

Datenanalyse und Statistik

K.Gerald van den Boogaart
<http://www.stat.boogaart.de/>

21. Oktober 2019

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)
- ▶ Probeklausuren (später auf OPAL)

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)
- ▶ Probeklausuren (später auf OPAL)
- ▶ Organisation

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zuzätzlich)
- ▶ Probeklausuren (später auf OPAL)
- ▶ Organisation
- ▶ Übungen

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)
- ▶ Probeklausuren (später auf OPAL)
- ▶ Organisation
- ▶ Übungen
- ▶ Klausur (Anmeldung 1+2, Hilfsmittel)

Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)
- ▶ Probeklausuren (später auf OPAL)
- ▶ Organisation
- ▶ Übungen
- ▶ Klausur (Anmeldung 1+2, Hilfsmittel)
- ▶ Vorlesung

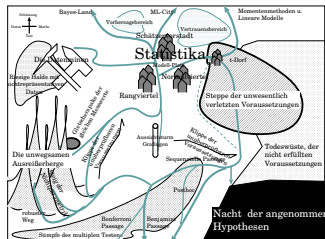
Organisation

- ▶ Webseite <http://stat.boogaart.de/DS>
- ▶ Bildungsserver <https://bildungsportal.sachsen.de/opal>
- ▶ Folien sind Folien (um Bemerkungen zur Vorlesung zu notieren)
- ▶ Skript (gibt es zusätzlich)
- ▶ Probeklausuren (später auf OPAL)
- ▶ Organisation
- ▶ Übungen
- ▶ Klausur (Anmeldung 1+2, Hilfsmittel)
- ▶ Vorlesung
- ▶ Wie bestehe ich? (Vorlesung, Lernen, Übungen)

Inhalt des Semesters

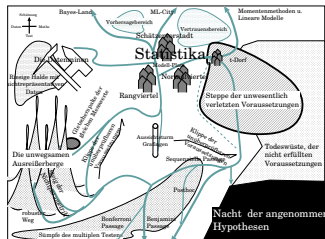
► Statistische Daten

Repräsentativität, Representation, Skala



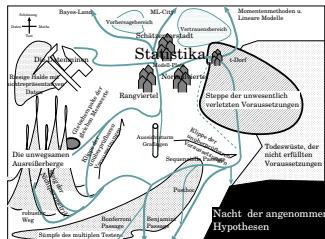
Inhalt des Semesters

- ▶ **Statistische Daten**
Repräsentativität, Representation, Skala
- ▶ **Statistische Graphik**
zur graphischen Analyse von Daten



Inhalt des Semesters

- ▶ **Statistische Daten**
Repräsentativität, Representation, Skala
- ▶ **Statistische Graphik**
zur graphischen Analyse von Daten
- ▶ **Statistische Test**
für empirische Nachweise in der Wissenschaft



Inhalt des Semesters

- ▶ **Statistische Daten**

Repräsentativität, Representation, Skala

- ▶ **Statistische Graphik**

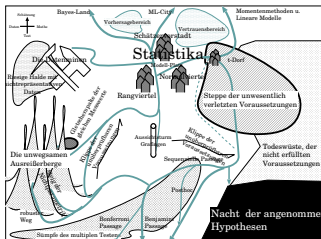
zur graphischen Analyse von Daten

- ▶ **Statistische Test**

für empirische Nachweise in der Wissenschaft

- ▶ **Lineare Modelle**

für das Modellieren von Einflüssen und die Vorhersage neuer Daten



Inhalt des Semesters

- ▶ **Statistische Daten**

Repräsentativität, Representation, Skala

- ▶ **Statistische Graphik**

zur graphischen Analyse von Daten

- ▶ **Statistische Test**

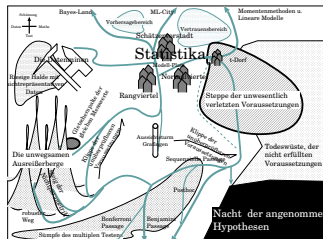
für empirische Nachweise in der Wissenschaft

- ▶ **Lineare Modelle**

für das Modellieren von Einflüssen und die Vorhersage neuer Daten

- ▶ **Ausblick auf weitere Möglichkeiten der Statistik**

andere Daten, andere Aufgaben, kausale Inferenz,...



Inhalt heute (Grundlagen)

- ▶ **Was ist Statistik?**
oder was heißt alles Statistik?

Inhalt heute (Grundlagen)

- ▶ **Was ist Statistik?**
- ▶ **Grundmodelle der Statistik**
Wie werden Daten erhoben und modelliert?

Inhalt heute (Grundlagen)

- ▶ **Was ist Statistik?**
- ▶ **Grundmodelle der Statistik**
- ▶ **Datenmatrix**
Wie werden Daten gespeichert?

Inhalt heute (Grundlagen)

- ▶ **Was ist Statistik?**
- ▶ **Grundmodelle der Statistik**
- ▶ **Datenmatrix**
- ▶ **Skala**
Wie werden Daten eingeteilt?

Inhalt heute (Grundlagen)

- ▶ **Was ist Statistik?**
- ▶ **Grundmodelle der Statistik**
- ▶ **Datenmatrix**
- ▶ **Skala**
- ▶ **Datentafel**
Wie kann man kategorielle Daten noch darstellen?

Was ist Statistik?

Wortwurzel: Aufstellungen (lat. stare)

Bedeutungen:

- ▶ **Datensammlung**
des Staats (ursprüngliche Bedeutung)

Was ist Statistik?

Wortwurzel: Aufstellungen (lat. stare)

Bedeutungen:

- ▶ **Datensammlung**
des Staats (ursprüngliche Bedeutung)
- ▶ **Wissenschaft**
von der Auswertung von Daten/vom Schließen aus Daten

Was ist Statistik?

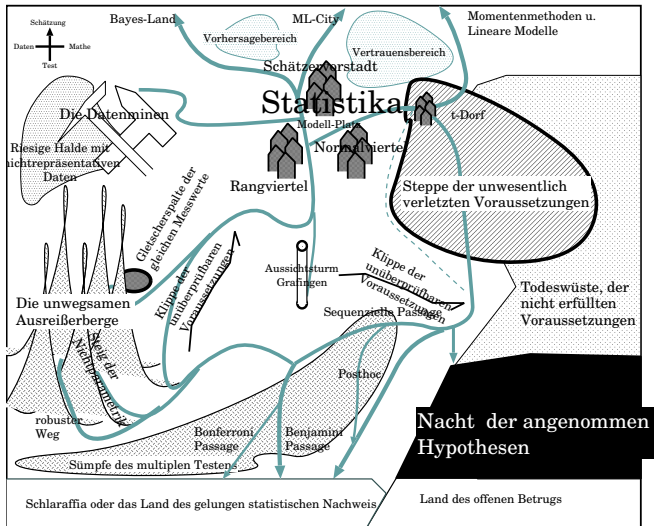
Wortwurzel: Aufstellungen (lat. stare)

Bedeutungen:

- ▶ **Datensammlung**
des Staats (ursprüngliche Bedeutung)
- ▶ **Wissenschaft**
von der Auswertung von Daten/vom Schließen aus Daten
- ▶ Aus beobachteten Zufallsvariablen berechnete weitere
Zufallsvariablen (z.B. der Mittelwert)

$$\bar{X} = \frac{1}{n} (X_1 + \dots + X_n)$$

Die Landkarte der Vorlesung



Nur mit repräsentativen Daten kann man statistische Schlüsse ziehen.

Nur mit repräsentativen Daten kann man statistische Schlüsse ziehen.

Repräsentativ heißt:

- ▶ *identisch verteilt*: alle Beobachtungen folgen dem gleichen (uns interessierenden) Zufallsgesetz.
- ▶ *stochastisch unabhängig*: jede Beobachtung ist neu nach diesem Zufallsgesetz zustande gekommen.

Es gibt zwei grundsätzlich verschiedene Wege zu repräsentativen Daten:

- ▶ **Zufällige und faire** Auswahl einer *Stichprobe* aus einer **Grundgesamtheit**.
- ▶ **Unabhängiger** Wiederholung **identischer** *Zufallsexperimente*.

Repräsentative Daten durch Stichprobe und Grundgesamtheit

Grundbegriffe

- ▶ **Grundgesamtheit**

Die Menge der Objekte für die man sich interessiert.
z.B. alle Wahlberechtigten.

Grundbegriffe

- ▶ **Grundgesamtheit**
- ▶ **statistisches Individuen**
Die Elemente der Grundgesamtheit.
z.B. ein Wahlberechtigter.

Grundbegriffe

- ▶ **Grundgesamtheit**
- ▶ **statistisches Individuen**
- ▶ **Stichprobe**

Die Menge der statistische Individuen zu den man tatsächlich Daten erhebt.

z.B. die Menge der Befragten Wahlberechtigten

Grundbegriffe

- ▶ **Grundgesamtheit**
- ▶ **statistisches Individuen**
- ▶ **Stichprobe**
- ▶ **repräsentativ**

Eine Stichprobe heißt repräsentativ, wenn jedes statistische Individuum aus der Grundgesamtheit unabhängig von allen anderen mit der gleichen Wahrscheinlichkeit in die Stichprobe gelangen kann.

Grundbegriffe

- ▶ **Grundgesamtheit**
- ▶ **statistisches Individuen**
- ▶ **Stichprobe**
- ▶ **repräsentativ**
- ▶ **Zufallsvariable**

das beim i -ten statistischen Individuum beobachtete Merkmal (z.B. bevorzugte Partei) wird durch eine Zufallsvariable X_i modelliert.

Grundbegriffe

- ▶ **Grundgesamtheit**
- ▶ **statistisches Individuen**
- ▶ **Stichprobe**
- ▶ **repräsentativ**
- ▶ **Zufallsvariable**
- ▶ **Realisierung der Zufallsvariable**

der tatsächlich beim i -ten Individuum beobachtete Wert x_i wird als Realisierung der Zufallsvariable X_i bezeichnet.

Beispiel: Bodenqualität

- ▶ **Grundgesamtheit:** Alle Punkte des Bodens im Untersuchungsgebiet.
- ▶ **Stichprobe:** Zufällig ausgewählte Untersuchungspunkte.
- ▶ **Zufallsvariablen:** Nährstoffgehalt in an diesen Stellen genommenen Bodenproben.
- ▶ **Realisierungen:** 5.34%, 7, 45%, ...

Beispiel: Werkstückprüfung

- ▶ **Grundgesamtheit:** Alle gefertigten Zahnräder der Teilenummer 45632N.
- ▶ **Stichprobe:** Zufällig zu Testzwecken entnommen Zahnräder.
- ▶ **Zufallsvariablen:** Betriebstunden im Testbetrieb bis Defekt.
- ▶ **Realisierungen:** $5343h, 7342h, \dots$

Vollerhebung

- ▶ Die **Vollerhebung** ist eine spezielle Art der Stichprobennahme.
- ▶ Bei Vollerhebung ist die Stichprobe gleich der Grundgesamtheit.
- ▶ Unabhängigkeit: alle kommen unabhängig von allen anderen sicher in die Stichprobe.
- ▶ gleiche Wahrscheinlichkeit: Wahrscheinlichkeit in die Stichprobe zu kommen ist 1.

Repräsentative Daten durch Zufallsexperimente

Grundbegriffe

▶ **Vorschrift für ein Zufallsexperiment**

Ein genaue Experimentiervorschrift für ein Experiment mit zufälligem Ausgang.

z.B. Würfelwurf

z.B. Messung mit Messfehler

z.B. Bruchversuch

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**

- ▶ **Zufallsexperiment**

Ein Experiment das nach dieser Experimentiervorschrift durchgeführt wurde

z.B. 3ter Wurf des Würfels

z.B. 2te Messung

z.B. 15ter Bruchversuch

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**
- ▶ **Zufallsexperiment**
- ▶ **identisch verteilt**

Zufallsexperimente heißen identisch verteilt, wenn sie nach der gleichen Vorschrift durchgeführt wurden, so dass jede Menge möglicher Ausgänge jeweils gleich wahrscheinlich sind.

z.B. gleiche Bruchversuche mit gleich produzierten Werkstücken.

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**
- ▶ **Zufallsexperiment**
- ▶ **identisch verteilt**
- ▶ **unabhängig**

Zufallsexperimente heißen unabhängig durchgeführt, wenn sich der Zufall der Experimente gegenseitig nicht beeinflussen kann.

z.B. nicht!!!, wenn die Werkstücke aus der gleichen Charge Rohmaterial hergestellt wurde, aber die Vorschrift nicht die Rohstoffcharge spezifiziert hat.

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**
- ▶ **Zufallsexperiment**
- ▶ **identisch verteilt**
- ▶ **unabhängig**
- ▶ **repräsentativ**

Zufallsexperimente heißt repräsentativ (für ihre Vorschrift), wenn sie unabhängig und identisch verteilt sind.

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**
- ▶ **Zufallsexperiment**
- ▶ **identisch verteilt**
- ▶ **unabhängig**
- ▶ **repräsentativ**
- ▶ **Zufallsvariable**

die beim i -ten Zufallsexperiment zu machende Beobachtung X_i .
z.B. die Bruchspannung.

Grundbegriffe

- ▶ **Vorschrift für ein Zufallsexperiment**
- ▶ **Zufallsexperiment**
- ▶ **identisch verteilt**
- ▶ **unabhängig**
- ▶ **repräsentativ**
- ▶ **Zufallsvariable**
- ▶ **Realisierung der Zufallsvariable**

der tatsächlich beim i -ten Individuum beobachtete Wert x_i wird als Realisierung der Zufallsvariable X_i bezeichnet.

Fadenbrüche

Anzahl Fadenbrüche bei verschiedenen Rahmenbedingungen:

```
> warpbreaks
```

```
  breaks wool tension
1      26    A      L
2      30    A      L
3      54    A      L
4      25    A      L
5      70    A      L
6      52    A      L
7      51    A      L
8      26    A      L
9      67    A      L
10     18    A      M
11     21    A      M
12     29    A      M
13     17    A      M
14     12    A      M
15     18    A      M
16     35    A      M
17     30    A      M
18     36    A      M
19     36    A      H
20     21    A      H
21     24    A      H
22     18    A      H
23     10    A      H
```

Beispiel: Lichtgeschwindigkeitsmessungen

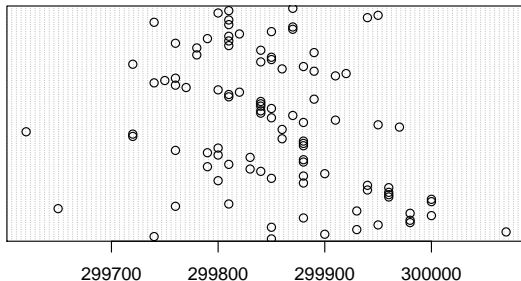
> *lightspeeds*

```
[1] 299850 299740 299900 300070 299930 299850 299950 299980
 [9] 299980 299880 300000 299980 299930 299650 299760 299810
[17] 300000 300000 299960 299960 299960 299940 299960 299940
[25] 299880 299800 299850 299880 299900 299840 299830 299790
[33] 299810 299880 299880 299830 299800 299790 299760 299800
[41] 299880 299880 299880 299860 299720 299720 299620 299860
[49] 299970 299950 299880 299910 299850 299870 299840 299840
[57] 299850 299840 299840 299840 299890 299810 299810 299820
[65] 299800 299770 299760 299740 299750 299760 299910 299920
[73] 299890 299860 299880 299720 299840 299850 299850 299780
[81] 299890 299840 299780 299810 299760 299810 299790 299810
[89] 299820 299850 299870 299870 299810 299740 299810 299940
[97] 299950 299800 299810 299870
```


Beispiel: Lichtgeschwindigkeitsmessungen

```
> dotchart(lightspeeds, main="Michelsons Lichtgeschwindigkeitsmessungen")
```

Michelsons Lichtgeschwindigkeitsmessungen



Repräsentativität

Allgemein (resultierende Zufallsvariablen)

- ▶ identisch verteilt
- ▶ stochastisch unabhängig

Stichproben (zufällige Auswahl)

- ▶ mit der gleichen Wahrscheinlichkeit
- ▶ unabhängig voneinander

Zufallsexperimente (Experiment mit zufälligem Ausgang)

- ▶ nach gleicher Vorschrift durchgeführt
- ▶ unabhängig voneinander

Mehrstichprobenmodell

Oft finden wir in einem Datensatz **zwei oder mehrer Gruppen** von Daten, die von unterschiedlichen

- ▶ Grundgesamtheit oder
- ▶ Zufallsexperimenten (Experimentiervorschriften)

herrühren.

Ein Datensatz kann also mehrer Stichproben enthalten.

Man spricht dann von einer **Zweistichproben- oder Mehrstichprobensituation**.

Zufälligkeit der Daten

Ein repräsentativer Datensatz ist grundsätzlich zufällig, da

- ▶ die Auswahl der Beobachtungen zufällig zustande gekommen ist, oder
- ▶ die Experimente zufällige Ergebnisse haben.

Wir interessieren uns aber nicht für die konkreten Daten, sondern für die dahinterstehenden Gesetze: z.B. für die Zahnräder, die tatsächlich ausgeliefert werden, was alle Deutschen wählen, oder welche Maschineneinstellung in Zukunft die besten Ergebnisse liefert.

Zufälligkeit der Kenngrößen

Das erste Ergebnis einer statistischen Analyse sind oft Kenngrößen, wie z.B. der Mittelwert.

► **Der Mittelwert als Zufallsvariable und Statistik**

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

Der Mittelwert ist selbst **zufällig!!!**.

Zufälligkeit der Kenngrößen

Das erste Ergebnis einer statistischen Analyse sind oft Kenngrößen, wie z.B. der Mittelwert.

- ▶ **Der Mittelwert als Zufallsvariable und Statistik**

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

Der Mittelwert ist selbst **zufällig!!!**.

- ▶ **Der Mittelwert als abstrakte Realisierung**

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

Zufälligkeit der Kenngrößen

Das erste Ergebnis einer statistischen Analyse sind oft Kenngrößen, wie z.B. der Mittelwert.

- ▶ **Der Mittelwert als Zufallsvariable und Statistik**

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

Der Mittelwert ist selbst **zufällig!!!**.

- ▶ **Der Mittelwert als abstrakte Realisierung**

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

- ▶ **Der realisierte Mittelwert**

[1] 299852.4

Repräsentation statistischer Daten

- ▶ Datenliste
- ▶ Datenmatrix
 - ▶ Fälle
 - ▶ Variablen
 - ▶ Skala (bestimmt die Auswertung!!!)
- ▶ Datentafel

Beispiel einer Datenliste

```
> lightspeeds
```

```
[1] 299850 299740 299900 300070 299930 299850 299950 299980  
[9] 299980 299880 300000 299980 299930 299650 299760 299810  
[17] 300000 300000 299960 299960 299960 299940 299960 299940  
[25] 299880 299800 299850 299880 299900 299840 299830 299790  
[33] 299810 299880 299880 299830 299800 299790 299760 299800  
[41] 299880 299880 299880 299860 299720 299720 299620 299860  
[49] 299970 299950 299880 299910 299850 299870 299840 299840  
[57] 299850 299840 299840 299840 299890 299810 299810 299820  
[65] 299800 299770 299760 299740 299750 299760 299910 299920  
[73] 299890 299860 299880 299720 299840 299850 299850 299780  
[81] 299890 299840 299780 299810 299760 299810 299790 299810  
[89] 299820 299850 299870 299870 299810 299740 299810 299940  
[97] 299950 299800 299810 299870
```

Beispiel Datenlisten

`$setosa`

```
[1] 5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0 4.4 4.9 5.4 4.8 4.8 4.3 5.8  
[16] 5.7 5.4 5.1 5.7 5.1 5.4 5.1 4.6 5.1 4.8 5.0 5.0 5.2 5.2 4.7  
[31] 4.8 5.4 5.2 5.5 4.9 5.0 5.5 4.9 4.4 5.1 5.0 4.5 4.4 5.0 5.1  
[46] 4.8 5.1 4.6 5.3 5.0
```

`$versicolor`

```
[1] 7.0 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 5.0 5.9 6.0 6.1 5.6  
[16] 6.7 5.6 5.8 6.2 5.6 5.9 6.1 6.3 6.1 6.4 6.6 6.8 6.7 6.0 5.7  
[31] 5.5 5.5 5.8 6.0 5.4 6.0 6.7 6.3 5.6 5.5 5.5 6.1 5.8 5.0 5.6  
[46] 5.7 5.7 6.2 5.1 5.7
```

`$virginica`

```
[1] 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 6.5 6.4 6.8 5.7 5.8  
[16] 6.4 6.5 7.7 7.7 6.0 6.9 5.6 7.7 6.3 6.7 7.2 6.2 6.1 6.4 7.2  
[31] 7.4 7.9 6.4 6.3 6.1 7.7 6.3 6.4 6.0 6.9 6.7 6.9 5.8 6.8 6.7  
[46] 6.7 6.3 6.5 6.2 5.9
```

Beispiel einer Datenmatrix

Ausschnitt eines Datensatzes:

```
> X
```

	Sepal.Length	Sepal.Width	Species
1	5.1	3.5	setosa
2	4.9	3.0	setosa
3	4.7	3.2	setosa
53	6.9	3.1	versicolor
54	5.5	2.3	versicolor
56	5.7	2.8	versicolor
58	4.9	2.4	versicolor

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.
- ▶ Jede Zeile X_i gehört zu einem statistischen Individuum

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.
- ▶ Jede Zeile $X_{i\cdot}$ gehört zu einem statistischen Individuum
- ▶ Jede Spalte $X_{\cdot j}$ gehört zu einem Merkmal

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.
- ▶ Jede Zeile $X_{i\cdot}$ gehört zu einem statistischen Individuum
- ▶ Jede Spalte $X_{\cdot j}$ gehört zu einem Merkmal
- ▶ Der Eintrag X_{ij} entspricht der Ausprägung des j -ten Merkmals am i -ten Individuum.

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.
- ▶ Jede Zeile $X_{i\cdot}$ gehört zu einem statistischen Individuum
- ▶ Jede Spalte $X_{\cdot j}$ gehört zu einem Merkmal
- ▶ Der Eintrag X_{ij} entspricht der Ausprägung des j -ten Merkmals am i -ten Individuum.
- ▶ Die Einträge einer Datenmatrix sind Zufallsvariablen bzw. ihre Realisierungen.

Die Datenmatrix

- ▶ $X_{ij}, i = 1, \dots, n, j = 1, \dots, m$ sind die Einträge einer Datenmatrix.
- ▶ Jede Zeile $X_{i\cdot}$ gehört zu einem statistischen Individuum
- ▶ Jede Spalte $X_{\cdot j}$ gehört zu einem Merkmal
- ▶ Der Eintrag X_{ij} entspricht der Ausprägung des j -ten Merkmals am i -ten Individuum.
- ▶ Die Einträge einer Datenmatrix sind Zufallsvariablen bzw. ihre Realisierungen.
- ▶ Die Einträge einer Datenmatrix sind nicht unbedingt reelle Zahlen!

Fälle

Die Zeile der Datenmatrix heißen Fälle. Sie entsprechen den statistischen Individuen.

```
> X
```

	Sepal.Length	Sepal.Width	Species
1	5.1	3.5	setosa
2	4.9	3.0	setosa
3	4.7	3.2	setosa
53	6.9	3.1	versicolor
54	5.5	2.3	versicolor
56	5.7	2.8	versicolor
58	4.9	2.4	versicolor

Der Begriff der Skala

Zu jeder Variable gehört eine **Skala**, also ein Wertebereich mit gewissen sinnvollen mathematischen Operationen.

Kriterien zur Bestimmung der Skala sind:

- ▶ Welche Werte sind möglich?
- ▶ Wieviele Werte sind möglich?
- ▶ Sind die möglichen Werte geordnet? (Fachabi < Abi?)
- ▶ Sind die Abstände der Werte vergleichbar?
- ▶ Ist die Differenz ein guter Unterschiedsbegriff?
- ▶ Ist das Verhältnis ein guter Unterschiedsbegriff?

Skalen

- ▶ **diskrete Skalen**

haben voneinander getrennte Werte

nominal ($()$), dichotom ($= \text{NOT}$), kategoriell ($=$), ordinal ($= <$),
Intervallskala ($= < -$), Anzahlen ($= < -*$)

- ▶ **stetige Skalen**

Anteil $< /$, positiv $< */$, reell $< *-$

- ▶ **spezielle Skalen**

z.B. Richtungen, Zusammensetzungen, Orientierungen,
Winkel, Zuordnungen, ...

Diskrete Skalen

- ▶ **nominal**

(lat. nomen = der Name) Namen, bezeichnet verschiedene Individuen, die prinzipiell gleichartig, aber individuell verschieden sind.

Diskrete Skalen

- ▶ **nominal**
- ▶ **kategoriiell**

(gr. katagorein = unterordnen, einordnen)

Kategorien, jedes statistische Individuum wird in eine von ein paar vor dem Experiment feststehende Kategorien eingeordnet.

Diskrete Skalen

- ▶ **nominal**
- ▶ **kategoriiell**
- ▶ **ordinal**

(angeordnet)

wie kategoriiell, nur dass die Kategorien in eine natürliche Reihenfolge gebracht werden können.

z.B. tot, krank, gesund

Diskrete Skalen

- ▶ **nominal**
- ▶ **kategorial**
- ▶ **ordinal**
- ▶ **dichotom**

(gr. dichys=zwei, tomein=schneiden)

wie kategorial, nur dass nur zwei Kategorien gibt.

z.B. männlich - weiblich, ja - nein, ...

Diskrete Skalen

- ▶ **nominal**
- ▶ **kategorial**
- ▶ **ordinal**
- ▶ **dichotom**
- ▶ **intervallskaliert**

(gr. intervall = Abstand)

wie ordinal, nur dass der Abstand der benachbarten Kategorien vergleichbar ist.

z.B. Schulnoten

Diskrete Skalen

- ▶ **nominal**
- ▶ **kategoriiell**
- ▶ **ordinal**
- ▶ **dichotom**
- ▶ **intervallskaliert**
- ▶ **Anzahl**

(von zählen)

Es wurde etwas gezählt. Wie Intervallskaliert nur, dass die Kategorien alle möglichen natürlichen Zahlen oder 0 sind.

Die diskreten Skalen

	Name	Geschlecht	Fach	Stufe	Note	Kinder
1	Maier	m	Chemie	Abi	4	0
2	Huber	w	Biologie	Vordiplom	1	1
3	Mueller	m	Geographie	Hauptdiplom	2	4

Stetige Skalen

- ▶ **reell**

(reelle Zahlen)

Jede beliebige reelle Zahl kann vorkommen. $+$, $-$, $*$ sind sinnvolle Operationen. Der Abstand von 10 zu 5 ist genauso groß wie der Abstand von 5 zu 0.

z.B. Temperaturänderung

Stetige Skalen

- ▶ **reell**
- ▶ **ratio / positiv reell /Verhältnisskala**
(ratio =Verhältnis)

Nur positive Zahlen können beobachtet werden. $*$, $/$ sind sinnvolle Operationen. Der Abstand von 10 zu 1 ist genauso groß, wie der Abstand von 1 zu 0.1.

z.B. Gewicht, Länge

Stetige Skalen

- ▶ **reell**
- ▶ **ratio / positiv reell / Verhältnisskala**
- ▶ **Anteilskala / Wahrscheinlichkeitskala**

(Anteil vom Ganzen)

Nur Werte zwischen 0 und 1 können beobachtet werden. Die Werte sind als Anteile interpretierbar.

Die stetigen Skalen

	Alkoholanteil	Menge	Temperatur
1	0.1	0.125	16
2	0.3	0.500	5
3	0.7	1.000	-20

Grobeinteilung der Skalen

Die Skala bestimmt welche statistischen Verfahren angewendet werden können. Oft genügt im ersten Schritt schon eine Grobeinteilung:

- ▶ **diskret**

Variablen mit diskreten Skalen heißen oft auch Faktor. Die Möglichen Werte heißen dann Stufen des Faktors.

- ▶ **stetig**

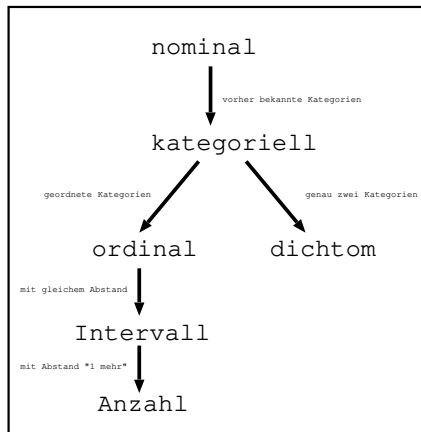
Variablen mit stetigen Skalen können ein unendlich viele verschiedene Zahlenwerte annehmen. Treten dabei der gleiche Wert mehrfach auf, so spricht man von **Bindungen**.

- ▶ **spezielle**

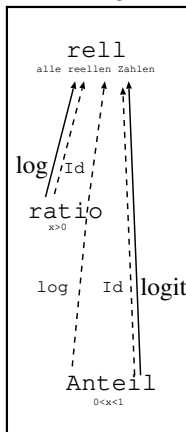
Variablen, die nicht ins Schema passen haben eine spezielle Skala.

Das feinste Skalenniveau

diskret



stetig



Versuchen wir es selbst

Ausschnitt des Iris Blüten Datensatzes:

```
> X
```

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
53	6.9	3.1	4.9	1.5	versicolor
54	5.5	2.3	4.0	1.3	versicolor
56	5.7	2.8	4.5	1.3	versicolor
58	4.9	2.4	3.3	1.0	versicolor

Welche Spalte hat welche Skala?

Wozu Skala?

- ▶ Die Skala bestimmt welche weiteren Verfahren angewendet werden sollten.
- ▶ Die Skala gibt Hinweise was in der weiteren Analyse beachtet werden sollte.
- ▶ Die Skala bestimmt, wie die Daten zusammengefaßt und beschrieben werden können.
- ▶ Die Bestimmung der Skala der Variablen ist daher der erste Schritt jeder Datenanalyse.

Datentafel

Die Datentafel ist eine alternative Darstellung zur Datenmatrix, wenn nur diskrete Skalen auftreten.

Datentafel (Beispiel)

```
> data(Titanic)
> ftable(Titanic,col.vars=c("Class","Survived"))
```

		Class		1st		2nd		3rd		Crew	
		Survived		No	Yes	No	Yes	No	Yes	No	Yes
Sex	Age										
	Male	Child		0	5	0	11	35	13	0	0
	Adult	118	57	154	14	387	75	670	192		
Female	Child	0		1	0	13	17	14	0	0	
	Adult	4	140	13	80	89	76	3	20		

Erklaerung Datentafel

		Class							
		1st		2nd		3rd		Crew	
Survived		No	Yes	No	Yes	No	Yes	No	Yes
Sex	Age								
	Child	0	5	0	11	35	13	0	0
Male	Adult	118	57	154	14	387	75	670	192
Female	Child	0	1	0	13	17	14	0	0
	Adult	4	140	13	80	89	76	3	20

Erklaerung Datentafel

		1st		2nd		3rd		Crew	
Sex	Age	Survived		Survived		Survived		Survived	
		No	Yes	No	Yes	No	Yes	No	Yes
Male	Child	0	5	0	11	35	13	0	0
	Adult	118	57	154	14	387	75	670	192
Female	Child	0	1	0	13	17	14	0	0
	Adult	4	140	13	80	89	76	3	20

Erklärung Datentafel

Sex	Age	Class 1st		2nd		3rd		Crew		
		Survived	No	Yes	No	Yes	No	Yes	No	Yes
Male	Child		0	5	0	11	35	13	0	0
	Adult		118	57	154	14	387	75	670	192
Female	Child		0	1	0	13	17	14	0	0
	Adult		4	140	13	80	89	76	3	20

Die Datentafel

- ▶ Jede Zelle der Datentafel enthält die Anzahl statistischer Individuen in der Stichprobe mit der gegebenen Faktorkombination.

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**

Datenmatrix, Datentafel, Unvorbereitet

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**

Wir lesen die Beschreibung!!!

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**

diskret: nominal, kategoriell, ordinal, Intervall, Anzahl

stetig: (Anzahl), reell, ratio, Anteil *speziell*: Irgendwie anders

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**
Eine Grundgesamtheit, zwei oder mehrere? bzw.
Eine Typ Zufallsexperiment, zwei oder mehrere?

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**
- ▶ **Was sind die Grundgesamtheiten?**

Welche wünschen wir uns? Für welche ist es wohl repräsentativ?

Erste Analyseschritte

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**
- ▶ **Was sind die Grundgesamtheiten?**
- ▶ **Sind die Daten für die Grundgesamtheit repräsentativ?**

Wie sind die Daten zustande gekommen?, unabhängig?,
identische Zufallsgesetze?

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**

Wir müssen die Daten gegebenenfalls in eine Datenmatrix übersetzen, um sie in Statistiksoftware einlesen zu können.

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**

Wir müssen den Datensatz verstehen, um ihn auswerten zu können.

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**

Die Skala bestimmt die Auswahl der Analyseverfahren.

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**

Stichproben müssen einzeln oder mit Mehrstichprobenverfahren ausgewertet werden.

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**
- ▶ **Was sind die Grundgesamtheiten?**

Alle Analyseergebnisse beziehen sich nur auf diese Grundgesamtheit.

Wozu die ersten Analyseschritte?

Eine Datenauswertung beginnt grundsätzlich mit den folgenden Analyseschritten:

- ▶ **Wie liegen die Daten vor?**
- ▶ **Welche Variablen gibt es und was bedeuten Sie?**
- ▶ **Welche Skala haben die einzelnen Variablen?**
- ▶ **Ein-, Zwei- oder Mehrstichprobensituation?**
- ▶ **Was sind die Grundgesamtheiten?**
- ▶ **Sind die Daten für die Grundgesamtheit repräsentativ?**
Nur repräsentative Daten führen zu richtigen Auswertungen!!!

Repräsentation statistischer Daten

- ▶ Datenliste
 - ▶ Nur ein Merkmal!!!
 - ▶ alle Skalen
- ▶ Datenmatrix
 - ▶ mehrere Variablen
 - ▶ alle Skalen
- ▶ Datentafel
 - ▶ mehrere Variablen
 - ▶ nur kategorielle Skalen

Zusammenfassung

- ▶ **Repräsentativität statistischer Daten**
- ▶ **Repräsentation statistischer Daten**
- ▶ **Skalen statistischer Daten**
- ▶ **Zufälligkeit statistischer Daten**

Zusammenfassung

- ▶ **Repräsentativität statistischer Daten**
Nur diese Daten erlauben Rückschlüsse.
- ▶ **Repräsentation statistischer Daten**
Nur diese Daten versteht jemand.
- ▶ **Skalen statistischer Daten**
Das bestimmt das Auswertungsverfahren.
- ▶ **Zufälligkeit statistischer Daten**
Das ist das Kernproblem bei der Auswertung.

Einordnung

