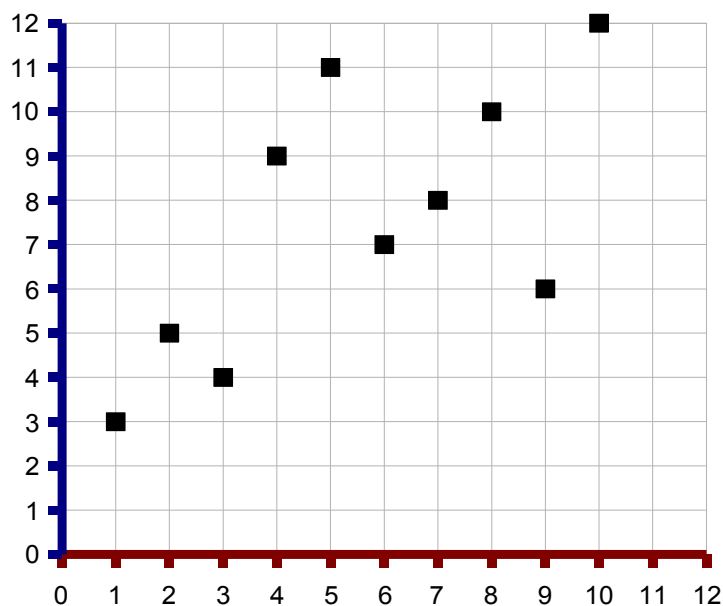


# Streudiagramme

für 2 metrisch skalierte Variablen

paarweise Messwerte (x,y)

|   |   |    |   |   |    |   |    |   |   |   |
|---|---|----|---|---|----|---|----|---|---|---|
| x | 3 | 10 | 6 | 4 | 8  | 2 | 5  | 7 | 9 | 1 |
| y | 4 | 12 | 7 | 9 | 10 | 5 | 11 | 8 | 6 | 3 |



## Aussagen zu Zusammenhängen

1.

### empirische Kovarianz

Standardabweichung der WertPAARE

$$s_{xy}$$

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

$(\text{Wert } x_1 - \text{Mittelwert aller } x)(\text{Wert } y_1 - \text{Mittelwert aller } y) + (\text{Wert } x_2 - \text{Mittelwert aller } x)(\text{Wert } y_2 - \text{Mittelwert aller } y) + \dots$   
Anzahl der WertPAARE

|   |   |    |   |   |    |   |    |   |   |   |                 |
|---|---|----|---|---|----|---|----|---|---|---|-----------------|
| x | 3 | 10 | 6 | 4 | 8  | 2 | 5  | 7 | 9 | 1 | $\bar{x} = 5,5$ |
| y | 4 | 12 | 7 | 9 | 10 | 5 | 11 | 8 | 6 | 3 | $\bar{y} = 7,5$ |

$$s_{xy} = \frac{1}{10} (3-5,5)(4-7,5) + (10-3,5)(12-7,5) + (4-3,5)(9-7,5) + \dots \approx 5,55$$

Erinnerung: Standardabweichung

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

$$s_y = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

## empirischer Korrelationskoeffizient Bravais-Pearson

"normierte" Standardabweichung der WertPAARE

$$r = r_{xy}$$



Der Teil  $1/n$   
aus den „Einzelformeln“  
kürzt sich weg

Normierung

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

$$\frac{\text{Abweichung zusammen}}{\text{Abweichung } x \cdot \text{Abweichung } y}$$

|   |   |    |   |   |    |   |    |   |   |   |                    |                       |
|---|---|----|---|---|----|---|----|---|---|---|--------------------|-----------------------|
| x | 3 | 10 | 6 | 4 | 8  | 2 | 5  | 7 | 9 | 1 | $s_x \approx 2,87$ | $s_{xy} \approx 5,55$ |
| y | 4 | 12 | 7 | 9 | 10 | 5 | 11 | 8 | 6 | 3 | $s_y \approx 2,87$ |                       |

$$r_{xy} = \frac{5,55}{2,87 \cdot 2,87} \approx 0,67$$

**r** kann **Werte zwischen -1 und +1** annehmen,  
und gibt damit Auskunft über **Stärke und Richtung des Zusammenhangs**

Interpretation:

|               |   |  |   |
|---------------|---|--|---|
| <b>r = -1</b> | vollständiger negativer Zusammenhang<br>(je mehr, desto weniger)<br><br><i>alle Punkte liegen auf einer fallenden Geraden</i> |  | <p><b>keine Aussage zu kausalen Beziehungen!</b></p> <p><b>Scheinkorrelationen möglich!</b></p> |
| <b>r = 1</b>  | vollständiger positiver Zusammenhang<br>(je mehr, desto mehr)<br><br><i>alle Punkte liegen auf einer steigenden Geraden</i>   |  |   |
| <b>r = 0</b>  | kein <b>LINEARER</b> Zusammenhang<br><br><i>...vielleicht aber ein nicht-linearer (exponentieller, u-förmiger oder ...)</i>   |  |   |

# lineare Regression

## Regressionsgerade

...auf der Suche nach  
einem berechenbaren Zusammenhang von  $x$  und  $y$   
(in Form einer linearen Gleichung /Funktion)

mit dem Ziel,  
anhand der Ausprägung eines (beliebigen)  $x$ -Wertes  
die Ausprägung des entsprechenden  $y$ -Wertes  
(ungefähr) berechnen zu können

Beziehung zwischen  
**U**nabhängiger **V**ariable + **A**bhängiger **V**ariable?

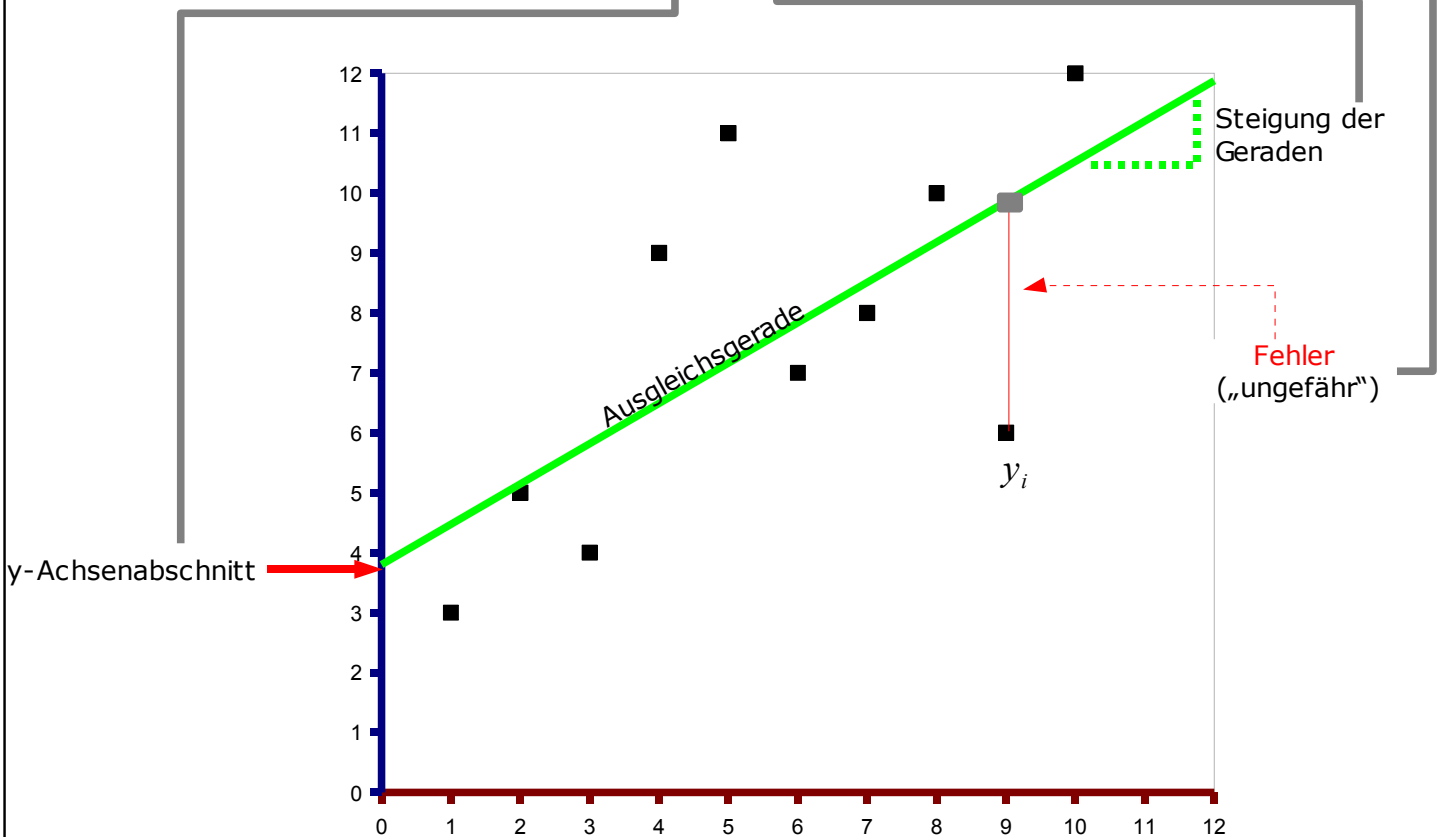
Gleichung

$$y_i = \alpha + \beta \cdot x_i + e_i$$

Funktion

$$f(x) = \alpha + \beta \cdot x$$

$$y_i = \alpha + \beta \cdot x_i + e_i$$



|  |                                 |
|--|---------------------------------|
| Berechnung der Ausgleichs-/Regressionsgeraden  | $f(x) = \alpha + \beta \cdot x$ |
| <h3>Kleinste-Quadrate-Schätzer</h3> <p>...auf der Suche nach<br/>einem Achsenabschnitt <math>a</math> und einer Steigung <math>\beta</math></p> <p>mit dem Ziel,<br/>den Gesamtfehler (<math>e_1+e_2+e_3+\dots</math>) so klein wie möglich zu halten</p> <p>-&gt; Minimierung<br/>der <b>durchschnittlichen Abweichung</b> der tatsächlichen Werte<br/>von den Werten auf der Ausgleichsgeraden</p> | $\hat{\alpha}$ $\hat{\beta}$    |

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

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

vergleiche:

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

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

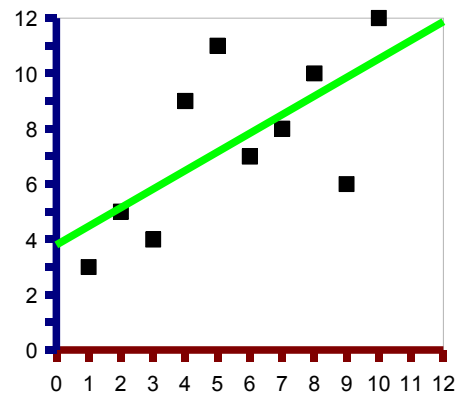
wegen irgendwelchen Wegkürzungen...

$\rightarrow s_{xx} \hat{=} s_x^2$

|   |   |    |   |   |    |   |    |   |   |   |               |                       |                      |
|---|---|----|---|---|----|---|----|---|---|---|---------------|-----------------------|----------------------|
| x | 3 | 10 | 6 | 4 | 8  | 2 | 5  | 7 | 9 | 1 | $\bar{x}=5,5$ | $s_{xy} \approx 5,55$ | $s_{xx} \approx 8,3$ |
| y | 4 | 12 | 7 | 9 | 10 | 5 | 11 | 8 | 6 | 3 | $\bar{y}=7,5$ |                       |                      |

$\hat{\beta} = \frac{5,5}{8,3} \approx 0,67$      $\rightarrow$  Sorry, entspricht in diesem Datensatz zufällig  $r$ ,  
weil  $s_x = s_y$  und damit  $s_x^2 = s_x \cdot s_y$

$$\hat{\alpha} = 7,5 - 0,67 \cdot 5,5 \approx 3,82$$



$$f(x) = 3,82 + 0,67x$$

Berechnung der einzelnen Fehler

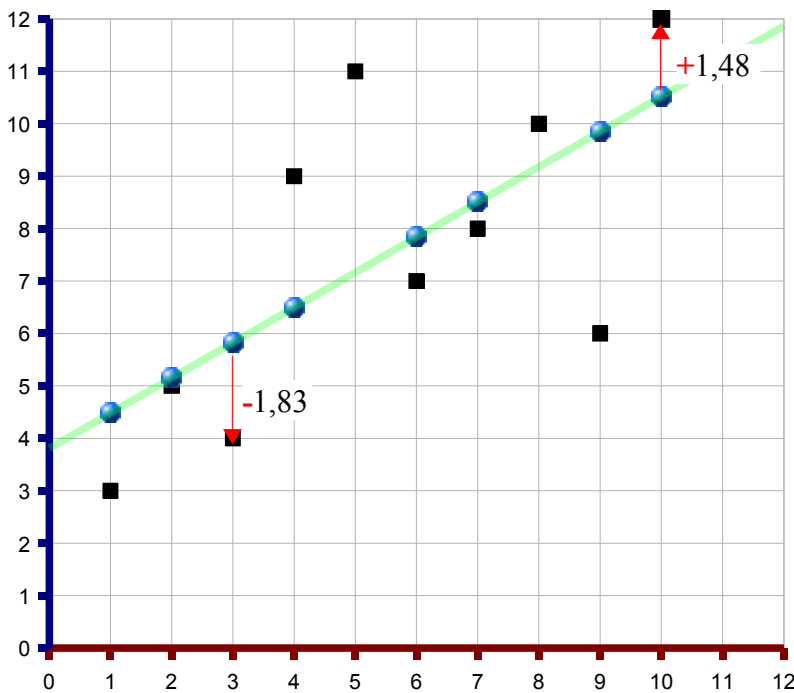
**Residuen**

$$\hat{e}_i$$

$$\hat{e}_i = y_i - \hat{y}_i$$

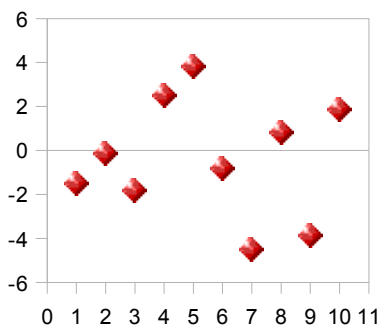
$$\hat{y}_i = \hat{\alpha} + \hat{\beta} \cdot x_i$$

|             |       |       |                              |       |      |      |       |       |      |       |                       |
|-------------|-------|-------|------------------------------|-------|------|------|-------|-------|------|-------|-----------------------|
| x           | 1     | 2     | 3                            | 4     | 5    | 6    | 7     | 8     | 9    | 10    |                       |
| y           | 3     | 5     | 4                            | 9     | 11   | 7    | 8     | 10    | 6    | 12    |                       |
| $\hat{y}_i$ | 4,49  | 5,16  | $3,82 + 0,67 \cdot 3 = 5,83$ | 6,5   | 7,17 | 7,84 | 8,51  | 9,18  | 9,85 | 10,52 | $\hat{y}_i + e_i = 0$ |
| $\hat{e}_i$ | -1,49 | -0,16 | $4 - 5,83 = -1,83$           | -1,83 | 2,5  | 3,83 | -0,84 | -4,51 | 0,82 | -3,85 |                       |



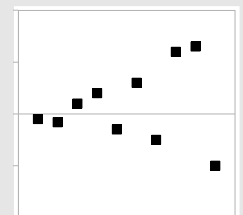
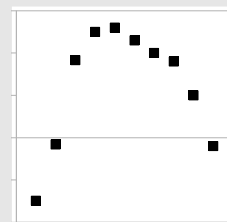
**Streudiagramm der Residuen/Residualanalyse**  
**Residualplot**

|             |       |       |       |     |      |       |       |      |       |      |
|-------------|-------|-------|-------|-----|------|-------|-------|------|-------|------|
| x           | 1     | 2     | 3     | 4   | 5    | 6     | 7     | 8    | 9     | 10   |
| $\hat{e}_i$ | -1,49 | -0,16 | -1,83 | 2,5 | 3,83 | -0,84 | -4,51 | 0,82 | -3,85 | 1,84 |



*ideal:*  
Die Residuen streuen **unsystematisch** und **möglichst nahe bei 0**

andere Residualplots weisen darauf hin, dass



der Datensatz **nicht die Voraussetzungen für ein lineares Regressionsmodell** erfüllt

alternative Plots (statt  $x, \hat{e}$ ) :  $\hat{y}, \hat{e}$  oder (mit standardisierten Residuen)  $\hat{y}, \hat{d}$  ...

## Berechnung der unterschiedlichen Varianzen/ Varianzanteile

**Streuungszerlegung**

$$SQT = SQE + SQR$$

**Sum of SQ**uares **T**otal:

Gesamtstreuung

Streuung  
der y-Werte  
um ihren Mittelwert  $\bar{y}$ 

$$SQT = \sum_{i=1}^n (y_i - \bar{y})^2$$

**Sum of SQ**uares **E**xplained:

erklärte Streuung


Streuung der  
von der Regressionsgeraden erfassten  
 $\hat{y}$  -Werte  
um den Mittelwert  $\bar{y}$ 

$$SQE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

**Sum of SQ**uares **R**esidual:Residualstreuung  
/ ReststreuungStreuung der  
von der Regressionsgeraden nicht erfassten  
y -Werte  
um ihren entsprechenden Wert  $\hat{y}$   
auf der Regressionsgeraden

$$SQR = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

*"Streuung der Fehler  $\hat{e}_i = y_i - \hat{y}_i$  "*

wäre im Beispiel ungefähr  $82,3 = 37,1 + 45,2$  (ziemlich großzügig passend gemacht ,  )

Wie genau lassen sich mit dem Modell (*hier: einfache lineare Regression*) die y-Werte vorhersagen?

**Bestimmtheitsmaß /  
Determinationskoeffizient**

$R^2$

$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2} = \frac{SQE}{SQT} \Rightarrow R^2 = r_{xy}^2$$

- Wertebereich **zwischen 0 und 1**

- **je näher an 0, desto schlechter** ist die Vorhersage

- im Fall  $R^2 = 1$  liegen bereits die Originaldaten auf einer Geraden

$$R^2 = \frac{37,1}{82,3} \approx 0,45 \quad \approx r^2 = 0,67^2$$

heißt: Das Modell eignet sich nicht so doll, um irgendwelche y-Werte vorauszusagen,

(z.B. "wie wäre der y-Wert zu x= 20?")

weil nur 45% der Gesamtstreuung von y durch die Regression erklärt werden

=  
Die Varianz der y-Werte ist nur zu 45% auf die Varianz der x-Werte zurückzuführen.