

LINEARE REGRESSION UND KORRELATION - EIN ELEMENTARER ZUGANG

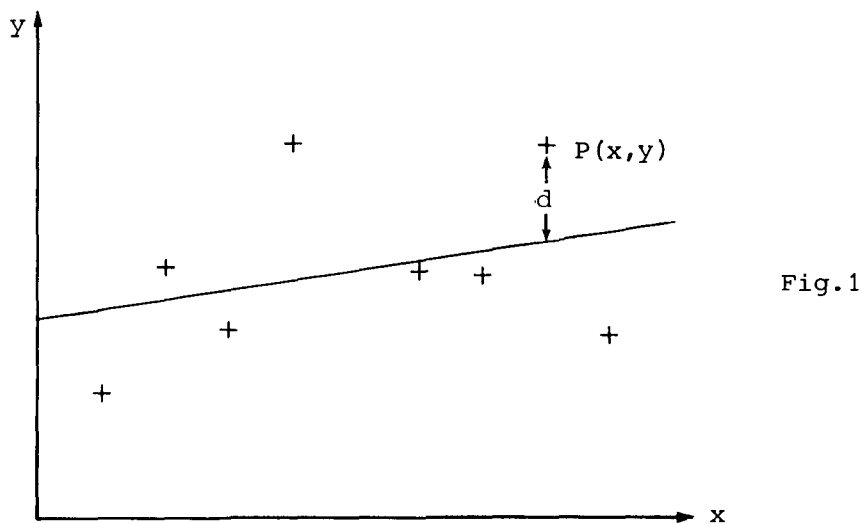
von S.M. Goode und E.J. Gold, Bristol

Originaltitel in "Teaching Statistics" Vol. 9 (1987) Nr. 2:
Linear Regression and Correlation - an enlightening approach

Übersetzung: Hans-Joachim Bentz, Osnabrück

Zusammenfassung: In diesem Artikel werden die Parameter der Regressionsgleichung sowie einige wesentliche mathematische Eigenschaften des Korrelationskoeffizienten auf *elementarem* Wege hergeleitet.

Gegeben seien die beiden Variablen x und y :



Die Regressionsgerade von y auf x ist so definiert, daß $\sum d_i^2$ minimiert wird. Nimmt man folgende Variablentransformation vor:

$$X = x - \bar{x}, \quad Y = y - \bar{y}, \quad \text{wobei} \quad \bar{x} = \frac{\sum x_i}{n}, \quad \bar{y} = \frac{\sum y_i}{n},$$

so erhält man folgendes Bild:

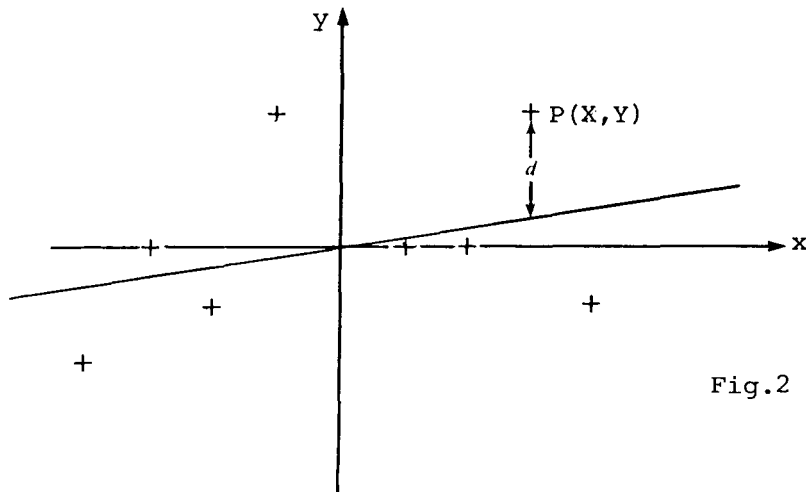


Fig.2

Bemerkung: In der deutschen Literatur ist manches Mal die Unterscheidung zwischen X und x üblich: X bezeichnet die Zufallsvariable, x einen konkreten Wert von X , etwa $x=14$. Die Bezeichnungen oben haben damit nichts gemeinsam.

In einigen Lehrbüchern werden die Themengebiete Regression und Korrelation sehr zu unserem Bedauern vollständig getrennt behandelt. In anderen wiederum wird vermieden, auf die Theorie einzugehen, wie man zu den Gleichungen für die Regressionsgerade bzw. zu einer Rechtfertigung für die Formel für den Produktmoment-Korrelationskoeffizienten kommt. Manchmal wird in Lehrbüchern auch auf partielle Differentiation zurückgegriffen, um die ersteren Gleichungen zu erhalten - ein Stoffgebiet, das in den meisten Lehrplänen für Reine Mathematik auf dem G/E Advanced Level nicht vorkommt. Dabei gibt es einen ganz elementaren und trotzdem vollständig verständlichen Zugang, dieser soll im folgenden dargestellt werden:

Die Regressionsgerade von Y auf X sei festgelegt durch:

$$Y = mX + c$$

Zu minimieren ist $\sum d_i^2$, d.h.

$$\sum d_i^2 = \sum (Y_i - mX_i - c)^2.$$

Der Übersicht wegen werden im folgenden die Indices i weglassen. Ausquadrieren liefert:

$$\begin{aligned} \sum d^2 &= \sum Y^2 + m^2 \sum X^2 + \sum c^2 - 2m \sum XY + 2mc \sum X - 2c \sum Y = \\ &= \sum Y^2 + m^2 \sum X^2 - 2m \sum XY + nc^2 \end{aligned}$$

Dies deswegen, weil $\sum X$ und $\sum Y$ gleich Null sind. Vervollständigen des Quadrats in m ergibt ferner:

$$\sum d^2 = (\sum X^2) \left[m - \frac{\sum XY}{\sum X^2} \right]^2 + \sum Y^2 - \frac{(\sum XY)^2}{\sum X^2} + nc^2$$

Nun sieht man unmittelbar, daß dieser Ausdruck minimal wird, wenn

$c = 0$ (d.h. die Regressionsgerade geht durch (\bar{x}, \bar{y}))

und

$$m = \frac{\sum XY}{\sum X^2} = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

ist.

Die Regressionsgerade von y auf x ist daher

$$Y = \frac{\sum XY}{\sum X^2} \cdot X \quad \text{bzw.} \quad y - \bar{y} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} (x - \bar{x}).$$

Die Regressionsgerade von y auf x kann man analog bestimmen.

Der minimale Wert der Abweichungsquadrate der Daten von der Regressionsgeraden $\sum d^2$ ist nun:

$$\sum d^2 = \sum Y^2 - \frac{(\sum XY)^2}{\sum X^2}$$

Dies ist eine Summe von Quadraten, daher folgt:

$$\sum Y^2 - \frac{(\sum XY)^2}{\sum X^2} \geq 0 \quad \text{d.h.} \quad \sum Y^2 \geq \frac{(\sum XY)^2}{\sum X^2} \quad \text{bzw.}$$

$$\frac{(\sum XY)^2}{\sum X^2 \sum Y^2} \leq 1, \quad \text{also:} \quad -1 \leq \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} \leq 1$$

Die Definition des Produkt-Moment-Korrelationskoeffizienten mit

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x) \text{var}(y)}}$$

folgt daraus in natürlicher Weise, wobei $r = \pm 1$ gilt, falls $\sum d^2 = 0$ (was wiederum bedeutet, daß die Punkte (x, y) alle auf einer Geraden liegen), und $r = 0$ gilt, falls die Kovarianz, $\text{cov}(x, y)$ gleich Null ist (d.h. die Punkte sind "gleichmäßig" verstreut in den vier Quadranten zum Punkt (\bar{x}, \bar{y})).