

Mathe III

Explorative Datenanalyse

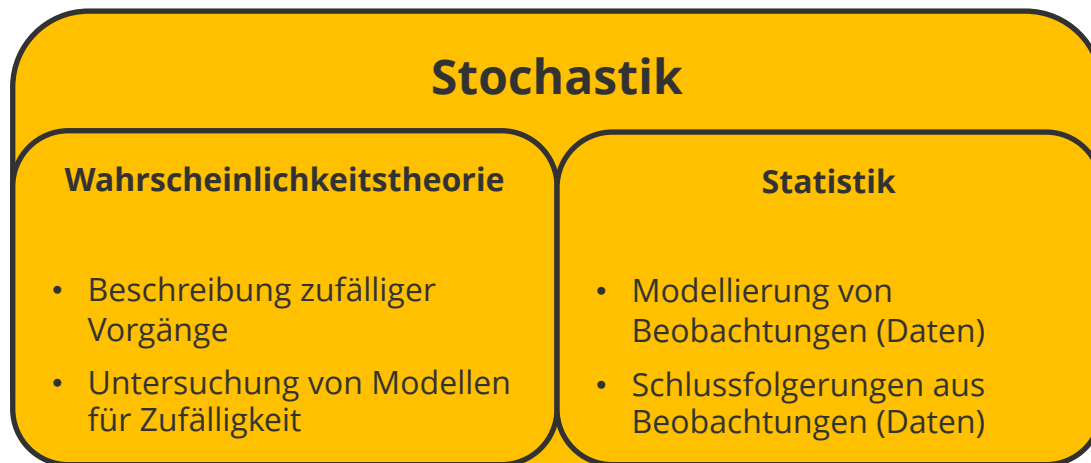
Ralf Herbrich

1. Empirische Daten
2. Visualisierung empirischer Daten

1. **Empirische Daten**
2. Visualisierung empirischer Daten

Wahrscheinlichkeitstheorie vs. Statistik

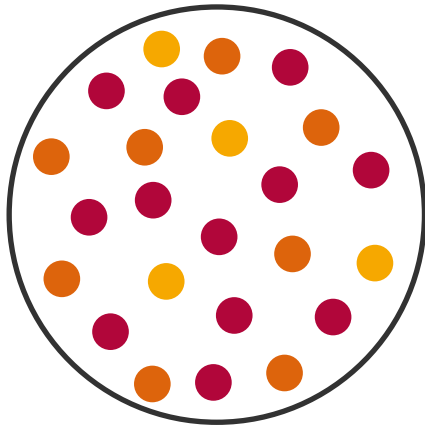
- **Wahrscheinlichkeitstheorie:** Wir haben uns bis jetzt hauptsächlich mit den Grundlagen der Wahrscheinlichkeitstheorie beschäftigt.
- **Statistik:** Statt nur Modelle zu beschreiben, wollen wir in der Statistik nun auch lernen, wie man reale Zufallsversuche, für die kein Modell bekannt ist, beschreiben kann.



Mathe III

Unit 8 –
Explorative Datenanalyse

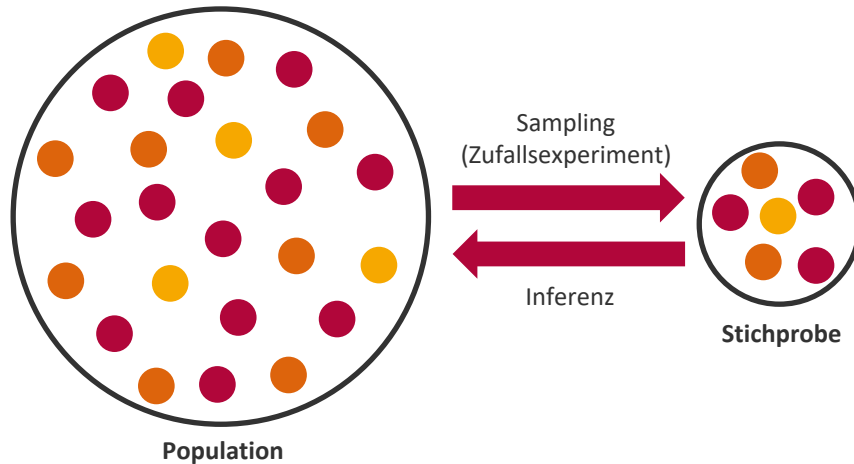
Population, Stichprobe und Modell



Population

- Eine Population beschreibt eine **Grundgesamtheit** von Objekten, welche **Merkmale** tragen.
- Wir wollen Aussagen über die Merkmalsausprägung in der Population treffen.
- In der Regel ist die Population zu groß, um sie vollständig zu vermessen.
- Daher möchte man zumindest eine „zuversichtliche“ Aussage treffen können.

Population, Stichprobe und Modell

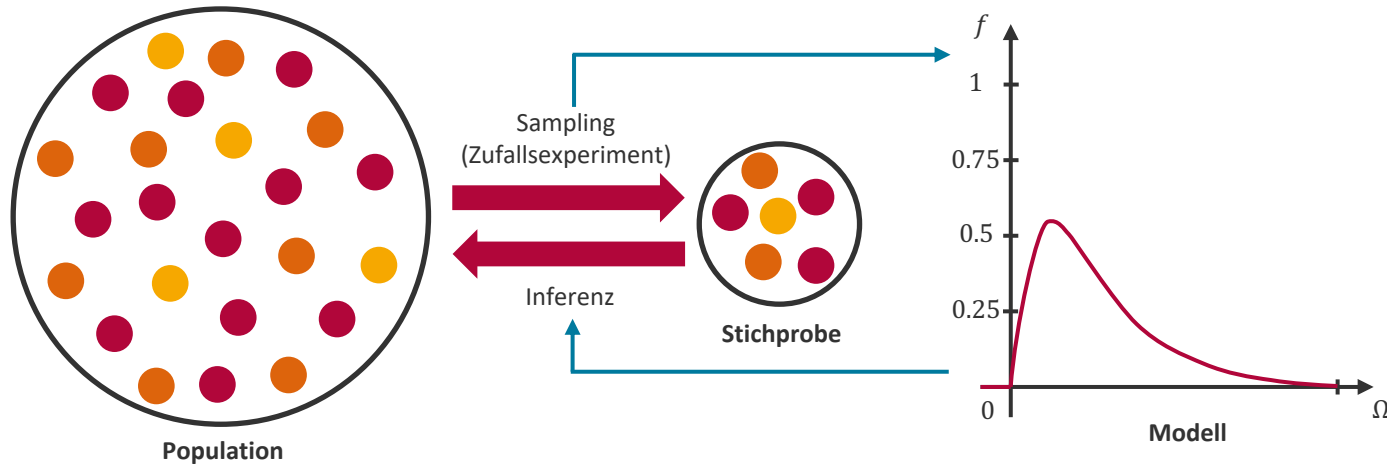


- Die **Stichprobe** ist eine Auswahl von Objekten aus einer Population.
- Die Auswahl der Stichprobe ist ein **Zufallsexperiment**.
- Wir können die Objekte der Stichprobe mit Kennzahlen charakterisieren (Deskriptive Statistik)
- Wir wollen damit Aussagen über die Population treffen (**Inferenz**).

Mathe III

Unit 8 –
Explorative Datenanalyse

Population, Stichprobe und Modell

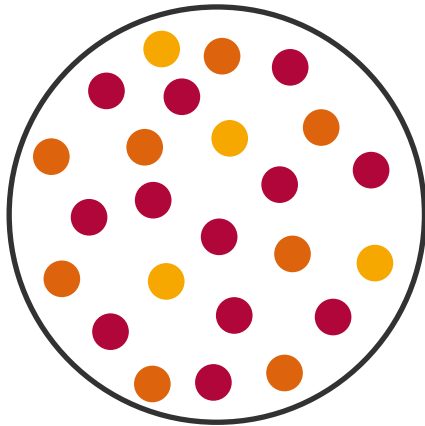


- Das **stochastische Modell** beschreibt das Zufallsexperiment der Stichprobenwahl.
- Wir nutzen dieses, um von Stichprobenmerkmalen auf die Populationsmerkmale zu schließen.
- Beim Aufstellen eines Modells trifft man Annahmen – stimmen diese nicht, sind die Ergebnisse hinfällig.

Mathe III

Unit 8 –
Explorative Datenanalyse

Probleme von Stichproben



Population



Repräsentative
Stichprobe



Nicht-repräsentative
Stichprobe

- Damit die Eigenschaften einer Stichprobe generalisierbar auf die Population sind, muss sie **repräsentativ** sein.
- Nicht-repräsentative Stichproben besitzen nur eine eingeschränkte Aussagekraft.

Mathe III

Unit 8 –
Explorative Datenanalyse

- **Definition (Mittelwert).** Für eine Sequenz x von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ ist der **Mittelwert** \bar{x} definiert als

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

- **Beispiel (Mittelwert).** Wir betrachten die Neupreise einer zufälligen Stichprobe von 4 Neuwagen:

x	Gezogenes Modell	Neupreis in 1000 €
x_1	BMW i5	72
x_2	smart #1	42
x_3	VW ID.3	40
x_4	Mercedes EQA	50



Mathe III

Unit 8 –
Explorative Datenanalyse

Hier beträgt der Mittelwert

$$\bar{x} = \frac{1}{4} \cdot (72 + 42 + 40 + 50) = \frac{204}{4} = 51$$

Statistiken: Quantile und Median

- **Definition (Quantil).** Für eine Sequenz x von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ und ein $p \in (0, 1)$ bezeichnen wir die Beobachtung an der Position $[n \cdot p]$ nach Sortieren der Werte x_1, \dots, x_n als p -Quantil.
- **Definition (Median).** Für eine Sequenz x von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ bezeichnen wir das $\frac{1}{2}$ -Quantil als Median \tilde{x} von x .
- **Beispiel (Median und Quantil).** Sei $x = (1, 2, 3, 4, 5)$.
 - Der Median \tilde{x} von x ist 3.
 - Das 0.25-Quantil von x ist 2.
 - Das 0.99-Quantil von x ist 5.

Statistiken: Quantile und Median

■ Bemerkungen (Quantile)

- Es existieren verschiedene Definitionen des Quantils, z. B. als arithmetisches Mittel zwischen dem größten Wert der $[n \cdot p]$ kleinsten Werte und dem kleinsten Wert der $[n \cdot (1 - p)]$ größten Werte.
- Die von uns verwendete Definition wird auch Untermedian genannt.
- Während der Mittelwert **ausreißer anfällig** ist, sind Quantile **robust**.

■ Beispiel (Robustheit des Medians). Sei $\mathbf{x} = (1, 2, 3, 4, 5)$ und $\mathbf{y} = (1, 2, 3, 4, 100)$

- Für die Mittelwerte der Stichproben gilt $\bar{x} = 3$ und $\bar{y} = 22$.
- Der Median für beide Stichproben hingegen ist unverändert 3.
- Der Ausreißer in \mathbf{y} hat demnach starken Einfluss auf den Mittelwert, jedoch kaum (bzw. hier keinen) Einfluss auf den Median.

Statistiken: *Median Absolute Deviation*

- **Definition (*Median Absolute Deviation – MAD*)**. Für eine Sequenz x von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ ist die *median absolute deviation* (MAD) definiert als

$$\text{MAD}(x) = \text{median}(|x_1 - \tilde{x}|, |x_2 - \tilde{x}|, \dots, |x_n - \tilde{x}|)$$

- **Beispiel (*Median Absolute Deviation*)**. Sei $x = (1, 2, 3, 4, 5)$.
 - Der Median \tilde{x} von x ist 3.
 - Die absoluten Abweichungen zum Median betragen (2, 1, 0, 1, 2).
 - Folglich gilt $\text{MAD}(x) = 1$.
- **Bemerkungen (*Median Absolute Deviation*)**
 - Der MAD gilt ebenso wie der Median als besonders robustes Maß einer Stichprobe.
 - Es gibt andere Definitionen für „MAD“, welche stattdessen den Mittelwert (Mean) nutzen – hier besteht Verwechslungsgefahr!

$$\overline{(|x_i - \bar{x}|)}$$

Statistiken: Empirische Varianz, Kovarianz und Korrelation

- **Definition (Empirische Varianz und Standardabweichung).** Für eine Sequenz \mathbf{x} von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ heißt

$$V[\mathbf{x}] = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{\mathbf{x}})^2$$

die empirische Varianz von \mathbf{x} , und $\sqrt{V[\mathbf{x}]}$ die empirische Standardabweichung von \mathbf{x} .

- **Definition (Empirische Kovarianz und Korrelation).** Für die Sequenzen \mathbf{x} und \mathbf{y} von Beobachtungen $x_1, \dots, x_n \in \mathbb{R}$ und $y_1, \dots, y_n \in \mathbb{R}$ heißt

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{\mathbf{x}}) \cdot (y_i - \bar{\mathbf{y}})$$

die empirische Kovarianz von \mathbf{x} und \mathbf{y} . Die empirische Korrelation ergibt sich als

$$\text{Corr}[\mathbf{x}, \mathbf{y}] = \frac{\text{Cov}[\mathbf{x}, \mathbf{y}]}{\sqrt{V[\mathbf{x}] \cdot V[\mathbf{y}]}}$$

Mathe III

Unit 8 –
Explorative Datenanalyse

- **Beispiel (Empirische Korrelation).** Wir betrachten die Neupreise und Reichweiten einer zufälligen Stichprobe \mathbf{x} und \mathbf{y} von 4 Neuwagen:

Gezogenes Modell	Neupreis in 1000 €	Reichweite in km	$(x_i - \bar{x})$	$(y_i - \bar{y})$
BMW i5	$x_1 = 72$	$y_1 = 570$	21	45
smart #1	$x_2 = 42$	$y_2 = 440$	-9	-85
VW ID.3	$x_3 = 40$	$y_3 = 560$	-11	35
Mercedes EQA	$x_4 = 50$	$y_4 = 530$	-1	5



- Es gilt $\bar{x} = \frac{1}{4} \cdot (72 + 42 + 40 + 50) = 51$ und $\bar{y} = \frac{1}{4} \cdot (570 + 440 + 560 + 530) = 525$.

- Daher gilt

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \frac{1}{4} \cdot (21 \cdot 45 + (-9) \cdot (-85) + (-11) \cdot 35 + (-1) \cdot 5) = 330$$

- Es gilt $V[\mathbf{x}] = \frac{1}{4} \cdot (21^2 + 9^2 + 11^2 + 1^2) = 161$ und $V[\mathbf{y}] = \frac{1}{4} \cdot (45^2 + 85^2 + 35^2 + 5^2) = 2625$.

- Daher gilt

$$\text{Corr}[\mathbf{x}, \mathbf{y}] = \frac{330}{\sqrt{161 \cdot 2625}} \approx 0.5$$

■ Bemerkungen (Empirische Kenngrößen)

- Die neu eingeführten Kenngrößen können genutzt werden, um gezogene Stichproben aus einer Population mit einer unbekannten Verteilung zu beschreiben.
- Die Kenngrößen sind die äquivalenten Momente von ein- und mehr-dimensionalen Zufallsvariablen wenn die **empirische Verteilung** benutzt wird mit

$$p_{\{x_1, \dots, x_n\}}(x) = \frac{|\{x_i \in \{x_1, \dots, x_n\} \mid x_i = x\}|}{n}$$

- Mithilfe dieser Kenngrößen können Rückschlüsse auf die tatsächliche Verteilung getroffen werden, z. B.
 - Mittelwert, Median → Erwartungswert
 - Empirische Varianz → Varianz
 - Empirische Standardabweichung, MAD → Standardabweichung
 - Empirische Kovarianz → Kovarianz
 - Empirische Korrelation → Korrelation

Mathe III

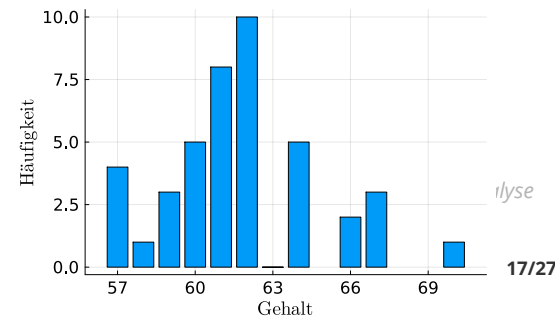
Unit 8 –
Explorative Datenanalyse

1. Empirische Daten
- 2. Visualisierung empirischer Daten**

Stichprobenvisualisierung

- **Problem:** Wie können wir einen **schnellen** Überblick über charakteristische Eigenschaften der Stichprobe bekommen?
 - Mittelwert
 - Streuung
 - Symmetrie oder Schiefe
 - Ausreißer
 - Verteilungsannahmen
- **Lösung:** Mittels Visualisierung der Stichprobe („ein Bild sagt mehr als tausend Worte“)
 - Schnelle Qualitätskontrollen für automatisierte Verfahren
 - Eigenschaften der Stichprobe an Endnutzer kommunizieren
- Während es sehr viele Visualisierungsarten gibt, werden wir nur eine Auswahl davon behandeln.

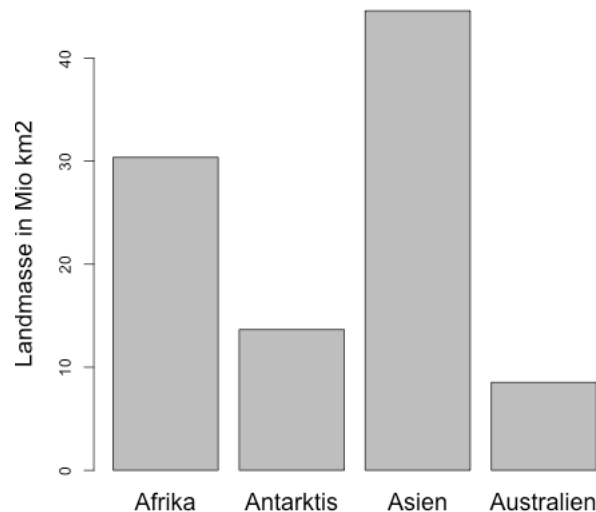
Gehalt (in Tausend €)	Häufigkeit
57	4
58	1
59	3
60	5
61	8
62	10
63	0
64	5
66	2
67	3
70	1



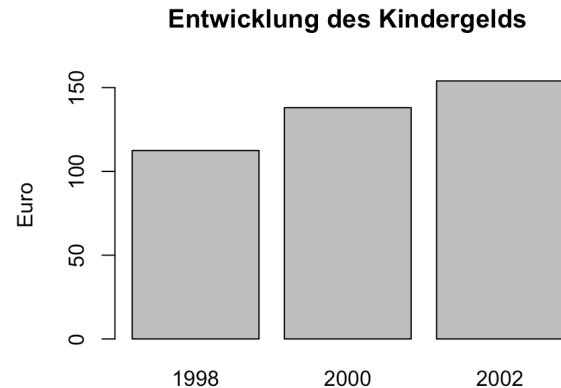
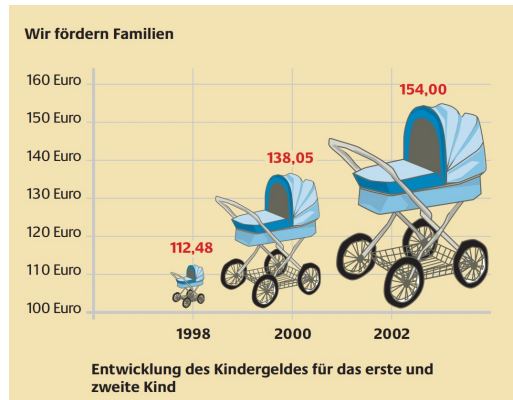
Säulendiagramm (*bar plot*)

- Ein **Säulendiagramm** ist die Darstellung einer kontinuierlichen Metrik für verschiedene Objekte oder Objektkategorien in der Stichprobe.
- Es eignet sich zum Vergleich der Objekte bzw. Objektkategorien in dieser Metrik.
- **Beispiel (Säulendiagramm)**

```
data <- data.frame(  
  continent = c("Afrika", "Antarktis", "Asien", "Australien"),  
  landmass = c(30.37, 13.66, 44.58, 8.53)  
)  
  
barplot(height = data$landmass, names = data$continent,  
        ylab = "Landmass in Mio km2")
```



- **Achtung:** Säulendiagramme können den Betrachter täuschen, wenn der Achsenschnittpunkt schlecht platziert ist oder die Säulen unterschiedlich breit sind.
- **Beispiel (Täuschung im Säulendiagramm)**



Mathe III

Unit 8 –
Explorative Datenanalyse

- Ein Säulendiagramm kommt schnell an seine Grenzen, falls

- die Objektkategorien aus einer zu großen Menge (>10 Kategorien) stammen.
- die Objektkategorien aus dem Raum der Zahlen kommen (d.h. Distanzen haben Bedeutung)

- Hierbei können Histogramme helfen.

- Es zählt, wie oft eine Größe vorkommt und bündelt die Ergebnisse in Klassen (*bins*).
- Ein Histogramm wird durch Ursprung x_0 und Klassenbreite (*bin width*) h bestimmt.
- Für alle $i \in \mathbb{Z}$ zählt die i -te Klasse die Vorkommen im Intervall

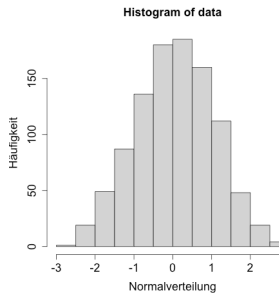
$$[x_0 + i \cdot h, x_0 + (i + 1) \cdot h)$$

- Jede Klasse ist eine Säule über ihrem Intervall mit der Häufigkeit als Höhe.

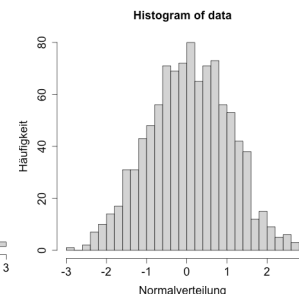
```
data <- rnorm(1000, mean = 0, sd = 1)
```

```
hist(data, xlab = "Normalverteilung",  
      ylab = "Häufigkeit", breaks = 10)
```

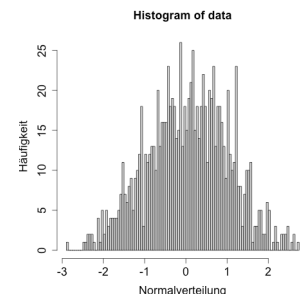
10 bins



30 bins



100 bins



Mathe III

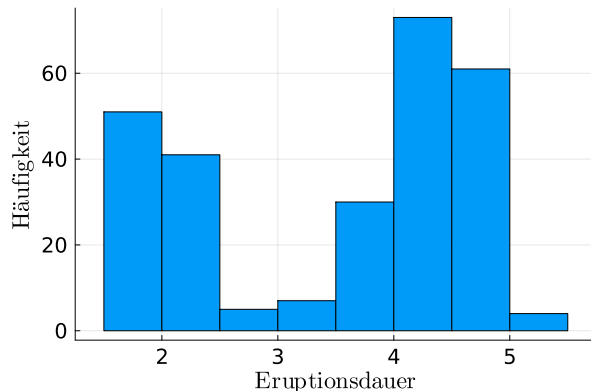
Unit 8 –
Explorative Datenanalyse

Histogramm

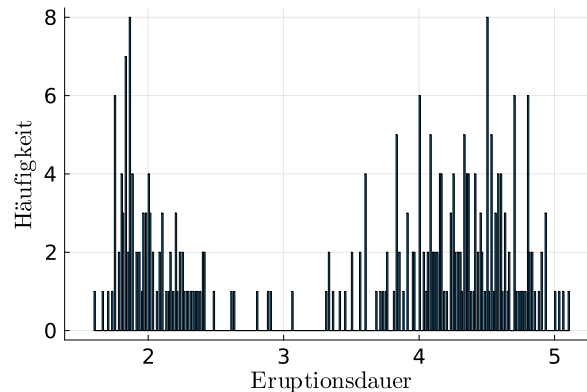
- **Achtung:** Ein Histogramm kann degenerieren, wenn
 1. die Klassenbreite zu schmal ist und viele Klassen kein oder wenige Vorkommen enthalten
 2. die Klassenbreite zu weit ist und eine zu starke Bündelung die Aussagekraft trübt.
- **Beispiel (Verschiedene Klassenbreiten im Histogramm).** Eruptionsdauer des Old Faithful Geysirs im Yellowstone-Nationalpark (in Minuten)



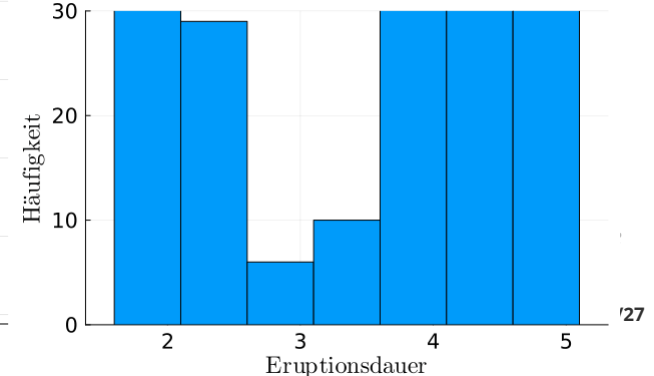
8 bins



500 bins

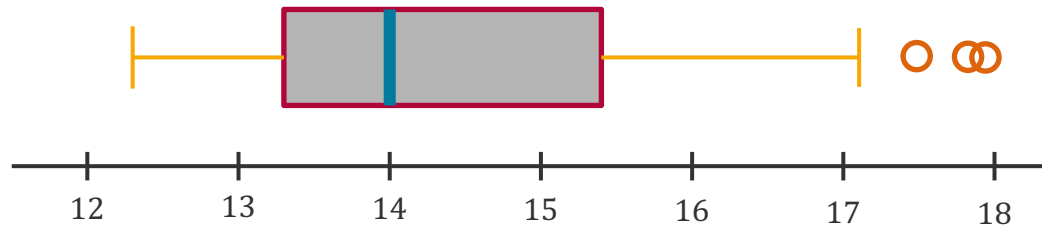


Bins = 8



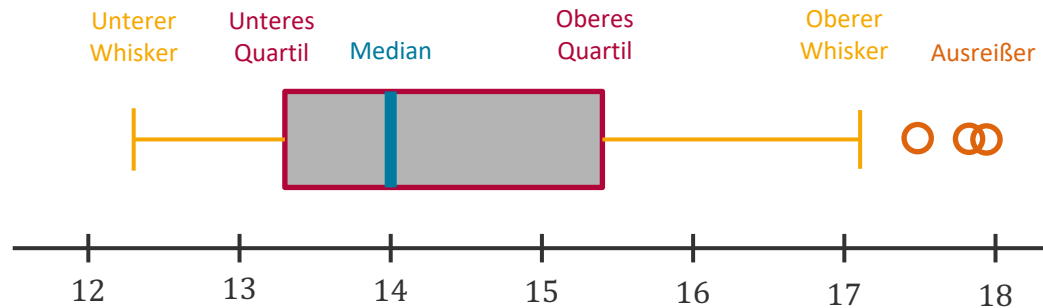
Boxplot

- Für eine besonders kompakte Darstellung der Verteilung von Werten kann ein Boxplot verwendet werden.
- Ein Boxplot enthält weniger Detailstufen als ein Histogramm, allerdings sind Kenngrößen wie Median oder Quantil auf einen Blick ersichtlich.
- Damit können mehrere Verteilungen effizient bezüglich dieser Kenngrößen verglichen werden.



Boxplot

- Das untere **Quartil** entspricht dem 0.25-Quantil, das obere dem 0.75-Quantil.
- Die **Whisker** sind im Allgemeinen so lang wie die Wertspanne, jedoch nicht länger als der anderthalbfache Interquartilabstand.
- Werte außerhalb der Whisker (weiter als der anderthalbfache Interquartilabstand vom Median entfernt) werden als **Ausreißer** gekennzeichnet.
- **Bemerkung.** Die Definition der Whisker ist in der Literatur nicht konsistent.



■ Beispiel (Vergleich von Daten mit Boxplots)

- Im Vergleich dargestellt sind drei Boxplots der Messwerte von Kelchblattlänge in Zentimetern für die Schwertlilien-Spezies *Iris setosa*, *Iris versicolor* und *Iris virginica*.
- Auch wenn die Wertebereiche nicht disjunkt sind, ist eine Tendenz klar erkennbar.

Iris setosa



Iris versicolor

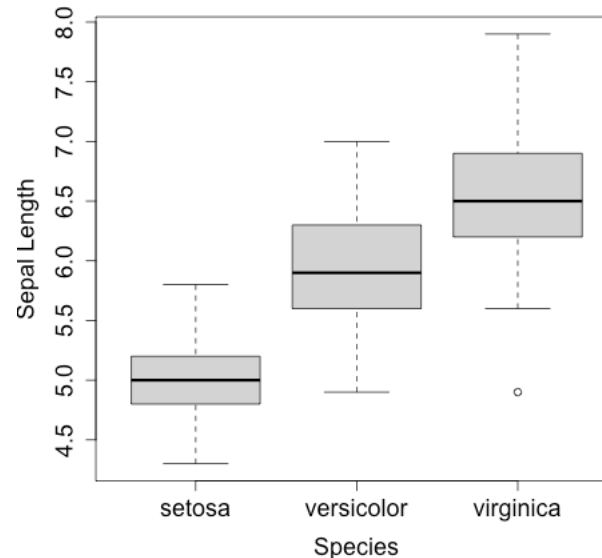


Iris virginica



```
data(iris)
```

```
boxplot(iris$Sepal.Length ~ iris$Species,  
        xlab = "Species",  
        ylab = "Sepal Length")
```



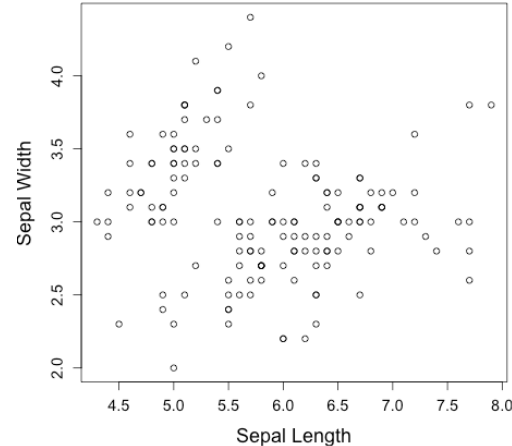
Mathe III

Unit 8 –
Explorative Datenanalyse

- Ein Streudiagramm (*scatter plot*) stellt zwei (oder mehr) kontinuierliche Eigenschaften einer Zufallsbeobachtung dar, welche gleichzeitig aufgetreten sind.
- Damit können Korrelationen gut sichtbar gemacht werden.

- **Beispiel (Streudiagramm)**

```
data(iris)
plot(iris$Sepal.Length, iris$Sepal.Width,
     xlab = "Sepal Length",
     ylab = "Sepal Width")
```



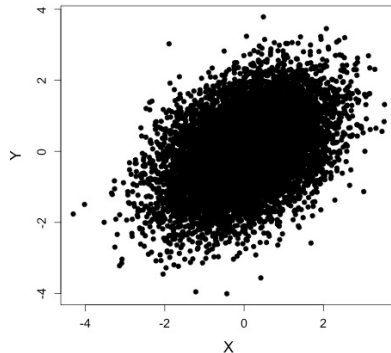
Mathe III

Unit 8 –
Explorative Datenanalyse

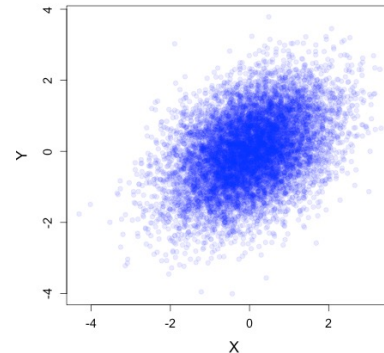
- **Achtung:** Überlagern sich mehrere Datenpunkte im Streudiagramm, kann die optische Wahrnehmung verfälscht werden.
 - Eine Darstellung der Messwert-Dichte statt der Einzelmesspunkte kann dies beheben.

- **Beispiel (Überlagerung in Streudiagrammen)**

```
plot(data[, 1], data[, 2],  
      xlab = "X", ylab = "Y", pch=19)
```



```
plot(data[, 1], data[, 2],  
      xlab = "X", ylab = "Y", pch=19,  
      col = rgb(0, 0, 1, 0.1))
```



Viel Spaß bis zur nächsten Vorlesung!