala | stetig | |
| i) | Fläche eines Bundeslands in km ² | Verhältnisskala | stetig | |
| j) | Temperatur in K | Verhältnisskala | stetig | 0 K ist ein natürlicher Nullpunkt |
| k) | Einwohnerzahl | Verhältnisskala | diskret | |
| l) | Pegelstand | Intervallskala | stetig | willkürlicher Nullpunkt |
| m) | Staatsangehörigkeit | Nominalskala | qualitativ | |
| n) | Interesse an Statistik (gering bis hoch) | Ordinalskala | qualitativ | |
| o) | Klausurnote | Ordinalskala | qualitativ | wird jedoch oft metrisch behandelt |
| p) | Bodentyp | Nominalskala | qualitativ | |
| q) | Entfernung zum Stadtzentrum in km | Verhältnisskala | stetig | |
| r) | Körpergröße | Verhältnisskala | stetig | |
| s) | Kleidergröße (S bis XXL) | Ordinalskala | qualitativ | |
| t) | Monatliches Nettoeinkommen | Verhältnisskala | stetig | oder diskret für Cent-Beträge |

2.4.7 Lösung 1-4

zur Aufgabenstellung

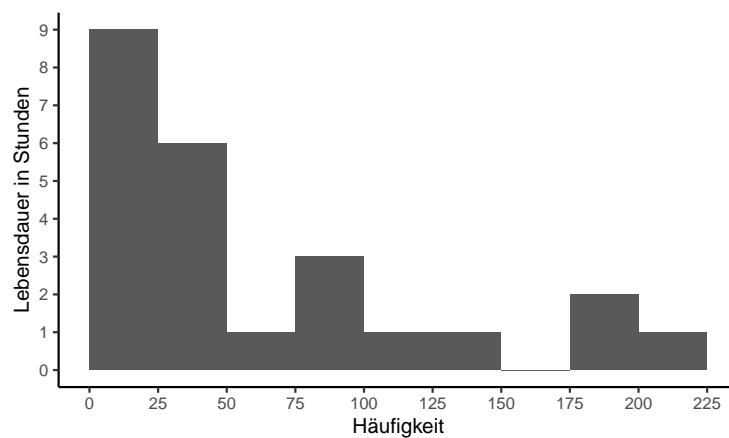
2.4.7.1 a)

Die Werte sind im Bereich zwischen 3 und 210 Stunden. Eine Klassengröße von 25 Stunden bietet sich an, es sind jedoch auch andere Größen denkbar. Da die Variable diskret zu sein scheint, können die Klassengrenzen als ganze Zahlen angegeben werden.

| Wert x_i | Häufigkeit f_i |
|------------------------------|------------------------------------|
| von 0 bis unter 25 h | 9 |
| von 25 bis unter 50 h | 5 |
| von 50 bis unter 75 h | 2 |
| von 75 bis unter 100 h | 3 |
| von 100 bis unter 125 h | 1 |
| von 125 bis unter 150 h | 1 |
| von 150 bis unter 175 h | 0 |
| von 175 bis unter 200 h | 2 |
| von 200 bis unter 225 h | 1 |

2.4.7.2 b)

Das Resultat sollte je nach gewählter Klassengröße in etwa so aussehen:

**2.4.7.3 c)**

Die Verteilung ist unregelmäßig abfallend.

2.4.8 Lösung 1-5

[zur Aufgabenstellung](#)

Sind die folgenden Aussagen wahr oder unwahr?

- a) wahr
- b) wahr
- c) unwahr
- d) wahr
- e) unwahr
- f) unwahr
- g) wahr
- h) wahr
- i) unwahr
- j) unwahr
- k) wahr
- l) wahr
- m) unwahr
- n) unwahr
- o) unwahr
- p) wahr
- q) wahr
- r) wahr

Quellenverzeichnis

- Bahrenberg, Gerhard, Ernst Giese und Josef Nipper. 2010. *Statistische Methoden in der Geographie*. Bd. 1. Univariate und bivariate Statistik. Stuttgart: Bornträger.
- Benninghaus, Hans. 2007. *Deskriptive Statistik. Eine Einführung für Sozialwissenschaftler*. Wiesbaden: VS Verlag.
- Bortz, Jürgen und Christof Schuster. 2010. *Statistik für Human- und Sozialwissenschaftler*. Berlin: Springer.
- Burt, James E. und Gerald M. Barber. 1996. *Elementary statistics for geographers*. 2nd ed. New York: Guilford Press.
- Haseloff, Otto W., Hans-Joachim Hoffmann, John H. Maindonald und W. John Braun. 1968. *Kleines Lehrbuch der Statistik DAAG. Data Analysis and Graphics Data and Functions*. Berlin: de Gruyter.
- Lange, Norbert de und Josef Nipper. 2018. *Quantitative Methodik in der Geographie*. UTB Geographie, Methoden, Statistische Verfahren 4933. Paderborn: Ferdinand Schöningh.
- R Core Team. 2018. R: A Language and Environment for Statistical Computing. Wien: R Foundation for Statistical Computing. <https://www.R-project.org/> (zugegriffen: 9. April 2021).
- Zimmermann-Janschitz, Susanne. 2014. *Statistik in der Geographie. Eine Exkursion durch die deskriptive Statistik*. Berlin: Springer.