

Korrelation ist ein Maß dafür, wie stark die Merkmalswerte (Daten) von dem gewählten Regressionsmodell (z.B. Regressionsgerade oder Regressionspolynom) abweichen.

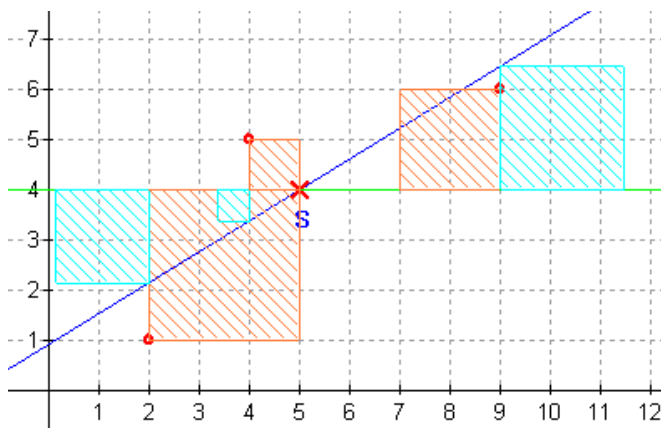
Man definiert den **Betrag der Korrelation r** als Verhältnis (Anteil) der Standardabweichung, die durch das Regressionsmodell festgelegt ist zur Standardabweichung, die durch die Daten gegeben ist :

Kurzschreibweise:

$$|r| = \frac{\sigma_{\text{Modell}}}{\sigma_{\text{Daten}}}$$

Was ist damit gemeint ?

Die folgende Grafik soll das erläutern (3 Punkte P1(2/1) , P2(4/5) , P3(9/6) seien gegeben):



Hier ist als Modell eine Gerade angenommen worden. Ihre Gleichung ist

$$y = \frac{8}{13}x + \frac{12}{13}$$

Für die beiden Standardabweichungen erhält man (Begründung siehe weiter unten):

$$\sigma_{\text{Daten}} = \sqrt{\frac{14}{3}} \quad \sigma_{\text{Modell}} = \sqrt{\frac{128}{13 \cdot 3}}$$

Das arithmetische Mittel der Daten heißt **Schwerpunkt** $S(\bar{x}; \bar{y})$, hier: **S(5 ; 4)** .

Es werden zunächst die roten Quadrate der vertikalen Abweichungen der Daten vom Schwerpunkt S betrachtet und ihre Flächeninhalte aufsummiert. $\sum_{i=1}^n (y_i - \bar{y})^2$

Teilt man diese Summe durch die Anzahl n der Daten (hier n = 3) und zieht dann die Wurzel,

so ergibt sich die Standardabweichung $\sigma_{\text{Daten}} = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n}}$

Entsprechend geht man vor bei der Bestimmung von σ_{Modell} . Dort werden die blauen Quadrate der vertikalen Abweichungen des Modells „Gerade“ vom Schwerpunkt S betrachtet .

Für die Standardabweichung erhält man: $\sigma_{\text{Modell}} = \sqrt{\frac{\sum_{i=1}^n (f(x_i) - \bar{y})^2}{n}}$

Begründung der beiden Ergebnisse für die Standardabweichungen :

Für die roten Quadrate kann man die Summe der Flächeninhalte am Grafikraster ablesen (14 Einheiten).

Für die blauen Quadrate rechnet man so:

$$\text{Summe der Flächeninhalte} = \left(\frac{8}{13} \cdot 2 + \frac{12}{13} - 4\right)^2 + \left(\frac{8}{13} \cdot 4 + \frac{12}{13} - 4\right)^2 + \left(\frac{8}{13} \cdot 9 + \frac{12}{13} - 4\right)^2 = \frac{128}{13} \approx 9,85$$

Teilt man jetzt σ_{Modell} durch σ_{Daten} , so erhält man $r = \sqrt{\frac{64}{91}} \approx \underline{\underline{0,8386}}$

Dies ist der gesuchte Korrelationskoeffizient für das obige Modell „Gerade“ .

Vereinfachung der Formel für den Korrelationskoeffizienten:

Bei der Division von σ_{Modell} durch σ_{Daten} kürzt sich n weg und man erhält die Formel:

$$|r| = \sqrt{\frac{\sum_{i=1}^n (f(x_i) - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad \text{Allgemeiner Korrelationskoeffizient ; } |r| \leq 1$$

Der Korrelationskoeffizient liegt also zwischen -1 und 1 .

Folgende Einteilung ist üblich:

$ r = 1$	$0,7 \leq r < 1$	$0,3 \leq r < 0,7$	$0 < r < 0,3$	$r = 0$
volle Korrelation	starke Korrelation	mittlere Korrelation	schwache Korrelation	keine Korrelation

Eine starke Korrelation bedeutet nicht zwangsläufig, dass zwischen den betrachteten Merkmalen ein kausaler Zusammenhang besteht.

Ob ein solcher Zusammenhang besteht, muss immer auch sachlich beurteilt werden !

Beispiel für ein quadratisches Modell mit $f(x) = 5,05x^2 - 0,21x + 0,2$:

Daten: (1 ; 5) (2 ; 20,1) (3 ; 44,9) (4 ; 80,2)

$$y_{\text{quer}} = (5 + 20,1 + 44,9 + 80,2) / 4 = 37,55$$

$$\text{sum}((y_i - y_{\text{quer}})^2) = ((5 - 37,55)^2 + (20,1 - 37,55)^2 + (44,9 - 37,55)^2 + (80,2 - 37,55)^2) = \underline{\underline{3237,05}}$$

$$\text{sum}((f(x_i) - y_{\text{quer}})^2) = ((5,04 - 37,55)^2 + (19,98 - 37,55)^2 + (45,02 - 37,55)^2 + (80,16 - 37,55)^2) = 1056,9 + 308,7 + 55,801 + 1815,6 = \underline{\underline{3237,018}}$$

$$r = \sqrt{(3237,018 / 3237,05)} = \underline{\underline{0,9999901...}}$$

Ersichtlich ist zur Berechnung des Allgemeinen Korrelationskoeffizienten r die Kenntnis der Funktionsgleichung der Modellannahme erforderlich.

Im speziellen Fall der linearen Funktion $f(x) = mx + b$ als Modell kommt man aber auch ohne die Kenntnis der Parameter m und b aus, wie die folgende Herleitung zeigt:

Korrelationskoeffizient für die Lineare Regression:

Ersetzt man in der allgemeinen Formel $f(x)$ durch $mx + b$, so folgt:

$$|r| = \frac{\sqrt{\sum_{i=1}^n (mx_i + b - \bar{y})^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad \bar{y} = m\bar{x} + b = \frac{\sqrt{\sum_{i=1}^n (mx_i + b - m\bar{x} - b)^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\sqrt{\sum_{i=1}^n m^2 \cdot (x_i - \bar{x})^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \sqrt{m^2 \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Jetzt setzt man den aus der Formel für die Lineare Regression bekannten Term für m ein:

$$|r| = \sqrt{\frac{\left[\frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^2 \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}}{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}}} \quad \text{Hier ist schon die Unabhängigkeit von } b, m \text{ zu sehen.}$$

Weitere Umformungen:

$$|r| = \sqrt{\frac{\left(\frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^2 \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}}{\left(\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \right)^2}} = \sqrt{\frac{\left(\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) \right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\left| \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) \right|}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Lässt man die Betragsstriche weg, so erhält man den vorzeichenbehafteten Wert von r :

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \quad \text{Linearer Korrelationskoeffizient}$$

Anmerkung:

Dividiert man den Term im Zähler durch die Anzahl n der Datenpaare, so erhält man die sog. **Kovarianz** .

Dividiert man die Faktoren im Nenner jeweils durch n , so erhält man die **Varianz der x-Daten** sowie die **Varianz der y-Daten**.

Hat man also die Kovarianz und die beiden anderen Varianzen berechnet, so lässt sich in einfacher Weise auch r berechnen, ohne dass m und b bekannt sein müssen !

Beispiel: $x_i = \{2, 4, 9\}$ $y_i = \{1, 5, 6\}$

$x_{\text{quer}} = 5$ $y_{\text{quer}} = 4$ (Schwerpunktkoordinaten: siehe oben)

$n \cdot \text{Kovarianz} = (2-5)(1-4) + (4-5)(5-4) + (9-5)(6-4) = 9 - 1 + 8 = 16$

$n \cdot \text{Varianz}_x = (2-5)^2 + (4-5)^2 + (9-5)^2 = 9 + 1 + 16 = 26$

$n \cdot \text{Varianz}_y = (1-4)^2 + (5-4)^2 + (6-4)^2 = 9 + 1 + 4 = 14$

$\rightarrow r = \frac{16}{\sqrt{26 \cdot 14}} = \frac{16}{\sqrt{364}} \approx \underline{\underline{0,8386}}$

Aufgabe zur (Lin.) Korrelation:

R.Doll untersuchte 1955 als erster systematisch den möglichen Zusammenhang zwischen Zigarettenkonsum und Erkrankungen an Lungenkrebs. Er trug folgende Daten zusammen:

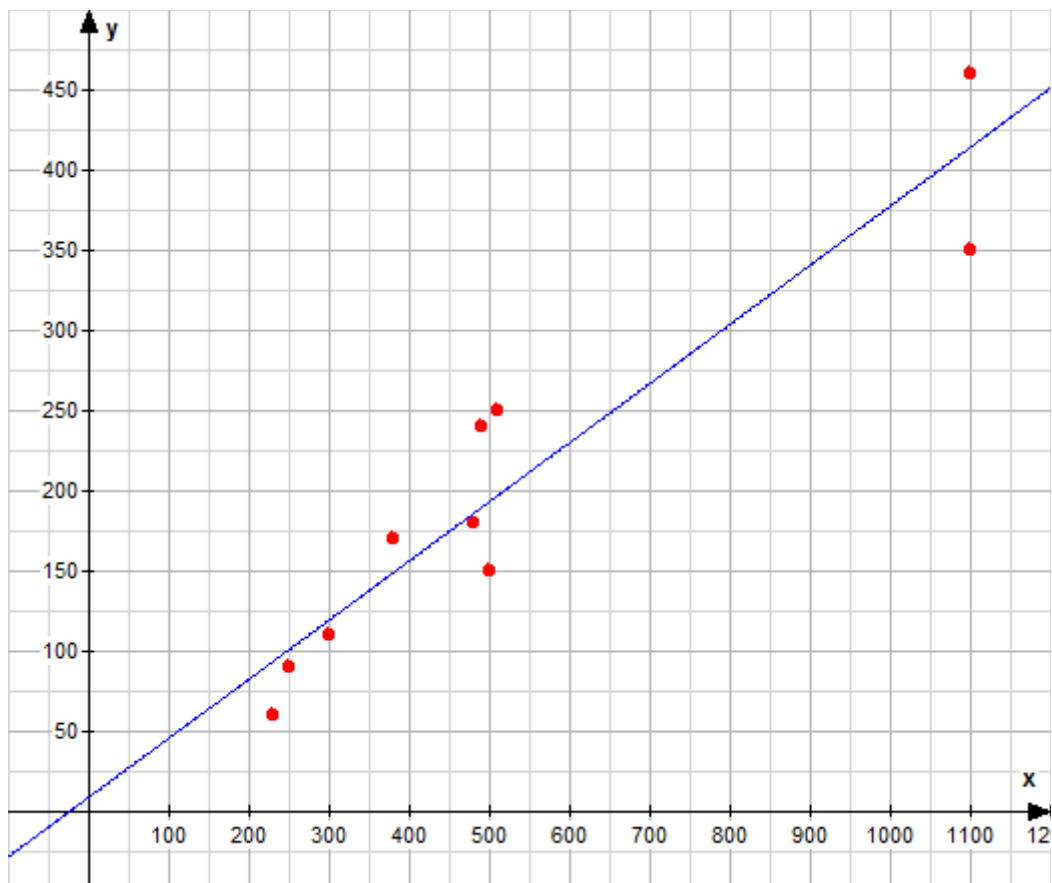
x = Zigarettenverbrauch pro Kopf 1930

y = Todesfälle an Lungenkrebs je Million 1950

z = Land, in dem die Untersuchung stattfand .

x	y	z
230	60	Island
250	90	Norwegen
300	110	Schweden
380	170	Dänemark
480	180	Australien
490	240	Niederlande
500	150	Kanada
510	250	Schweiz
1100	350	Finnland
1100	460	England

Lösung : $y = 0,368653x + 9,139335$ $r = 0,942763$



Es liegt eine starke Korrelation vor.

Dennoch ist es gewagt, von einem kausalen (ursächlichen) Zusammenhang zu sprechen.

Anmerkung:

Inzwischen ist dieser kausale Zusammenhang aufgrund weiterer Untersuchungen bewiesen.