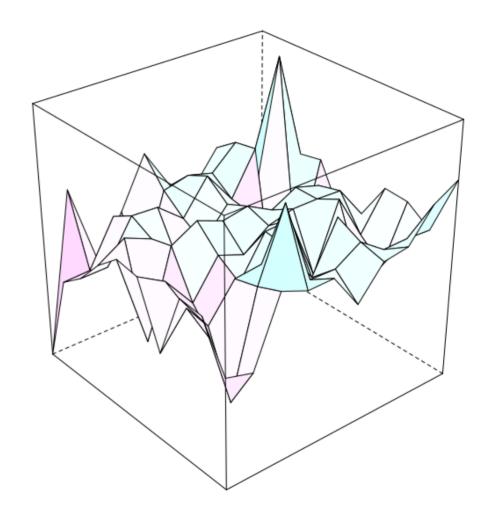
## **Datenanalyse mit R**Prof. Dr. Dennis Klinkhammer



#### (I) GRUNDLAGEN



- Grundlage für seriöse Forschung
- Beziehen sich auf mögliche <u>Schwachstellen</u>
  - des Datensatzes
  - der Erhebung des Datensatzes
  - der Forscherinnen und Forscher selbst

- Ermöglichen Berücksichtigung von Schwachstellen in der <u>Ergebnisinterpretation</u>
- Unterscheidung in <u>drei Hauptarten</u>
  - Objektivität
  - Reliabilität
  - Validität

- Objektivität bedeutet Unabhängigkeit
- Der Forschungsprozess wird nicht für eigene oder "Wunschergebnisse" Dritter <u>beeinflusst</u>
- Insbesondere die <u>Finanzierung</u> von Forscherinnen und Forschern durch Dritte bietet Anlass zur Frage nach deren Objektivität

- Reliabilität bedeutet Reproduzierbarkeit unter Einsatz geeigneter Erhebungsinstrumente
- Zentrale Aspekte
  - Stabilität und Genauigkeit einer Messung
  - Berücksichtigung und Ausweisung der Rahmenbedingungen einer Messung

- Validität bezieht sich auf die <u>Präzision</u> des Erhebungsinstruments
- Drei Unterarten werden unterschieden
  - Konstruktvalidität (Womit wird gemessen?)
  - Kriteriumsvalidität (Was wird gemessen?)
  - prognostische Validität (Sind Schlüsse möglich?)

- Validität bezieht sich auch auf die Rahmenbedingungen einer Erhebung
- Ergebnisinterpretationen erfordern eine Ausweisung der zugrundeliegenden Validität
  - Interne Validität (kontrollierte Bedingungen)
  - Externe Validität (natürliche Bedingungen)

- Validität setzt Objektivität und Reliabilität voraus
- Forschung ohne entsprechende Ausweisung der wissenschaftlichen Gütekriterien ist keine seriöse Forschung

# Klausurrelevant (!)

### Obligatorischer Exkurs

- Simpson-Paradoxon
  - Die <u>Bewertung verschiedener Gruppen</u> fällt <u>unterschiedlich</u> aus, je nachdem ob man die Ergebnisse der Gruppen kombiniert oder nicht
  - Zusatzaufgabe: Recherche eines Beispiels

#### (I) GRUNDLAGEN



- R ist eine <u>Programmiersprache</u> für statistische Berechnungen und Grafiken
- RStudio ist eine grafische Benutzeroberfläche für die Programmiersprache R
- Die in der Vorlesung genannten Beispiele und Übungsaufgaben erfordern R und RStudio

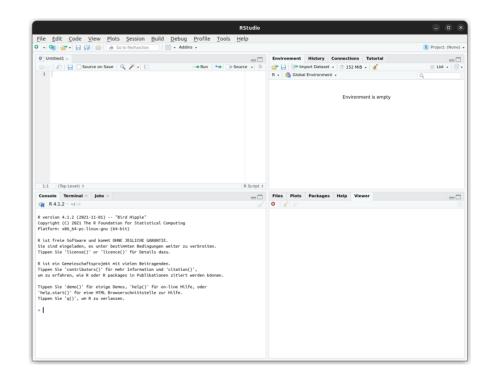
- R ist eine <u>objektbasierte Programmiersprache</u>
- Alles was <u>existiert</u> ist ein Objekt
  - Datensätze
  - Variablen
  - Merkmalsausprägungen

- Alles was <u>passiert</u> ist eine Funktion
  - Funktionen bearbeiten Objekte
  - Funktionen führen zu neuen Objekten
  - Funktionen berücksichtigen Argumente
- Die <u>Syntax</u> bildet Objekte, Funktionen und Argumente ab und ermöglicht Dritten die Reproduktion

- Erläuterungen werden über einen # eingegeben und dadurch von R nicht interpretiert
- Der Zuweisungspfeil legt neue Objekte an:

```
# Erläuterung mit Syntaxbeispiel
Kuchen <- backen(Milch, Mehl, Eier, Zucker)
```

- RStudio: 4 Fenster
  - Syntax (links oben)
  - Konsole (links unten)
  - Historie (rechts oben)
  - Grafiken (rechts unten)



 Die Funktion <u>plot(...)</u> kann auf den Datensatz <u>swiss</u> direkt angewendet oder als neues Objekt zwischengespeichert werden:

```
# Erste Beispiele
plot(swiss)
neues_objekt <- plot(swiss)
```

- Packages erweitern den Funktionsumfang
- Eine Korrelation kann ohne erweiterten Funktionsumfang über die Funktion cor(...) aufgerufen werden:

# Ohne Package CORRPLOT cor(swiss)

 Die Funktionen <u>install.packages(...)</u> und <u>library(...)</u> installieren Packages und aktivieren diese in RStudio:

```
# Package CORRPLOT installieren und aktivieren install.packages(corrplot) library("corrplot")
```

 Das neu installierte und aktivierte Package bietet zusätzlich die Funktion <u>corrplot(...)</u> und Abwandlungen derselben:

```
# Mit Package CORRPLOT
noch_ein_neues_objekt <- cor(swiss)
corrplot.mixed(noch_ein_neues_objekt)
```

 Für viele Objekte und Funktionen stehen über die Funktion <u>help(...)</u> weiterführende Informationen bereit:

```
# Hilfe zum Package CORRPLOT help(corrplot)
```

#### (I) GRUNDLAGEN



- Der TREES Datensatz ermöglicht einen ersten Einblick in die Datenanalyse mit R
- Ziele einer Datenanalyse
  - Beschreiben
  - Erklären
  - Vorhersagen

- Die verschiedenen Ziele der Datenanalyse werden mit unterschiedlichen <u>Teilbereichen der</u> <u>Statistik</u> umgesetzt
  - Univariate Statistik
  - Bivariate Statistik
  - Multivariate Statistik

- Ziel: Ein allgemeines Modell zur <u>Vorhersage</u> des Volumens von Bäumen aufstellen
- Variablen des TREES Datensatzes
  - (X<sub>1</sub>) Girth
  - (x<sub>2</sub>) Height
  - (Y) Volume

• Die Funktionen <u>summary(...)</u> und <u>boxplot(...)</u> sind die gängigsten Funktionen der univariaten Statistik und beschreiben jeweils <u>eine Variable</u>:

# Univariate Statistik summary(trees) boxplot(trees)

- Die Korrelation mit der Funktion <u>cor(...)</u> gehört zum Teilbereich der bivariaten Statistik
- Hier werden die Zusammenhänge zwischen zwei Variablen erklärt:

# Bivariate Statistik cor(trees)

 Die Zusammenhänge zwischen zwei Variablen können zusätzlich über die Funktionen plot(...) und abline(...) visualisiert werden:

```
# Bivariate Statistik
plot(trees$Volume~trees$Girth)
abline(lm(trees$Volume~trees$Girth))
```

- In der multivariaten Statistik können <u>mehrere</u> unabhängige Variablen zur <u>Vorhersage einer abhängigen Variable</u> herangezogen werden
- Dabei kann zusätzlich die <u>Stärke des Einflusses</u> der unabhängigen <u>Variablen</u> auf die abhängige Variable analysiert werden

 Der lineare Zusammenhang ermöglicht die Funktion <u>Im(...)</u> in der multivariaten Statistik:

```
# Multivariate Statistik
regression_model <- lm(Volume~., data=trees)
regression_model
summary(regression_model)
```

#### (I) GRUNDLAGEN

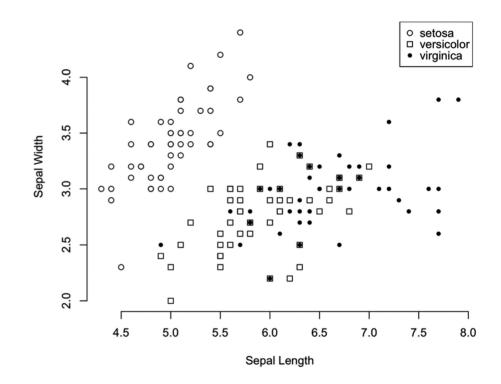


- Der IRIS Datensatz beinhaltet Variablen zu den Schwertlilienarten <u>setosa</u>, <u>versicolor</u> und <u>virginica</u>
- Ziel: Eine manuelle Identifikation der Schwertlilienarten über die unabhängigen Variablen und die Funktion <u>subset(...)</u>

- Variablen im IRIS Datensatz
  - $-(x_1)$  Sepal.Length
  - (x<sub>2</sub>) Sepal.Width
  - (X<sub>3</sub>) Petal.Length
  - (X<sub>4</sub>) Petal.Width
  - (Y) Species

Herausforderung:

 Schwertlilienarten
 können in Länge und
 Breite ihrer Blätter
 übereinstimmen



 Die Schwertlilienart setosa kann bspw. wie folgt identifiziert werden:

```
identification <- subset(iris, Petal.Length<="2" &
Petal.Width<="1", select=c(Species))
summary(identification)</pre>
```

- Alle erforderlichen <u>Funktionen</u>, <u>Objekte</u> und <u>Argumente</u> sind in der Übungsaufgabe hervorgehoben und können entsprechend in die eigene Syntax überführt werden
- Entsprechend sind noch versicolor und virginica zu identifizieren

#### (I) GRUNDLAGEN



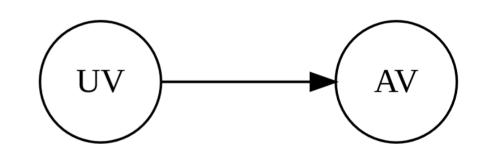
- Analysemodelle visualisieren die theoretisch fundierten Zusammenhänge zwischen allen unabhängigen und abhängigen Variablen
- Dadurch werden <u>Haupt- und Nebenhypothesen</u> abgebildet und mögliche <u>Interdependenzen</u> offengelegt

- Mit dem Package <u>DiagrammeR</u> können Analysemodelle in R visualisiert werden
- Dazu muss DiagrammeR zunächst über <u>install.packages(...)</u> und <u>library(...)</u> installiert und aktiviert werden
- Diese Installation erfordert <u>Dependencies</u>

- Funktionen aus Packages lassen sich auch über :: aufrufen
- Es können <u>label</u>
   vergeben und mit <u>-></u>
   verbunden werden

```
DiagrammeR::grViz("
digraph {graph [layout = circo]
node [shape = circle]
A [label = 'UV']
B [label = 'AV']
edge []
A -> B} ")
```

- Hypothese: Je mehr (weniger) UV, desto mehr (weniger) AV
- Dies ist ein Beispiel für eine spezifische Hypothese

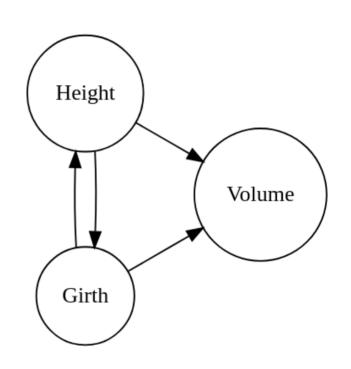


- Rückblick: Der TREES Datensatz beinhaltet die Variablen <u>Girth</u>, <u>Height</u> und <u>Volume</u>
- In der vorherigen Analyse konnte festgestellt werden, dass alle drei Variablen voneinander abhängen

- Entsprechende Hypothesen könnten lauten
  - Je mehr Girth [Height], desto mehr Height [Girth]
  - Je mehr Height, desto mehr Volume
  - Je mehr Girth, desto mehr Volume
- Für Girth [Height] und Height [Girth] liegt eine unspezifische Hypothese vor

- Demnach beeinflussen sich die unabhängigen Variablen gegenseitig; Dies nennt man Interdependenzen
- Über die multivariate Statistik konnte ermittelt werden, dass Girth einen stärkeren Einfluss auf Volume nimmt

```
DiagrammeR::grViz("
digraph {graph [layout = circo]
node [shape = circle]
A [label = 'Girth']
B [label = 'Height']
C [label = 'Volume']
edge []
A \rightarrow C
A \rightarrow B
B \rightarrow A
B \to C ")
```



- Es empfiehlt sich, <u>zu jedem neuen Datensatz</u> zunächst <u>ein Analysemodell</u> zu skizzieren
- Analysemodelle <u>vereinfachen die Interpretation</u> der Befunde aus der bivariaten und multivariaten Statistik

#### (I) GRUNDLAGEN

#### **SWISS Datensatz**



 Während der TREES Datensatz einen Einblick in die Teilbereiche der Statistik ermöglicht, bietet der SWISS Datensatz darüber hinaus einen Einblick in die <u>Moderationseffekte der</u> <u>unabhängigen Variablen</u>

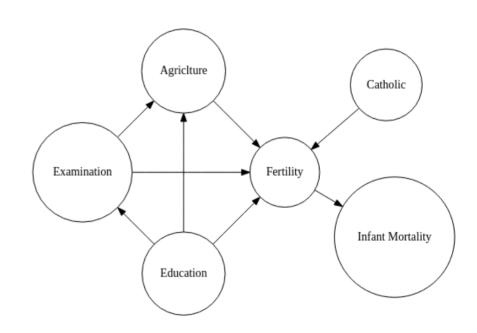
 Die Funktionen <u>head(...)</u> und <u>help(...)</u> zeigen erste Details zum SWISS Datensatz:

```
# Datensatz einsehen
head(swiss)
help(swiss)
```

- Unabhängige <u>Variablen</u>
  - (X<sub>1</sub>) Agriculture
  - (X<sub>2</sub>) Examination
  - (X<sub>3</sub>) Education
  - (X<sub>4</sub>) Catholic
  - (X<sub>5</sub>) Infant.Mortality

- Abhängige <u>Variable</u>
  - (Y) Fertility
- Die Variablen ergeben ein <u>komplexes</u> <u>Analysemodell</u>

- Variablenmoderation
  - Agriculture
  - Examination
  - Education
- Auswirkungen auf die multivariate Statistik



 Für die univariate Statstik stehen erneut die Funktionen <u>summary(...)</u> und <u>boxplot(...)</u> zur Verfügung:

```
# Univariate Statistik 
summary(swiss) 
boxplot(swiss)
```

 Moderationseffekte lassen sich über einen Vergleich zwischen bivariater und multivariater Statistik erkennen # Bivariate Statistik cor(swiss)

# Multivariate Statistik Im(Fertility~., data=swiss)

- Mögliche <u>Folgen der Moderationseffekte</u>
  - Veränderung bei den Effektstärken
  - Veränderung bei den Vorzeichen
- Die korrekte Interpretation von Moderationseffekten erfordert eine theoretische Fundierung

#### (I) GRUNDLAGEN



- Für die 2. Übungsaufgabe steht der MTCARS Datensatz zur Verfügung
- Dieser beinhaltet Variablen zu den technischen <u>Eigenschaften</u> von 32 verschiedenen <u>Automobilen</u> aus dem Jahr 1974

- Unabhängige <u>Variablen</u>
  - $-(X_1)$  cyl  $-(X_6)$  qsec
  - $(X_2) disp (X_7) vs$
  - $(X_3) hp (x_8) am$
  - $-(X_4)$  drat  $-(x_9)$  gear
  - $(X_5)$  wt  $(x_{10})$  carb

- Abhängige <u>Variablen</u>
  - (Y<sub>1</sub>) mpg
  - $(X_6) \rightarrow (Y_2)$  qsec
- Auch im MTCARS
   Datensatz ergeben die
   Variablen ein komplexes
   Analysemodell

#### Variablenmoderation

- vs

- cyl

- hp

- disp

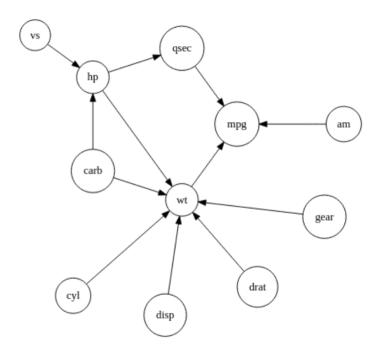
- carb

- drat

- wt

- gear

- qsec



 Zur Analyse von naturwissenschaftlichen Datensätzen kann im Rahmen der multivariaten Statistik auf die Funktion step(...) zurückgegriffen werden:

```
# Multivariate Statistik step(lm(data=mtcars, mpg~.), trace=0, steps=11)
```

- Alle erforderlichen <u>Funktionen</u>, <u>Objekte</u> und <u>Argumente</u> sind wieder in der Übungsaufgabe hervorgehoben und können entsprechend in die eigene Syntax überführt werden
- Anstatt alle unabhängigen Variablen manuell einzugeben, kann auch ein <u>.</u> gesetzt werden

#### (I) GRUNDLAGEN



- Häufig stellen <u>Befragungen</u> die Grundlage für zu analysierende Datensätze dar
- Tools zur Erstellung von Onlinebefragungen und zur Berechnung der Stichprobengröße
  - https://www.soscisurvey.de/
  - https://www.surveymonkey.de/mp/sample-size-calculator/

 Wenn Daten mittels einer Befragung und / oder Onlinebefragung gewonnen werden sollen, so sind insgesamt acht <u>Herausforderungen in der</u> <u>Interpretation dieser Daten</u> zu berücksichtigen

#### Reaktivität

- Wenn Teilnehmerinnen und Teilnehmer den Forschungsgegenstand kennen, dann kann sich daraufhin ihr (Antwort-)verhalten <u>verändern</u>
- Beispiel: Gesundheitsbezogenes Verhalten

- Soziale Erwünschtheit
  - Eine mögliche Verzerrung des Antwortverhaltens aufgrund der Annahme der Teilnehmerinnen und Teilnehmer über mit der Antwort in Verbindung stehende <u>Normen</u> und / oder <u>Erwartungen</u>
  - Beispiel: Deviantes Verhalten

- Schweigeverzerrung
  - Teilnehmerinnen und Teilnehmer können ein anderes Antwortverhalten aufweisen als nicht-Teilnehmerinnen und nicht-Teilnehmer
  - Beispiel: Produktbewertungen im Internet

- Selektive Aufmerksamkeit
  - Menschliche Wahrnehmung ist ein selektiver Prozess, orientiert an <u>vertrauten Mustern und</u> <u>Strukturen</u> und somit hinsichtlich der zu verarbeitenden Informationsmenge begrenzt
  - Beispiel: Eltern sehen mehr Gefahrenquellen

- Tendenz zur Mitte
  - Bei mehrstufigen Antwortskalen bewirkt die Tendenz zur Mitte ein <u>Antwortverhalten in der Mitte</u>
  - Beispiel: Unreflektierte Befragungsteilnahme

- Tendenz zur Milde / Härte
  - Bei mehrstufigen Antwortskalen erfolgt an Stelle eines objektiven Antwortverhaltens eine subjektive Unter- / Überbewertung eines realen Phänomens
  - Beispiel: Sympathie bei Lehrevaluationen

- Retrospektionseffekt
  - Erlebnisse und Ereignisse k\u00f6nnen im R\u00fcckblick, bspw. am n\u00e4chsten Tag, positiver oder negativer bewertet werden, als in der Situation selbst
  - Beispiel: Verkehrsunfall

- Rückschaufehler
  - Unzutreffende Erinnerungen, wenn Menschen, nachdem Sie den tatsächlichen Ausgang eines Ereignisses erfahren haben, sich falsch an ihre frühere Vorhersage des Ausgangs erinnern
  - Beispiel: Ausgang von politischen Wahlen

#### (II) FORMELSAMMLUNG

# Formeln und Verteilungstabellen



- Zusammenfassung der grundlegenden Formeln und Verteilungstabellen für die gängigsten statistischen Analyseverfahren
- Dabei <u>bauen</u> die Formeln der univariaten, bivariaten und multivariaten Statistik <u>aufeinander auf</u>

- Die <u>korrigierte Stichprobenvarianz</u> wird bspw. für die <u>Standardabweichung</u> benötigt
- Tipp: Die Standardabweichung lässt sich mit einer Taste auf dem Taschenrechner aufrufen

- Gleichermaßen erfordert der <u>Korrelations-</u> <u>koeffizient</u> die <u>Standardabweichung</u>
- Tipp: Vorherige Ergebnisse können einfach in die nächste Formel überführt werden

- Schließlich führen die <u>Standardabweichung</u> und der <u>Korrelationskoeffizient</u> zu den <u>linearen</u> <u>Regressionsgewichten</u>
- Tipp: Dabei können auch die Ergebnisse aus mehreren Formeln kombiniert werden

- Für manche Formeln sind darüber hinaus <u>Verteilungstabellen</u> mit statistischen Referenzwerten erforderlich
  - Chi-Quadrat-Test
  - t-Test

# (III) DATENANALYSE



- Skalenniveaus werden häufig auch als <u>Messniveaus</u> bezeichnet
- Sie beziehen sich auf den <u>Detailgrad der</u> <u>Merkmalsausprägungen</u> einer Variable

- Geeignete Skalenniveaus sind in der Lage, die unterschiedlichen <u>Facetten eines realen</u> <u>Phänomens ausreichend abzubilden</u>
- Für eine zielführende Datenanalyse mit R ist die Differenzierung von <u>vier unterschiedlichen</u> <u>Skalenniveaus</u> erforderlich

- Von der Nominalskala zur Verhältnisskala <u>nimmt</u> der Detailgrad der Merkmalsausprägungen zu
  - Nominalskala
  - Ordinalskala
  - Intervallskala
  - Verhältnisskala

- Die <u>Nominalskala</u> differenziert zwischen Merkmalsausprägungen <u>ohne vorgegebene Rangfolge</u>
  - Geschlecht (männlich, weiblich, divers) → 1 ≠ 2 ≠ 3
  - Postleitzahl (50939, 53225) → 1 ≠ 2
- Ein Spezialfall sind die dichotomen Variablen
  - Postleitzahl (50939 für Köln, 53225 für Bonn) → 0 ≠ 1

- Mit der <u>Ordinalskala</u> kann eine Rangfolge unter den Merkmalsausprägungen abgebildet werden, <u>ohne diese Rangfolge im Detail</u> <u>interpretieren zu können</u>
  - Einkommen (niedrig, hoch) → niedrig < hoch</li>
  - Schulnoten (sehr gut, gut) → sehr gut > gut

- Bei der <u>Intervallskala</u> lassen sich die unterschiedlichen <u>Merkmalsausprägungen mittels</u> <u>Zahlen</u> abbilden und interpretieren
  - Intelligenz (110, 100)  $\rightarrow$  110 10 = 100
  - Temperatur (30°C, 35°C)  $\rightarrow$  30°C + 5°C = 35°C
- Bei der Intervallskala gibt es keinen Nullpunkt

- Die <u>Verhältnisskala</u> erlaubt die Berechnung der <u>exakten Beziehung</u> zwischen den Merkmalsausprägungen
  - Einkommen (3000, 1500)  $\rightarrow$  3000 : 2 = 1500
  - Temperatur (20°K, 40°K)  $\rightarrow$  20°K \* 2 = 40°K
- Die Verhältnisskala besitzt einen Nullpunkt

- Höhere Skalenniveaus können in niedrigere Skaleniveaus überführt werden
- Umgekehrt ist dies jedoch nicht möglich
  - Rauchverhalten (ja, nein) → Anzahl an Zigaretten?
  - Einkommen (hoch) → Betrag in Euro?

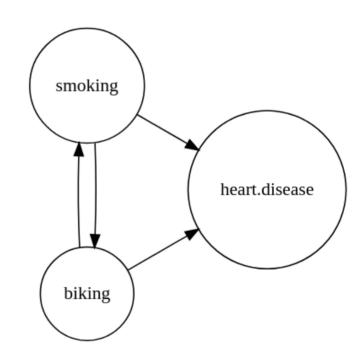
# (III) DATENANALYSE



 Für die univariate Statistik wird auf den externen HEART Datensatz zurückgegriffen, der über die Funktion <u>read.csv(...)</u> eingelesen werden kann:

```
# Datensatz einlesen
heart.data <-
read.csv("https://raw.githubusercontent.com/statistical-
thinking/datasets/main/heart.data.csv")
```

- Dieser beinhaltet zwei unabhängige und eine abhängige Variable
  - (X<sub>1</sub>) biking
  - (X<sub>2</sub>) smoking
  - (Y) heart.disease



 Wie bisher geben die Funktionen dim(...), summary(...) und boxplot(...) einen ersten Überblick über die Variablen:

```
# Univariate Statistik
dim(heart.data)
summary(heart.data)
boxplot(heart.data)
```

 Neben dem <u>Minimum</u> und <u>Maximum</u> der Merkmalsausprägungen werden zusätzlich <u>Quartile</u> und <u>Lagemaße</u> (hier: der Median und das arithmetische Mittel) ausgewiesen

- Quartile geben über <u>Schwellenwerte</u> (1. Quartil, Median, 3. Quartil) Auskunft zu vier Gruppen
  - ≤ 25% der Fälle
  - ≤ 50% der Fälle
  - ≤ 75% der Fälle
  - ≤ 100% der Fälle

- <u>Lagemaße</u> sollen die <u>zentrale Tendenz</u> eines Datensatzes zum Ausdruck bringen
  - Arithmetisches Mittel
  - Median
  - Modus

- Arithmetisches Mittel
  - Auch als <u>Durchschnittswert</u> oder <u>Mean</u> bekannt
  - Summe der Merkmalsausprägungen dividiert durch die Anzahl der Fälle

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + x_3 + \dots + x_n)$$

#### Median

- Auch als **Zentralwert** bekannt
- Die Merkmalsausprägung in der Mitte einer aufsteigend sortierten Liste aller Merkmalsausprägungen

$$\tilde{x}_{ungerade} = x_{\frac{n+1}{2}}$$
 bzw.  $\tilde{x}_{gerade} = \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$ 

#### • Modus

- Auch als Modalwert bekannt
- Bezeichnet die häufigste Merkmalsausprägung innerhalb einer Variable

$$\bar{x}_d = H\ddot{a}ufigster\ Beobachtungswert$$

- Sowohl die Quartile als auch die Lagemaße geben <u>Auskunft über die Verteilung der</u> <u>Merkmalsausprägungen</u> einer Variable
- Allerdings lässt sich aus diesen Funktionen nicht ableiten, wie häufig einzelne Merkmalsausprägungen in einer Variable vertreten sind

Auskunft über die <u>Häufigkeit einzelner</u>
 <u>Merkmalsausprägungen</u> einer Variable liefert
 ein <u>Histogramm</u> in Verbindung mit der
 dazugehörigen <u>Dichtefunktion</u>

• Ein Histogramm kann zunächst über die Funktionen <u>hist(...)</u> aufgerufen werden:

```
# Verteilungsfunktion und Lagemaße hist(heart.data$heart.disease, freq=FALSE, breaks=40, ylim=c(0,0.10), xlim=c(-5,26))
```

 Die dazugehörige Dichtefunktion ergibt sich daraufhin aus der Funktion <u>curve(...)</u>:

```
curve(dnorm(x,mean=mean(heart.data$heart.disease),
sd=sd(heart.data$heart.disease)), add=TRUE, lwd=5)
```

 Anschließend lassen sich über die Funktion abline(...) die Lagemaße ergänzen:

```
abline(v=10.17, col="red") # Mean
abline(v=10.38, col="green") # Median
abline(v=6.75, col="blue") # Modus
```

• Für eine passende Legende wird auf die Funktion <u>legend(...)</u> zurückgegriffen:

```
legend(19, 0.1, legend=c("Mean", "Median", "Modus"), col=c("red", "green", "blue"), lty=1)
```

- Die Dichtefunktion der Variable <u>heart.disease</u> entspricht der sogenannten <u>Normalverteilung</u>
  - glockenförmig
  - asymptotisch
  - symmetrisch
  - unimodal

- Dabei gilt: Das Integral der Dichtefunktion f(x) ist die sogenannte <u>Verteilungsfunktion F(x)</u>
- Die Verteilungsfunktion gibt an, wie groß die <u>Wahrscheinlichkeit</u> ist, dass die <u>Merkmals-</u> <u>ausprägung ≤ x</u> ist

- Neben den Lagemaßen lassen sich auch die <u>Streuungsmaße</u> in den Plot überführen
- Zunächst die Funktionen var(...) und sd(...):

```
# Varianz und Standardabweichung (Teil 1) var(heart.data$heart.disease) sd(heart.data$heart.disease)
```

- Streuungsmaße beziehen sich auf die Streubreite einzelner Merkmalsausprägungen um ausgewiesene Lagemaße (häufig das arithmetische Mittel)
  - Korrigierte Stichprobenvarianz
  - Standardabweichung

- Korrigierte Stichprobenvarianz
  - In vielen Fällen einfach Varianz genannt
  - Mittlere quadratische Abweichung vom arithmetischen Mittel

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

- Standardabweichung
  - Wurzel aus der korrigierten Stichprobenvarianz
  - Durchschnittliche Entfernung aller Merkmalsausprägungen vom arithmetischen Mittel

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

- Mittels der Standardabweichung kann auf die 68-95-Regel zurückgegriffen werden
  - 68% der Fälle sind ± eine Standardabweichung vom arithmetischen Mittel entfernt
  - 95% der Fälle sind ± zwei Standardabweichungen vom arithmetischen Mittel entfernt

 Zunächst wird für den Plot erneut auf die Funktion <u>hist(...)</u> zurückgegriffen:

```
# Varianz und Standardabweichung (Teil 2)
hist(heart.data$heart.disease, freq=FALSE, breaks=40,
ylim=c(0,0.10), xlim=c(-5,26))
```

 Zur besseren Veranschaulichung wird auch die Dichtefunktion über die Funktion <u>curve(...)</u> wieder ergänzt:

```
curve(dnorm(x,mean=mean(heart.data$heart.disease),
sd=sd(heart.data$heart.disease)), add=TRUE, lwd=5)
```

 Die Funktion <u>abline(...)</u> ermöglicht daraufhin die Anwendung der <u>68-95-Regel</u>:

```
abline(v=5.6, col="red") # Untere Grenze (68 %)
abline(v=14.74, col="red") # Obere Grenze (68 %)
abline(v=1.03, col="green") # Untere Grenze (95 %)
abline(v=19.31, col="green") # Obere Grenze (95 %)
```

 Schließlich verdeutlicht die Funktion <u>legend(...)</u> die Bereiche, in dem sich <u>68% bzw. 95% der</u> <u>Merkmalsausprägungen</u> befinden:

```
legend(20, 0.1, legend=c("68 %", "95 %"), col=c("red", "green"), lty=1)
```

# (III) DATENANALYSE



 Der z-Wert ist die Differenz eines Rohwertes vom arithmetischen Mittel dividiert durch die Standardabweichung:

$$z = \frac{x - x}{S_x}$$

- Mit der zugrundeliegenden z-Transformation werden Rohwerte in Normwerte überführt
- Normwerte ermöglichen gegenüber Rohwerten einen <u>standardisierten Vergleich</u>

- Zwei Schüler haben in unterschiedlichen Fächern an der <u>PISA-Studie</u> teilgenommen
  - Moritz in Mathematik
  - Fritz in Deutsch
- Beide haben 620 Punkte erzielt
  - Wer hat besser abgeschnitten?

- Für die z-Transformation benötigte Angaben
  - Durchschnittliche Punktezahl in Mathematik (490)
  - Standardabweichung in Mathematik (99,4)
    - -----
  - Durchschnittliche Punktezahl in Deutsch (484)
  - Standardabweichung in Deutsch (110,9)

- Moritz in Mathematik
  - $-z_{M} = (620 490)/99,4$
  - $-z_{\rm M}=1.31$

- Fritz in Deutsch
  - $-z_F = (620 484)/110,9$
  - $-z_F = 1.22$

- Fritz hat in der Differenz zur durchschnittlichen Punktezahl in Deutsch mehr Punkte als Moritz in seiner Vergleichsgruppe erzielt
- Die unterschiedlichen Standardabweichungen in den beiden Vergleichsgruppen bedingen aber einen höheren z-Wert zugunsten von Moritz

• Demnach hat Moritz ( $z_M = 1,31$ ) besser abgeschnitten als Fritz ( $z_F = 1,22$ )

- Die <u>z-Transformation einer Variable</u> führt zur Standardisierung derselben, so dass aus einer Normalverteilung eine <u>Standardnormalverteilung</u> mit  $\mu$  = 0 und  $\sigma$  = 1 wird
- Tipp: Die Funktion scale(...) führt zu Normwerten

# (III) DATENANALYSE



- Die bivariate Statistik ermöglicht das <u>Erklären</u> der <u>Zusammenhänge</u> zwischen <u>zwei Variablen</u>
- Gängige Verfahren der bivariaten Statistik
  - Korrelation (und Kovarianz)
  - Chi-Quadrat-Test
  - t-Test

- Die <u>Korrelation</u> misst die <u>Stärke des (linearen)</u>
   <u>Zusammenhangs</u> von zwei Variablen
- Der Korrelationskoeffizient ist <u>ungerichtet</u>, d.h. die Richtung der Stärke des Zusammenhangs muss theoretisch begründet werden
- Korrelationen sind kein Beweis für Kausalität

 Formal setzt sich der Korrelationskoeffizient aus der <u>Kovarianz</u> und den dazugehörigen <u>Standardabweichungen</u> zusammen

$$r_{xy} = \frac{\hat{\sigma}_{xy}}{s_x * s_y} = \frac{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2 * \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (y_i - \bar{y})^2}}}$$

• Die zugrundeliegende <u>Kovarianz</u> ist eine Erweiterung der empirischen <u>Stichproben-</u> <u>varianz</u> um eine <u>zweite Variable</u>

$$\hat{\sigma}_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

- Demnach ist die <u>Kovarianz</u> primär dazu geeignet, einen <u>positiven oder negativen</u> <u>Zusammenhang</u> zu identifizieren
- Eine <u>standardisierte Interpretation</u> dieses Zusammenhangs ist allerdings nur mit der <u>Korrelation</u> möglich

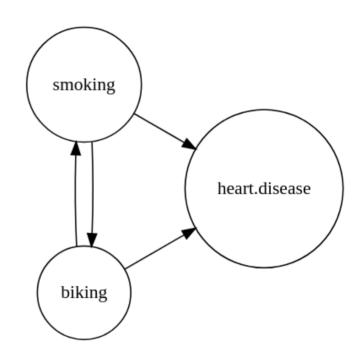
- Die Standardisierung bedingt die Ausprägungen des Korrelationskoeffizienten
  - perfekt negativer Zusammenhang (- 1,00)
  - perfekt positiver Zusammenhang (+ 1,00)
  - kein Zusammenhang (± 0,00)

- Der <u>Blick für die Stärke des Zusammenhangs</u> einer Korrelation zwischen zwei Variablen lässt sich natürlich auch ganz anschaulich trainieren
  - http://www.guessthecorrelation.com/

 Für die Umsetzung in R kann wieder auf den externen HEART Datensatz über die Funktion read.csv(...) zurückgegriffen werden:

```
# Datensatz einlesen
heart.data <-
read.csv("https://raw.githubusercontent.com/statistical-
thinking/datasets/main/heart.data.csv")
```

- Diesmal geht es um den <u>Zusammenhang</u> der <u>Variablen</u>
  - (X<sub>1</sub>) biking
  - (X<sub>2</sub>) smoking
  - (Y) heart.disease



 Bei überschaubaren Datensätzen lassen sich die Korrelationskoeffizienten in der <u>Korrelations-</u> <u>matrix</u> über die Funktion <u>cor(...)</u> aufrufen:

```
# Korrelation (Teil 1) cor(heart.data)
```

 Der Zusammenhang zwischen zwei Variablen lässt sich darüber hinaus mit den Funktionen plot(...) und abline(...) visualisieren:

```
# Korrelation (Teil 2)
plot(heart.data$heart.disease~heart.data$biking)
abline(lm(heart.data$heart.disease~heart.data$biking),
col="green")
```

- Demnach setzt die (lineare) Korrelation mindestens intervall- oder verhältnisskalierte Variablen voraus
- Liegen hingegen zwei <u>nominalskalierte</u>
   <u>Variablen</u> vor, kann auf den <u>Chi-Quadrat-Test</u>
   zurückgegriffen werden

 Weil der HEART Datensatz keine nominalskalierten Variablen beinhaltet, werden diese zunächst über die Funktion ifelse(...) generiert:

```
# Chi-Quadrat-Test
nominal_y <- ifelse(heart.data$heart.disease > 10.17, 1, 0)
nominal_x <- ifelse(heart.data$smoking > 15.43, 1, 0)
```

Nominalskalierte
 Variablen lassen sich
 über den table(...)
 Befehl als <u>Tabelle</u>
 aufrufen:

table(nominal\_y, nominal\_x)

 Schließlich wird der <u>Chi-Quadrat-Test</u> mittels der Funktion <u>chisq.test(...)</u> durch-geführt:

chisq.test(nominal\_y,
nominal\_x)

 Die zugrundeliegende Formel des Chi-Quadrat-Tests vergleicht <u>beobachtete Werte</u> (n<sub>j</sub>) mit <u>erwarteten Werten</u> (n<sub>j0</sub>)

$$\chi^2 = \sum_{j=1}^m \frac{(n_j - n_{j0})^2}{n_{j0}}$$

- Die <u>beobachteten Werte</u> (n<sub>j</sub>) befinden sich in der zuvor generierten <u>Tabelle mit den nominal-</u> <u>skalierten Variablen</u>
- Ausgehend von den beobachteten Werten (nj) soll die <u>Hypothese</u> (H) untersucht werden, dass ein <u>Zusammenhang zwischen x und y</u> vorliegt

- Beim Chi-Quadrat-Test wird die Hypothese (H) allerdings nicht direkt überprüft, sondern indirekt über die Nullhypothese (H<sub>0</sub>)
- Die Nullhypothese (H<sub>0</sub>) geht von <u>keinem</u>
   Zusammenhang aus und bedingt dadurch die erwarteten Werte (n<sub>j0</sub>)

• Beobachtete Werte (n<sub>j</sub>)

• Erwartete Werte (n<sub>j0</sub>)

 $\begin{array}{ccc} & \underline{x_0} & \underline{x_1} \\ y_0 & 140 & 101 \\ y_1 & 105 & 152 \end{array}$ 

• H: Rauchen (x) <u>führt</u> zu Herzerkrankungen (y) ?

H<sub>0</sub>: Rauchen (x) <u>führt nicht</u>
 zu Herzerkrankungen (y)

- Wahrscheinlichkeit für eine Herzerkrankung
  - $P(y_1) = (105 + 152) / 498$
  - $P(y_1) = 51,6\%$
- Wahrscheinlichkeit für keine Herzerkrankung
  - $P(y_0) = 100\% 51,6\%$
  - $P(y_0) = 48,4\%$

- Wenn die Nullhypothese (H<sub>0</sub>) gilt, dann müssten diese Wahrscheinlichkeiten unabhängig vom Rauchverhalten (x<sub>0</sub> und x<sub>1</sub>) sein
  - 51,6% mit und 48,4% ohne Herzerkrankung bei x<sub>0</sub>
  - 51,6% mit und 48,4% ohne Herzerkrankung bei x<sub>1</sub>

- Die Wahrscheinlichkeiten sind also gleich verteilt und ergeben die erwarteten Werte (n<sub>j0</sub>)
  - Für  $x_0$ : 245 \* 51,6% = 126
  - Für  $x_0$ : 245 \* 48,4% = 119
  - Für  $x_1$ : 253 \* 51,6% = 131
  - Für  $x_1$ : 253 \* 48,4% = 122

• Beobachtete Werte (n<sub>j</sub>)

• Erwartete Werte (n<sub>j0</sub>)

 $x_0$   $x_1$   $x_2$  140 101  $x_2$  105 152

 $\frac{x_0}{y_0}$   $\frac{x_1}{122}$   $\frac{x_1}{y_1}$   $\frac{x_2}{126}$   $\frac{x_1}{131}$ 

• H: Rauchen (x) <u>führt</u> zu Herzerkrankungen (y)

 H₀: Rauchen (x) <u>führt nicht</u> zu Herzerkrankungen (y)

Der berechnete Chi-Quadrat-Wert beträgt

$$-\chi^{2} \approx (140 - 119)^{2} / 119 + (101 - 122)^{2} / 122 + (105 - 126)^{2} / 126 + (152 - 131)^{2} / 131 -\chi^{2} \approx 15$$

Der Chi-Quadrat-Wert wird mit den <u>Referenz-werten der Chi-Quadrat-Wert-Verteilungstabelle</u> verglichen

 Referenzwerte in Abhängigkeit von <u>Freiheits-</u> graden (dF) und <u>Irrtumswahrscheinlichkeiten</u> (α)

	1 - α						
Freiheitsgrade	00,85	00,90	00,95	00,975	00,99	00,995	
1	02,07	02,71	03,84	05,02	06,63	07,88	
2	03,79	04,61	05,99	07,38	09,21	10,60	
3	05,32	06,25	07,81	09,35	11,34	12,84	
4	06,74	07,78	09,49	11,14	13,28	14,86	
5	08,12	09,24	11,07	12,83	15,09	16,75	
()	()	()	()	()	()	()	

- Berechnung der <u>Freiheitsgrade</u> (dF)
  - dF = (Zeilen 1) \* (Spalten 1)
  - dF = (2-1) \* (2-1)
  - dF = 1
- Die <u>Irrtumswahrscheinlichkeit</u> steht für die Wahrscheinlichkeit, dass <u>H<sub>0</sub> irrtümlich abgelehnt</u> wird

- Ist der Chi-Quadrat-Wert größer als der Referenzwert, wird die Nullhypothese (H<sub>0</sub>) verworfen
- Dies bedeutet im vorliegenden Beispiel
  - 15 > 7,88 (Referenzwert)
  - 99,5% Sicherheit in der Ablehnung von H<sub>0</sub>
  - H: Rauchen (x) führt zu Herzerkrankungen (y)

- Ein ähnliches Verfahren unter Rückgriff auf Referenzwerte zeigt sich beim <u>t-Test</u>
- Dabei geht es um den Zusammenhang zwischen einer mindestens intervallskalierten abhängigen Variable und einer nominalskalierten unabhängigen Variable

- Die nominalskalierte unabhängige Variable ermöglicht dem t-Test u.a. <u>Gruppenvergleiche</u>
  - Arithmetisches Mittel der Gruppe x<sub>1</sub>
  - Arithmetisches Mittel der Gruppe x<sub>2</sub>
- Tipp: Die Gruppen können auch nach den Merkmalsausprägungen 0 oder 1 benannt sein

Formal ist dies durch den <u>standardisierten</u>
 <u>Vergleich</u> beider arithmetischer Mittel in Abhängigkeit von der <u>Anzahl der Fälle</u> (N) möglich

$$t = \frac{x_1 - x_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}$$

• Beim t-Test beinhaltet die <u>t-Wert-Verteilungs-</u> <u>tabelle</u> die entsprechen <u>Referenzwerte</u>

	1 - α						
Freiheitsgrade	00,85	00,90	00,95	00,975	00,99	00,995	
20	01,06	01,32	01,72	02,09	02,53	02,86	
40	01,05	01,30	01,68	02,02	02,42	02,70	
60	01,04	01,29	01,67	02,00	02,39	02,66	
80	01,04	01,29	01,66	01,99	02,37	02,63	
100	01,04	01,29	01,66	01,98	02,36	02,62	
()	()	()	()	()	()	()	

- Anwendung und Interpretation erfolgen analog dem für den Chi-Quadrat-Test vorgestellten Beispiel
- Tipp: Der hier vorgestellte t-Test ist eigentlich der Welch-Test, bei dem die intervallskalierte abhängige Variable nicht normalverteilt sein muss

 Zur Visualisierung eines Gruppenvergleichs werden im HEART Datensatz Gruppen über die Funktion <u>subset(....)</u> angelegt:

```
# t-Test (Teil 1)
low_smoking_areas <- subset(heart.data,
heart.data$smoking < 15.43)
high_smoking_areas <- subset(heart.data,
heart.data$smoking > 15.43)
```

• Über die Funktion <u>par(...)</u> werden die beiden Funktionen <u>boxplot(...)</u> zusammengefasst:

```
par(mfrow=c(1,2))
boxplot(low_smoking_areas[c(3)], ylim=c(0, 22),
main="Low Smoking Areas (0)")
boxplot(high_smoking_areas[c(3)], ylim=c(0, 22),
main="High Smoking Areas (1)")
```

 Schließlich lassen sich die arithmetischen Mittel der beiden Gruppen auch direkt über die Funktion <u>t.test(...)</u> vergleichen:

```
# t-Test (Teil 2)
t.test(heart.data$heart.disease~nominal_x)
```

# Klausurrelevant (!)

## Obligatorischer Exkurs

• Exponentialfunktion mit negativem Exponenten:

```
x <- seq(0,50,1)

y <- runif(1,5,5)*exp(-runif(1,0.1,0.1)*x)+rnorm(51,0,0.5)

a\_start <- 10

b\_start <- 2*log(2)/a\_start

m <- nls(y~a*exp(-b*x), start=list(a=a\_start, b=b\_start))

plot(x,y)

lines(x, predict(m), col="blue", lty=2,lwd=3)
```

# Klausurrelevant (!)

## Obligatorischer Exkurs

Michaelis-Menten-Gleichung:

```
x <- seq(0,50)

y <- (runif(1,10,20)*x)/(runif(1,0,10)+x)+rnorm(51,0,1)

a_start <- 10

b_start <- 2*log(2)/a_start

m <- nls(y\sim a*x/(b+x), start=list(a=a_start, b=b_start))

plot(x,y)

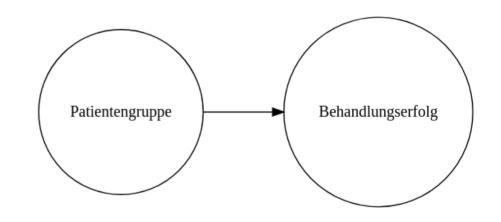
lines(x, predict(m), lty=2, col="blue", lwd=3)
```

## (III) DATENANALYSE



- In dieser Übungsaufgaben stehen <u>zwei</u> nominalskalierte Variablen zur Verfügung
  - Patientinnen und Patienten in Gruppe A oder B
  - Behandlung erfolgreich bzw. nicht erfolgreich
- Diesmal ist der Datensatz nicht in R enthalten

- Nominalskalierte Variablen
  - Patientengruppe (X)
  - Behandlungserfolg (Y)
- Hypothese: Behandlungserfolg in Gruppe A größer als in Gruppe B



 Über die Funktionen <u>cbind(...)</u> und <u>rbind(...)</u> kann der Datensatz aus der Übungsaufgabe in R überführt werden:

```
# Zusätzliche Syntax für R (Teil 1)
erfolgreich <- cbind(70, 55)
nicht_erfolgreich <- cbind(30, 45)
matrix <- rbind(erfolgreich, nicht_erfolgreich)
```

 Die Funktion <u>chisq.test(...)</u> ermittelt in R das exakte <u>Signifikanzniveau</u>:

```
# Zusätzliche Syntax für R (Teil 2) chisq.test(matrix)
```

- Ziele dieser Übungsaufgabe (ohne R)
  - Nullhypothese aufsetzen
  - Erwartete Werte bestimmen
  - Chi-Quadrat-Wert bestimmen
  - Nullhypothese mit <u>Referenzwert</u> pr

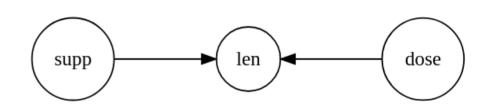
    üfen

## (III) DATENANALYSE



- Die Übungsaufgabe zum TOOTHGROWTH Datensatz kann mit dem <u>t-Test</u> gelöst werden
- Hierbei geht es um den Einfluss von <u>Orangensaft</u> oder purem <u>Vitamin C</u> als Supplement (Variable: supp) auf das Wachstum der Zähne (Variable: len) von Meerschweinchen

- Insgesamt stehen drei <u>Variablen</u> zur Verfügung
  - supp (X<sub>1</sub>)
  - dose (X<sub>2</sub>)
  - len (Y)



 Dabei sollen die Unterschiede zwischen den Gruppen der Meerschweinchen zunächst über die Funktionen <u>subset(...)</u> und <u>boxplot(...)</u> visualisiert werden:

```
# Vitamin C Meerschweinchen vc <- subset(ToothGrowth, supp=="VC") boxplot(vc[c(1,3)], main="VC")
```

## (III) DATENANALYSE

## Multivariate Statistik



- Das Prinzip der multivariaten Statistik, also der <u>Vorhersage</u> einer abhängigen Variable <u>über</u> <u>mehrere unabhängige Variablen</u>, kann über zwei Verfahren veranschaulicht werden
  - Lineare Regression
  - Logistische Regression

• Bei der <u>(einfachen) linearen Regression</u> wird über den <u>Korrelationskoeffizienten</u> und die dazugehörigen <u>Standardabweichungen</u> das <u>Regressionsgewicht</u> (b<sub>1</sub>) der unabhängigen Variable bestimmt

$$b_1 = r_{xy} * \frac{s_y}{s_x}$$

• Mit dem Regressionsgewicht (b<sub>1</sub>) und dem arithmetischen Mittel für y und x<sub>1</sub> lässt sich daraufhin der <u>Schnittpunkt der Regressionsgeraden mit der y-Achse</u> (a) bestimmen

$$a = \bar{y} - b_1 \bar{x}_1$$

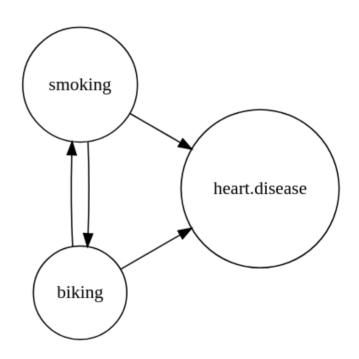
- Regressionsanalysen geben zusätzlich Auskunft über den Grad an Varianzaufklärung
- Der <u>Determinationskoeffizient</u> (R²) basiert bei der (einfachen) linearen Regression auf dem Korrelationskoeffizienten

$$R^2 = (r_{xy})^2$$

 Um dies zu veranschaulichen wird ein letztes Mal auf den externen HEART Datensatz über die Funktion <u>read.csv(...)</u> zurückgegriffen:

```
# Datensatz einlesen
heart.data <-
read.csv("https://raw.githubusercontent.com/statistical-
thinking/datasets/main/heart.data.csv")
```

- Diesmal geht es um die <u>Vorhersage</u> von heart.disease
  - (X<sub>1</sub>) biking
  - (X<sub>2</sub>) smoking
  - (Y) heart.disease



• Eine (einfache) lineare Regression kann über die Funktion <u>lm(...)</u> aufgerufen werden:

```
# Multivariate Statistik (Teil 1) (Im(heart.disease~biking, data=heart.data))
```

• Tipp: Über die äußere (...) lässt sich das Regressionsmodell bei Bedarf mit anderen Funktionen kombinieren, bspw. mit <u>summary(...)</u>

 Die Kennwerte dieses Regressionsmodells lassen sich mit den Funktionen plot(...) und abline(...) visualisieren:

plot(heart.data\$heart.disease~heart.data\$biking) abline(lm(heart.data\$heart.disease~heart.data\$biking), col="green")

- Das <u>Regressionsgewicht</u> (b<sub>1</sub>) entspricht dabei der <u>Steigung der linearen Funktion</u>
- Der <u>Schnittpunkt mit der y-Achse</u> (a) kann als theoretisch begründbarer Ausgangspunkt der linearen Funktion betrachtet werden

• Die Kennwerte der (einfachen) linearen Regression für die <u>unabhängige Variable smoking</u> lassen sich ebenfalls über die Funktion <u>lm(...)</u> aufrufen:

(Im(heart.disease~biking+smoking, data=heart.data))

• Die Funktion <u>lm(...)</u> ermöglicht auch die lineare Regression <u>mit beiden unabhängigen Variablen</u>:

```
# Multivariate Statistik (Teil 2)
reg1 <- (lm(heart.disease~., data=heart.data))
summary(reg1)
```

• Die Funktion <u>summary(...)</u> gibt dabei eine zusätzliche Modellzusammenfassung aus

- In der Modellzusammenfassung sind insbesondere die <u>Residuen</u> von Interesse
- Dabei geht es um die <u>"Fehler" in der Vorher-sagegenauigkeit</u> des Regressionsmodells

- Die <u>Residuen</u> sollten bei der linearen Regression <u>normalverteilt</u> sein
- Dies lässt sich über die Modellzusammenfassung mit der Funktion <u>summary(...)</u> überprüfen

- Darüber hinaus sollten die <u>Residuen keine</u> <u>Muster beinhalten</u>, die auf eine nicht-lineare Vorhersagbarkeit schließen lassen
- Dies ist bspw. bei <u>Heteroskedastizität</u> der Fall: Die <u>Varianz der Residuen variiert mit zunehmenden bzw. abnehmenden Werten</u> des Prädiktors

 Die Residuen eines Regressionsmodells können über die Funkion plot(....) auf solche Muster hin überprüft werden:

plot(reg1\$residuals)

 Wenn die <u>abhängige Variable</u> lediglich als <u>nominalskalierte Variable</u> vorliegt, so kann auf die <u>logistische Regression</u> zurückgegriffen werden

 Zur Veranschaulichung wird über die Funktion ifelse(...) eine nominalskalierte abhängige Variable generiert:

```
# Multivariate Statistik (Teil 3)
nominal y <- ifelse(heart.data$heart.disease > 10.17, 1, 0)
```

 Die Herausforderung in der Interpretation des Zusammenhangs zeigt sich in der Visualisierung über die Funktion plot(...):

plot(nominal\_y~heart.data\$biking)

- Einer nominalskalierten abhängigen Variablen fehlen kleinschrittige Merkmalsausprägungen, so dass eine lineare Funktion nicht zur Vorhersage geeignet ist
- In diesem Fall wird auf eine <u>logistische Funktion</u> mit <u>S-förmigen Verlauf</u> zurückgegriffen

• Dafür wird in der Funktion glm(...) das Argument <u>family=binomial</u> spezifiziert:

```
reg2 <- glm(nominal_y~heart.data$biking, family=binomial) summary(reg2)
```

 Achtung: Die Estimates sind nicht wie lineare Regressionsgewichte zu interpretieren!

- Die nominalskalierte abhängige Variable wird in der logistischen Regression mit <u>logarithmierten</u> <u>Odds Ratios als Regressionsgewichte</u> vorhergesagt
- Dazu ein Beispiel: Die logarithmierte Odds Ratio ein nerdiger <u>Star Wars Fan</u> zu sein

- 7 von 10 Nerds sind Star Wars Fans
- 4 von 10 Normalos sind Star Wars Fans
  - 7 / 10 = 70% als <u>Wahrscheinlichkeit für Nerds</u>
  - 4 / 10 = 40% als Wahrscheinlichkeit für Normalos
  - -70% / (1 70%) = 2,33 als Odds für Nerds
  - -40% / (1 40%) = 0,66 als Odds für Normalos

- Die Odds von 2,33 für einen Nerd und 0,66 für KEINEN Nerd ergeben demnach...
  - 2,33 / 0,66 = 3,5 als <u>Odds Ratio</u>
  - In(3,5) = 1,25 als <u>logarithmierte Odds Ratio</u>
- Dieser Wert wird auch <u>Logit</u> genannt
  - Logit = In(Odds Ratio)

### Obligatorischer Exkurs

Logarithmierte Odds Ratio eines nerdigen Star Wars Fans:

### (III) DATENANALYSE



- Datensätze können <u>viele Variablen und Fälle</u> beinhalten, die sich dennoch (theoriegeleitet) <u>zusammenfassen</u> lassen
  - Variablen der Bildung vs. des Einkommens
  - Fälle im Schwimmverein vs. Basketballverein

- Bekannte <u>Verfahren der Komplexitätsreduktion</u>
  - Faktorenanalyse (bzgl. Variablen)
  - <u>Clusteranalyse</u> (bzgl. Fälle)
- Zunächst wird eine Einführung in die Faktorenanalyse gegeben

 Die Packages PSYCH und CORRPLOT beinhalten Funktionen für die <u>Faktorenanalyse</u> und können nach erfolgreicher Installation über die Funktion <u>library(...)</u> aufgerufen werden:

# Funktionsumfang von R erweitern library(psych) library(corrplot)

- Als Beispiel wird der BFI Datensatz des Packages PSYCH herangezogen
- Erste Details zum BFI Datensatz verraten die Funktionen dim(...) und summary(...):

```
# Deskriptive Statistik dim(bfi) summary(bfi)
```

 Fehlwerte (NA) werden über die Funktion <u>na.omit(...)</u> ausgeschlossen und <u>relevante</u> <u>Variablen</u> in den Spalten 11 bis 20 fokussiert:

```
data <- na.omit(bfi[c(11:20)])
dim(data)
summary(data)
```

- Variablen der <u>Extraversion</u>
  - (X<sub>1</sub>) E1
  - (X<sub>2</sub>) E2
  - (X<sub>3</sub>) E3
  - $-(X_4) E4$
  - (X<sub>5</sub>) E5

- Variablen des <u>Neurotizismus</u>
  - $-(X_6) N1$
  - $-(X_7) N2$
  - (X<sub>8</sub>) N3
  - $-(X_9) N4$
  - $-(X_{10})$  N5

• Fasst man die Variablen von <u>Extraversion</u> zusammen, lässt sich das arithmetische Mittel des Summenscores über die Funktion <u>mean(...)</u> aufrufen:

```
# Summenscores vergleichen
extraversion_sum <-
(data$E1+data$E2+data$E3+data$E4+data$E5)/5
mean(extraversion_sum)
```

 Analog ist dies für <u>Neurotizismus</u> über die Funktion <u>mean(...)</u> möglich:

```
neuroticism_sum <-
(data$N1+data$N2+data$N3+data$N4+data$N5)/5
mean(neuroticism_sum)
```

Die arithmetischen Mittel unterscheiden sich

 Zusätzlich bestätigt die Funktion cor(...), dass die Variablen stärker innerhalb ihrer theoretischen Konstrukte korrelieren:

```
# Bivariate Statistik cor_matrix <- cor(data) corrplot(cor_matrix)
```

 Schließlich führen die Funktionen <u>fa.parallel(...)</u> und <u>fa(...)</u> zur Faktorenanalyse:

```
# Faktorenanalyse
fa.parallel(data, fa="both")
factors <- fa(data, nfactors=2, rotate="varimax")
factors
```

- Die Funktion <u>plot(...)</u> visualisiert die Befunde:
   plot(factors)
- <u>Faktorenanzahl</u> und <u>Rotationsverfahren</u> lassen sich über die Funktion <u>fa(...)</u> variieren, sollten aber theoretisch begründet werden können:

fa(data, nfactors=3, rotate="varimax")

- Auswahl an Rotationsverfahren
  - Varimax (Reduktion der Varianz)
  - Promax (wie Varimax mit Standardisierung)
- Rotationsverfahren <u>drehen das zugrunde-</u> <u>liegende Koordinatensystem</u>, bis ein <u>Kriterium</u> erfüllt wird (bspw. Anzahl an Faktoren)

 Tipp: Ein Abgleich der Faktorenanalyse mit den Befunden der ursprünglichen Studie ist über die Funktion help(...) möglich:

help(bfi)

 Für die <u>Clusteranalyse</u> steht das Package CLUSTER zur Verfügung, welches nach erfolgreicher Installation über die Funktion <u>library(cluster)</u> aufgerufen wird:

# Funktionsumfang von R erweitern library(cluster)

 Hierbei wird auf den bereits bekannten IRIS Datensatz und die Variablen zu den Schwertlilienarten setosa, versicolor und virginica zurückgegriffen

- Variablen im IRIS Datensatz
  - (x<sub>1</sub>) Sepal.Length
  - (x<sub>2</sub>) Sepal.Width
  - (X<sub>3</sub>) Petal.Length
  - (X<sub>4</sub>) Petal.Width
  - (Y) Species < - Cluster (!)

 Die Funktionen <u>dim(...)</u> und <u>summary(...)</u> beschreiben den IRIS Datensatz:

```
# Deskriptive Statistik dim(iris) summary(iris)
```

 Mittels <u>z-Transformation</u> und der Funktion <u>scale(...)</u> werden die <u>unabhängigen Variablen</u> standardisiert:

```
# Relevante Variablen z-transformieren 
cluster_data <- scale(iris[c(1:4)]) 
summary(cluster_data)
```

 Danach werden über die Funktion <u>sum(...)</u> die <u>Within Cluster Sum of Squares</u> (WSS) herangezogen, um die Anzahl an Clustern zu ermitteln:

```
# Within Cluster Sum of Squares (WSS) ermitteln
wss <- (nrow(cluster_data)-1)*sum(apply(cluster_data,2,var))
for (i in 2:10) wss[i] <- sum(kmeans(cluster_data, centers=i)
$withinss)
```

 Die <u>WSS</u> lassen sich über die Funktion <u>plot(...)</u> wieder entsprechend visualisieren:

```
# Grafische Darstellung der Anzahl an Clustern (WSS) plot(1:10, wss, type="c", xlab="Cluster", ylab="WSS", main="Clusteranzahl in Abhängigkeit von WSS")
```

 Mit der entsprechenden Anzahl an Clustern wird die <u>Clusteranalyse</u> bspw. über die Funktion <u>kmeans(...)</u> aufgerufen:

```
# K-Means Algorithmus anwenden
k_means_cluster <- kmeans(iris[,-5], 3, nstart=30)
clusplot(iris, k_means_cluster$cluster, color=TRUE,
shade=TRUE, lines=0)
```

## Komplexitätsreduktion

• Die <u>Präzision</u> der Clusteranalyse lässt sich über die Funktion <u>table(...)</u> überprüfen:

```
# Präzision des K-Means Algorithmus table(iris$Species, k_means_cluster$cluster)
```

#### (III) DATENANALYSE

# 5. Übungsaufgabe



# 5. Übungsaufgabe

- Diese Übungsaufgabe orientiert sich am Beispiel zum <u>BFI Datensatz</u>
- Diesmal sind alle <u>fünf Faktoren mittels</u>
   <u>Faktorenanalyse zu identifizieren</u>

# 5. Übungsaufgabe

 Alle benötigten Informationen sind bereits in der Übungsaufgabe enthalten und können ggfls. über die Funktion help(...) spezifiziert werden: help(bfi)

 Tipp: Die Packages PSYCH und CORRPLOT setzen Dependencies voraus (!)

#### (IV) MACHINE LEARNING



- Die in einem Datensatz enthaltenen <u>Informationen</u> werden für <u>automatisierte Verallgemeinerungen</u> herangezogen
- Dazu werden die im Datensatz enthaltenen Informationen in <u>erlernbare Beispiele</u> überführt

- Algorithmen sollen dabei <u>Muster und Gesetz-mäßigkeiten</u> über die erlernbaren Beispiele hinaus <u>identifizieren</u> können
- Algorithmen sind <u>Handlungsvorschriften</u>, oftmals unter Rückgriff auf statistische Formeln

- Beim Machine Learning sind <u>zwei Vorgehens-</u> weisen voneinander zu unterscheiden
  - Supervised Machine Learning
  - Unsupervised Machine Learning

- Beim <u>Supervised Machine Learning</u> werden <u>Input und Output definiert</u> und mittels Algorithmen <u>in ein Modell überführt</u>
- Dieses Modell soll den Output vorhersagen
  - <u>Prediction</u> (bspw. Lineare Regression)
  - <u>Classification</u> (bspw. Logistische Regression)

- Beim <u>Unsupervised Machine Learning</u> versuchen die Algorithmen zunächst <u>Muster</u> <u>und Gesetzesmäßigkeiten</u> zu erfassen und daraufhin <u>in ein Modell zu überführen</u>
  - <u>Dimensionality Reduction</u> (bspw. Faktorenanalyse)
  - <u>Clustering</u> (bspw. Clusteranalyse)

- Mögliche Anwendungsgebiete
  - Bilderkennung (bspw. auf dem Smartphone)
  - Sentiment Analysis (bspw. auf Social Media)
  - eMail Klassifizierung (bspw. SPAM identifizieren)

#### (IV) MACHINE LEARNING

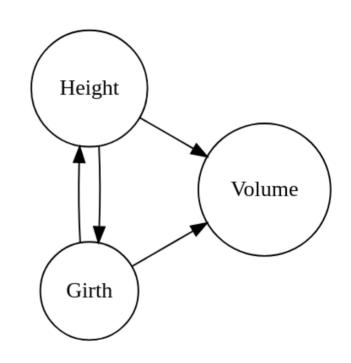
#### Praktische Einführung



- Für die praktische Einführung wird erneut auf den bereits bekannten TREES Datensatz zurückgegriffen
- Die Funktion dim(...) gibt dabei Auskunft über die Anzahl der enthaltenen Fälle:

dim(trees)

 Ein Algorithmus soll ausgehend von aus- gewählten Fällen ein Modell zur Vorher- sage des Volumens ermöglichen



• Über die Funktion <u>subset(...)</u> werden die Fälle 1 bis 5 sowie 26 bis 30 ausgewählt:

```
# Trainingsdatensatz anlegen
subset1 <- subset(trees[1:5,])
subset2 <- subset(trees[26:30,])
```

 Die beiden Subsets werden über die Funktion <u>rbind(...)</u> zu einem Trainingsdatensatz zusammengefasst:

```
training_data <- rbind(subset1, subset2) training_data
```

• Der <u>Algorithmus lernt am Trainingsdatensatz</u>

 Die verbleibenden Fälle 6 bis 25 sowie 31 werden über die Funktion <u>subset(...)</u> für den Validierungsdatensatz ausgewählt:

```
# Validierungsdatensatz anlegen
subset3 <- subset(trees[6:25,])
subset4 <- subset(trees[31,])
```

 Der über die Funktion <u>rbind(...)</u> zusammengefasste Validierungsdatensatz ist somit kleiner als der Trainingsdatensatz:

```
validation_data <- rbind(subset3, subset4)
validation_data</pre>
```

- Der Trainingsdatensatz ist in der Regel größer als der Validierungsdatensatz
  - 70 % Trainingsdatensatz
  - <u>30 % Validierungsdatensatz</u>
- Oftmals wird zusätzlich ein <u>unabhängiger</u> <u>Testdatensatz</u> zur Überprüfung eingesetzt

- Anhand der <u>Trainingsdaten</u> und unter Rückgriff auf Algorithmen wird ein Modell <u>aufgesetzt</u>
  - <u>Algorithmus</u> (bspw. eine lineare Regression)
  - <u>Modell</u> (bspw. Befunde einer linearen Regression)
- Anhand des <u>Validierungsdatensatzes</u> lässt sich das Modell <u>überprüfen</u>

 Für den TREES Datensatz eignet sich eine lineare Regression über die Funktion <u>Im(...)</u> zur Vorhersage der Variable Volume:

```
# Lineares Regressionsmodell als Algorithmus algorithm <- lm(data=training_data, Volume~Girth+Height)
```

 Eine Überprüfung ermöglicht die Funktion <u>predict(...)</u>, welche die <u>Befunde der linearen</u> <u>Regression</u> auf die <u>Merkmalsausprägungen des</u> <u>Validierungsdatensatzes</u> anwendet:

# Algorithmus auf Validierungsdatensatz anwenden validation <- predict(algorithm, validation\_data)

 Die Unterschiede zwischen den tatsächlichen Merkmalsausprägungen und den vorhergesagten Werten verdeutlicht die <u>durchschnittliche</u> <u>Differenz</u> über die Funktion <u>mean(...)</u>:

difference <- validation\_data-validation
mean(difference\$Volume)</pre>

 Die <u>Präzision</u> lässt sich über einen Vergleich der <u>durchschnittlichen Differenz</u> mit der <u>Standardab-</u> <u>weichung</u> über die Funktion <u>sd(...)</u> abschätzen: <u>sd(trees\$Volume)</u>

 Weniger Abweichungen als bei der Standardabweichung sind ein <u>präzises Modell</u>

#### (IV) MACHINE LEARNING

# **CARET Package**



- Das Package CARET stellt Algorithmen zum <u>Classification and Regression Training</u> in RStudio bereit
- Die Installation erfordert <u>Dependencies</u> und nimmt etwas mehr Zeit in Anspruch (!)

• Die Installation und Aktivierung erfolgt über die Funktionen <u>install.packages(...)</u> und <u>library(...)</u>:

# Zusatzprogramm CARET installieren und aktivieren install.packages("caret", dependencies=TRUE) library(caret)

 Im nächsten Schritt werden über den <u>Zuwei-</u> <u>sungspfeil</u> die unabhängigen (x) und abhängigen Variablen (y) als <u>Features</u> definiert:

```
# Features anlegen
x <- iris[,1:4]
y <- iris[,5]
```

• Die Funktion <u>featurePlot(...)</u> visualisiert die Merkmalsausprägungen der unabhängigen Variablen in Bezug auf die Schwertlilienarten:

```
# Grafische Darstellung der Features featurePlot(x=x, y=y, plot="box")
```

• Über die Funktion <u>shapes(...)</u> kann den Schwertlilienarten ein Symbol zur besseren Unterscheidung zugewiesen werden:

```
# Grafische Darstellung: SEPAL.WIDTH und SEPAL.LENGTH shapes=c(1,0,20) shapes <- shapes[as.numeric(iris$Species)]
```

• Die Funktionen <u>plot(...)</u> und <u>legend(...)</u> visualisieren das Zusammenspiel der unabhängigen Variablen:

```
plot(x=iris$Sepal.Length, y=iris$Sepal.Width, frame=FALSE, xlab="Sepal Length", ylab="Sepal Width", pch=shapes) legend("topright", legend=levels(iris$Species), pch=c(1,0,20))
```

 Das Package CARET ermöglicht über die Funktion <u>createDataPartition(...)</u> die Erstellung eines Trainings- und Validierungsdatensatzes:

```
# Training Data und Validation Data anlegen validation_index <- createDataPartition(iris$Species, p=0.80, list=FALSE) validation <- iris[-validation_index,] training <- iris[validation_index,]
```

 Trainings- und Validierungsdatensatz können über die Funktion <u>summary(...)</u> inspiziert werden:

# Training Data und Validation Data einsehen summary(validation) summary(training)

- Im nächsten Schritt erfordert das Package CARET die Spezifizierung einer Validierungsmethode
- Bei der Kreuzvalidierung wird die <u>Vorhersage-</u> genauigkeit eines <u>Modells</u> über <u>Teilmengen des</u> <u>Datensatzes</u> in mehreren Schritten <u>überprüft</u>

 Eine Kreuzvalidierung über zehn Schritte wird mit der Funktion trainControl(...) spezifiziert:

```
# Validierung festlegen
control <- trainControl(method="cv", number=10)
metric <- "Accuracy"
```

- Schließlich werden verschiedene Machine Learning <u>Algorithmen</u> aufgerufen
  - Linear Discriminant Analysis (LDA)
  - K-Nearest-Neighbor (KNN)
  - Random Forest (RF)

 Der <u>LDA-Algorithmus</u> versucht die aus den Merkmalsausprägungen zwei unabhängiger Variablen bestehenden <u>Datenpunkte mittels</u> einer linearen Funktion voneinander zu trennen

 Über die Funktion <u>train(...)</u> wird der LDA-Algorithmus aufgerufen:

```
# Linear Discriminant Analysis trainieren fit.lda <- train(Species~., data=training, method="lda", metric=metric, trControl=control)
```

 Der KNN-Algorithms verbindet naheliegende Datenpunkte ausgehend von ihrem Schwerpunkt der Verteilung und trennt diese von weiter entfernt liegenden Datenpunkten

• Über die Funktion <u>train(...)</u> wird der KNN-Algorithmus aufgerufen:

```
# K-Nearest Neighbor trainieren fit.knn <- train(Species~., data=training, method="knn", metric=metric, trControl=control)
```

 Der <u>RF-Algorithmus</u> kombiniert die Ergebnisse verschiedener <u>Entscheidungsbäume</u> durch <u>Variation einzelner Bedingungen unter den</u> <u>Merkmalsausprägungen</u>, um bestmögliche Entscheidungen zu treffen

 Schließlich wird auch der RF-Algorithmus über die Funktion train(...) aufgerufen:

```
# Random Forest trainieren fit.rf <- train(Species~., data=training, method="rf", metric=metric, trControl=control)
```

 Das Abschneiden der Algorithmen kann zunächst über die Funktion <u>resample(...)</u> aufgerufen und zwischengespeichert werden:

```
# Vergleich der Machine Learning Algorithmen results <- resamples(list(lda=fit.lda, knn=fit.knn, rf=fit.rf))
```

 Über die Funktionen <u>summary(...)</u> und <u>dotplot(...)</u> ist ein Vergleich der Algorithmen möglich:

summary(results) dotplot(results)

 Ist die Entscheidung für einen Algorithmus gefallen, kann dieser über die Funktion print(...) im Detail analysiert werden:

# Fokus auf besten Algorithmus print(fit.lda)

 Die <u>Validitätsüberprüfung</u> erfolgt schließlich über die Funktionen <u>predict(...)</u> und <u>confusionMatrix(...)</u>:

# Algorithmus auf Validation Data anwenden predictions <- predict(fit.lda, validation) confusionMatrix(predictions, validation\$Species)

## Obligatorischer Exkurs

- Weitere Datensätze, viele davon für Machine Learning Algorthmen geeignet, sind online zur kostenlosen Nutzung hinterlegt
  - https://www.kaggle.com/datasets/
  - https://openpsychometrics.org/\_rawdata/

#### (IV) MACHINE LEARNING

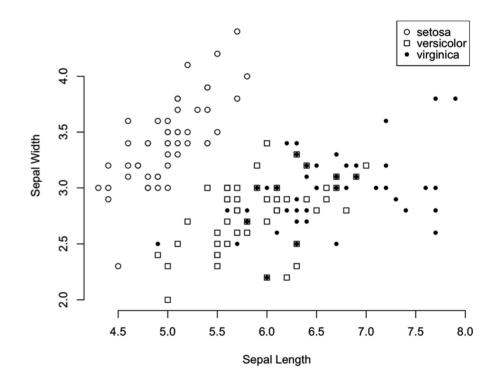
# 6. Übungsaufgabe



- Diese Übungsaufgabe reproduziert die zuvor vorgestellten Schritte zur <u>Anwendung von</u> <u>Machine Learning Algorithmen</u> am IRIS Datensatz
- Ziel: Der beste Machine Learning Algorithmus ist zu bestimmen

- Variablen im IRIS Datensatz
  - $-(x_1)$  Sepal.Length < --- Feature (!)
  - $-(x_2)$  Sepal.Width < --- Feature (!)
  - $(X_3)$  Petal.Length < --- Feature (!)
  - $-(X_4)$  Petal.Width < --- Feature (!)
  - (Y) Species

 Erinnerung an die Herausforderung: Schwertlilienarten können in Länge und Breite ihrer Blätter übereinstimmen



 Zusatzaufgabe: Wie schneiden die Machine Learning Algorithmen im Vergleich zu der Lösung aus der 1. Übungsaufgabe ab?

#### (IV) MACHINE LEARNING



- Machine Learning Algorithmen können nur im Rahmen der vorgegebenen Datensätze (oftmals <u>Stichproben</u>) trainiert werden
- Ein <u>Bias in den Daten</u> führt demnach zu einem <u>Bias des Modells</u>

- Mit <u>Resampling</u> werden Verfahren bezeichnet, die auf Basis einer Stichprobe <u>wiederholt kleine</u> <u>Stichproben</u> generieren
- Demnach werden <u>statistische Kennwerte</u> nicht einmalig berechnet, sondern für die Anzahl der wiederholten Stichproben <u>gemittelt</u>

- Bewährte Resamplingverfahren sind die Permutationstests und das Bootstrapping
  - Permutationstest: Variation der Reihenfolge aller
     Fälle einer Stichprobe zur Abschätzung der p-Werte
  - Bootstrapping: Wiederholte Ziehung einzelner Fälle einer Stichprobe mit Zurücklegen zur Abschätzung der weiteren statistischen Kennwerte

- Für <u>Machine Learning</u> Algorithmen sind insbesondere die weiteren <u>statistischen</u> <u>Kennwerte</u> von Bedeutung
- Bootstrapping: Beispiel zum <u>Regressions-gewicht</u> bei der linearen Regression

 Bestimmung des Regressionsgewichtes im MTCARS Datensatz über die Funktion <u>Im(...)</u>:

```
# Lineare Regression 
Im(mtcars$mpg~mtcars$wt)
```

• Regressionsgewicht (b<sub>Regression</sub>) einmalig ermittelt

 Zur Vorbereitung des Resamplingverfahrens und zur Reproduzierbarkeit der Befunde wird auf die Funktion <u>set.seed(...)</u> zurückgegriffen:

```
# Resampling set.seed(1701)
```

• Das Resampling (R) soll <u>1.000-fach wiederholt</u> werden und als <u>Regressionsgewicht</u> (b<sub>Bootstrap</sub>) angelegt werden:

```
R <- 1000
b <- vector(length=R)
```

Eine entsprechende Bootstrapping-Funktion:

```
for(i in 1:R){
boot.data <-
mtcars[sample(1:nrow(mtcars),size=100,replace=T),]
boot.result <-
summary(lm(boot.data$mpg~boot.data$wt))
b[i] <- boot.result$coefficients[2,1]
}
```

 Die Summe aller Regressionsgewichte (b<sub>Bootstrap</sub>) kann daraufhin über die Funktion <u>summary(...)</u> aufgerufen werden:

summary(b)

 Die Funktionen plot(...) und abline(...) stellen b<sub>Reg</sub> und b<sub>Bootstrap</sub> gegenüber:

```
# Verteilung der Regressionsgewichte plot(density(b), main="Stichprobenwiederholung der Regressionsgewichte", ylim=c(0.05,1.25), xlim=c(-7,-3.5)) abline(0,0,0,-5.344, col="red")
```

 Mit der Funktion <u>legend(...)</u> wird eine entsprechende Legende hinzugefügt:

```
legend("topright", inset=.01, legend=c("b-coefficient
(original) = -5.344"), col=c("red"), box.lty=0, lty=1:2,
cex=1.0)
```

 Tipp: Resamplingverfahren tragen nicht nur im Rahmen des Machine Learnings zur <u>besseren</u> <u>Interpretation statistischer Kennwerte</u> bei

## Obligatorischer Exkurs

- Tipp: Die <u>Bootstrapping-Funktion</u> lässt sich entsprechend adaptieren:
  - 1. Zusatzaufgabe: <u>p-Werte</u> für
     t.test(ToothGrowth\$len~ToothGrowth\$supp)
  - Zusatzaufgabe: <u>Chi-Quadrat-Werte</u> für mtcars\$binary\_mpg <- ifelse(mtcars\$mpg<=20,1,0) chisq.test(mtcars\$binary\_mpg,mtcars\$vs)

#### (IV) MACHINE LEARNING



 Die von den Machine Learning Algorithmen verwendeten <u>Muster und Gesetzmäßigkeiten</u> entsprechen <u>nicht immer einer theoretischen</u> <u>Fundierung</u> oder basieren auf einem <u>realen</u> <u>Phänomen</u>

- Schließlich gilt: Korrelation ≠ Kausalität
- Herausforderungen beim Machine Learning
  - Unzureichende <u>Datenqualität</u>
  - Übermäßige Optimierung
  - Limitationen bei der Durchführung
  - Fehlende <u>Mathematik- und Statistikkenntnisse</u>

- Unzureichende Datenqualität
  - Statistische Kennwerte und Machine Learning Algorithmen sind nur unter Berücksichtigung der Datenqualität <u>angemessen interpretierbar</u>
  - Fokus auf die <u>univariate Statistik</u> in Bezug auf die relevanten Variablen (Verteilung, Fehlwerte, etc.)

- Übermäßige Optimierung
  - Weniger Präzision der Machine Learning
     Algorithmen in Bezug auf die Trainings- und
     Validierungsdatensätze kann von Vorteil sein
  - Es gilt: Ein flexibler Dietrich ist manchmal besser als ein starrer Schlüssel

- <u>Limitationen</u> auf einen Algorithmus / Trainingsund Validierungsdatensatz
  - Der <u>Vergleich verschiedener Algorithmen</u> beinhaltet wertvolle Informationen zur Interpretation
  - Resampling vermag die <u>Varianz eines realen</u>
     <u>Phänomens</u> besser abzubilden als ein einmaliger
     Datensatz / eine einmalige Stichprobe

## Herausforderungen

- Fehlende Mathematik- / Statistikkenntnisse
  - Die zugrundeliegenden Analyseverfahren, bspw. die <u>Faktorenanalyse</u>, sind primär der <u>mathematischen</u> <u>Statistik</u> zuzuordnen
  - Bei Herausforderungen in der Anwendung empfiehlt sich ein Verständnis der zugrundeliegenden Funktionsweise, bspw. dem <u>Rechnen mit Vektoren</u>

# Herausforderungen

- Tipp: In vielen Fällen führen die <u>klassischen</u> <u>Analyseverfahren</u> aus der Statistik, bspw. die lineare Regression, zu <u>verlässlichen Befunden</u>
- <u>Machine Learning</u> Algorithmen sind demnach <u>nicht immer erforderlich</u>

## Herausforderungen

- Schließlich lauten zwei der wichtigsten Regeln für eine zielführende Datenanalyse mit R
  - Prediction is only as good as the sum of its parts
  - KISS: Keep It Short & Simple (!)



# Vicht klausurrelevant (!)

## Fakultativer Exkurs

- Tipp: Weiterführende Kursempfehlungen
  - https://www.coursera.org/
    - Mathematics for Machine Learning (Imperial College London)
    - <u>Data Science Specialization</u>
       (Johns Hopkins University)
    - Machine Learning (Stanford Online)

#### • <u>Simpson-Paradoxon</u>:

```
x <- c(3, 4, 5, 7, 8, 9)
y <- c(8, 9, 10, 2, 3, 4)
simpson < - data.frame(x, y)
plot(simpson, ylim=c(1, 11), xlim=c(1, 11))
abline(Im(data=simpson, y\sim x))
abline(lm(data=simpson[c(1, 2, 3),], y\sim x), col="red")
abline(lm(data=simpson[c(4, 5, 6),], y\sim x), col="green")
legend(8, 10.5, legend=c("Gruppe A", "Gruppe B"),
col=c("red", "green"), lty=1)
```

Parallelen zwischen t-Test und linearer Regression:

t.test(ToothGrowth\$len~ToothGrowth\$supp) Im(ToothGrowth\$len~ToothGrowth\$supp)

- Mittelwert der Gruppe 0 entspricht dem Schnittpunkt mit der y-Achse (a)
- Differenz zwischen Gruppe 0 und Gruppe 1 entspricht der Steigung der Regressionsgeraden (b<sub>1</sub>)

Moderationseffekt (Teil 1):

```
set.seed(1701)
data <- data.frame(
fan = rep(c("Star Wars", "Star Trek"), each = 30),
nerdindex = rep(c("1", "2", "3"), each = 10, times = 2),
tvindex = c(runif(10, -3, 3), runif(10, 0, 5), runif(10, 4, 6),
runif(10, -4, 2), runif(10, 0, 3), runif(10, 5, 8)))
```

Moderationseffekt (Teil 2):

```
interaction.plot(x.factor = data$nerdindex,
trace.factor = data$fan, response = data$tvindex,
fun = median, main = "Moderationseffekt (Schnittstelle)",
ylab = "TV-Index (Stunden pro Tag)",
xlab = "Nerd-Index (1=klein, 2=mittel, 3=groß)",
col = c("red", "blue"), lty = 2, lwd = 2,
trace.label = "Franchise")
```

- Moderationseffekt (Teil 3):
  - Schnittpunkte in der Interaktionsgrafik deuten an, dass zwischen den Variablen <u>nerdindex</u> und <u>fan</u> ein Moderationseffekt vorliegt
  - Dies wird in der Funktion <u>lm(...)</u> berücksichtigt:
     summary(lm(data=data, tvindex~nerdindex\*fan))

# Nicht klausurrelevant (!)

### Fakultativer Exkurs

Plot der ersten Seite (Cover):

```
library(mgcv)
library(lattice)
x <- rnorm(100)
y <- rnorm(100)
tab <- data.frame(x,y,z)
mod <- gam(z~te(x,y), data=tab)
z <- matrix(fitted(mod), ncol=10)
wireframe(z, drape=TRUE, colorkey=TRUE)</pre>
```

The End