

Statistik

HA 20.4.2016: Aufgaben 34-40

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, 14.00-15.30	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	J3.19	22.03.2016	
Offener Statistikraum	Etschberger/Tutoren	Mo, 14.00-17.45	B3.05	04.04.16	
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. desk. 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



- Berechnung eines linearen Modells der Bundesligadaten
- dabei: Punkte $\hat{=}$ y und Etat $\hat{=}$ x:

\bar{x}	33,83
\bar{y}	46,89
$\sum x_i^2$	25209
$\sum x_i y_i$	31474
n	18

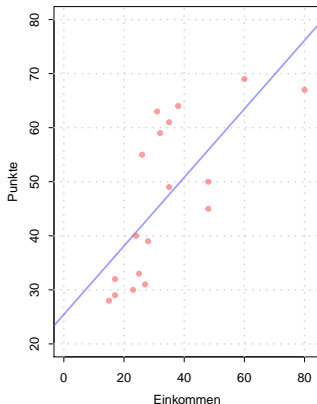
$$\Rightarrow \hat{b} = \frac{31474 - 18 \cdot 33,83 \cdot 46,89}{25209 - 18 \cdot 33,83^2}$$

$$\approx 0,634$$

$$\Rightarrow \hat{a} = 46,89 - \hat{b} \cdot 33,83$$

$$\approx 25,443$$

- Modell: $\hat{y} = 25,443 + 0,634 \cdot x$



- Prognosewert für Etat = 30:

$$\hat{y}(30) = 25,443 + 0,634 \cdot 30$$

$$\approx 44,463$$

1. Einführung

2. Deskriptive Statistik

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

Lineare Regression

3. W-Theorie

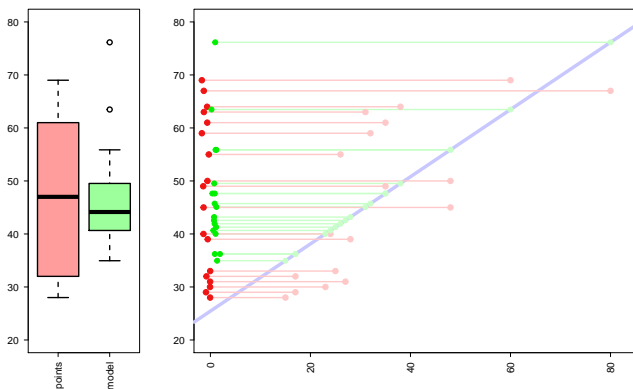
4. Induktive Statistik

Quellen

Tabellen



- **Varianz** der Daten in abhängiger Variablen y_i als Repräsentant des **Informationsgehalts**
- Ein Bruchteil davon kann in Modellwerten \hat{y}_i abgebildet werden



- Empirische Varianz (mittlere quadratische Abweichung) für „rot“ bzw. „grün“ ergibt jeweils

$$\frac{1}{18} \sum_{i=1}^{18} (y_i - \bar{y})^2 \approx 200,77 \quad \text{bzw.} \quad \frac{1}{18} \sum_{i=1}^{18} (\hat{y}_i - \bar{y})^2 \approx 102,78$$

1. Einführung

2. Deskriptive Statistik

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

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen



- ▶ Gütemaß für die Regression: **Determinationskoeffizient** (Bestimmtheitskoeffizient):

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n \hat{y}_i^2 - n\bar{y}^2}{\sum_{i=1}^n y_i^2 - n\bar{y}^2} = r^2 \in [0; 1]$$

- ▶ Mögliche Interpretation von R^2 :
Durch die Regression erklärter Anteil der Varianz
- ▶ $R^2 = 0$ wird erreicht wenn X, Y unkorreliert
 $R^2 = 1$ wird erreicht wenn $\hat{y}_i = y_i \forall i$ (alle Punkte auf Regressionsgerade)
- ▶ Im (Bundesliga-)Beispiel:

$$R^2 = \frac{\sum_{i=1}^{18} (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^{18} (y_i - \bar{y})^2} \approx \frac{102,78}{200,77} \approx 51,19\%$$

1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

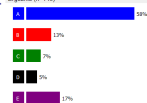
Tabellen

Regression: 4 eindimensionale Beispiele

Umfrage: Ich habe

- A) Zahlen und Grafik richtig
- B) Gerade falsch eingezeichnet, sonst alles richtig
- C) Fehler in den Zahlen (Grafik konsequent)
- D) Fehler in Zahlen und Grafik
- E) alles falsch
- F) nichts gemacht

Ergebnis (n=143)



► Berühmte Daten aus den 1970er Jahren:

i	x_{1i}	x_{2i}	x_{3i}	x_{4i}	y_{1i}	y_{2i}	y_{3i}	y_{4i}
1	10	10	10	8	8,04	9,14	7,46	6,58
2	8	8	8	8	6,95	8,14	6,77	5,76
3	13	13	13	8	7,58	8,74	12,74	7,71
4	9	9	9	8	8,81	8,77	7,11	8,84
5	11	11	11	8	8,33	9,26	7,81	8,47
6	14	14	14	8	9,96	8,10	8,84	7,04
7	6	6	6	8	7,24	6,13	6,08	5,25
8	4	4	4	19	4,26	3,10	5,39	12,50
9	12	12	12	8	10,84	9,13	8,15	5,56
10	7	7	7	8	4,82	7,26	6,42	7,91
11	5	5	5	8	5,68	4,74	5,73	6,89

Aufgaben:

- Modellparameter (a , b)
- R^2
- Streuplot mit eingezeichneter Gerade

Lage und Streuung

ion
male
s
gression
orie

4. Induktive Statistik

Quellen

Tabellen

(Quelle: Anscombe (1973))



- ▶ In folgender Tabelle: Jeweils Ergebnisse der linearen Regressionsanalyse
- ▶ dabei: x_k unabhängige Variable und y_k abhängige Variable
- ▶ Modell jeweils: $y_k = a_k + b_k x_k$

k	\hat{a}_k	\hat{b}_k	R_k^2
1	3,0001	0,5001	0,6665
2	3,0010	0,5000	0,6662
3	3,0025	0,4997	0,6663
4	3,0017	0,4999	0,6667

1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

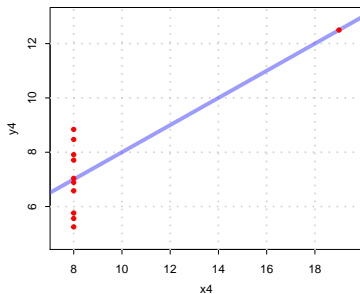
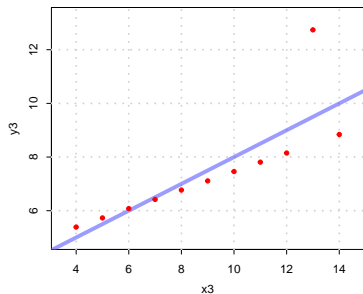
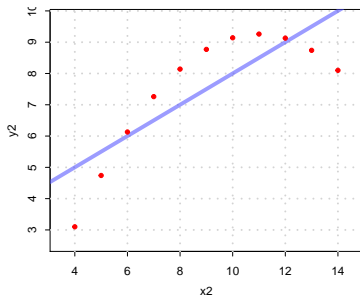
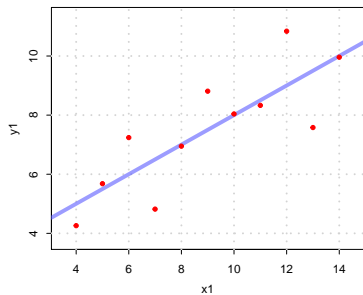
3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

Plot der Anscombe-Daten



1. Einführung

2. Deskriptive Statistik

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

3. W-Theorie

4. Induktive Statistik

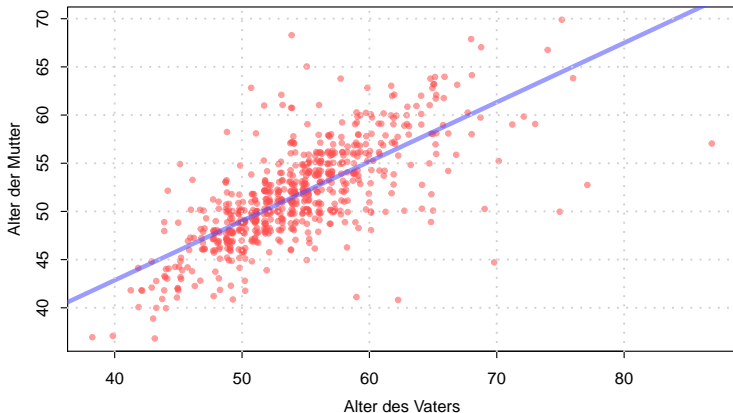
Quellen

Tabellen

```
meineRegression = lm(AlterM ~ AlterV)  
meineRegression
```

```
plot(AlterV, AlterM,  
     xlab="Alter des Vaters",  
     ylab="Alter der Mutter")  
abline(meineRegression)
```

```
##  
## Call:  
## lm(formula = AlterM ~ AlterV)  
##  
## Coefficients:  
## (Intercept)      AlterV  
##      18.2234      0.6159
```



1. Einführung

2. Deskriptive Statistik

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

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen



- ▶ Oft Kritisch: Einzelne Punkte, die Modell stark beeinflussen
- ▶ Idee: Was würde sich ändern, wenn solche Punkte weggelassen würden?
- ▶ **Cook-Distanz**: Misst den Effekt eines gelöschten Objekts
- ▶ Formel für ein lineares Modell mit einem unabh. Merkmal:

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(\text{ohne } i)})^2}{\text{MSE}}$$

- ▶ Dabei bedeutet:
 - \hat{y}_j : Prognosewert des kompletten Modells für das j-te Objekt
 - $\hat{y}_{j(\text{ohne } i)}$: Prognosewert des Modells ohne Objekt i für das j-te Objekt
 - $\text{MSE} = \frac{1}{n} \cdot \sum (\hat{y}_i - y_i)^2$: Normierender Term (Schätzwert für Fehlerstreuung)

1. Einführung

2. Deskriptive Statistik

Häufigkeiten
Lage und Streuung
Konzentration
Zwei Merkmale
Korrelation
Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

► Anscombe-Daten: Regressionsmodell Nr. 3



1. Einführung

2. Deskriptive Statistik

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

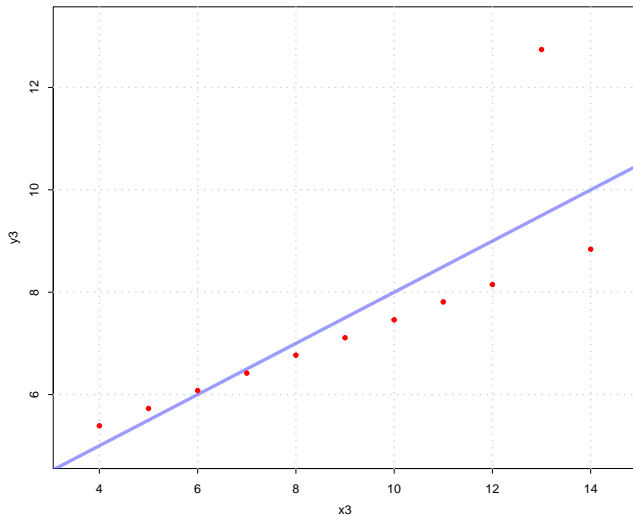
Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen



- ▶ Anscombe-Daten: Regressionsmodell Nr. 3
- ▶ Darstellung der Cook-Distanz neben Punkten
- ▶ Faustformel: Werte über 1 sollten genau untersucht werden



1. Einführung

2. Deskriptive Statistik

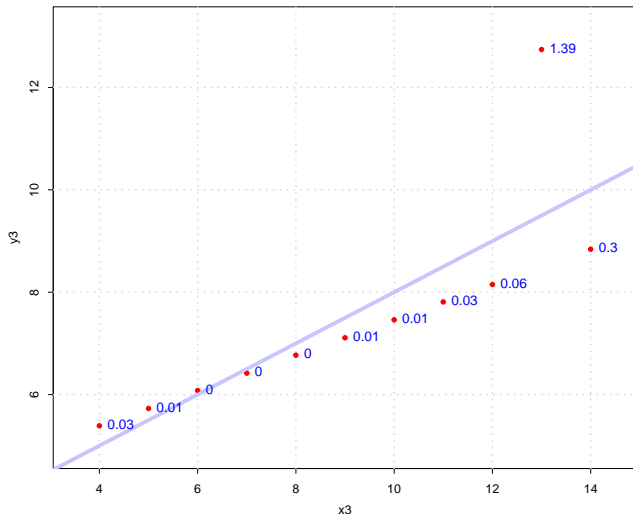
Häufigkeiten
Lage und Streuung
Konzentration
Zwei Merkmale
Korrelation
Preisindizes
Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen



- ▶ Oft aufschlussreich: Verteilung der **Residuen** e_i
- ▶ Verbreitet: Graphische Darstellungen der Residuen
- ▶ Z.B.: e_i über \hat{y}_i



1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

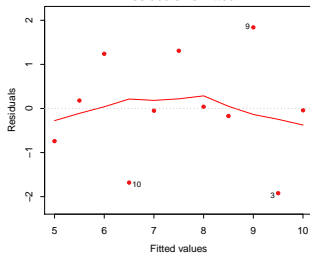
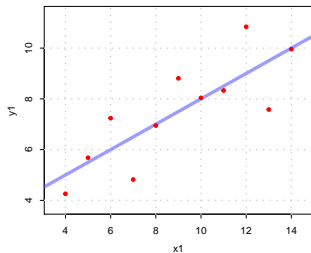
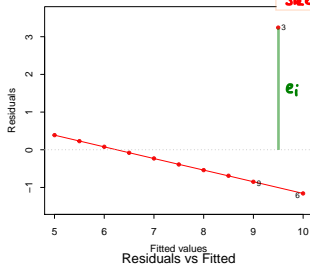
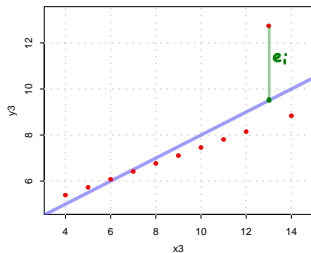
Tabellen

- ▶ Oft aufschlussreich: Verteilung der **Residuen** e_i
- ▶ Verbreitet: Graphische Darstellungen der Residuen
- ▶ Z.B.: e_i über \hat{y}_i

auch schlecht wäre:
steigende, fallende
oder ungleichmäßig
verteilte Streuung
in Residuenplots

gut
(homoskedastisch)

schlecht
(heteroskedastisch)



schlecht!
Tendenz
oder
nicht
gleichmäßig

2. Deskriptive Statistik

- Häufigkeiten
- Lage und Streuung
- Konzentration
- Zwei Merkmale
- Korrelation
- Preisindizes
- Lineare Regression
- 3W-Theorie

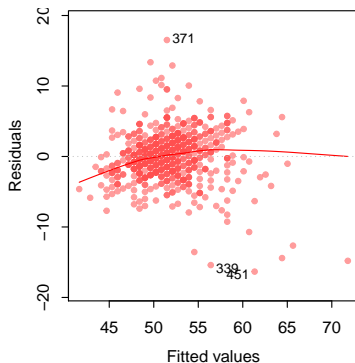
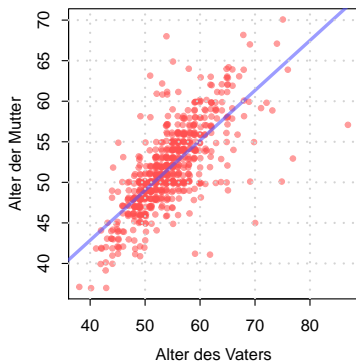
4. Induktive Statistik

- Quellen
- Tabellen



Wichtige Eigenschaften der Residuenverteilung

- ▶ Möglichst **keine systematischen Muster**
- ▶ Keine Änderung der Varianz in Abhängigkeit von \hat{y}_i (**Homoskedastizität**)
- ▶ Nötig für inferentielle Analysen: Näherungsweise **Normalverteilung** der Residuen (q-q-plots)



1. Einführung

2. Deskriptive Statistik

Häufigkeiten
Lage und Streuung
Konzentration
Zwei Merkmale
Korrelation
Preisindizes

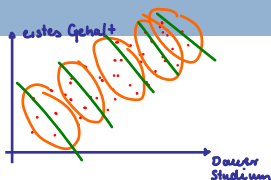
Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen



Exkurs: Kausalität vs. Korrelation

- ▶ Meist wichtig für sinnvolle Regressionsanalysen:
- ▶ **Kausale Verbindung** zwischen unabhängigem und abhängigem Merkmal
- ▶ Sonst bei Änderung der unabhängigen Variablen keine sinnvollen Prognosen möglich
- ▶ Oft: **Latente Variablen** im Hintergrund



1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

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



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

Kombinatorik

Permutationen $\hat{=}$ Anordnungen von Objekten

Beispiel: Außendienstmitarbeiter, 50 Stationen bei Kunden; gesucht: Anzahl möglicher Routen

$$50 \cdot 49 \cdot 48 \cdot \dots \cdot 3 \cdot 2 \cdot 1 = 50! \approx 3 \cdot 10^{64}$$

allgemein: n (unterschiedliche) Objekte kann man auf

$n!$

verschiedene Arten anordnen

(!): factorial(n)

Beispiel: M I I I I S S S S P P

ges: Anzahl der (unterscheidbaren) Anordnungen

$$\left. \begin{array}{l} \text{M I S S I S S I P P I} \\ \text{M I S S I S S I P P I} \\ \vdots \end{array} \right\} 4! \cdot 4! \cdot 2! \hat{=} \text{Anzahl der Permutationen pro "Wort", also die nicht unterscheidbaren Varianten}$$

$$\left. \begin{array}{l} \downarrow \\ \text{Vert. des "I"} \\ \downarrow \\ \text{Vert. des "S"} \\ \downarrow \\ \text{Vert. des "P"} \end{array} \right\}$$

\Rightarrow unterscheidbare Buchstabenanordnungen

$$\frac{11!}{4! \cdot 4! \cdot 2!} = 34650$$

allgemein: n Objekte, in k Gruppen zu jeweils n_1, \dots, n_k nicht unterscheidbaren Elementen

\Rightarrow Anzahl der (unterscheidbaren) Permutationen

$$\frac{n!}{\prod_{i=1}^k (n_i!)}$$

Kombinationen $\hat{=}$ Auswahl von k aus n Elementen

Beispiel: Anzahl der Möglichkeiten für EC-Karten-PIN

$$10 \cdot 10 \cdot 10 \cdot 10 = 10000 = 10^4$$

allgemein: n^k mit Wiederholung mit Reihenfolge

Beispiel: Top-3-Liste von 100 Mitarbeitern (Besten, Zweitbesten, 3-Beste(r))

$$100 \cdot 99 \cdot 98 = \frac{100!}{97!} = \frac{100!}{(100-3)!} = 970200$$

allgemein: $\frac{n!}{(n-k)!}$ ohne WH mit Reihenfolge

$$[TR: \frac{100!}{(100-3)!} = 100 \cdot 99 \cdot 98 = 970200]$$

Beispiel: Wahl eines 3-köpfigen, gleichberechtigten Teams aus 250 Leuten

$$\frac{250 \cdot 249 \cdot 248}{3!} = \frac{250!}{(250-3)! \cdot 3!} = \binom{250}{3} = 2.573.000$$

allgemein: $\binom{n}{k}$ ohne WH
ohne Reihenfolge

$$[\text{im TR: } \binom{250}{3} = 250 \cdot nCr 3 = 2.573.000]$$

Beispiel: B, S, W (3 Sorten Getränke)
10 Flaschen sollen gekauft werden

B	B	B		S	S	S	S		W	W	W
B	B	B	B	B	B	B	B	B			
	S	S	S	S	S	S	S	S		W	
B		S	S		W	W	W	W	W	W	
F	F	F		F	F	F	F		F	F	F

hier: $n=3$
 $k=10$

$$\text{insgesamt: } \frac{12 \cdot 11}{2 \cdot 1} = 66 = \frac{(10+3-1)!}{2! \cdot 10!}$$

$$\text{allgemein: } \frac{(k+n-1)!}{(n-1)! \cdot k!} = \binom{n+k-1}{k} \text{ mit WH ohne R.F.}$$



2-mal Würfeln, das heißt Auswahl von $k = 2$ aus $n = 6$ Zahlen.



(1,1)	(1,2)	(1,3)	(1,4)	(1,5)	(1,6)
(2,1)	(2,2)	(2,3)	(2,4)	(2,5)	(2,6)
(3,1)	(3,2)	(3,3)	(3,4)	(3,5)	(3,6)
(4,1)	(4,2)	(4,3)	(4,4)	(4,5)	(4,6)
(5,1)	(5,2)	(5,3)	(5,4)	(5,5)	(5,6)
(6,1)	(6,2)	(6,3)	(6,4)	(6,5)	(6,6)

- ▶ mit WH, mit RF: alle Möglichkeiten, $6^2 = 36$
- ▶ ohne WH, mit RF: Diagonale entfällt, $36 - 6 = 30 = 6 \cdot 5 = \frac{6!}{(6-2)!}$

- ▶ ohne WH, ohne RF: Hälfte des letzten Ergebnisses: $\frac{30}{2} = 15 = \frac{6!}{4!2!} = \binom{6}{2}$
- ▶ mit WH, ohne RF: Letztes Ergebnis plus Diagonale, $15 + 6 = 21 = \binom{7}{2}$

Auswahl von k aus n Dingen

	mit Wiederholung	ohne Wiederholung
mit Reihenfolge	n^k	$\frac{n!}{(n-k)!}$
ohne Reihenfolge	$\binom{n+k-1}{k}$	$\binom{n}{k}$



1. Einführung

2. Deskriptive Statistik

3. W-Theorie

Kombinatorik

Zufall und Wahrscheinlichkeit

Zufallsvariablen und Verteilungen

Verteilungsparameter

4. Induktive Statistik

Quellen

Tabellen



- ▶ **Zufallsvorgang:** Geschehen mit ungewissem Ausgang, z.B. Münzwurf
- ▶ **Elementarereignis** ω : Ein möglicher Ausgang, z.B. „Kopf“
Elementarereignisse schließen sich gegenseitig aus („Kopf“ oder „Zahl“)!
kleines omega
- ▶ **Ergebnismenge** Ω : Menge aller ω
- ▶ **Beispiel:** Werfen zweier Würfel:

$$\Omega : \left\{ \begin{array}{cccc} (1,1) & (1,2) & \dots & (1,6) \\ (2,1) & (2,2) & \dots & (2,6) \\ \vdots & \vdots & \ddots & \vdots \\ (6,1) & (6,2) & \dots & (6,6) \end{array} \right\}$$

$$\Rightarrow \Omega = \{(x_1, x_2) : x_1, x_2 \in \{1, \dots, 6\}\}$$

1. Einführung

2. Deskriptive Statistik

3. W-Theorie

Kombinatorik

Zufall und Wahrscheinlichkeit

Zufallsvariablen und
Verteilungen

Verteilungsparameter

4. Induktive Statistik

Quellen

Tabellen



► **Ereignis** A: Folgeerscheinung eines Elementarereignisses

► Formal:

$$A \subset \Omega$$

$$\begin{aligned} \Omega &= \{1, 2, 3\} \\ \omega = 1 &\Rightarrow \omega \in \Omega \\ A &= \{x \in \Omega : x < 2\} = \{1\} \\ A &\subset \Omega \\ A &= \{\omega\} \subset \Omega \end{aligned}$$

► Ereignisse schließen sich nicht gegenseitig aus!

► **Beispiel:** Werfen zweier Würfel:

Ereignis	verbal	formal
A	Augensumme = 4	$\{(1,3), (2,2), (3,1)\}$
B	Erste Zahl = 2	$\{(2,1), (2,2), \dots, (2,6)\}$

► **Wahrscheinlichkeit** $P(A)$: Chance für das Eintreten von A

► **Laplace-Wahrscheinlichkeit:**

$$P(A) = \frac{|A|}{|\Omega|} = \frac{\text{Anzahl der für A günstigen Fälle}}{\text{Anzahl aller möglichen Fälle}}$$

- 1. Einführung
- 2. Deskriptive Statistik
- 3. W-Theorie
 - Kombinatorik
 - Zufall und Wahrscheinlichkeit
 - Zufallsvariablen und Verteilungen
 - Verteilungsparameter
- 4. Induktive Statistik
- Quellen
- Tabellen