

# Deskriptive Statistik (eindimensional)

Peter Büchel

HSLU I

ASTAT: Block 02

- Daten und Statistiken bestimmen immer mehr unser Leben
- Zeitung: Prognose zur nächsten Abstimmung oder Wahlen  
→ Befragung
- Googeln: Wie „weiss“ Google so genau, was man suchen will?  
→ Google wertet Suchanfragen aus
- Passkontrolle am Flughafen: Wie erkennt die Software Gesichter?  
→ Gesichter werden charakterisiert
- Wetterbericht: Wie kommt die Vorhersage zustande?  
→ Modell aufgrund früherer Wetterdaten (und Theorie)
- Börsenkurse: Wie lässt sich aus dem Börsenverlauf der letzten paar Tage, der Kurs für die nächsten paar Tage vorhersagen?  
→ Modellierung aus alten Daten

# Datensätze (eindimensional)

- *Liste*: Einfachste Variante eines Datensatzes
- Bsp: Körpergrössen von 5 Personen

1.75,      1.80,      1.72,      1.65,      1.54

- Solche Listen heissen: *Eindimensionale Datensätze* oder *Messreihen*

# Datensätze (zweidimensional)

- Häufigste Form von Datensätzen: *Tabellen* oder *zweidimensionale Datensätze*
- Bsp:

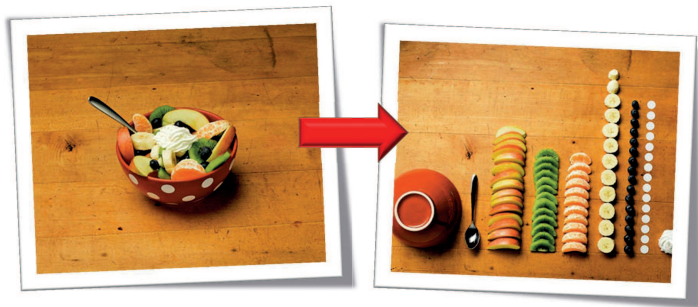
Person	Grösse	Gewicht	Geschlecht	Nationalität
A	1.82	72	m	CH
B	1.75	82	w	D
C	1.61	70	w	CH
D	1.80	83	m	A
E	1.89	95	w	FL

- Grösse und Gewicht: *Quantitative* Daten, also (gemessene) Zahlen
- Können, zumindest theoretisch, jeden beliebigen Zahlwert in einem Bereich annehmen
- Geschlecht und Nationalität: *Qualitative* Daten
- Nehmen nur bestimmte Anzahl Werte an (müssen keine Zahlen sein)

- Deskriptive Statistik: Darstellung von Datensätzen
- Datensätze
  - ▶ durch gewisse Zahlen charakterisieren (z.B. Mittelwert)
  - ▶ *und* graphisch darstellen
- Zunächst *eindimensionale* Daten: *Eine* Messgrösse wird an einem Untersuchungsobjekt ermittelt (zweidimensionale später)

# Ziele der Deskriptiven Statistik

- Daten *zusammenfassen* durch *numerische Kennwerte*
- *Graphische Darstellung* der Daten



# Beispieldatensatz

- Messungen Körpergewicht
- Erfahrung: Steht am Morgen auf Waage und merkt sich Gewicht
- Steht nochmals auf die Waage und erhält leicht anderes Resultat
- Wir wollen es genauer wissen
- Nehmen 80 Kilogramm schweren Metallblock, der geeicht ist, d.h. er hat mit sehr grosser Genauigkeit 80 kg
- Gewicht dieses Metallblocks wird mehrere Male mit zwei Waagen  $A$  und  $B$  gemessen
- Zwei *Datensätze* mit Gewichten (in kg; auf 10 g genau gemessen)

- Tabelle:

Waage A	79.98	80.04	80.02	80.04	80.03	80.03	80.04	79.97	80.05
Waage A	80.03	80.02	80.00	80.02					
Waage B	80.02	79.94	79.98	79.97	79.97	80.03	79.95	79.97	

- Frage: Warum führen verschiedene Messungen, die am gleichen Objekt stattfinden zu unterschiedliche Resultaten?
- Messungen finden nie unter *exakt* denselben Bedingungen statt
- Scheinbar genaue Angaben sind nur *ungefähre* Angaben
  - ▶ Kalorienzahl auf einer Packung Schokolade
  - ▶ Inhalt 500 ml Pet-Flasche: Keine zwei Pet-Flaschen sind *absolut* gleich
  - ▶ Gesichtserkennung am Flughafen: Sie haben *nie* denselben Gesichtsausdruck



# Zurück zum Beispiel der Waagen

- Messungen wurden mit grösstmöglichen Sorgfalt durchgeführt
- Trotzdem variieren die Messwerte innerhalb beider Waagen
- Es stellen sich hier nun die folgenden Fragen:
  - ▶ Gibt es einen Unterschied zwischen der Waage *A* und der Waage *B*?
  - ▶ Falls ja, wie können wir diesen Unterschied ermitteln?
- Es fällt auf:
  - ▶ Beide Waagen: Messwerte um 80 herum liegen (sollte auch so sein)
  - ▶ Waage *A*: Nur 2 Werte von 13 *unter* 80
  - ▶ Waage *B*: Nur 2 von 8 Werten *über* 80 liegen
  - ▶ Werte der Waage *A* sind also *eher* grösser als die der Waage *B*

- Was heisst hier aber „eher“?
- Wie kann man die beiden Messreihen miteinander vergleichen?
- Ziel: Messreihen irgendwie *zusammenzufassen*, um die beiden Waagen miteinander vergleichen zu können
- *Deskriptive Statistik* beschäftigt sich damit, auf welche Weisen Daten organisiert und zusammengefasst werden können
- Ziel: Interpretation und darauffolgende statistische Analyse dieser Daten vereinfachen

- Kennzahlen sollen Daten numerisch zusammenfassen und grob charakterisieren
- Bei statistischen Analysen ist es sehr wichtig, nicht einfach blind ein Modell anzupassen oder *ein* statistisches Verfahren anzuwenden
- Daten sollten *immer* mit Hilfe von geeigneten graphischen Mitteln *und* den Kennzahlen dargestellt werden
- Nur auf diese Weise kann man (teils unerwartete) Strukturen und Besonderheiten entdecken

- Aber:

### **Warnung!!!**

Wann immer ein Datensatz „reduziert“ wird (durch Kennzahlen oder Graphiken), geht *Information verloren!*

- Bsp: Noten einer Schulklasse mit 24 Lernenden an einer Prüfung:  
4.2, 2.3, 5.6, 4.5, 4.8, 3.9, 5.9, 2.4, 5.9, 6, 4, 3.7, 5, 5.2, 4.5, 3.6, 5, 6, 2.8, 3.3, 5.5, 4.2, 4.9, 5.1
- Notendurchschnitt ist 4.51
- Dieser Wert sagt über Klasse als *Ganzes* etwas aus, aber nichts mehr über die einzelnen Lernenden
- Kennen nur Zahl 4.51: Keine Information mehr, wie die einzelnen Lernenden abgeschnitten haben
- Wissen nicht einmal, wieviele Lernende in der Klasse sind

# Bezeichnungen

- Standardbezeichnung von Daten mit

$$x_1, x_2, \dots, x_n$$

- $n$ : *Umfang* der Messreihe (Daten, Datensatz)
- Beispiel: Messreihe der Waage  $A$  hat Umfang  $n = 13$ :

$$x_1 = 79.98, \quad x_2 = 80.04, \quad \dots, \quad x_{13} = 80.02$$

# Kennzahlen

70082950	0.25365536	0.96081524	0.85124829	0.05214026	0.63052716	0.36719205	0.10695418	0.35956536	0.8341956	0.49614412	0.76273099	0.43001
25980996	0.37021603	0.07884733	0.71977404	0.07237495	0.68020504	0.48657579	0.53165132	0.59685485	0.78909487	0.93854889	0.95425422	0.5002
74579848	0.30692408	0.05351679	0.2853162	0.39888676	0.39349628	0.61886139	0.73188697	0.42457447	0.31000296	0.156226	0.50062453	0.4875
82994033	0.83220426	0.9372354	0.73133803	0.96199504	0.55862717	0.32692428	0.61868638	0.56245289	0.71896155	0.34543829	0.75111871	0.1583
92944405	0.64783158	0.60979875	0.52364734	0.26584028	0.40918689	0.16443477	0.25090652	0.04425809	0.06631721	0.45026614	0.96015307	0.5999
1.3322061	0.87182226	0.22334968	0.45692102	0.38131123	0.91921094	0.56080453	0.42412237	0.79812259	0.12081416	0.18896155	0.24489878	0.4241
97712468	0.50452793	0.57458309	0.02272522	0.12008212	0.68844427	0.93512611	0.35232595	0.54222107	0.74300188	0.10069157	0.22498337	0.6473
57467084	0.16038595	0.20683896	0.58934436	0.55401355	0.78000419	0.67956489	0.09056988	0.68952151	0.00707904	0.26790229	0.42494747	0.6355
72574951	0.60798922	0.00653834	0.80803689	0.88663097	0.14771898	0.75301527	0.48470291	0.54921568	0.04009414	0.8453546	0.67167616	0.8958
12893952	0.7431223	0.4202211	0.53911787	0.24420123	0.78464218	0.78235327	0.30197733	0.38276003	0.63617851	0.72978276	0.90730678	0.5484
50684686	0.14058675	0.07426667	0.6377913	0.44437689	0.32789424	0.38075527	0.28287319	0.55515924	0.17444947	0.44069165	0.35637294	0.2464
72021194	0.52889677	0.51331006	0.20434876	0.5249763	0.71545814	0.61285279	0.87822767	0.53536095	0.28884442	0.69949788	0.84420515	0.7418
47268391	0.3610854	0.310148								0.399793	0.71514861	0.55
04257944	0.09101231	0.10635								0.782089	0.04599336	0.9347
33114474	0.80847503	0.589571								0.395522	0.613164	0.0035
17245673	0.67983345	0.231912								0.171166	0.25283066	0.3387
40573334	0.59170081	0.718914								0.88086	0.64948237	0.2252
00561757	0.02425735	0.973367								0.089384	0.00563944	0.31221
82481867	0.18901555	0.627044								0.409241	0.29417144	0.4912
42911629	0.89390795	0.820254								0.6370891	0.15453231	0.8502
15493105	0.51554705	0.81666845	0.33193235	0.110345	0.35500368	0.75014733	0.50944245	0.60935806	0.62794021	0.58346955	0.47319041	0.6518
18653266	0.37671214	0.09282844	0.734327	0.79912816	0.67877946	0.22687246	0.40043241	0.61701288	0.49018961	0.03681597	0.2230552	0.9720
38415242	0.04575544	0.18294704	0.0735783	0.49763891	0.15634616	0.47553336	0.39954434	0.49785766	0.19208229	0.03939701	0.5054817	0.1786
07747484	0.7417904	0.48776921	0.34229175	0.65785054	0.77978943	0.20129577	0.62714576	0.46987345	0.69996167	0.48786104	0.99177657	0.6729
71427139	0.83346645	0.50236663	0.59062007	0.29268677	0.67964115	0.09614286	0.14222698	0.66263698	0.42537685	0.64928539	0.5648649	0.2613
96293853	0.6974188	0.85632265	0.45947964	0.00242453	0.68051404	0.20703925	0.87558209	0.679752	0.45999782	0.8722821	0.04547348	0.8243
04080904	0.5989028	0.87059205	0.12444579	0.26178908	0.8533065	0.20800837	0.90760418	0.06746495	0.61181415	0.37402957	0.36137753	0.8349
1.5616472	0.78210485	0.26718637	0.74856241	0.93690527	0.51338037	0.94582627	0.60380999	0.19747357	0.34424067	0.05237252	0.91349594	0.8796
71333452	0.28822987	0.65203382	0.49709346	0.70379359	0.27200958	0.85341908	0.15968767	0.34960955	0.6796046	0.34255204	0.62727145	0.9353
33192659	0.72932196	0.70036634	0.31364757	0.31615678	0.62072333	0.68964657	0.47503972	0.80823875	0.9708966	0.32082118	0.11199293	0.2306
91696324	0.46608963	0.38554788	0.09440939	0.18995497	0.19254922	0.8299711	0.63238203	0.87524562	0.38170458	0.40120436	0.12882023	0.0850
1.8707509	0.06485663	0.22943682	0.41974316	0.9098332	0.86713599	0.88315761	0.31558244	0.63788522	0.48528904	0.17606219	0.17009773	0.4134
06291977	0.05277628	0.48101212	0.1043349	0.30497809	0.0559275	0.64358846	0.19723847	0.74347764	0.6704249	0.26325428	0.04458277	0.4040
22521559	0.30987268	0.99622375	0.94174692	0.28813039	0.20353298	0.84322955	0.54332297	0.34110065	0.68044315	0.87158643	0.41122531	0.8023

$$\bar{x} = 0.53$$

# Überblick über die Kennzahlen

- Bekannt:  $n$  beobachtete Datenpunkte (Messungen)

$$x_1, x_2, \dots, x_n$$

(z.B. Verkehrsaufkommen an  $n$  verschiedenen Tagen)

- Unterscheidung zwischen Lage- und Streuungsparametern
- *Lageparameter* („Wo liegen die Beobachtungen auf der Mess-Skala?“)
  - ▶ Arithmetisches Mittel („Durchschnitt“)
  - ▶ Median
  - ▶ Quantile
- *Streuungsparameter* („Wie streuen die Daten um ihre mittlere Lage?“)
  - ▶ Empirische Varianz / Standardabweichung
  - ▶ Quartilsdifferenz

# Arithmetisches Mittel

- Umgangsprachlich: *Durchschnitt*
- Addiert alle Daten und teilt Summe durch Anzahl Daten (Umfang)
- Definition:

## Arithmetisches Mittel

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Sprechweise: „x quer“
- Beispiel Waage A: Arithmetische Mittel der  $n = 13$  Messungen

$$\bar{x} = \frac{79.98 + 80.04 + \dots + 80.03 + 80.02 + 80.00 + 80.02}{13} = 80.020\,77$$



- R-Befehl für arithmetisches Mittel `mean(...)`:

```
waageA <- c(79.98, 80.04, 80.02, 80.04, 80.03, 80.03, 80.04,  
            79.97, 80.05, 80.03, 80.02, 80.00, 80.02)  
  
mean(waageA)  
## [1] 80.02077
```

- Arithmetisches Mittel: Anschaulich



# Vergleich der Waagen

- Waagen mit dem arithmetischen Mittel miteinander vergleichen:

```
waageB <- c(80.02, 79.94, 79.98, 79.97, 79.97, 80.03, 79.95, 79.97)

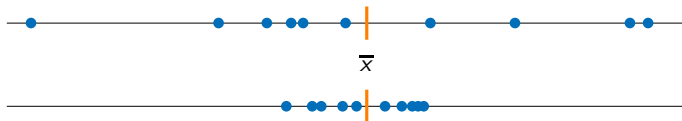
mean(waageB)

## [1] 79.97875
```

- Waage *B* hat durchschnittlich die tieferen Werte als die Waage *A*
- Waagen vergleichbar, obwohl verschiedener Umfang

# Streuung graphisch

- Arithmetisches Mittel sagt schon einiges über einen Datensatz aus
- Arithmetisches Mittel : „Wo ist die „Mitte“ der Daten?“
- Aber: Arithmetisches Mittel sagt nicht alles über Messreihe aus
- Graphisches Beispiel:



- Beide Datenreihen haben dasselbe arithmetische Mittel
- Punkte 2. Datenreihe: Durchschnittlich viel näher beim Mittelwert  $\bar{x}$  als Punkte der 1. Datenreihe
- Verschiedene *Streuung* der Daten um den Durchschnitt

# Streuung numerisch

- Beispiel von (fiktiven) Schulnoten:

2, 6, 3, 5      und      4, 4, 4, 4

- Beide Mittelwert 4, aber Verteilung der Daten um Mittelwert ziemlich unterschiedlich
  - ▶ 1. Fall: Zwei gute und zwei schlechte Lernende
  - ▶ 2. Fall: Alle Lernende gleich gut
- Datensätze haben eine verschiedene *Streuung* um den Mittelwert
- 2. Klasse: Keine Streuung (Streuung 0)
- Wie kann man diese Streuung *mathematisch* beschreiben?
- Suchen also Mass für den Unterschied zum Mittelwert

- 1. Idee: Durchschnitt der *Unterschiede zum Mittelwert*

- 1. Fall:

$$\frac{(2 - 4) + (6 - 4) + (3 - 4) + (5 - 4)}{4} = \frac{-2 + 2 - 1 + 1}{4} = \frac{0}{4} = 0$$

- Zweiter Fall auch 0:

$$\frac{(4 - 4) + (4 - 4) + (4 - 4) + (4 - 4)}{4} = \frac{0 + 0 + 0 + 0}{4} = \frac{0}{4} = 0$$

- Dieser Ausdruck macht keine Aussage über unterschiedliche Streuung
- Problem: Unterschiede können *negativ* werden, können sich aufheben

- Nächste Idee: Unterschiede durch die *Absolutwerte* ersetzen

- 1. Fall:

$$\frac{|(2 - 4)| + |(6 - 4)| + |(3 - 4)| + |(5 - 4)|}{4} = \frac{2 + 2 + 1 + 1}{4} = 1.5$$

- D.h.: Noten weichen im Schnitt 1.5 vom Mittelwert ab

- 2. Fall: Dieser Wert natürlich auch 0

$$\frac{|(4 - 4)| + |(4 - 4)| + |(4 - 4)| + |(4 - 4)|}{4} = \frac{0 + 0 + 0 + 0}{4} = 0$$

- Je grösser dieser Wert (immer grösser gleich 0) , desto mehr unterscheiden sich die Daten bei gleichem Mittelwert voneinander

- Dieser Wert für die Streuung: *Mittlere absolute Abweichung*

- Aber: Theoretische Nachteile

# Empirische Varianz und Standardabweichung

- Besser: *Empirische Varianz* und *empirische Standardabweichung*
- Mass für *Variabilität* oder *Streuung* der Messwerte
- Definition:

## Empirische Varianz $\text{var}(x)$ und Standardabweichung $s_x$

$$\text{Var}(x) = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1} = \frac{1}{n - 1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$s_x = \sqrt{\text{Var}(x)} = \sqrt{\frac{1}{n - 1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

# Eigenschaften der Varianz

- Bei Varianz: Abweichungen  $x_i - \bar{x}$  quadrieren, damit sich Abweichungen nicht gegenseitig aufheben können
- Nenner  $n - 1$ , anstelle von  $n$ : Mathematisch begründet
- In einigen Büchern:  $n$  anstatt  $n - 1$ : Spielt für grosse  $n$  keine Rolle (siehe Jupyter Notebook [varianz\\_n.ipynb](#))
- Standardabweichung ist die Wurzel der Varianz
- Durch Wurzelziehen wieder dieselbe Einheit wie bei den Daten selbst:
  - ▶ Daten in cm
  - ▶ Varianz wegen der Quadrierung der Daten in  $\text{cm}^2$
  - ▶ Einheit der ursprünglichen Daten: Aus Varianz Wurzel ziehen



- Ist empirische Varianz (und damit die Standardabweichung) gross, so ist die Streuung der Messwerte um das arithmetische Mittel gross
- Wert der empirischen Varianz hat keine physikalische Bedeutung
- Man weiss nur, je grösser der Wert umso grösser die Streuung
- Wichtig: *Nur die Standardabweichung  $s_x$  lässt sich konkret interpretieren*
- Für normalverteilte Daten hat die Standardabweichung noch eine schöne geometrische Interpretation (siehe später)

## Beispiel: Waage A

- Arith. Mittel der  $n = 13$  Messungen:  $\bar{x} = 80.02$  (siehe Slide 16)
- Empirische Varianz:

$$\begin{aligned}\text{Var}(x) &= \frac{(79.98 - 80.02)^2 + (80.04 - 80.02)^2 + \dots + (80.00 - 80.02)^2 + (80.02 - 80.02)^2}{13 - 1} \\ &= 0.000574\end{aligned}$$

- Empirische Standardabweichung:

$$s_x = \sqrt{0.000574} = 0.024$$

- D.h.: „Mittlere“ Abweichung vom Mittelwert 80.02 kg ist 0.024 kg

- Von Hand sehr mühsam

- Mit R:

```
var(waageA)
## [1] 0.000574359
sd(waageA)
## [1] 0.02396579
```

- **sd**: Standard deviation

# Median

- Ein weiteres Lagemass für die „Mitte“: *Median*
- Sehr vereinfacht: Wert, bei dem die Hälfte der Messwerte unter oder gleich diesem Wert sind
- Andere Hälfte ist gleich diesem Messwert oder darüber
- Beispiel: Prüfung in der Schule ist Median 4.6
  - ▶ D.h.: Hälfte der Klasse hat *diese Note oder ist schlechter*
  - ▶ Umgekehrt: Andere Hälfte der Klasse hat *diese Note oder ist besser*
- Obige Interpretation: Medians sehr vereinfacht dargestellt
- Exakte Definition folgt nun

# Geordnete Strichprobe

- Datensatz in aufsteigender Reihenfolge *ordnen*
- Bezeichnung der *geordneten Daten* mit  $x_{(i)}$ :

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

- Runde Klammern im Index: Daten geordnet
- Bestimmung *Median*: Daten zuerst der Grösse nach ordnen

# Beispiel Waage A

- Daten Waage A geordnet:

79.97, 79.98, 80.00; 80.02, 80.02, 80.02, 80.03, 80.03, 80.03, 80.04, 80.04, 80.04, 80.05

- Median ist nun sehr einfach zu bestimmen
- Unter diesen 13 Messungen: Wert der mittleren Beobachtung
- Wert der 7. Beobachtung:

79.97, 79.98, 80.00; 80.02, 80.02, 80.02, 80.03, 80.03, 80.03, 80.04, 80.04, 80.04, 80.05

- Median des Datensatzes der Waage  $A$  ist 80.03
- D.h.: Knapp die Hälfte der Messwerte, nämlich 6 Beobachtungen sind kleiner oder gleich 80.03
- Ebenso sind 6 Messwerte grösser oder gleich dem Median
- Ungerade Anzahl Messungen: *Genau* eine mittlere Messung

## Beispiel Waage B

- Vorher: Anzahl der Daten ungerade und damit ist die mittlere Beobachtung eindeutig bestimmt
- Anzahl Daten gerade: *Keine* mittlere Beobachtung
- *Definition* Median: Mittelwert der beiden mittleren Beobachtungen
- Beispiel: Datensatz der Waage B hat 8 Beobachtungen
- Ordnen den Datensatz: Median Durchschnitt von der 4. und 5. Beobachtung

79.94, 79.95, 79.97, 79.97, 79.97, 79.94, 80.02, 80.03

$$\frac{79.97 + 79.97}{2} = 79.97$$



- R-Befehl:

```
median(waageA)
## [1] 80.03
waageB <- c(80.02, 79.94, 79.98, 79.97, 79.97, 80.03, 79.95, 79.97)
median(waageB)
## [1] 79.97
```

- Als Median kann Wert auftreten, der in Messreihe nicht ist
- Beispiel: Median Schulklasse:

```
noten <- c(4.2, 2.3, 5.6, 4.5, 4.8, 3.9, 5.9, 2.4, 5.9, 6, 4, 3.7,
           5, 5.2, 4.5, 3.6, 5, 6, 2.8, 3.3, 5.5, 4.2, 4.9, 5.1)
median(noten)
## [1] 4.65
```

- Aber: Noten auf Zehntelnoten genau

# Median vs. arithmetisches Mittel

- Zwei Lagemasse für die Mitte eines Datensatzes
- Welches ist nun „besser“?
- Dies kann man so nicht sagen:
  - ▶ Kommt auf die jeweilige Problemstellung an
  - ▶ Am besten werden beide Masse gleichzeitig betrachtet
- Eigenschaft des Medians: *Robustheit*
- Das heisst: Wird viel weniger stark durch extreme Beobachtungen beeinflusst als das arithmetische Mittel

# Median vs. arithmetisches Mittel

- Beispiel: Bei der grössten Beobachtung ( $x_9 = 80.05$ ) ist ein Tippfehler passiert und  $x_9^* = 800.5$  eingegeben worden
- Arithmetische Mittel ist dann anstatt 80.02:

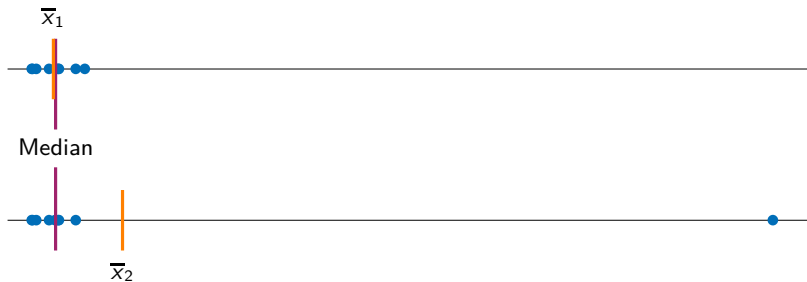
$$\bar{x}^* = 135.44$$

- Median ist aber nach wie vor

$$x_{(7)}^* = x_{(7)} = 80.03$$

- Arithmetisches Mittel: Durch Veränderung *einer* Beobachtung sehr stark beeinflusst
- Median bleibt hier gleich: Robust

# Graphisch



# Beispiel

- Untersuchen typisches Haushaltseinkommen von Vororten von Seattle um Lake Washington
- Durchschnittliches Einkommen von Medina und Windermere wird sehr unterschiedlich sein
- Grund: Bill Gates lebt in Medina
- Grundsätzlich: Für das mittlere Einkommen wird praktisch immer der Median genommen und nicht der Durchschnitt, da dies gerechter ist
- Zügelt Bill Gates von Medina nach Windermere, so ändert sich das Durchschnittseinkommen, aber niemand hat was davon

# Bemerkung

- Median: Auch *Zentralwert* oder *mittlerer Wert* (nicht zu verwechseln mit dem Mittelwert) genannt
- Exakte Interpretation des Medians ist noch erstaunlich schwierig
- Für uns ausreichend: Hälfte der Werte kleiner oder gleich und die andere Hälfte grösser oder gleich dem Median

# Quartile

- Median: Wert, wo die Hälfte der Beobachtungen kleiner (oder gleich) wie dieser Wert sind
- Analoge Überlegung: Unteres und oberes Quartil
- Unteres Quartil: Wert, wo 25 % aller Beobachtungen kleiner oder gleich und 75 % grösser oder gleich sind wie dieser Wert
- Oberes Quartil: Wert, wo 75 % aller Beobachtungen kleiner oder gleich und 25 % grösser oder gleich wie dieser Wert sind
- Achtung: Meist gibt es nicht *exakt* 25 % der Beobachtungen
- Man *definiert* Wert für das untere Quartil bzw. obere Quartil

## Beispiel: Waage

- Waage A hat  $n = 13$  Messpunkte: 25 % davon ist 3.25
- Man *wählt* nächstgrösseren Wert  $x_{(4)}$  als unteres Quartil:

79.97, 79.98, 80.00, 80.02, 80.02, 80.02, 80.03, 80.03, 80.03, 80.04, 80.04, 80.04, 80.05

- Unteres Quartil ist 80.02
- Knapp ein Viertel der Messwerte ist gleich oder kleiner 80.02
- Oberes Quartil: Wählen  $x_{(10)}$ , da für  $0.75 \cdot 13 = 9.75$  die Zahl 10 der nächsthöhere Wert ist

79.97, 79.98, 80.00, 80.02, 80.02, 80.02, 80.03, 80.03, 80.03, 80.04, 80.04, 80.04, 80.05

- Knapp drei Viertel der Messwerte sind kleiner oder gleich 80.04



- Waage B: 25 % der Werte 2: ganze Zahl
- Man *wählt* Durchschnitt von  $x_{(2)}$  und  $x_{(3)}$  als unteres Quartil
- Dann sind 2 Beobachtungen kleiner und 6 Beobachtungen grösser als dieser Wert

79.94, 79.95, 79.97, 79.97, 79.97, 79.94, 80.02, 80.03

$$\frac{79.95 + 79.97}{2} = 79.96$$

- Das untere Quartil der Methode B ist also 79.96

# Bemerkungen

- Hier jeweils aufgerundet, falls 25 % bzw. 75 % der Anzahl Beobachtungen nicht ganz ist
- Hätten auch *abrunden* können: Andere Werte für die Quartile
- Für grosse Datensätze: Spielt praktisch keine Rolle, ob aufgerundet, abgerundet oder gerundet wird
- Aber: Es gibt keine einheitliche Definition für die Quartile

- Software **R** kennt keine eigenen Befehle für die Quartile
- Allgemeinerer Befehl **quantile** (Quantile kommen gleich)
- **R**: Quartile nach unserer Definition: Option **type=2**

```
# Syntax für das untere Quartil: p=0.25
```

```
quantile(waageA, p = 0.25, type = 2)
```

```
##    25%
```

```
## 80.02
```

```
quantile(waageB, p = 0.25, type = 2)
```

```
##    25%
```

```
## 79.96
```

```
# Syntax für das obere Quartil: p=0.75
```

```
quantile(waageA, p = 0.75, type = 2)
```

```
##    75%
```

```
## 80.04
```

# Quartilsdifferenz

- *Quartilsdifferenz* ist ein Streuungsmass für die Daten  
oberes Quartil – unteres Quartil
- Es misst die Länge des Intervalls, das etwa die Hälfte der „mittleren“ Beobachtungen enthält
- Je kleiner dieses Mass, umso näher liegt die Hälfte aller Werte um den Median und umso kleiner ist die Streuung
- Dieses Streuungsmass ist robust
- Quartilsdifferenz der Methode A

$$80.04 - 80.02 = 0.02$$

- R: Quartilsdifferenz (interquartile range) **IQR**

```
IQR(waageA, type=2)  
## [1] 0.02
```

- Die (oder ungefähr die) Hälfte der Messwerte liegt also in einem Bereich der Länge 0.02
- Beim Boxplot: Anschauliche Interpretation der Quartilsdifferenz

# Beispiel

- Schulklasse mit 24 Lernenden: Noten an einer Prüfung:

4.2, 2.3, 5.6, 4.5, 4.8, 3.9, 5.9, 2.4, 5.9, 6, 4, 3.7, 5, 5.2, 4.5, 3.6, 5, 6, 2.8, 3.3, 5.5, 4.2, 4.9, 5.1

- Berechnen mit R die Quartile und die Quartilsdifferenz:

```
quantile(noten, p = c(0.25, 0.75), type = 2)
## 25% 75%
## 3.80 5.35
IQR(noten, type = 2)
## [1] 1.55
```

- Hälfte der Lernenden liegen innerhalb von 1.55 Noten, nämlich zwischen 3.8 und 5.35
- 25 % der Klasse 3.8 oder weniger; rund 25 % der Klasse 5.35 und mehr

# Quantile

- Quartile auf jede andere Prozentzahl verallgemeinern: Quantile
- 10 %-Quantil: Wert, wo 10 % der Werte kleiner oder gleich und 90 % der Werte grösser oder gleich diesem Wert sind
- Definition analog wie bei Quartilen
- Median ist 50 %-Quantil
- 25 %-Quantil ist unteres Quartil
- 75 %-Quantil ist oberes Quartil

- R: 10 %- und 70 %-Quantil der Waage A:

```
quantile(waageA, p = .1, type = 2)
##    10%
## 79.98

quantile(waageA, p = .7, type = 2)
##    70%
## 80.04
```

- Knapp 10 % der Messwerte sind kleiner oder gleich 79.97
- Entsprechend: Knapp 70 % der Messwerte kleiner oder gleich 80.04



# Beispiel

- Noten an Prüfung in Schulklasse mit 24 Lernenden:

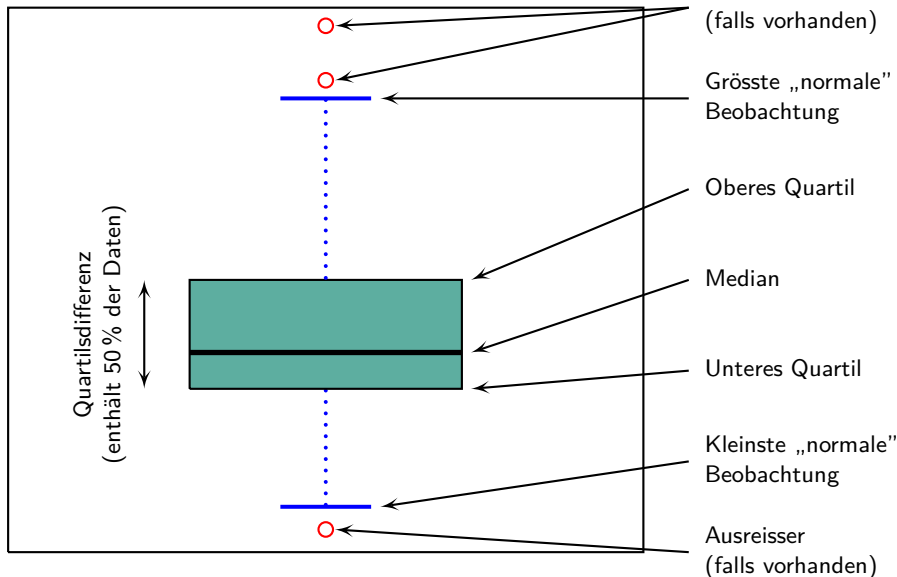
4.2, 2.3, 5.6, 4.5, 4.8, 3.9, 5.9, 2.4, 5.9, 6, 4, 3.7, 5, 5.2, 4.5, 3.6, 5, 6, 2.8, 3.3, 5.5, 4.2, 4.9, 5.1

- Verschiedene Quantile mit R:

```
quantile(noten, p = seq(from = 0.2, to = 1, by = 0.2), type = 2)
## 20% 40% 60% 80% 100%
## 3.6 4.2 5.0 5.6 6.0
```

- D.h.: Knapp 20 % der Lernenden sind schlechter als 3.6
- Genau 20 % der Lernenden nicht möglich: 4.8 Lernende
- 60 %-Quantil: (Knapp) diese Anzahl Prozent der Lernenden waren schlechter oder gleich einer 5

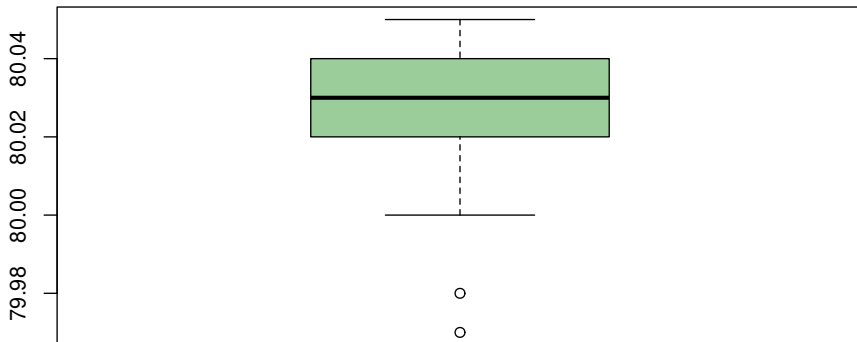
# Boxplot: Schematischer Aufbau



# Boxplot: Schematischer Aufbau

- Rechteck, dessen Höhe vom empirischen 25 %- und vom 75 %-Quantil begrenzt wird: Box
- Horizontaler Strich für den Median in Box (schwarz)
- Linien, die von diesem Rechteck bis zum kleinsten- bzw. grössten „normalen“ Wert führen (blau eingezeichnet)
  - ▶ Definition: „Normaler“ Wert höchstens 1.5 mal die Quartilsdifferenz von einem der beiden Quartile entfernt
- Ausreisser: Kleine Kreise (rot)

```
boxplot(waageA,  
        col = "darkseagreen3"  
)
```

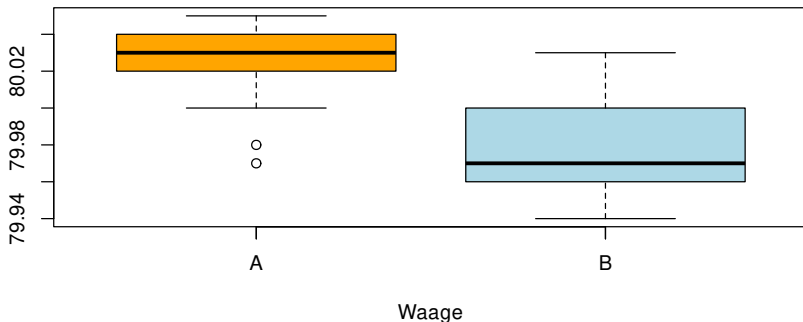


- Hälfte der Beobachtungen befindet sich zwischen dem oberen Quartil 80.04 und dem unteren Quartil 80.02, mit Quartilsdifferenz 0.02
- Median liegt bei 80.03
- Der „normale“ Bereich der Werte liegt zwischen 80.00 und 80.05
- Zwei Ausreisser: 79.97 und 79.98
- Erste beiden Punkte: Früher berechnet
- Boxplot: Graphische Darstellung von Median und Quartilen

# Vergleich von Datensätzen

- Boxplot: Darstellungen von verschiedenen Gruppen

```
boxplot(waageA, waageB,  
        xlab = "Waage",  
        col = c("orange", "lightblue")  
)  
axis(side = 1, at = c(1, 2), labels = c("A", "B"))
```



- Waage  $A$  grössere Werte als Waage  $B$ : Median von  $A$  grösser ist
- Daten von Waage  $A$  haben weniger Streuung als die Daten von Waage  $B$ : Rechteck weniger hoch (Quartilsdifferenz!)
- R-Code `axis(...)`: Siehe Skript