

## Anmerkungen zur Vorlesung Statistik vom 08.11.2016

### Inhalt:

- (einfache lineare ) Regressionsanalyse
  - Regressionsmodell
  - Regressionskoeffizienten, Regressionsgerade
  - Gütemaß: Determinationskoeffizient/Bestimmtheitsmaß  $R^2$
  - R: Umsetzung, Darstellung (Plot) und Outputinterpretation
  - Residualanalyse
- Klausur SS2016: Aufgabe 2 (WDH Korrelation / Regression)

### Aufgabensammlung Statistik:

|              |  |    |
|--------------|--|----|
| 25.10.2016   | Aufgabe 23: Rangkorrelation . . . . .        | 31 |
| Hausaufgabe  | Aufgabe 24: Lage Korrelation . . . . .       | 32 |
| A23-A31 ohne | Aufgabe 25: Kontingenzkoeffizient . . . . .  | 33 |
| Regressions- | Aufgabe 26: Kontingenzkoeffizient . . . . .  | 34 |
| aspekte      | Aufgabe 27: Korrelation Regression . . . . . | 35 |
| 08.11.2016   | Aufgabe 28: Korrelation Regression . . . . . | 36 |
| Hausaufgabe  | Aufgabe 29: Korrelation Regression . . . . . | 37 |
| A27-A31      | Aufgabe 30: Korrelation Regression . . . . . | 38 |
| Regressions- | Aufgabe 31: Korrelation Regression . . . . . | 39 |
| aspekte      | Aufgabe 32: Regression . . . . .             | 40 |
| A32-A33      | Aufgabe 33: Regression . . . . .             | 41 |

Bundesliga 2008/2009

- ▶ Gegeben: Daten zu den 18 Vereinen der ersten Bundesliga in der Saison 2008/09
- ▶ Merkmale: **Vereinssetat** für Saison (nur direkte Gehälter und Spielergehälter) und **Ergebnispunkte** in Tabelle am Ende der Saison

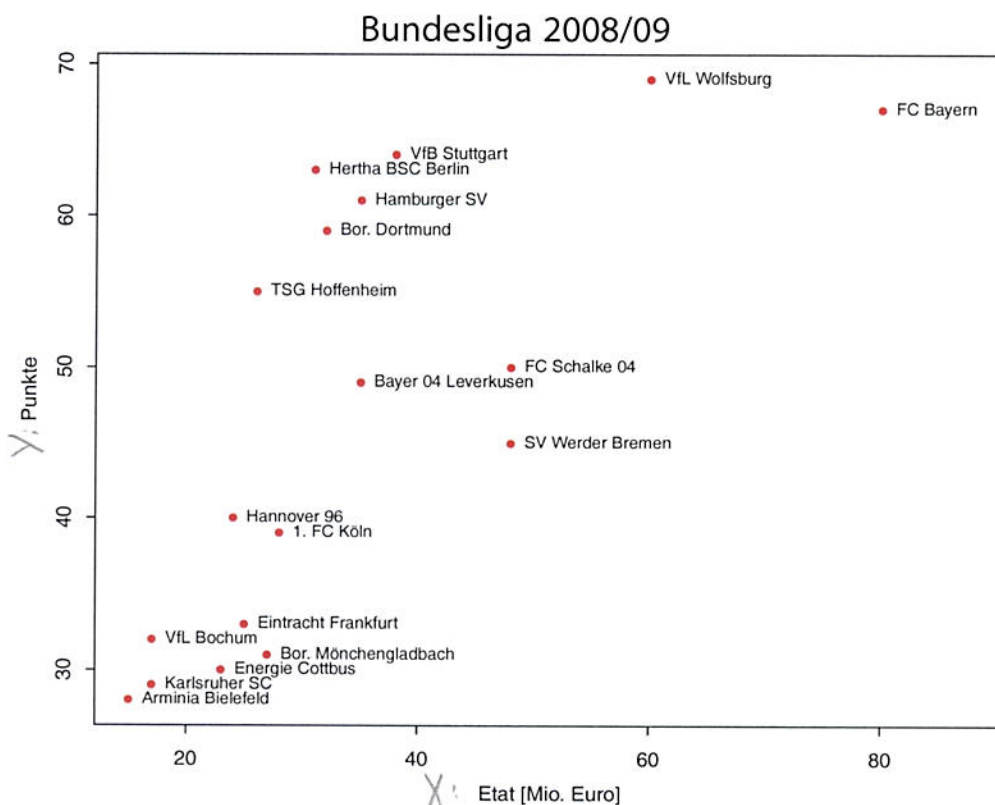
|                      | X    | Y      |
|----------------------|------|--------|
|                      | Etat | Punkte |
| FC Bayern            | 80   | 67     |
| VfL Wolfsburg        | 60   | 69     |
| SV Werder Bremen     | 48   | 45     |
| FC Schalke 04        | 48   | 50     |
| VfB Stuttgart        | 38   | 64     |
| Hamburger SV         | 35   | 61     |
| Bayer 04 Leverkusen  | 35   | 49     |
| Bor. Dortmund        | 32   | 59     |
| Hertha BSC Berlin    | 31   | 63     |
| 1. FC Köln           | 28   | 39     |
| Bor. Mönchengladbach | 27   | 31     |
| TSG Hoffenheim       | 26   | 55     |
| Eintracht Frankfurt  | 25   | 33     |
| Hannover 96          | 24   | 40     |
| Energie Cottbus      | 23   | 30     |
| VfL Bochum           | 17   | 32     |
| Karlsruher SC        | 17   | 29     |
| Arminia Bielefeld    | 15   | 28     |

(Quelle: Welt)



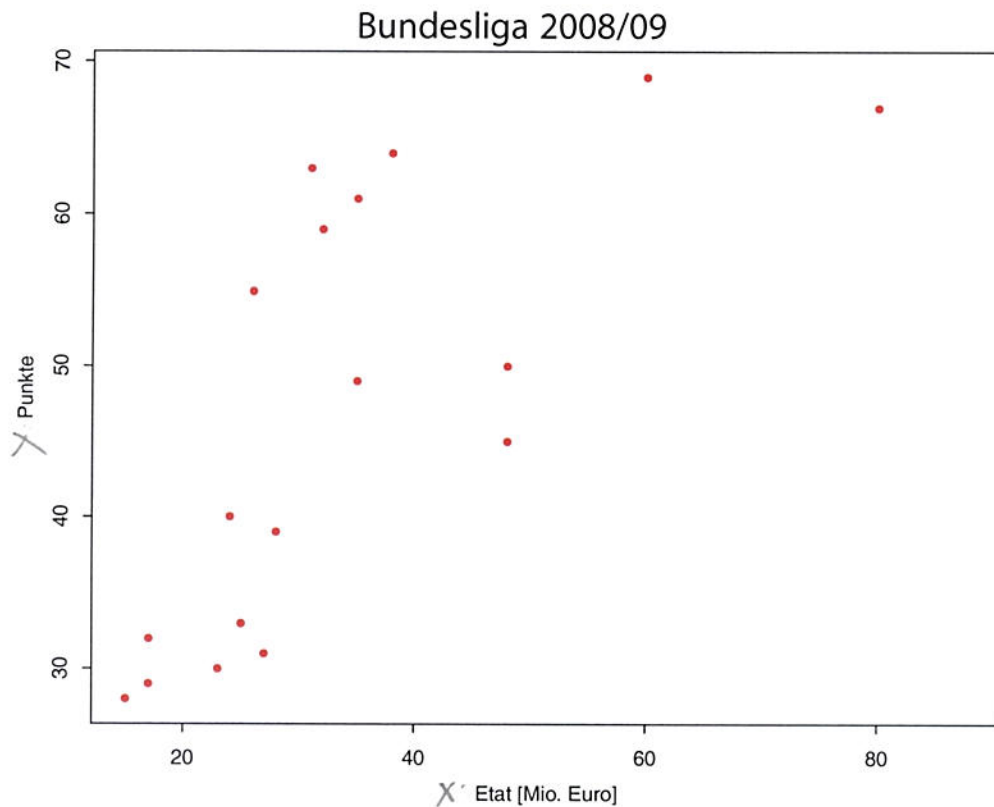
- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

Darstellung der Daten in Streuplot



- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

## Darstellung der Daten in Streuplot



Statistik



1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

99

## Trend als lineares Modell

- ▶ Kann man die **Tabellenpunkte** näherungsweise über einfache Funktion **in Abhängigkeit des Vereinsetats** darstellen?
- ▶ Allgemein: Darstellung einer Variablen  $Y$  als Funktion von  $X$ :

$$y = f(x)$$

- ▶ Dabei:
  - $X$  heißt **Regressor** bzw. **unabhängige Variable**
  - $Y$  heißt **Regressand** bzw. **abhängige Variable**
- ▶ Wichtiger (und einfachster) Spezialfall:  $f$  beschreibt einen linearen Trend:

$$y = a + b x$$

- ▶ Dabei anhand der Daten zu schätzen:  $a$  (Achsenabschnitt) und  $b$  (Steigung)
- ▶ Schätzung von  $a$  und  $b$ : **Lineare Regression**

Statistik



1. Einführung

2. Deskriptive Statistik

Häufigkeiten

Lage und Streuung

Konzentration

Zwei Merkmale

Korrelation

Preisindizes

Lineare Regression

3. W-Theorie

4. Induktive Statistik

Quellen

Tabellen

100



- ▶ Pro Datenpunkt gilt mit Regressionsmodell:

$$y_i = a + bx_i + \epsilon_i$$

- ▶ Dabei:  $\epsilon_i$  ist jeweils Fehler (der Grundgesamtheit),
- ▶ mit  $e_i = y_i - (\hat{a} + \hat{b}x_i)$ : Abweichung (**Residuen**) zwischen gegebenen Daten der Stichprobe und durch Modell geschätzten Werten
- ▶ Modell gut wenn alle Residuen  $e_i$  zusammen möglichst klein
- ▶ Einfache Summe aber nicht möglich, denn  $e_i$  positiv oder negativ
- ▶ Deswegen: Summe der Quadrate von  $e_i$
- ▶ **Prinzip der kleinsten Quadrate**: Wähle  $a$  und  $b$  so, dass

$$Q(a, b) = \sum_{i=1}^n [y_i - (a + bx_i)]^2 \rightarrow \min$$

$y_i$  tats. Beobachtung "Prognosewert"

Notwendige Bedingung für Minimum:

$$\frac{\partial Q}{\partial a} \stackrel{!}{=} 0 \quad \text{und} \quad \frac{\partial Q}{\partial b} \stackrel{!}{=} 0$$

- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen



- ▶ Beste und eindeutige Lösung:

$$\hat{b} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x}$$



$(\bar{x}, \bar{y})$  ist Datenschwerpunkt und liegt immer auf der Regressionsgerade!

- ▶ **Regressionsgerade**:

$$\hat{y} = \hat{a} + \hat{b}x$$

- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

## Bundesligabeispiel

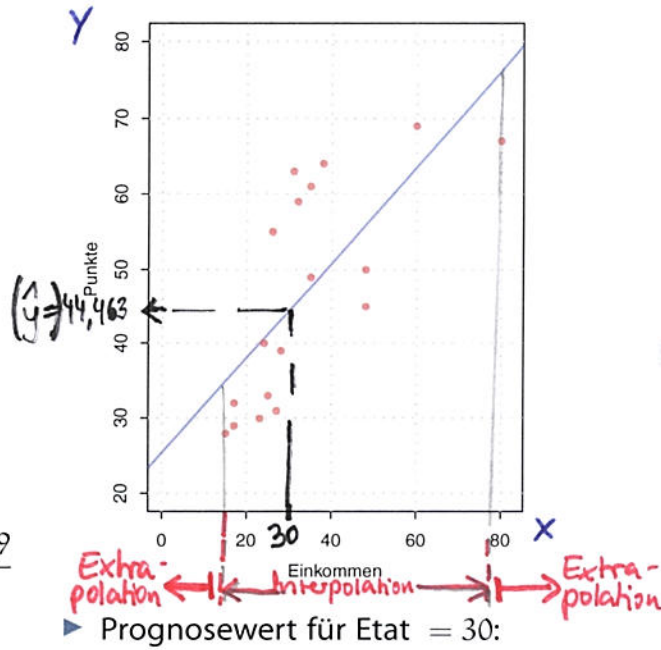
- ▶ Berechnung eines linearen Modells der Bundesligadaten
- ▶ dabei: Punkte  $\hat{=}$  y und Etat  $\hat{=}$  x:

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

$$\Rightarrow \hat{b} = \frac{31474 - 18 \cdot 33,83 \cdot 46,89}{25209 - 18 \cdot 33,83^2} \approx 0,634$$

$$\Rightarrow \hat{a} = 46,89 - \hat{b} \cdot 33,83 \approx 25,443$$

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



- ▶ Prognosewert für Etat = 30:

$$\hat{y}(30) = 25,443 + 0,634 \cdot 30 \approx 44,463$$

Statistik

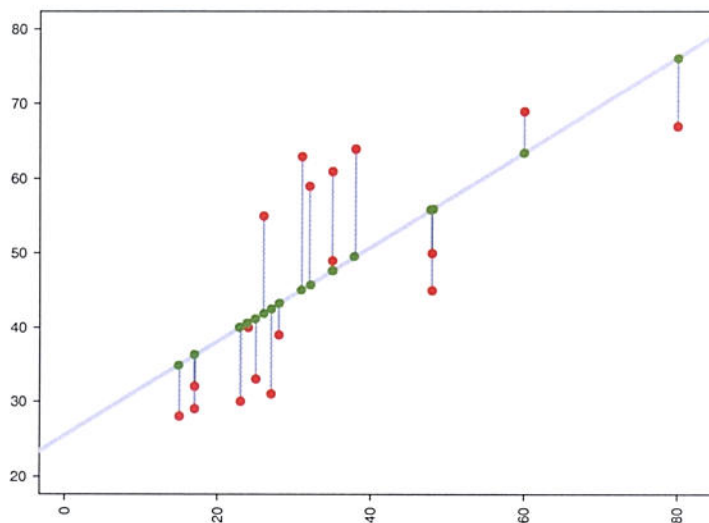


1. Einführung
  2. Deskriptive Statistik
    - Häufigkeiten
    - Lage und Streuung
    - Konzentration
    - Zwei Merkmale
    - Korrelation
    - Preisindizes
    - Lineare Regression
  3. W-Theorie
  4. Induktive Statistik
- Quellen  
 Tabellen

103

## Varianz und Information

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



Statistik

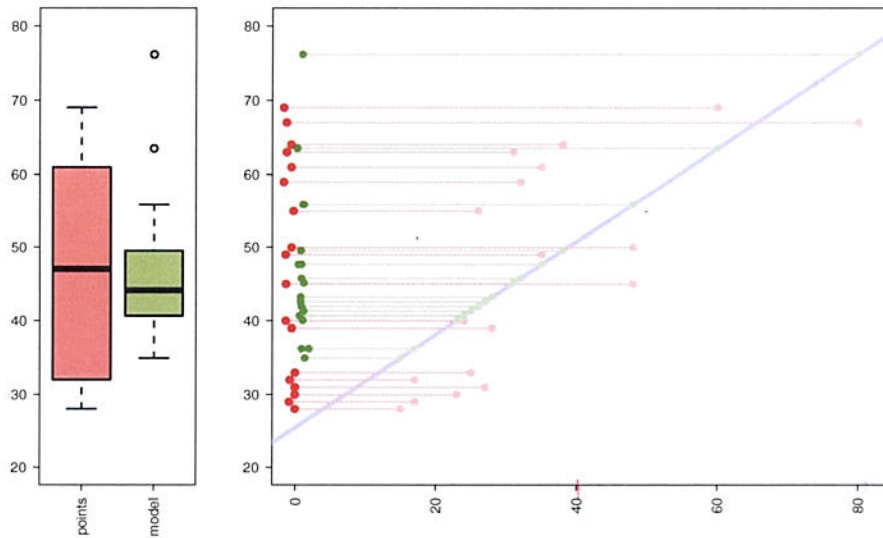


1. Einführung
  2. Deskriptive Statistik
    - Häufigkeiten
    - Lage und Streuung
    - Konzentration
    - Zwei Merkmale
    - Korrelation
    - Preisindizes
    - Lineare Regression
  3. W-Theorie
  4. Induktive Statistik
- Quellen  
 Tabellen

104



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

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

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

# Determinationskoeffizient



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

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

*Prognosewerte* (above  $\hat{y}_i$ )  
*tatsächliche Werte* (below  $y_i$ )  
*Bravais-Pearson Korrelationskoeff. zum Quadrat!* (next to  $r^2$ )

- ▶ 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:

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

*d.h. ca. 51% der in der abh. Variable y vorhandenen Information kann über die Regression erklärt werden (Rest: nicht berücksichtigte abh. Variablen Effekte)*

- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen



**OPTN** **▼** **1** (Addition) **1** bis **8**

**Anzahl der Stichproben:**  $n^*$  / **Mittelwert:**  $\bar{x}^*$ ,  $\bar{y}$  / **Varianz der**

**Grundgesamtheit:**  $\sigma_x^2$ ,  $\sigma_y^2$  / **Grundgesamtheits-Standardabweichung:**

$\sigma_x^*$ ,  $\sigma_y$  / **Stichprobenvarianz:**  $s_x^2$ ,  $s_y^2$  / **Stichproben-Standardabweichung:**  
 $S_x^*$ ,  $S_y$

**OPTN** **▼** **2** (Variable) **1** bis **8**, **▼** **1** bis **▼** **3**

**Minimalwert:**  $\min(x)^*$ ,  $\min(y)$  / **Maximalwert:**  $\max(x)^*$ ,  $\max(y)$

Wenn die statistische Berechnung mit Einzelvariable ausgewählt ist:

**OPTN** **▼** **3** (Minimum/Maximum) **1**, **5**

Wenn die statistische Berechnung mit Variablenpaar ausgewählt ist:

**OPTN** **▼** **3** (Minimum/Maximum) **1** bis **4**

**Erstes Quartil:**  $Q_1^*$  / **Median:**  $Med^*$  / **Drittes Quartil:**  $Q_3^*$  (Nur bei statistischen Berechnungen mit Einzelvariable)

**OPTN** **▼** **3** (Minimum/Maximum) **2** bis **4**

**Regressionskoeffizienten:**  $a$ ,  $b$  / **Korrelationskoeffizient:**  $r$  / **Schätzwerte:**  
 $\hat{x}$ ,  $\hat{y}$

**OPTN** **▼** **4** (Regressionen) **1** bis **5**

**Regressionskoeffizienten für quadratische Regression:**  $a$ ,  $b$ ,  $c$  /

**Schätzwerte:**  $\hat{x}_1$ ,  $\hat{x}_2$ ,  $\hat{y}$

**OPTN** **▼** **4** (Regressionen) **1** bis **6**

•  $\hat{x}$ ,  $\hat{x}_1$ ,  $\hat{x}_2$  und  $\hat{y}$  sind Befehle mit einem Argument unmittelbar davor.

**Bsp. 2:** Geben Sie die Daten  $x = \{1; 2; 2; 3; 3; 3; 4; 4; 5\}$  für eine einzelne Variable ein, verwenden Sie dabei die Freq-Spalte, um die Anzahl der Wiederholungen für jedes Element anzugeben  $\{x_{i,j}; \text{freq}_{i,j}\} = \{1;1, 2;2, 3;3, 4;2, 5;1\}$  und berechnen Sie den Mittelwert.

**SHIFT** **MENU** (SETUP) **▼** **2** <sup>\*1</sup> oder **▼** **3** <sup>\*2</sup> (Statistik) **1** (Ein)

\*1: fx-87DE X \*2: fx-991DE X

**OPTN** **1** (Typ auswählen) **1** (1 Variable)

1 **≡** 2 **≡** 3 **≡** 4 **≡** 5 **≡** **▼** **▶**  
1 **≡** 2 **≡** 3 **≡** 2 **≡**

|   | x | Freq |
|---|---|------|
| 2 | 2 | 2    |
| 3 | 3 | 3    |
| 4 | 4 | 2    |
| 5 | 5 | 1    |

**AC** **OPTN** **▼** **2** (Variable) **1** ( $\bar{x}$ ) **≡**

3

**Bsp. 3:** Berechnen Sie die Korrelationskoeffizienten für die logarithmische Regression für folgende Variablenpaar-Daten und bestimmen Sie die Regressionsformel:  $(x; y) = (20; 3150), (110; 7310), (200; 8800), (290; 9310)$ . Legen Sie Fix 3 (drei Dezimalstellen) für die Ergebnisse fest.

**SHIFT** **MENU** (SETUP) **▼** **2** <sup>\*1</sup> oder **▼** **3** <sup>\*2</sup> (Statistik) **2** (Aus)

\*1: fx-87DE X \*2: fx-991DE X

**SHIFT** **MENU** (SETUP) **3** (Zahlenformat) **1** (Fix) **3**

**OPTN** **1** (Typ auswählen) **4** ( $y=a+b \cdot \ln(x)$ )

20 **≡** 110 **≡** 200 **≡** 290 **≡** **▼** **▶**  
3150 **≡** 7310 **≡** 8800 **≡** 9310 **≡**

|   | x   | y    |
|---|-----|------|
| 2 | 110 | 7310 |
| 3 | 200 | 8800 |
| 4 | 290 | 9310 |

**AC** **OPTN** **▼** **4** (Regressionen) **3** (r) **≡**

0,998

**AC** **OPTN** **▼** **4** (Regressionen) **1** (a) **≡**

-3857,984

**AC** **OPTN** **▼** **4** (Regressionen) **2** (b) **≡**

2357,532

## Schätzwerte berechnen

Anhand der mit einer statistischen Berechnung mit Variablenpaar erhaltenen Regressionsformel kann der Schätzwert von  $y$  für einen gegebenen  $x$ -Wert berechnet werden. Der entsprechende  $x$ -Wert (zwei



Beispiel:

Taschenrechnerfunktion

vgl. F 71

| $i$ | $x_i$ | $y_i$ |
|-----|-------|-------|
| 1   | 2     | 4     |
| 2   | 4     | 3     |
| 3   | 3     | 6     |
| 4   | 9     | 7     |
| 5   | 7     | 8     |

Gesucht: i)  $\hat{a}$ ,  $\hat{b}$ ,  $R^2$  ( $=r^2$ )

ii) Regressionsgerade

iii) Prognosewert für  $x = 6$

i)  $\hat{a} = 3,1$     $\hat{b} = 0,5$     $r = 0,7030$   
 $\rightarrow R^2 = 0,4942$

ii)  $\hat{y} = 3,1 + 0,5 \cdot x$

iii)  $\hat{y}(6) = 3,1 + 0,5 \cdot 6 = 6,1$

**fx-991 DEX:**

MENU  $\rightarrow$  Statistik  $\rightarrow$   $\boxed{\equiv}$   $\rightarrow$   $\boxed{2}$

$\swarrow$   $\boxed{6}$   $\searrow$

< Dateneingabe >

[ mittels  $\boxed{\equiv}$  abschließen

$\boxed{\wedge}$   $\boxed{\vee}$   $\boxed{>}$   $\boxed{<}$  Pfeile

nutzen]

$\boxed{\text{OPTN}}$   $\rightarrow$   $\boxed{4}$

$\leadsto a = 3,1$

$b = 0,5$

$r = 0,7030$

# ① Regressionsanalyse mit R am

Beispiel

| $x_i$ | $y_i$ |
|-------|-------|
| 2     | 4     |
| 4     | 3     |
| 3     | 6     |
| 9     | 7     |
| 7     | 8     |

Umsetzung  
mit R

R-Code

1) Regression

2) Plot mit Regressionsgerade

# ② Anscombe - Daten (F.106-108)

# ③ Regressionsanalyse mit Falldaten

(F.109)

1) summary (meine Regression)

↳ Outputinterpretation

i)  $\hat{a}$ ,  $\hat{b}$ ,  $R^2$  ?

ii) Regressionsgerade ?

iii) Prognosewert für Alter  $V = 50$

iv)  $r$  ?

2) plot (meine Regression)

↳ Residualanalyse zur Modellprüfung

Output -  
interpretation

# Regression: 4 eindimensionale Beispiele



► Berühmte Daten aus den 1970er Jahren:

|    | ①               | ②               | ③               | ④               | ①               | ②               | ③               | ④               |
|----|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| i  | x <sub>1i</sub> | x <sub>2i</sub> | x <sub>3i</sub> | x <sub>4i</sub> | y <sub>1i</sub> | y <sub>2i</sub> | y <sub>3i</sub> | y <sub>4i</sub> |
| 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            |

- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

(Quelle: Anscombe (1973))

in R: `data(anscombe)`  
`anscombe[order(anscombe[, x1]), ]`  
 # Daten sortieren w) Erster Eindruck<sup>106</sup> der Datenlage anhand der Datentabelle

# Regression: 4 eindimensionale Beispiele



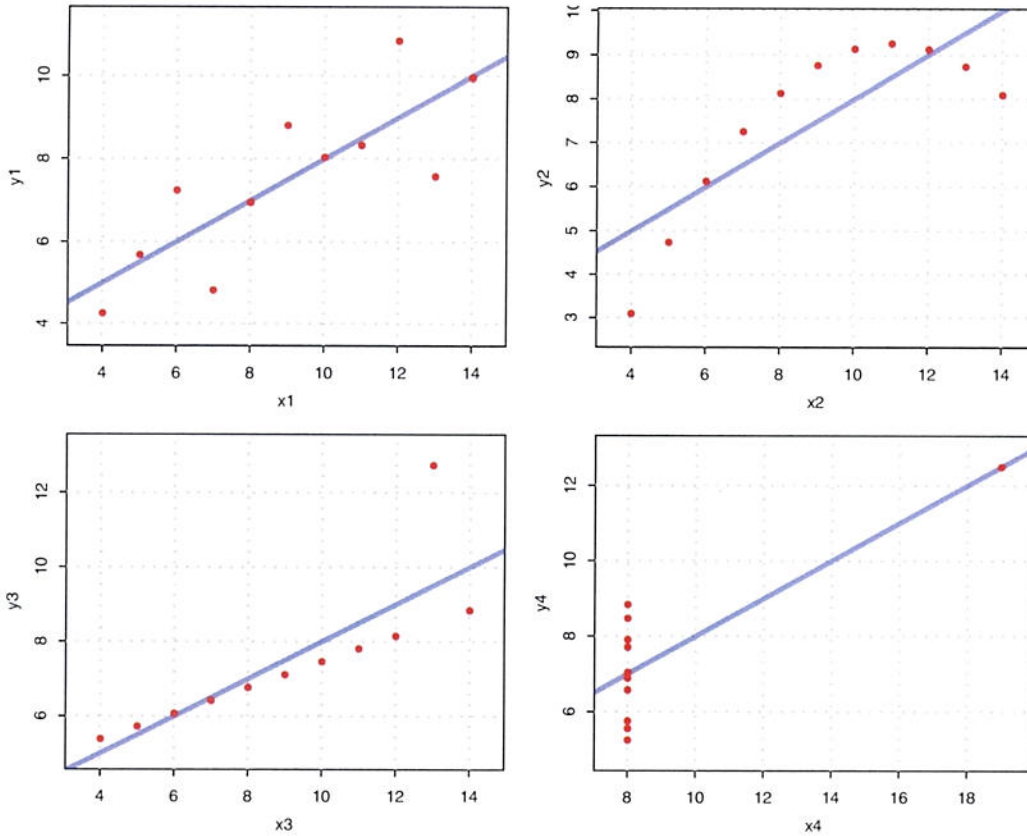
- In folgender Tabelle: Jeweils Ergebnisse der linearen Regressionsanalyse
- dabei: x<sub>k</sub> unabhängige Variable und y<sub>k</sub> abhängige Variable
- Modell jeweils:  $y_k = a_k + b_k x_k$

| 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

[?] Wie lautet r?  
 1.) selbes Vorzeichen wie  $\hat{b}$  (Steigungsparameter)  
 2.)  $r = \sqrt{R^2}$  (hier alle positiv.)

# Plot der Anscombe-Daten



Statistik



- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

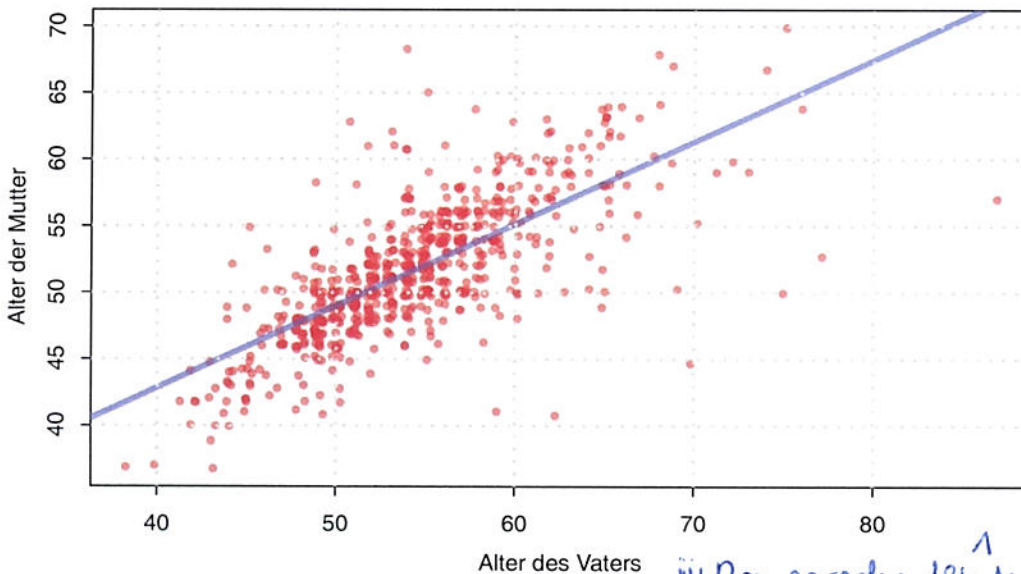
108

# Beispieldaten

```
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
```



Statistik



- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

aktuelle Daten:

i)  $\hat{\alpha} = 18,2309$   
 $\hat{\beta} = 0,6153$   
 $R^2 = 0,5211$

ii)  $\rightarrow r = +\sqrt{0,5211} = \dots$   
 iii) Reg. gerade:  $\text{AlterM} = 18,2309 + 0,6153 \cdot \text{AlterV}$

iv) Prognose:  $\hat{\text{AlterM}}(50) = 18,2309 + 0,6153 \cdot 50$   
 $\approx 48,9959$

### ③ Regressionsanalyse mit aktuellen Falldaten (R-Outputinterpretation)

```
> meineRegression=lm(AlterM~AlterV)
>
> plot(AlterV, AlterM, xlab="Alter des Vaters", ylab="Alter der Mutter", col=rgb(1,0,0,0.7), pch=16)
> abline(meineRegression, col="blue", lwd=2)
>
> meineRegression
```

```
Call:
lm(formula = AlterM ~ AlterV)
```

```
Coefficients:
(Intercept)      AlterV
    18.2309         0.6154
```

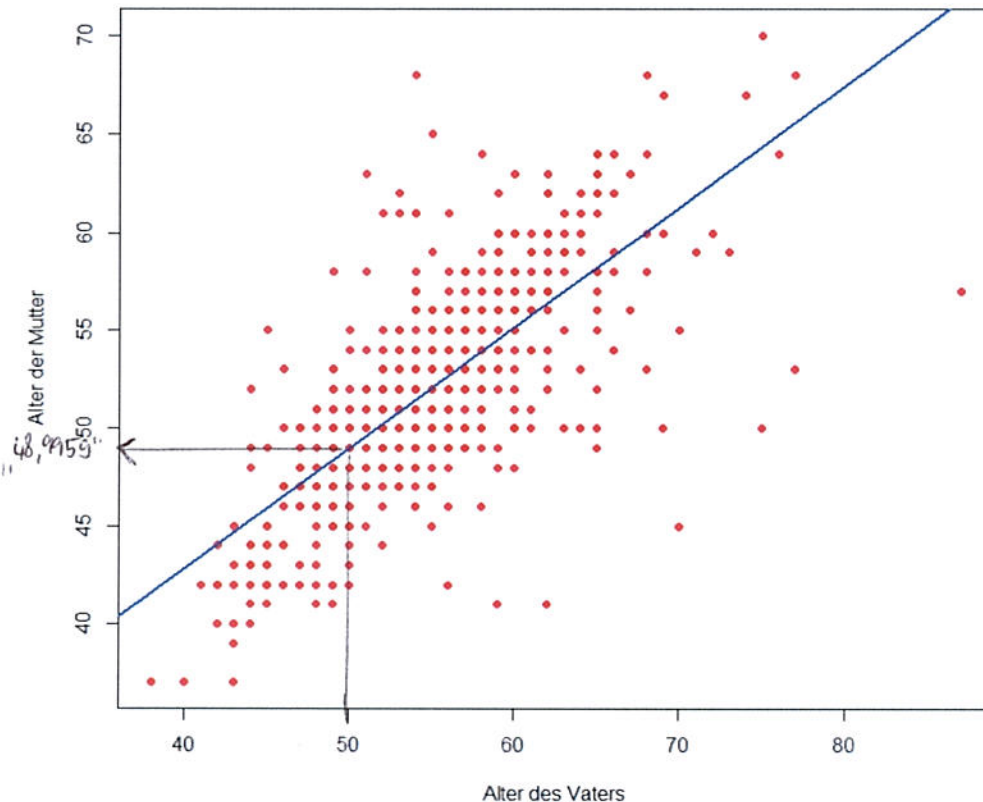
```
> summary(meineRegression)
```

```
Call:
lm(formula = AlterM ~ AlterV)
```

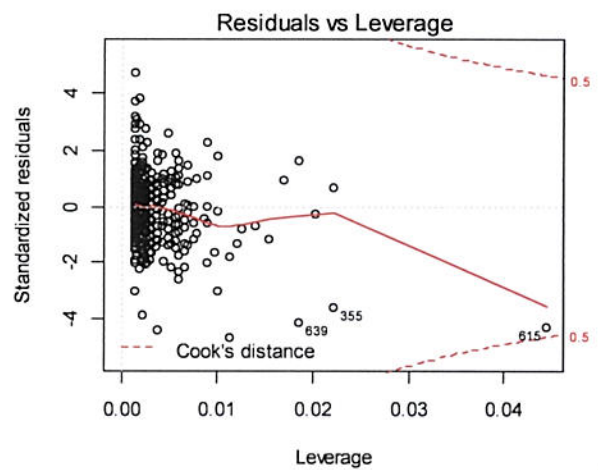
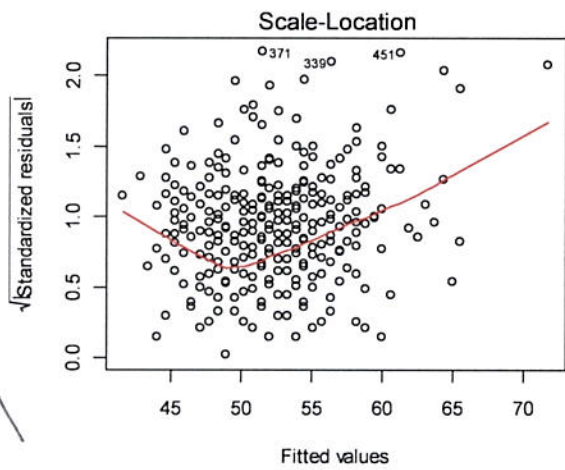
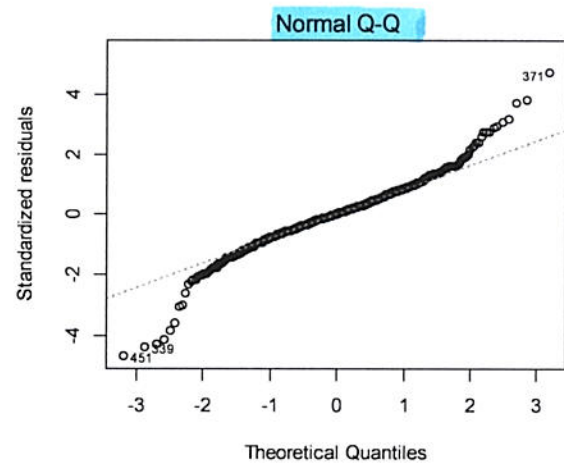
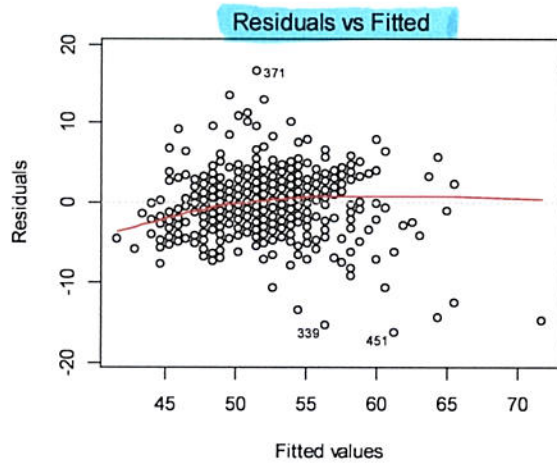
```
Residuals:
    Min       1Q   Median       3Q      Max
-16.3055  -1.9213   0.0015   1.9247  16.5401
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    18.2309     1.2190   14.96 <2e-16 ***
AlterV          0.6153     0.0223   27.60 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 3.51 on 700 degrees of freedom
(4 observations deleted due to missingness)
Multiple R-squared:  0.5211,    Adjusted R-squared:  0.5204
F-statistic: 761.7 on 1 and 700 DF,  p-value: < 2.2e-16
```



```
> par(mfrow=c(2,2))
> plot(meineRegression)
```





- ▶ 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:

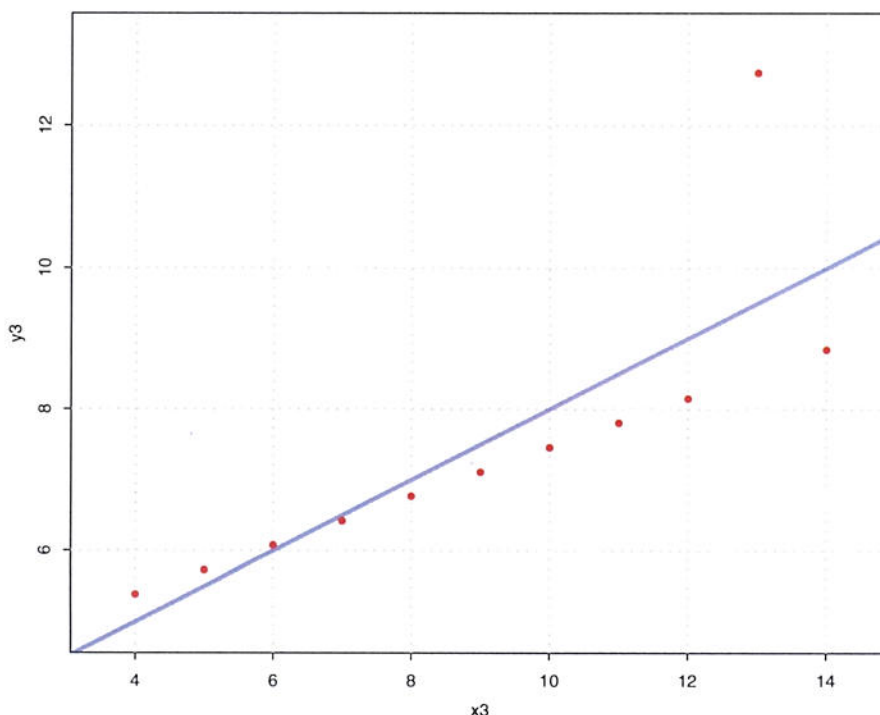
- $\hat{y}_j$ : Prognosewert des kompletten Modells für das j-te Objekt
- $\hat{y}_{j(\text{ohne } i)}$ : Prognosewert des Modells ohne Objekt i für das j-te Objekt
- $\text{MSE} = \frac{1}{n} \cdot \sum (\hat{y}_i - y_i)^2$ : Normierender Term (Schätzwert für Fehlerstreuung)

! ohne F. 110-111

- 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



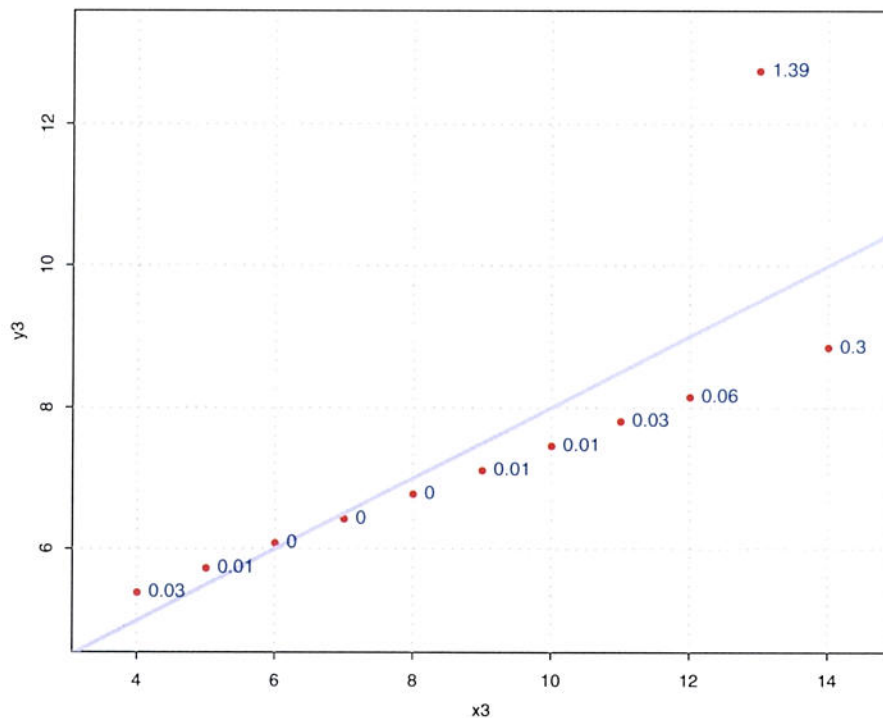
- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen



- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

## Ausreißer?

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



## Residualanalyse

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

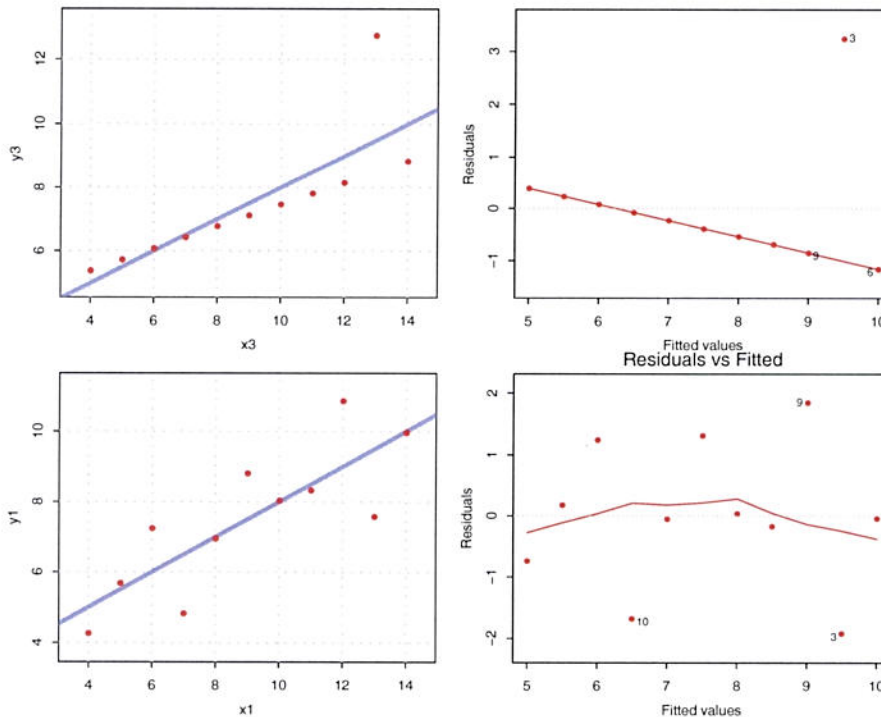


- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen





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

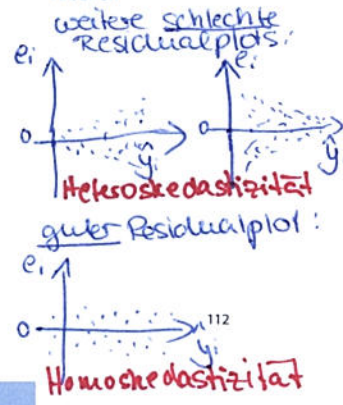


*schlecht!  
Trend / Trend  
in Residuen*

1. Einführung
2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
3. W-Theorie
4. Induktive Statistik

Quellen

Tabellen

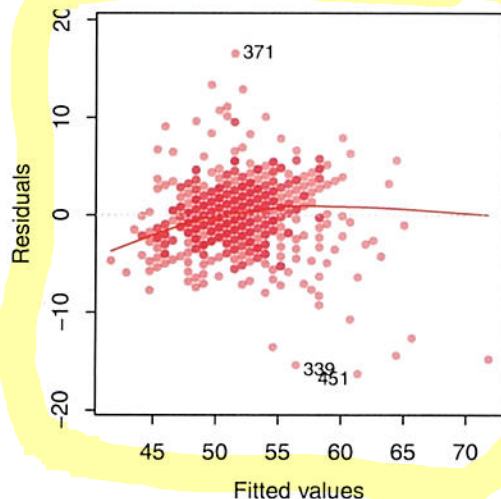
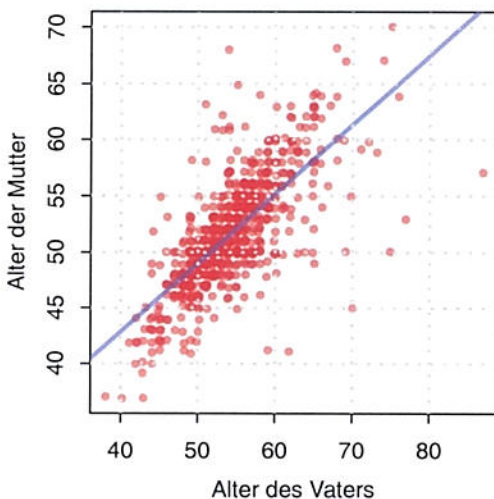


*R: plot (meine Regression)*

## 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. Plot*  
*2. Plot nicht abgebildet*



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

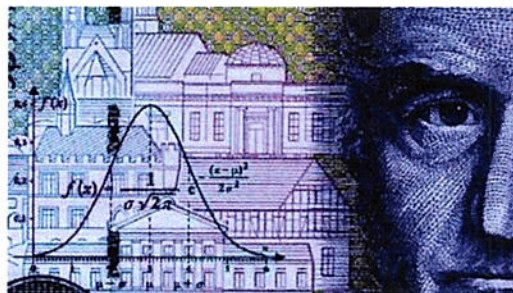
- ▶ Meist wichtig für sinnvolle Regressionsanalysen:
- ▶ **Kausale Verbindung** zwischen unabhängigem und abhängigem Merkmal (Wirkungsbeziehung)
- ▶ Sonst bei Änderung der unabhängigen Variablen keine sinnvollen Prognosen möglich
- ▶ Oft: **Latente Variablen** im Hintergrund



- 1. Einführung
- 2. Deskriptive Statistik
  - Häufigkeiten
  - Lage und Streuung
  - Konzentration
  - Zwei Merkmale
  - Korrelation
  - Preisindizes
  - Lineare Regression
- 3. W-Theorie
- 4. Induktive Statistik
- Quellen
- Tabellen

## Statistik: Table of Contents

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



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

## Aufgabe 2

18 Punkte

Stefan Jumper nimmt seit 8 Jahren am gleichen Marathon teil und hat jedes Jahr seine Trainings- und Ergebnisdaten dokumentiert. Er hat dazu pro Jahr jeweils die Ergebniszeit im Marathon (Merkmal *Ergebnis Marathon*, in Minuten) sowie die durchschnittliche Anzahl gelaufener Trainingskilometer in den 8 bzw. 16 Wochen vor dem Marathon (Merkmal *8-Wochen-Trainings-Durchschnitt* bzw. *16-Wochen-Trainings-Durchschnitt*, jeweils in km pro Woche) in der folgenden Tabelle festgehalten:

| Jahr                             | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   |
|----------------------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Ergebnis Marathon                | 172 | 161 | 156 | 147 | 152 | 167 | 157 | 168 |
| 8-Wochen-Trainings-Durchschnitt  | 85  | 105 | 125 | 150 | 130 | 90  | 110 | 95  |
| 16-Wochen-Trainings-Durchschnitt | 70  | 80  | 120 | 125 | 125 | 100 | 105 | 75  |

Sie haben Stefan laufend von den interessanten Themen der Statistik-Vorlesung berichtet, die Sie in diesem Semester besucht haben. Stefan wünscht sich daraufhin von Ihnen eine statistische Auswertung seiner Ergebnisse.

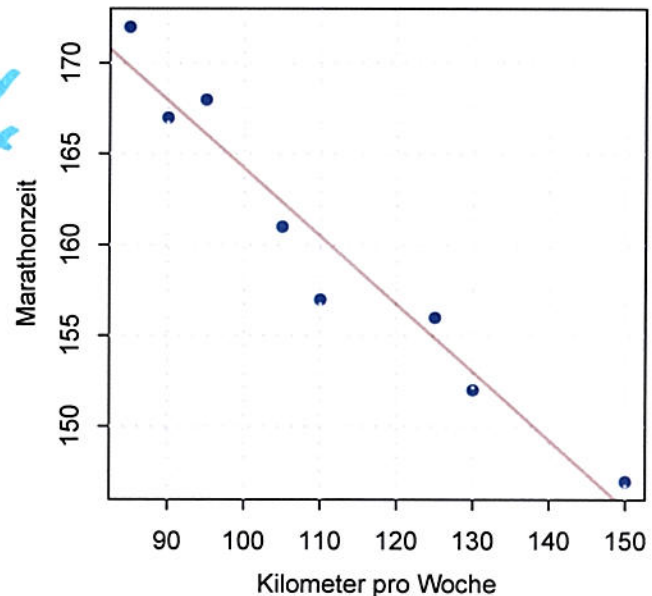
- a) Stefan vermutet, dass der 8-Wochen-Trainings-Durchschnitt eine bessere Prognose für seine gelaufene Marathonzeit darstellt als der 16-Wochen-Trainings-Durchschnitt.

Berechnen Sie jeweils einen geeigneten Korrelationskoeffizienten, vergleichen Sie die beiden Ergebnisse und geben Sie an, ob Stefans Vermutung durch das Ergebnis gestützt wird. ?

- b) Bestimmen Sie den Funktionsterm eines linearen Regressionsmodells für das Ergebnis im Marathon in Abhängigkeit vom 8-Wochen-Trainings-Durchschnitt.

- c) Zeichnen Sie die Datenpunkte sowie die Regressionsgerade in das nebenstehende Koordinatensystem ein. Beschriften Sie dazu auch die Achsen geeignet. (Hinweis: Der Schnittpunkt der Koordinatenachsen soll nicht bei (0, 0) liegen)

- d) Welche Marathonzeit erwartet Stefan auf Grundlage des linearen Modells für einen 8-Wochen-Trainings-Durchschnitt von 175 km pro Woche?



- e) Korrigieren Sie den folgenden R-Code mit dem ein wenig begabter Student einen Teil der obigen Aufgabe lösen wollte:

a) Bravais-Pearson-Korrelationskoeffizient

$$\begin{aligned} 8\text{-Wo./Mar. } r &\approx -0,9728 \\ 16\text{-Wo./Mar. } r &\approx -0,8805 \end{aligned}$$

$$r = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{\sum x_i^2 - n \bar{x}^2} \cdot \sqrt{\sum y_i^2 - n \bar{y}^2}}$$

Zusammenhang zw. 8-Wo. und Marathonzeit

b)  $\hat{\text{Marathon}} = 201,6290 - 0,3742 \cdot 8\text{-Wo.}$



d)  $\hat{\text{Marathon}}(175) = 201,6290 - 0,3742 \cdot 175 \approx 136,144$

lm(y~x)  
plot(x,y)

```
Zeit = c[172:161:156:147:152:167:157:168] ( ) statt [ ] ; " statt ;"  
8Wochendurchschnitt = seq(85,105,125,150,130,90,110,95) "e" statt "seq"; keine Zahlen,  
cor(Zeit, 8Wochendurchschnitt, method = "spearman") "e" statt "seq"; keine Zahlen,  
z <- LM(8Wochendurchschnitt~Zeit) "pearson" sollte zeichnen  
plotter(Zeit,8Wochendurchschnitt, color = "blau", blue Beginn von Objekt  
ylab = "Kilometer pro Woche", ylab = "Zeit")  
abline(z)
```

### Lösungshinweis:

```
Marathon = c(172, 161, 156, 147, 152, 167, 157, 168)  
WochenKM.8 = c(85, 105, 125, 150, 130, 90, 110, 95)  
WochenKM.16 = c(70, 80, 120, 125, 125, 100, 105, 75)
```

```
cor(Marathon, WochenKM.8, method = "pearson")  
cor(Marathon, WochenKM.16, method = "pearson")
```

```
Marathon.Regression = lm(Marathon ~ WochenKM.8)  
Marathon.Regression
```

```
plot(WochenKM.8, Marathon, col = "blue", pch=16, cex=1.5,  
      xlab = "Kilometer pro Woche", ylab = "Marathonzeit")  
grid()  
abline(Marathon.Regression, lwd=2, col="#90000050")
```

zu c) 2 geeignete Punkte der Regressionsgerade :

$$(\bar{x} | \bar{y}) = (111,25 | 160)$$

$$(90 | 167,951)$$

$$(140 | 149,241)$$