

# Statistik

für Betriebswirtschaft, Internationales Management,  
Wirtschaftsinformatik und Informatik

## Sommersemester 2016

Veranstaltungen zur Statistik für BW/IM Sommersemester 2016				
Was?	Wer?	Wann?	Wo?	ab wann?
Vorlesung Statistik	Etschberger	Mi, 14.00-17.00	B2.14	16.03.2016
Vorlesung Statistik PLUS	Etschberger/Jansen	- Blocktermin -	?	?
Übung Statistik	Etschberger	Mi, 17.00-18.30	A1.10	30.03.2016
Übung Statistik	Jansen	Di, 11.30-13.00	W1.06	22.03.2016
Übung Statistik	Jansen	Di, 14.00-15.30	W2.14	22.03.2016
Übung Statistik	Jansen	Mi, 11.30-13.00	W2.11	30.03.2016
Übung Statistik	Jansen	Do, 14.00-15.30	W2.14	31.03.2016
Übung Statistik	Schneller	Do, 15.30-17.00	W3.03	31.03.2016
Übung Statistik	Schneller	Do, 15.30-17.00	W3.03	31.03.2016
Übung Statistik	Wins	Di, 14.00-15.30	J3.19	22.03.2016
Übung Statistik	Wins	Di, 15.30-17.00	J3.13	22.03.2016
Offener Statistikraum	Etschberger/Tutoren	?	?	?
Veranstaltungen für Teilnehmer der WiMa-Klausur im Juli 2016				
Was?	Wer?	Wann?	Wo?	ab wann?
Tutorium Mathematik	Burkart	Do 13.30-15.00	W1.06	07.04.2016
Tutorium Mathematik	Burkart	Do 15.00-16.15	W1.06	07.04.2016
Offener Matheraum	Jansen/Tutoren	?	?	?

HSA Statistik SS 2016 Sessionlist		
Datum	Statistik für IM/BW	Nr.
Mittwoch, 16. März 2016	Einführung, R Installation, Rstudio Einführung, Skalen	1
Mittwoch, 23. März 2016	univ. deskr. Stat., Quantile, Plots	2
Mittwoch, 30. März 2016	Streuung, Konzentrationsmaße	3
Mittwoch, 6. April 2016	Kontingenztabellen, Mosaikplots, Korrelation	4
Mittwoch, 13. April 2016	Preisindizes, lineare Regression	5
Mittwoch, 20. April 2016	Kombinatorik, Wahrscheinlichkeit	6
Mittwoch, 27. April 2016	Wahrscheinlichkeit, diskrete Zufallsvariablen	7
Mittwoch, 4. Mai 2016	Pyramid	
Mittwoch, 11. Mai 2016	Binomial-, Hypergeom.-, Poisson-Verteilung	8
Mittwoch, 18. Mai 2016	Stetige ZV, Gleichverteilung	9
Mittwoch, 25. Mai 2016	Normalverteilung, Verteilungsparameter	10
Mittwoch, 1. Juni 2016	Schätzfunktionen und Punktschätzer	11
Mittwoch, 8. Juni 2016	Konfidenzintervalle	12
Mittwoch, 15. Juni 2016	Tests	13
Mittwoch, 22. Juni 2016	Puffer, WH, Fragen zur Probekl.	14
Mittwoch, 29. Juni 2016	AW Prüfungswoche	

Prof. Dr. Stefan Etschberger  
Hochschule Augsburg



## 1 Statistik: Einführung

- Berühmte Leute zur Statistik
- Wie lügt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

## 2 Deskriptive Statistik

- Häufigkeiten
- Lage und Streuung
- Konzentration
- Zwei Merkmale
- Korrelation
- Preisindizes
- Lineare Regression

## 3 Wahrscheinlichkeitstheorie

- Kombinatorik
- Zufall und Wahrscheinlichkeit
- Zufallsvariablen und Verteilungen
- Verteilungsparameter

## 4 Induktive Statistik

- Grundlagen
- Punkt-Schätzung
- Intervall-Schätzung
- Signifikanztests

1. Einführung

2. Deskriptive Statistik

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

## Kursmaterial:

- ▶ Aufgabensatz (beinhaltet Aufgaben zu R)
- ▶ Handout der Folien
- ▶ Alle Folien inklusive Anmerkungen (nach der jeweiligen Vorlesung)
- ▶ Beispieldaten
- ▶ Alle Auswertungen als **R**-Datei






1. Einführung
2. Deskriptive Statistik
3. W-Theorie
4. Induktive Statistik

Quellen

Tabellen

## Literatur:

-  Bamberg, Günter, Franz Baur und Michael Krapp (2011). **Statistik**. 16. Aufl. München: Oldenbourg Verlag. ISBN: 3486702580.
-  Dalgaard, Peter (2002). **Introductory Statistics with R**. New York: Springer.
-  Fahrmeir, Ludwig, Rita Künstler, Iris Pigeot und Gerhard Tutz (2009). **Statistik: Der Weg zur Datenanalyse**. 7. Aufl. Berlin, Heidelberg: Springer. ISBN: 3642019382.

## Klausur:

- ▶ **Klausur** am Ende des Semesters
- ▶ Bearbeitungszeit: **90 Minuten**
- ▶ Erreichbare Punktzahl: 90
- ▶ R ist prüfungsrelevant: Siehe Anmerkungen in Übungsaufgaben!
- ▶ Hilfsmittel:
  - **Schreibzeug**,
  - **Taschenrechner**, der nicht 70! berechnen kann,
  - **ein** Blatt (DIN-A4, vorne und hinten beschrieben) mit handgeschriebenen Notizen (keine Kopien oder Ausdrücke),
- ▶ Danach (optional): Für Teilnehmer der **Statistik-Plus** Vorlesung noch eine 30-minütige Teilklausur über zusätzliche Inhalte (2 Wahlfachcredits zusätzlich möglich; Hilfsmittel TR und **ein** Blatt)  
<sub>2</sub>



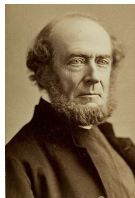
- 1 Statistik: Einführung
- 2 Deskriptive Statistik
- 3 Wahrscheinlichkeitstheorie
- 4 Induktive Statistik



- 1 **Statistik: Einführung**
  - Berühmte Leute zur Statistik
  - Wie lügt man mit Statistik?
  - Gute und schlechte Grafiken
  - Begriff Statistik
  - Grundbegriffe der Datenerhebung
  - R und RStudio

▶ **Leonard Henry Courteney (1832-1918):**

„There are three kinds of lies: lies, damned lies and statistics.“



▶ **Winston Churchill (1874-1965) angeblich:**

„Ich glaube nur den Statistiken, die ich selbst gefälscht habe.“



▶ **Andrew Lang (1844-1912):**

„Wir benutzen die Statistik wie ein Betrunkener einen Laternenpfahl: Vor allem zur Stütze unseres Standpunktes und weniger zum Beleuchten eines Sachverhalts.“



DRAWN BY BLAKE MERRIFIELD.

ENGRAVED BY J. F. KIRKING.



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen

## Morgens in Zeitung: Mehr Statistiken als Goethe und Schiller im ganzen Leben gesehen haben:

- ▶ Arbeitslosenzahlen wachsen
- ▶ Vogelgrippe breitet sich aus 14,28%
- ▶ 78,643% der Deutschen unzufrieden mit Löw
- ▶ Bundesbürger verzehrt 5,8 Liter Speiseeis pro Jahr
- ▶ Musiker leben länger als andere Leute
- ▶ Tennisspieler B hat noch nie gegen einen brilletragenden Linkshänder verloren, der jünger ist als er
- ▶ in New York schläft man am sichersten im Central Park

**Viele dieser Statistiken: Falsch, bewußt manipuliert oder unpassend ausgesucht.**

Fehlerquellen:

- ▶ Zahlenmanipulation
- ▶ irreführende Darstellung der Zahlen
- ▶ ungenügendes Wissen



### 1. Einführung

Berühmte Leute zur Statistik?

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen



## 1. Frage:

„Finden Sie, dass in einem Betrieb alle Arbeiter in der Gewerkschaft sein sollten?“

### Resultat:

- ▶ Dafür: 44%
- ▶ Dagegen: 20%
- ▶ Unentschieden: 36%

## 2. Frage:

„Finden Sie, dass in einem Betrieb alle Arbeiter in der Gewerkschaft sein sollten oder muss man es jedem einzelnen überlassen, ob er in der Gewerkschaft sein will oder nicht?“

### Resultat:

- ▶ Dafür: 24%
- ▶ Dagegen: 70%
- ▶ Unentschieden: 6%



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

### Quellen

### Tabellen

## Laut einem „Bericht zur Bekämpfung des Analphabetismus in Deutschland“:

- ▶ Heute gibt es in Deutschland ca. 7 Millionen Analphabeten
- ▶ Zu Kaiser Wilhelms Zeiten gab es weniger als 10 000

Was leiten Sie daraus ab?

### BILDUNG

#### 7,5 Millionen Deutsche sind Analphabeten

Ein Siebtel der erwerbsfähigen Bevölkerung kann laut einer Studie kaum lesen und schreiben – doppelt so viel wie bisher gedacht. Bildungsministerin Schavan will reagieren [\[weiter...\]](#)



### ANALPHABETISMUS

#### Ein Land verliert das Lesen

Studenten verstehen abstrakte Texte nicht mehr, ein Schulbuchverlag kürzt Klassiker, Banker besuchen Lesekurse: Viele Deutsche haben keine Lust mehr zu lesen. [\[weiter...\]](#)

### ANALPHABETISMUS

#### Buchstäblich resigniert

Mehr als sieben Millionen Deutsche können kaum lesen und schreiben. Erst jetzt hat die Politik das Problem erkannt. Aber es gibt zu wenig Geld für Kurse. Von M. Spiewak [\[weiter...\]](#)

Quelle: Zeit.de

## Definition

### Zu Kaiser Wilhelms Zeiten:

„Analphabet ist, wer seinen Namen nicht schreiben kann.“

### Definition heute:

„Ein Analphabet ist eine Person, die sich nicht beteiligen kann an all den zielgerichteten Aktivitäten ihrer Gruppe und ihrer Gemeinschaft, bei denen Lesen, Schreiben und Rechnen erforderlich ist und an der weiteren Nutzung dieser Kulturtechniken für ihre weitere Entwicklung und die der Gesellschaft“.



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen



## Aussage des Vertriebsleiters:

„Unser Umsatz stieg vor einem Jahr um 1%. Dieses Jahr stieg das Umsatzwachstum um 50%!“

## Im Klartext:

- ▶ Basisjahr: Umsatz 100
- ▶ Dann: Wachstum auf 101
- ▶ Dieses Jahr: Wachstum des Wachstums um 50% bedeutet 1,5% Wachstum. Also Umsatz dann 102,5049

### 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

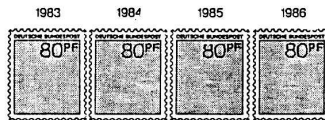
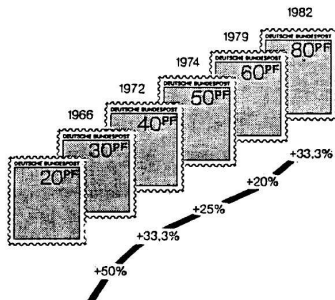
### Quellen

### Tabellen

## Seit 1983 stabile Gebühren

Sie, lieber Postkunde, sehen es selbst anhand unserer Zeichnung: Seit 1983 sind die Gebühren für Briefe, Päckchen und Pakete nicht mehr gestiegen. Und Sie bleiben auch 1986 stabil.

Das heißt: eine Legislaturperiode ohne Portonerhöhung. Und das seit 20 Jahren zum erstenmal wieder!



Diese erfreuliche Tatsache ist der konsequenten Stabilitätspolitik der Post seit 1983 zu verdanken. **1983-1986 +0%**

Quelle Kramer (2011)



### 1. Einführung

- Berühmte Leute zur Statistik?
- Wie lügt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

### Quellen

### Tabellen

## Grafik aussagekräftig?



Quelle: Bach u. a. (2006)



### 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

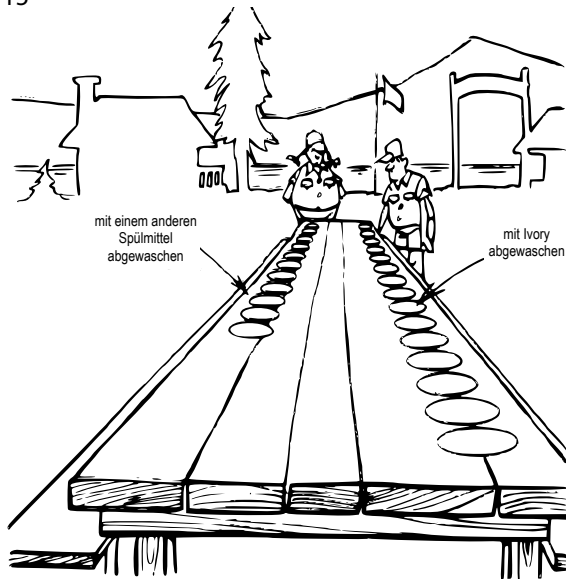
### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen

11 zu 15



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen

- ▶ Ein Einzelhändler bezieht ein Produkt zu 100 € und verkauft es für 200 €. Hat er eine Gewinnspanne von 50% oder 100%?
- ▶ Bahn: 9 Tote pro 10 Mio Passagieren je Kilometer  
Flugzeug: 3 Tote pro 10 Mio Passagieren je Kilometer  
Bahn: 7 pro 10 Mio Passagiere je Stunde  
Flugzeug: 24 pro 10 Mio Passagiere je Stunde
- ▶ Nur 40 % aller durch Autounfälle Gestorbenen hatten keinen Sicherheitsgurt angelegt  
Also: Keinen Gurt anlegen ist sicherer
- ▶ Die Hälfte der Todesfälle ereignen sich in Krankenhäusern  
Also: Krankenhäuser sind lebensgefährlich
- ▶ Zwei Drittel aller alkoholabhängigen Personen sind verheiratet  
Also: die Ehe führt zum Alkohol



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

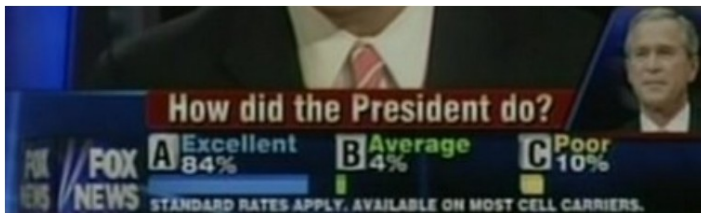
## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen

## Fernsehumfragen



- ▶ Kostenpflichtige Telefonabstimmung nach regierungsfreundlichem Bericht im Fernsehen
- ▶ In den meisten Umfragen erreichte Bush zu diesem Zeitpunkt nur 30 % Zustimmung



### 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen



## Challenger-Katastrophe



Am 28. Januar 1986, 73 Sekunden nach dem Start der Mission STS-51-L, brach die Raumfähre in etwa 15 Kilometer Höhe auseinander. Dabei starben alle sieben Astronauten. Es war der bis dahin schwerste Unfall in der Raumfahrtgeschichte der USA.



### 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

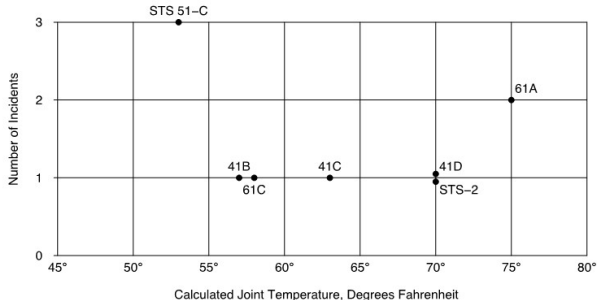
### 4. Induktive Statistik

Quellen

Tabellen

- ▶ Grund für Explosion: 2 Gummidichtungsringe waren undicht
- ▶ Die Temperatur der Dichtungsringe: Unter 20° F (ca. -6,7° C).
- ▶ Probleme mit Dichtungsringen bei Start der vorigen Föhre: Umgebungstemperatur 53° F (ca. 11,7° C).
- ▶ Frage: Ist der Dichtungsfehler durch die Umgebungstemperatur zu prognostizieren?

## O-Ring Failure Data



### 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen

- ▶ Fehler in Analyse: Starts ohne Fehler wurden nicht berücksichtigt
- ▶ Korrekte Modellierung mittels **logistischer Regression** liefert:



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

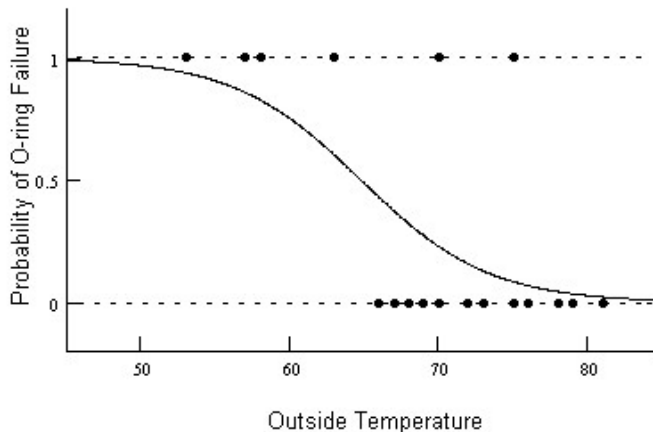
## 2. Deskriptive Statistik

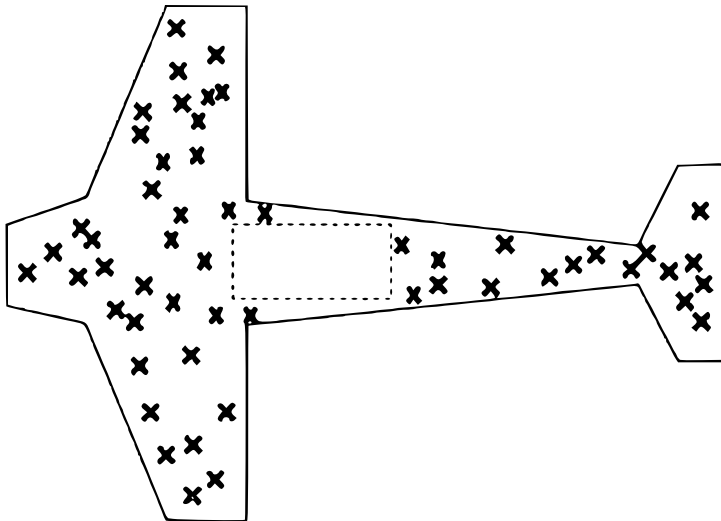
## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen





## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der

Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

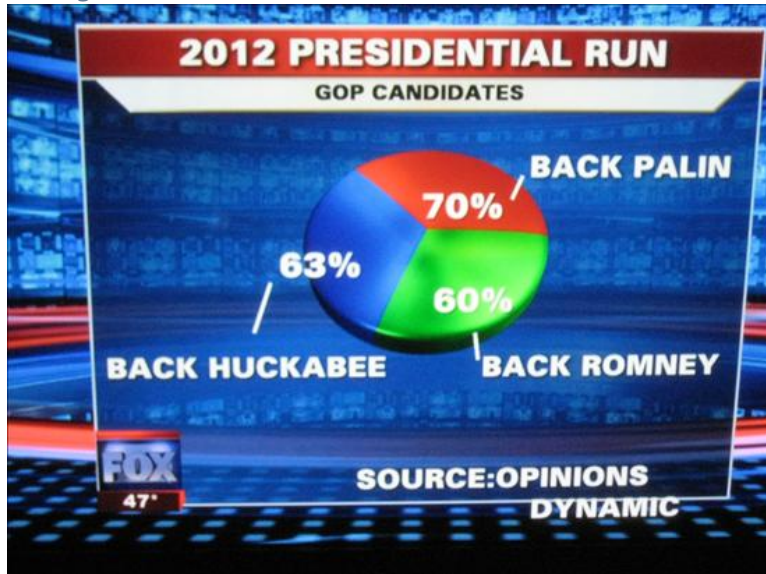
## 3. W-Theorie

## 4. Induktive Statistik

## Quellen

## Tabellen

## Aussage?



### 1. Einführung

- Berühmte Leute zur Statistik?
- Wie lügt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

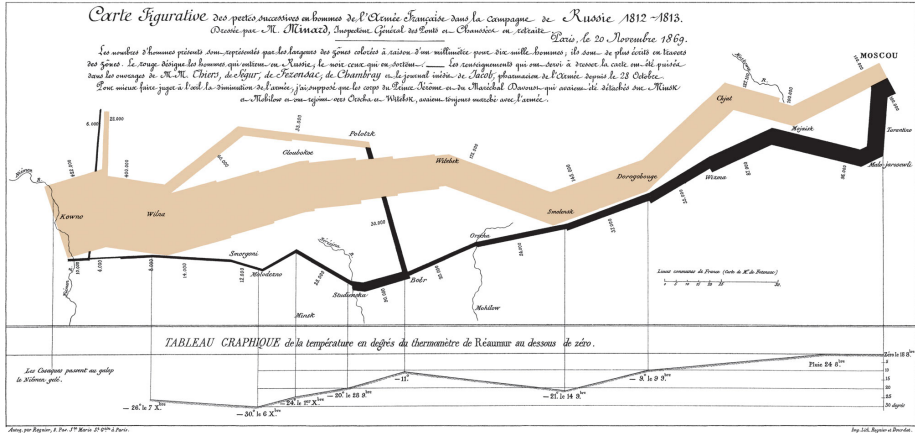
### 4. Induktive Statistik

### Quellen

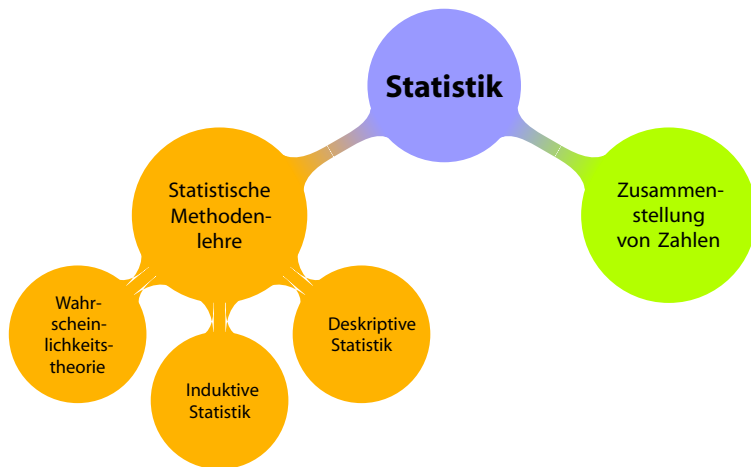
### Tabellen



## Minards Grafik von 1869 über Napoleons Rußlandfeldzug



Quelle: Wikimedia Commons, <http://goo.gl/T7ZNme>, Stand November 2014



## 1. Einführung

Berühmte Leute zur Statistik  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken

### Begriff Statistik

Grundbegriffe der Datenerhebung  
R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

### Quellen

### Tabellen

## Beispiel

12 Beschäftigte werden nach der Entfernung zum Arbeitsplatz (in km) befragt.

Antworten: 4, 11, 1, 3, 5, 4, 20, 4, 6, 16, 10, 6

► deskriptiv:

- Durchschnittliche Entfernung: 7,5
- Klassenbildung:

Klasse	[0;5)	[5;15)	[15;30)
Häufigkeit	5	5	2

► induktiv:

- Schätze die mittlere Entfernung **aller** Beschäftigten.
- Prüfe, ob die mittlere Entfernung geringer als 10 km ist.



### 1. Einführung

Berühmte Leute zur Statistik  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken

#### Begriff Statistik

Grundbegriffe der  
Datenerhebung  
R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

#### Quellen

#### Tabellen



*Objekt, cases*

- ▶ **Merkmalsträger**: Untersuchte statistische Einheit
- ▶ **Merkmal**: Interessierende Eigenschaft des Merkmalsträgers
- ▶ (Merkmals-) **Ausprägung**: Konkret beobachteter Wert des Merkmals
- ▶ **Grundgesamtheit**: Menge aller relevanten Merkmalsträger
- ▶ **Typen** von Merkmalen:
  - a) qualitativ – quantitativ
    - qualitativ: z.B. Geschlecht
    - quantitativ: z.B. Schuhgröße
    - Qualitative Merkmale sind quantifizierbar (z.B.: weiblich 1, männlich 0)
  - b) diskret – stetig
    - **diskret**: Abzählbar viele unterschiedliche Ausprägungen
    - **stetig**: Alle Zwischenwerte realisierbar



## 1. Einführung

Berühmte Leute zur Statistik  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken  
Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen

## Nominalskala:

- ▶ Zahlen haben nur Bezeichnungsfunktion
- ▶ z.B. Artikelnummern

## Ordinalskala:

- ▶ zusätzlich Rangbildung möglich
- ▶ z.B. Schulnoten
- ▶ Differenzen sind aber **nicht** interpretierbar!
  - Addition usw. ist unzulässig.

*metrische Merkmale*

## Kardinalskala:

- ▶ zusätzlich Differenzbildung sinnvoll
- ▶ z.B. Gewinn
- ▶ Noch feinere Unterscheidung in: **Absolutskala**, **Verhältnisskala**, **Intervallskala**



### 1. Einführung

Berühmte Leute zur Statistik  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken  
Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen

Ziel der Skalierung: Gegebene Information angemessen abbilden, möglichst ohne Über- bzw. Unterschätzungen

Es gilt:

- ▶ Grundsätzlich können alle Merkmale nominal skaliert werden.
- ▶ Grundsätzlich kann jedes metrische Merkmal ordinal skaliert werden.

Das nennt man **Skalendegression**. Dabei: **Informationsverlust**

Aber:

- ▶ Nominale Merkmale dürfen **nicht** ordinal- oder metrisch skaliert werden.
- ▶ Ordinale Merkmale dürfen **nicht** metrisch skaliert werden.

Das nennt man **Skalenprogression**. Dabei: Interpretation von **mehr Informationen** in die Merkmale, als inhaltlich vertretbar.  
(Gefahr der **Fehlinterpretation**)



## 1. Einführung

Berühmte Leute zur Statistik?  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken  
Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

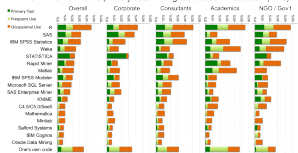
Quellen

Tabellen

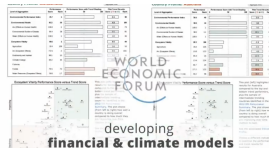
- ▶ R ist ein **freies** Softwarepaket zu Statistik und Datenanalyse
- ▶ R ist sehr mächtig und **weit verbreitet** in Wissenschaft und Industrie (sogar von mehr Leuten benutzt als z.B. SPSS)



- The average data miner reports using 4 software tools.
- R is used by the most data miners (47%).
- STATISTICA is the primary data mining tool chosen most often (17%).



source: <http://goo.gl/axhGhh>



graphics source: <http://goo.gl/W70kms>

## 1. Einführung

- Berühmte Leute zur Statistik
- Wie lügt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

## 2. Deskriptive Statistik

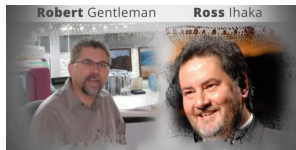
## 3. W-Theorie

## 4. Induktive Statistik

## Quellen

## Tabellen

- ▶ R ist ein **freies** Softwarepaket zu Statistik und Datenanalyse
- ▶ R ist sehr mächtig und **weit verbreitet** in Wissenschaft und Industrie (sogar von mehr Leuten benutzt als z.B. SPSS)
- ▶ Ursprung von R: **1993** an der Universität Auckland von Ross Ihaka and Robert Gentleman entwickelt
- ▶ Seitdem: Viele Leute haben R verbessert mit **tausenden von Paketen** für viele Anwendungen



## 1. Einführung

Berühmte Leute zur Statistik

Wie lögt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

## Quellen

## Tabellen

- ▶ R ist ein **freies** Softwarepaket zu Statistik und Datenanalyse
- ▶ R ist sehr mächtig und **weit verbreitet** in Wissenschaft und Industrie (sogar von mehr Leuten benutzt als z.B. SPSS)
- ▶ Ursprung von R: **1993** an der Universität Auckland von Ross Ihaka and Robert Gentleman entwickelt
- ▶ Seitdem: Viele Leute haben R verbessert mit **tausenden von Paketen** für viele Anwendungen
- ▶ Nachteil (auf den ersten Blick): Kein point und click tool

```
> summary(diamonds$price)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  326     950    2401   3933   5324   18820
> aveSize <- round(mean(diamonds$carat), 4)
> clarity <- levels(diamonds$clarity)
> p <- qplot(carat, price,
+           data=diamonds, color=clarity,
+           xlab="Carat", ylab="Price",
+           main="Diamond Pricing")
>
> format.plot(p, size=24)
> |
```



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen

- ▶ R ist ein **freies** Softwarepaket zu Statistik und Datenanalyse
- ▶ R ist sehr mächtig und **weit verbreitet** in Wissenschaft und Industrie (sogar von mehr Leuten benutzt als z.B. SPSS)
- ▶ Ursprung von R: **1993** an der Universität Auckland von Ross Ihaka and Robert Gentleman entwickelt
- ▶ Seitdem: Viele Leute haben R verbessert mit **tausenden von Paketen** für viele Anwendungen
- ▶ Nachteil (auf den ersten Blick): Kein point und click tool
- ▶ Großer Vorteil (auf den zweiten Blick): Kein point und click tool

```
> summary(diamonds$price)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  326     950    2401   3933   5324   18820
> aveSize <- round(mean(diamonds$carat), 4)
> clarity <- levels(diamonds$clarity)
> p <- qplot(carat, price,
+           data=diamonds, color=clarity,
+           xlab="Carat", ylab="Price",
+           main="Diamond Pricing")
>
> format.plot(p, size=24)
> |
```

**Download: [R-project.org](http://R-project.org)**



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

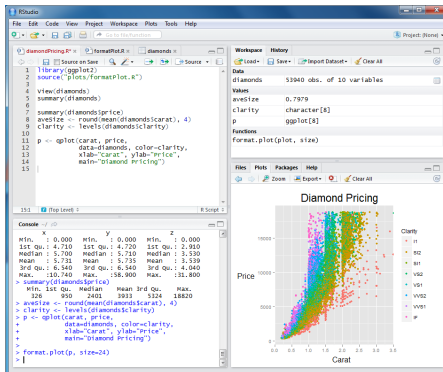
Tabellen

# Was ist RStudio?

- ▶ RStudio ist ein **Integrated Development Environment (IDE)** um R leichter benutzen zu können.
- ▶ Gibt's für OSX, Linux und Windows
- ▶ Ist auch frei
- ▶ Trotzdem: Sie müssen Kommandos schreiben
- ▶ Aber: RStudio unterstützt Sie dabei
- ▶ **Download: RStudio.com**



Free & Open-Source IDE for R



## 1. Einführung

- Berühmte Leute zur Statistik
- Wie legt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

## Quellen

## Tabellen





## RStudio Kennenlernen

- ▶ Code
- ▶ Console
- ▶ Workspace
- ▶ History
- ▶ Files
- ▶ Plots
- ▶ Packages
- ▶ Help
- ▶ Auto-Completion
- ▶ Data Import

The screenshot shows the RStudio environment with the following components:

- Source Editor:** Contains R code for loading the 'diamonds' dataset, summarizing it, and creating a scatter plot of Price vs. Carat, faceted by Clarity.
- Console:** Shows the execution of the code, including summary statistics for 'diamonds' and 'diamonds\$price'.
- Workspace:** Lists the loaded objects: 'diamonds' (53940 obs. of 10 variables), 'avesize', 'clarity', 'p', and 'format\_plot'.
- Plots Panel:** Displays a scatter plot titled 'Diamond Pricing' showing Price (y-axis, 0 to 15000) versus Carat (x-axis, 0.0 to 3.5). Points are colored by Clarity (I1, SI2, SI1, VS2, VS1, VVS2, VVS1, IF).

```

1 library(ggplot2)
2 source("plots/formatPlot.R")
3
4 view(diamonds)
5 summary(diamonds)
6
7 summary(diamonds$price)
8 avesize <- round(mean(diamonds$carat), 4)
9 clarity <- levels(diamonds$clarity)
10
11
12 p <- ggplot(carat, price,
13             data=diamonds, color=clarity,
14             xlab="carat", ylab="price",
15             main="Diamond Pricing")
    
```

```

> summary(diamonds)
  Min.   0.000  Min.   0.000  Min.   0.000
 1st Qu. 4.710  1st Qu. 4.720  1st Qu. 2.910
  Median 5.700  Median 5.710  Median 3.530
  Mean   5.731  Mean   5.735  Mean   3.539
 3rd Qu. 6.540  3rd Qu. 6.540  3rd Qu. 4.040
  Max.  110.740  Max.  158.900  Max.  131.800

> summary(diamonds$price)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 326    950    2403    3933    5324   18820

> avesize <- round(mean(diamonds$carat), 4)
> clarity <- levels(diamonds$clarity)
> p <- ggplot(carat, price,
+            data=diamonds, color=clarity,
+            xlab="carat", ylab="price",
+            main="Diamond Pricing")
>
> format_plot(p, size=24)
> |
    
```

### 1. Einführung

- Berühmte Leute zur Statistik
- Wie lügt man mit Statistik?
- Gute und schlechte Grafiken
- Begriff Statistik
- Grundbegriffe der Datenerhebung
- R und RStudio

### 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

### Quellen

### Tabellen



```
# Arbeitsverzeichnis setzen (alternativ über Menü)
setwd("C:/ste/work/vorlesungen/2015SS_HSA_Statistik")

# Daten einlesen aus einer csv-Datei (Excel)
MyData = read.csv2(file="../genericFiles/Daten/Umfrage_HSA_2015_03.csv", header=TRUE)
```

```
# inspect structure of data
str(MyData)

## 'data.frame': 377 obs. of 17 variables:
## $ Jahrgang : int 2015 2015 2015 2015 2015 2015 2015 2015 2015 2015 ...
## $ Alter : int 20 25 19 21 25 20 25 20 23 21 ...
## $ Groesse : int 174 157 163 185 178 170 165 175 180 161 ...
## $ Geschlecht : Factor w/ 2 levels "Frau","Mann": 1 1 1 2 2 1 1 2 2 1 ...
## $ AlterV : int 55 54 51 52 60 50 60 52 56 70 ...
## $ AlterM : int 53 61 49 50 63 55 60 49 50 55 ...
## $ GroesseV : int 187 185 178 183 170 183 185 175 175 180 ...
## $ GroesseM : int 169 160 168 165 160 160 170 169 170 165 ...
## $ Geschwister : int 3 1 1 4 2 2 4 1 1 2 ...
## $ Farbe : Factor w/ 6 levels "blau","gelb",...: 4 6 4 4 1 6 1 6 4 4 ...
## $ AusgKomm : num 240 119 270 40 550 ...
## $ AnzSchuhe : int 25 30 25 6 5 65 10 7 10 22 ...
## $ AusgSchuhe : int 450 300 100 100 80 250 150 400 150 300 ...
## $ Essgewohnheiten: Factor w/ 5 levels "carnivor","fruktarisch",...: 1 1 1 1 1 1 5 1 1 1 ...
## $ Raucher : Factor w/ 2 levels "ja","nein": NA 2 2 2 1 2 2 2 1 ...
## $ NoteMathe : num 2.3 3.3 1.7 2 4 4 3.3 2.7 3.7 3.3 ...
## $ MatheZufr : Ord.factor w/ 4 levels "unzufrieden"<.: 2 2 2 2 2 2 2 2 2 2 ...
```

## 1. Einführung

Berühmte Leute zur Statistik  
Wie liigt man mit Statistik?  
Gute und schlechte Grafiken  
Begriff Statistik  
Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen



```
# Erste Zeilen in Datentabelle
```

```
head(MyData, 6)
```

```
##   Jahrgang Alter Groesse Geschlecht AlterV AlterM GroesseV GroesseM Geschwister Farbe AusgKomm
## 1   2015    20   174      Frau      55    53    187    169      3 schwarz  240.0
## 2   2015    25   157      Frau      54    61    185    160      1 weiss   119.4
## 3   2015    19   163      Frau      51    49    178    168      1 schwarz  270.0
## 4   2015    21   185      Mann      52    50    183    165      4 schwarz  40.0
## 5   2015    25   178      Mann      60    63    170    160      2 blau    550.0
## 6   2015    20   170      Frau      50    55    183    160      2 weiss   420.0
##   AnzSchuhe AusgSchuhe Essgewohnheiten Raucher NoteMathe MatheZufr
## 1         25         450      carnivore <NA>      2.3 geht so
## 2         30         300      carnivore nein      3.3 geht so
## 3         25         100      carnivore nein      1.7 geht so
## 4          6         100      carnivore nein      2.0 geht so
## 5          5          80      carnivore ja       4.0 geht so
## 6         65         250      carnivore nein      4.0 geht so
```

```
# lege MyData als den "Standard"-Datensatz fest
```

```
attach(MyData)
```

```
# Wie Viele Objekte gibt's im Datensatz?
```

```
nrow(MyData)
```

```
## [1] 377
```

```
# Wie Viele Merkmale?
```

```
ncol(MyData)
```

```
## [1] 17
```

## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

## 3. W-Theorie

## 4. Induktive Statistik

Quellen

Tabellen



```
# Auswahl spezieller Objekte und Merkmale über [Zeile, Spalte]
MyData[1:3, 2:5]
```

```
##   Alter Groesse Geschlecht AlterV
## 1   20    174      Frau    55
## 2   25    157      Frau    54
## 3   19    163      Frau    51
```

```
# Auswahl von Objekten über logische Ausdrücke
Auswahl = (MyData$Geschlecht=="Mann" & MyData$Alter < 19)
# zeige die ersten Einträge
head(Auswahl, 30)
```

```
## [1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [17] FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
```

```
# Ausgabe der Auswahl: Alter, Alter des Vaters und der Mutter
MyData[Auswahl, # Objektauswahl
       c("Alter", "AlterM", "AlterV")] # Welche Merkmale?
```

```
##   Alter AlterM AlterV
## 23    18     44     48
## 268   18     46     52
```

## 1. Einführung

Berühmte Leute zur Statistik  
Wie lügt man mit Statistik?  
Gute und schlechte Grafiken  
Begriff Statistik  
Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen

```
# Zeige die Männer, die mehr als 1000 Euro für Schuhe  
# und Mobilfunk zusammen ausgegeben haben  
MyData[MyData$Geschlecht=="Mann" & MyData$AusgSchuhe + MyData$AusgKomm > 1000,  
c("Alter", "Geschwister", "Farbe", "AusgSchuhe", "AusgKomm")]
```

##	Alter	Geschwister	Farbe	AusgSchuhe	AusgKomm
## 19	20	0	weiss	300	924
## 23	18	1	silber	300	1000
## 33	25	1	schwarz	300	1000
## 42	24	1	schwarz	1000	600
## 81	25	2	silber	200	1900
## 106	21	1	schwarz	200	860
## 121	22	0	silber	300	1100
## 142	20	2	schwarz	290	1570
## 161	19	1	schwarz	600	800
## 168	21	2	blau	600	505
## 179	21	0	silber	300	825
## 211	23	1	schwarz	450	630
## 223	20	1	rot	400	815
## 227	20	1	schwarz	200	1250
## 249	20	1	blau	1000	350
## 256	25	0	schwarz	280	1200
## 272	24	1	schwarz	300	900
## 281	19	2	schwarz	500	720
## 315	21	1	weiss	200	1300
## 353	20	0	schwarz	400	950



## 1. Einführung

Berühmte Leute zur Statistik

Wie lügt man mit Statistik?

Gute und schlechte Grafiken

Begriff Statistik

Grundbegriffe der  
Datenerhebung

R und RStudio

## 2. Deskriptive Statistik

### 3. W-Theorie

### 4. Induktive Statistik

Quellen

Tabellen