



Allgemeines Lineares Modell

BSc Psychologie SoSe 2022

Prof. Dr. Dirk Ostwald

(2) Korrelation

Grundlagen

Korrelation und Bestimmtheitsmaß

Korrelation und lineare Abhängigkeit

Korrelation und Regression

Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Grundlagen

Korrelation und Bestimmtheitsmaß

Korrelation und lineare Abhängigkeit

Korrelation und Regression

Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Anwendungsszenario

Psychotherapie



Mehr Therapiestunden

⇒ Höhere Wirksamkeit?

Unabhängige Variable

- Anzahl Therapiestunden

Abhängige Variable

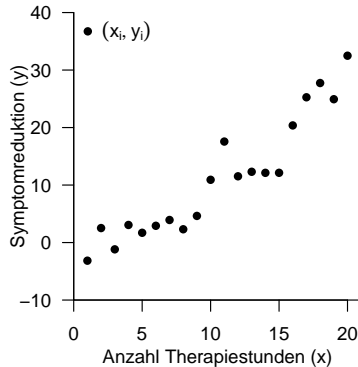
- Symptomreduktion

Beispieldatensatz

$i = 1, \dots, 20$ Patient:innen, y_i Symptomreduktion bei Patient:in i , x_i Anzahl Therapiestunden von Patient:in i

y_i	x_i
-3.15	1
2.52	2
-1.18	3
3.06	4
1.70	5
2.91	6
3.92	7
2.31	8
4.63	9
10.91	10
17.56	11
11.52	12
12.31	13
12.12	14
12.13	15
20.37	16
25.26	17
27.75	18
24.93	19
32.49	20

Beispieldatensatz



Wie stark hängen Anzahl Therapiestunden und Symptomreduktion zusammen?

Definition (Korrelation)

Die *Korrelation* zweier Zufallsvariablen X und Y ist definiert als

$$\rho(X, Y) := \frac{\mathbb{C}(X, Y)}{\mathbb{S}(X)\mathbb{S}(Y)} \quad (1)$$

wobei $\mathbb{C}(X, Y)$ die Kovarianz von X und Y und $\mathbb{V}(X)$ und $\mathbb{V}(Y)$ die Varianzen von X und Y , respektive, bezeichnen.

Bemerkungen

- $\rho(X, Y)$ wird auch *Korrelationskoeffizient* von X und Y genannt.
- Wir haben bereits gesehen, dass $-1 \leq \rho(X, Y) \leq 1$ gilt.
- Wenn $\rho(X, Y) = 0$ ist, werden X und Y *unkorreliert* genannt.
- Wir haben bereits gesehen, dass aus der Unabhängigkeit von X und Y , folgt dass $\rho(X, Y) = 0$.
- Aus $\rho(X, Y) = 0$ folgt aber wie bereits gesehen die Unabhängigkeit von X und Y im Allgemeinen nicht.

Definition (Stichprobenkorrelation)

$\{(x, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}$ sei eine Wertemenge. Weiterhin seien:

- Die Stichprobenmittel der x_i und y_i definiert als

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i \text{ und } \bar{y} := \frac{1}{n} \sum_{i=1}^n y_i. \quad (2)$$

- Die Stichprobenstandardabweichungen x_i und y_i definiert als

$$s_x := \sqrt{\frac{1}{n-1} (x_i - \bar{x})^2} \text{ und } s_y := \sqrt{\frac{1}{n-1} (y_i - \bar{y})^2}. \quad (3)$$

- Die Stichprobenkovarianz der $(x, y_1), \dots, (x_n, y_n)$ definiert als

$$c_{xy} := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n). \quad (4)$$

Dann ist die *Stichprobenkorrelation* der $(x, y_1), \dots, (x_n, y_n)$ definiert als

$$r_{xy} := \frac{c_{xy}}{s_x s_y} \quad (5)$$

und wird auch *Stichprobenkorrelationskoeffizient* genannt.

Beispiel

```
# Laden des Beispieldatensatzes
fname = file.path(getwd(), "2_Daten", "2_Korrelation_Beispieldatensatz.csv") # Dateipfad
D      = read.table(fname, sep = ",", header = TRUE)                       # Laden als Dataframe
x_i    = D$x_i                                                             # x_i Werte
y_i    = D$y_i                                                             # y_i Werte
n      = length(x_i)                                                       # n

# "Manuelle" Berechnung der Stichprobenkorrelation
x_bar  = (1/n)*sum(x_i)                                                     # \bar{x}
y_bar  = (1/n)*sum(y_i)                                                     # \bar{y}
s_x    = sqrt(1/(n-1)*sum((x_i - x_bar)^2))                                # s_x
s_y    = sqrt(1/(n-1)*sum((y_i - y_bar)^2))                                # s_y
c_xy   = 1/(n-1) * sum((x_i - x_bar) * (y_i - y_bar))                      # c_{xy}
r_xy   = c_xy/(s_x * s_y)                                                   # r_{xy}
print(r_xy)                                                                # Ausgabe

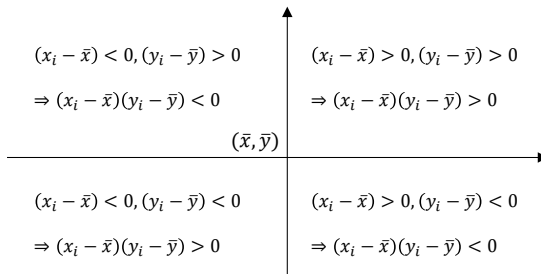
> [1] 0.938

# Automatische Berechnung mit cor()
r_xy   = cor(x_i,y_i)                                                       # r_{xy}
print(r_xy)                                                                # Ausgabe

> [1] 0.938
```

⇒ Anzahl Therapiestunden und Symptomreduktion sind hochkorreliert.

Mechanik der Kovariationsterme

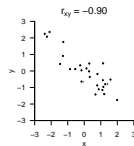
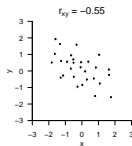
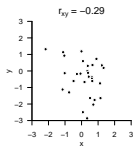
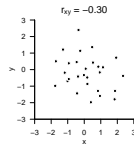
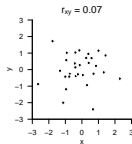
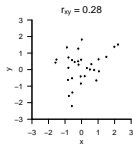
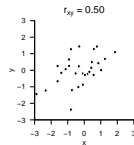
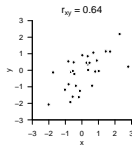
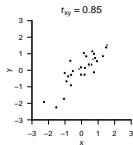


Häufige richtungsgleiche Abweichung der x_i und y_i von ihren Mittelwerten \Rightarrow Positive Korrelation

Häufige richtungsungleiche Abweichung der x_i und y_i von ihren Mittelwerten \Rightarrow Negative Korrelation

Keine häufigen richtungsgleichen oder -entgegengesetzten Abweichungen \Rightarrow Keine Korrelation

Beispiele



Theorem (Stichprobenkorrelation bei linear-affinen Transformationen)

Für eine Wertemenge $\{(x_i, y_i)\}_{i=1, \dots, n} \subset \mathbb{R}^2$ sei $\{(\tilde{x}_i, \tilde{y}_i)\}_{i=1, \dots, n} \subset \mathbb{R}^2$ eine linear-affin transformierte Wertemenge mit

$$(\tilde{x}_i, \tilde{y}_i) = (a_x x_i + b_x, a_y y_i + b_y), a_x, a_y \neq 0. \quad (6)$$

Dann gilt

$$|r_{\tilde{x}\tilde{y}}| = |r_{xy}|. \quad (7)$$

Bemerkungen

- Der Betrag der Stichprobenkorrelation ändert sich bei linear-affiner Datentransformation nicht.
- Man sagt, dass die Stichprobenkorrelation im Gegensatz zur Stichprobenkovarianz *maßstabsunabhängig* ist.

Grundlagen

Beweis

Es gilt

$$\begin{aligned} r_{\tilde{x}\tilde{y}} &:= \frac{\frac{1}{n-1} \sum_{i=1}^n (\tilde{x}_i - \bar{\tilde{x}})(\tilde{y}_i - \bar{\tilde{y}})}{\sqrt{\frac{1}{n-1} \left(\sum_{i=1}^n \tilde{x}_i - \bar{\tilde{x}} \right)^2} \sqrt{\frac{1}{n-1} \left(\sum_{i=1}^n \tilde{y}_i - \bar{\tilde{y}} \right)^2}} \\ &= \frac{\sum_{i=1}^n (a_x x_i + b_x - (a_x \bar{x} + b_x))(a_y y_i + b_y - (a_y \bar{y} + b_y))}{\sqrt{\sum_{i=1}^n (a_x x_i + b_x - (a_x \bar{x} + b_x))^2} \sqrt{\sum_{i=1}^n (a_y y_i + b_y - (a_y \bar{y} + b_y))^2}} \\ &= \frac{a_x a_y \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{a_x^2 \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{a_y^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (8) \\ &= \frac{a_x a_y}{|a_x| |a_y|} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{a_x a_y}{|a_x| |a_y|} \frac{c_{xy}}{s_x s_y} \\ &= \frac{a_x a_y}{|a_x| |a_y|} r_{xy}. \end{aligned}$$

Also folgt, durch Durchspielen aller möglichen Vorzeichenfälle, dass

$$|r_{\tilde{x}\tilde{y}}| = |r_{xy}|. \quad (9)$$

□

Grundlagen

Korrelation und Bestimmtheitsmaß

Korrelation und lineare Abhängigkeit

Korrelation und Regression

Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Überblick

Das sogenannte Bestimmtheitsmaß R^2 ist eine beliebte Statistik.

Numerisch ist R^2 das Quadrat des Stichprobenkorrelationskoeffizienten.

Ist die Stichprobenkorrelation $r_{xy} = 0.5$, dann ist $R^2 = 0.25$, ist $r_{xy} = -0.5$, dann ist $R^2 = 0.25$.

⇒ R^2 enthält also weniger Information über die Rohdaten als r_{xy} , da das Vorzeichen wegfällt.

⇒ *Perse* ist die Angabe von R^2 anstelle von r_{xy} im Kontext der Korrelation zweier Variablen wenig sinnvoll.

Ein tieferes Verständnis von R^2 erlaubt jedoch

- (1) Einen Einstieg in das Konzept von Quadratsummenzerlegungen, einem wichtigen ALM Evaluationsprinzip.
- (2) Einen Einstieg in das Verständnis der Zusammenhänge von Ausgleichsgerade und Stichprobenkorrelation.
- (3) Einen ersten Einblick in die Tatsache, dass Korrelationen (nur) linear-affine Zusammenhänge quantifizieren.

Definition (Erklärte Werte und Residuen einer Ausgleichsgerade)

Gegeben seien eine Wertemenge $\{(x, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ und die zu dieser Wertemenge gehörende Ausgleichsgerade

$$f_{\hat{\beta}} : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto f_{\hat{\beta}}(x) := \hat{\beta}_0 + \hat{\beta}_1 x \quad (10)$$

Dann werden für $i = 1, \dots, n$

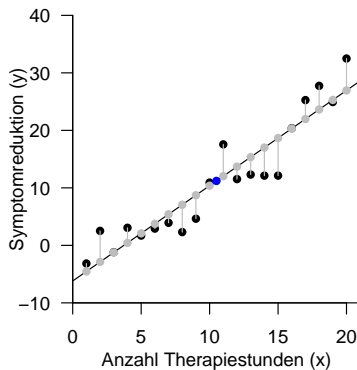
$$\hat{y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (11)$$

die durch die Ausgleichsgerade *erklärten Werte* genannt und

$$\hat{\varepsilon}_i := y_i - \hat{y}_i \quad (12)$$

die *Residuen* der Ausgleichsgerade genannt.

Erklärte Werte und Residuen



$\bullet (x_i, y_i)$ $\bullet (\bar{x}, \bar{y})$ $— f_{\hat{\beta}}(x)$ $\bullet \hat{y}_i$ $— \hat{\varepsilon}_i$ $i = 1, \dots, n$

Theorem (Quadratsummenzerlegung bei Ausgleichsgerade)

Für eine Wertemenge $\{(x, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ und ihre zugehörige Ausgleichsgerade $f_{\hat{\beta}}$ seien für

$$\bar{y} := \frac{1}{n} \sum_{i=1}^n y_i \text{ und } \hat{y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i, \text{ für } i = 1, \dots, n \quad (13)$$

das Stichprobenmittel der y -Werte und die durch die Ausgleichsgerade erklärten Werte, respektive. Weiterhin seien

$$\text{SQT} := \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{die Total Sum of Squares}$$

$$\text{SQE} := \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad \text{die Explained Sum of Squares}$$

$$\text{SQR} := \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{die Residual Sum of Squares}$$

Dann gilt

$$\text{SQT} = \text{SQE} + \text{SQR} \quad (14)$$

Bemerkungen

- SQT repräsentiert die Gesamtstreuung der y_i -Werte um ihren Mittelwert \bar{y} .
- SQE repräsentiert die Streuung der erklärten Werte \hat{y}_i um ihren Mittelwert
 - ⇒ Große Werte von SQE repräsentieren eine große absolute Steigung der y_i mit den x_i
 - ⇒ Kleine Werte von SQE repräsentieren eine kleine absolute Steigung der y_i mit den x_i
- SQE ist also ein Maß für die Stärke des linearen Zusammenhangs der x - und y -Werte
- SQR ist die Summe der quadrierten Residuen, es gilt

$$\text{SQR} := \sum_{i=1}^n (y_i - \hat{y}_i)^2 := \sum_{i=1}^n \hat{\epsilon}_i^2 \quad (15)$$

- ⇒ Große Werte von SQR repräsentieren große Abweichungen der erklärten von den beobachteten y -Werten
- ⇒ Kleine Werte von SQR repräsentieren geringe Abweichungen der erklärten von den beobachteten y -Werten
- SQR ist also ein Maß für die Güte der Beschreibung der Datenmenge durch die Ausgleichsgerade.

Beweis

$$\begin{aligned}\text{SQT} &= \sum_{i=1}^n (y_i - \bar{y})^2 \\&= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\&= \sum_{i=1}^n ((y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}))^2 \\&= \sum_{i=1}^n \left((y_i - \hat{y}_i)^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + (\hat{y}_i - \bar{y})^2 \right) \\&= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\&= \text{SQE} + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \text{SQR} \\&= \text{SQE} + \text{SQR}\end{aligned}\tag{16}$$

Beweis (fortgeführt)

Dabei ergibt sich die letzte Gleichung mit

$$\bar{\hat{y}} := \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i) = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x} = \bar{y} \quad (17)$$

und damit auch

$$\bar{\hat{y}} = \bar{y} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n y_i \Leftrightarrow \sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i \Leftrightarrow \bar{y} \sum_{i=1}^n \hat{y}_i = \bar{y} \sum_{i=1}^n y_i \quad (18)$$

sowie

$$\bar{\hat{y}} = \bar{y} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n y_i \Leftrightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i \Leftrightarrow \sum_{i=1}^n y_i \hat{y}_i = \sum_{i=1}^n \hat{y}_i \hat{y}_i \quad (19)$$

aus

Beweis (fortgeführt)

$$\begin{aligned}\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n (y_i \hat{y}_i - y_i \bar{y} - \hat{y}_i \hat{y}_i + \hat{y}_i \bar{y}) \\&= \sum_{i=1}^n y_i \hat{y}_i - \sum_{i=1}^n y_i \bar{y} - \sum_{i=1}^n \hat{y}_i \hat{y}_i + \sum_{i=1}^n \hat{y}_i \bar{y} \\&= \sum_{i=1}^n y_i \hat{y}_i - \sum_{i=1}^n \hat{y}_i \hat{y}_i + \bar{y} \sum_{i=1}^n \hat{y}_i - \bar{y} \sum_{i=1}^n y_i \\&= 0 + 0 \\&= 0\end{aligned}\tag{20}$$

□

Definition (Bestimmtheitsmaß R^2)

Für eine Wertemenge $\{(x, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ und ihre zugehörige Ausgleichsgerade $f_{\hat{\beta}}$ sowie die zugehörigen Explained Sum of Squares SQE und Total Sum of Squares SQT heißt

$$R^2 := \frac{\text{SQE}}{\text{SQT}} \quad (21)$$

Bestimmtheitsmaß oder Determinationskoeffizient.

Theorem (Stichprobenkorrelation und Bestimmtheitsmaß)

Für eine Wertemenge $\{(x, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ sei R^2 das Bestimmtheitsmaß und r_{xy} sei die Stichprobenkorrelation. Dann gilt

$$R^2 = r_{xy}^2. \quad (22)$$

Bemerkungen

- Mit $-1 \leq r_{xy} \leq 1$ folgt aus dem Theorem direkt, dass $0 \leq R^2 \leq 1$.
- Es gilt $R^2 = 0$ genau dann, wenn $SQE = 0$ ist
 - \Rightarrow Für $R^2 = 0$ ist die erklärte Streuung der Daten durch die Ausgleichsgerade gleich null.
 - $\Rightarrow R^2 = 0$ beschreibt also den Fall einer denkbar schlechten Erklärung der Daten durch die Ausgleichsgerade.
- Es gilt $R^2 = 1$ genau dann, wenn $SQE = SQT$ ist.
 - \Rightarrow Für $R^2 = 1$ ist also die Gesamtstreuung gleich der durch die Ausgleichsgerade erklärten Streuung.
 - $\Rightarrow R^2 = 1$ beschreibt also den Fall das sämtliche Datenvariabilität durch die Ausgleichsgerade erklärt wird.

Korrelation und Bestimmtheitsmaß

Beweis

Wir halten zunächst fest, dass mit

$$\bar{\hat{y}} := \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i) = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x} = \bar{y} \quad (23)$$

folgt, dass

$$\begin{aligned} \text{SQE} &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 \\ &= \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \hat{\beta}_0 - \hat{\beta}_1 \bar{x})^2 \\ &= \sum_{i=1}^n (\hat{\beta}_1 (x_i - \bar{x}))^2 \\ &= \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned} \quad (24)$$

Korrelation und Bestimmtheitsmaß

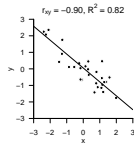
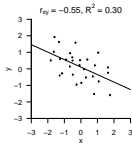
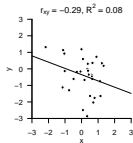
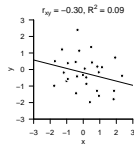
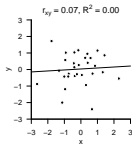
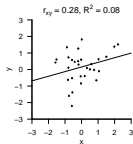
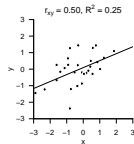
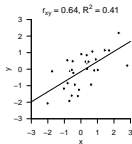
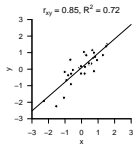
Beweis

Damit ergibt sich dann

$$\begin{aligned} R^2 &= \frac{SQE}{SQT} \\ &= \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\ &= \hat{\beta}_1^2 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\ &= \frac{c_{xy}^2 \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}{s_x^4 \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} \\ &= \frac{c_{xy}^2}{s_x^4} \frac{s_x^2}{s_y^2} \\ &= \frac{c_{xy}^2}{s_x^2 s_y^2} \\ &= \left(\frac{c_{xy}}{s_x s_y} \right)^2 \\ &= r_{xy}^2. \end{aligned} \tag{25}$$

□

Beispiele



Grundlagen

Korrelation und Bestimmtheitsmaß

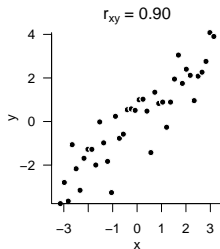
Korrelation und lineare Abhängigkeit

Korrelation und Regression

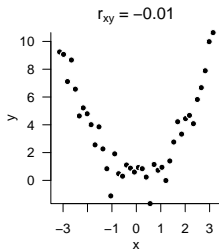
Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Funktionale Abhängigkeiten und Stichprobenkorrelation

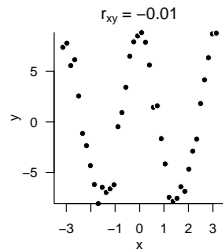


$$y_i = x_i + \varepsilon_i$$



$$y_i = x_i^2 + \varepsilon_i$$

$$\varepsilon_i \sim N(0, 1)$$



$$y_i = 8 \cos(2x_i) + \varepsilon_i$$

Theorem (Korrelation und linear-affine Abhängigkeit)

X und Y seien zwei Zufallsvariablen mit positiver Varianz. Dann besteht genau dann eine lineare-affine Abhängigkeit der Form

$$Y = \beta_0 + \beta_1 X \text{ mit } \beta_0, \beta_1 \in \mathbb{R} \quad (26)$$

zwischen X und Y , wenn

$$\rho(X, Y) = 1 \text{ oder } \rho(X, Y) = -1 \quad (27)$$

gilt.

Bemerkungen

- Die lineare Abhängigkeit $Y = \beta_0 + \beta_1 X$ impliziert eine lineare Abhängigkeit $X = \tilde{\beta}_0 + \tilde{\beta}_1 Y$, denn

$$Y = \beta_0 + \beta_1 X \Leftrightarrow -\beta_0 + Y = \beta_1 X \Leftrightarrow X = -\frac{\beta_0}{\beta_1} + \frac{1}{\beta_1} Y \Leftrightarrow X = \tilde{\beta}_0 + \tilde{\beta}_1 Y \quad (28)$$

mit

$$\tilde{\beta}_0 = -\frac{\beta_0}{\beta_1} \text{ und } \tilde{\beta}_1 = \frac{1}{\beta_1}. \quad (29)$$

Korrelation und lineare Abhängigkeit

Beweis

Wir beschränken uns auf den Beweis der Aussage, dass aus $Y = \beta_0 + \beta_1 X$ folgt, dass $\rho(X, Y) = \pm 1$ ist. Dazu halten wir zunächst fest, dass mit den Theoremen zu den Eigenschaften von Erwartungswert und Varianz gilt, dass

$$\mathbb{E}(Y) = \beta_0 + \beta_1 \mathbb{E}(X) \text{ und } \mathbb{V}(Y) = \beta_1^2 \mathbb{V}(X). \quad (30)$$

Wegen $\mathbb{V}(X) > 0$ und $\mathbb{V}(Y) > 0$ gilt damit $\beta_1 \neq 0$. Es folgt dann

$$\beta_1 > 0 \Rightarrow \mathbb{S}(Y) = \beta_1 \mathbb{S}(X) > 0 \text{ und } \beta_1 < 0 \Rightarrow \mathbb{S}(Y) = -\beta_1 \mathbb{S}(X) > 0. \quad (31)$$

Weiterhin gilt

$$\begin{aligned} Y - \mathbb{E}(Y) &= \beta_0 + \beta_1 X - \mathbb{E}(Y) \\ &= \beta_0 + \beta_1 X - \beta_0 - \beta_1 \mathbb{E}(X) \\ &= \beta_1 X - \beta_1 \mathbb{E}(X) \\ &= \beta_1 (X - \mathbb{E}(X)). \end{aligned} \quad (32)$$

Für die Kovarianz von X und Y ergibt sich also

$$\begin{aligned} \mathbb{C}(X, Y) &= \mathbb{E}((Y - \mathbb{E}(Y))(X - \mathbb{E}(X))) \\ &= \mathbb{E}(\beta_1 (X - \mathbb{E}(X))(X - \mathbb{E}(X))) \\ &= \beta_1 \mathbb{E}((X - \mathbb{E}(X))^2) \\ &= \beta_1 \mathbb{V}(X). \end{aligned} \quad (33)$$

Damit ergibt für die Korrelation von X und Y

$$\rho(X, Y) = \frac{\mathbb{C}(X, Y)}{\mathbb{S}(X)\mathbb{S}(Y)} = \pm \frac{\beta_1 \mathbb{V}(X)}{\mathbb{S}(X)\beta_1 \mathbb{S}(X)} = \pm \frac{\beta_1 \mathbb{V}(X)}{\beta_1 \mathbb{V}(X)} = \pm 1. \quad (34)$$

Grundlagen

Korrelation und lineare Abhängigkeit

Korrelation und Regression

Korrelation und Bestimmtheitsmaß

Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Überblick

Der fundamentale Unterschied zwischen “Korrelation” und “Regression” ist, dass

- bei Korrelation sowohl die UV (die x 's) als auch die AV (die y 's) als Zufallsvariablen modelliert werden,
- bei Regression dagegen lediglich die AV als Zufallsvariable modelliert wird und die UV als vorgegeben gilt.

Dieser Tatsache unbenommen, kann man auf gegebene Daten prinzipiell natürlich sowohl “Korrelation” als auch “Regression” anwenden. Das Ergebnis einer Regressionsanalyse lässt sich in das Ergebnis einer Korrelationsanalyse umrechnen. Die zusätzlich Durchführung einer Korrelationsanalyse bei durchgeführter Regressionsanalyse erzeugt kein mehr an Information oder Verständnis über den Zusammenhang von UV und AV.

Für ein tieferes Verständnis dieser Zusammenhänge ist ein Regressionsmodell nötig, indem auch die UV eine Zufallsvariable ist. In Abgrenzung zum Modell der einfachen linearen Regression, in dem die UV keine Zufallsvariable ist, bezeichnen wir dieses Modell als *Regression*. Letztlich gerät die Terminologie hier an eine Grenze und es muss jeweils geprüft bzw. geschlossen werden, welches Modell Datenanalysten nun tatsächlich vorschwebt.

Weiterhin treffen wir die Annahme, dass sowohl die UV als auch die AV im Regressionsmodell normalverteilt sind. Diese Annahme ist nicht zwingend nötig, da Aussagen zur Regression zweier Zufallsvariablen im Wesentlichen die Erwartungswerte, Varianzen, und Kovarianzen berühren. Allerdings wird die Normalverteilungsannahme für UV und AV im Anwendungskontext häufig getroffen und bereitet didaktisch sinnvoll auf das Konzept multivariater Normalverteilungen vor, das für die ALM Theorie zentral ist.

Theorem (Bivariate Normalverteilung)

Z_1 und Z_2 seien zwei standardnormalverteilte Zufallsvariablen und $\mu_1 \in \mathbb{R}, \mu_2 \in \mathbb{R}, \sigma_1 > 0, \sigma_2 > 0$ und $\rho \in]-1, 1[$ seien Konstanten. Weiterhin seien x und z zwei Zufallsvariablen definiert als

$$\begin{aligned}x &:= \sigma_1 Z_1 + \mu_1 \\z &:= \sigma_2 \left(\rho Z_1 + (1 - \rho^2)^{\frac{1}{2}} Z_2 \right) + \mu_2\end{aligned}\tag{35}$$

Dann hat die gemeinsame WDF von x und z die Form

$$\begin{aligned}p : \mathbb{R}^2 \rightarrow \mathbb{R}_{>0}, \begin{pmatrix} x \\ z \end{pmatrix} \mapsto p \left(\begin{pmatrix} x \\ z \end{pmatrix} \right) &= \frac{1}{2\pi(1 - \rho^2)^{1/2}\sigma_1\sigma_2} \\&\times \exp \left(-\frac{1}{2(1 - \rho^2)} \left(\left(\frac{x - \mu}{\sigma_1} \right)^2 - 2\rho \left(\left(\frac{x - \mu_1}{\sigma_1} \right) \left(\frac{z - \mu_2}{\sigma_2} \right) \right) + \left(\frac{z - \mu}{\sigma_2} \right)^2 \right) \right)\end{aligned}\tag{36}$$

Definition (Korrelationsmatrix eines Zufallsvektors)

y sei ein n -dimensionaler Zufallsvektor. Die *Korrelationsmatrix* von y ist definiert als

$$\mathbb{R}(y) = \left(\rho(y_i, y_j) \right)_{1 \leq i, j \leq n} = \begin{pmatrix} \rho(y_1, y_1) & \rho(y_1, y_2) & \cdots & \rho(y_1, y_n) \\ \rho(y_2, y_1) & \rho(y_2, y_2) & \cdots & \rho(y_2, y_n) \\ \vdots & \vdots & \ddots & \vdots \\ \rho(y_n, y_1) & \rho(y_n, y_2) & \cdots & \rho(y_n, y_n) \end{pmatrix}. \quad (37)$$

Bemerkung

- Die Korrelationsmatrix $\mathbb{R}(y)$ ist also die Matrix der Korrelationen der Komponenten von y .

Theorem (Korrelationsmatrix und Kovarianzmatrix)

y sei ein n -dimensionaler Zufallsvektor mit Kovarianzmatrix $\mathbb{C}(y)$. Weiterhin sei

$$S_y := \text{diag} \left(\sqrt{\mathbb{C}(y_1, y_1)}, \sqrt{\mathbb{C}(y_2, y_2)}, \dots, \sqrt{\mathbb{C}(y_n, y_n)} \right) \quad (38)$$

die Diagonalmatrix der Standardabweichungen der Komponenten von y . Dann gelten

$$\mathbb{R}(y) = S_y^{-1} \mathbb{C}(y) S_y^{-1} \quad (39)$$

und

$$\mathbb{C}(y) = S_y \mathbb{R}(y) S_y \quad (40)$$

Bemerkungen

- $\mathbb{C}(y)$ kann mithilfe ihrer Diagonalelemente in $\mathbb{R}(y)$ umgerechnet werden.
- $\mathbb{R}(y)$ kann mithilfe der Diagonalelemente von $\mathbb{C}(y)$ in $\mathbb{C}(y)$ umgerechnet werden.
- $\mathbb{C}(y)$ enthält also mehr Information als $\mathbb{R}(y)$.

Beweis

Wir halten zunächst fest, dass

$$S_y^{-1} = \text{diag} \left(\frac{1}{\sqrt{\mathbb{C}(y_1, y_1)}}, \frac{1}{\sqrt{\mathbb{C}(y_2, y_2)}}, \dots, \frac{1}{\sqrt{\mathbb{C}(y_n, y_n)}} \right) \quad (41)$$

Es ergibt sich dann

$$\begin{aligned} S_y^{-1} \mathbb{C}(y) S_y^{-1} &= \left(\frac{\mathbb{C}(y_i, y_j)}{\sqrt{\mathbb{C}(y_i, y_i)}} \right)_{1 \leq i, j \leq n} S_y^{-1} \\ &= \left(\frac{\mathbb{C}(y_i, y_j)}{\sqrt{\mathbb{C}(y_i, y_i)} \sqrt{\mathbb{C}(y_j, y_j)}} \right)_{1 \leq i, j \leq n} \\ &= \left(\rho(y_i, y_j) \right)_{1 \leq i, j \leq n} \\ &=: \mathbb{R}(y) \end{aligned} \quad (42)$$

Beweis

Analog ergibt sich

$$\begin{aligned} S_y \mathbb{R}(y) S_y &= \left(\frac{\sqrt{\mathbb{C}(y_i, y_i) \mathbb{C}(y_i, y_j)}}{\sqrt{\mathbb{C}(y_i, y_i)} \sqrt{\mathbb{C}(y_j, y_j)}} \right)_{1 \leq i, j \leq n} S_y \\ &= \left(\frac{\mathbb{C}(y_i, y_j) \sqrt{\mathbb{C}(y_j, y_j)}}{\sqrt{\mathbb{C}(y_j, y_j)}} \right)_{1 \leq i, j \leq n} \\ &= \left(\mathbb{C}(y_i, y_j) \right)_{1 \leq i, j \leq n} \\ &=: \mathbb{C}(y). \end{aligned} \tag{43}$$

Theorem (Korrelationsmatrix eines normalverteilten Zufallsvektors)

$y \sim N(\mu, \Sigma)$ sei ein multivariat normalverteilter Zufallsvektor mit Erwartungswertparameter $\mu \in \mathbb{R}^n$ und Kovarianzmatrixparameter

$$\Sigma := \left(\sigma_{ij}^2 \right)_{1 \leq i, j \leq n} = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \cdots & \sigma_{1n}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \cdots & \sigma_{2n}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1}^2 & \sigma_{n2}^2 & \cdots & \sigma_{nn}^2 \end{pmatrix} \in \mathbb{R}^{n \times n} \text{ p.d..} \quad (44)$$

Dann gilt

$$\mathbb{R}(y) = \left(\frac{\sigma_{ij}^2}{\sigma_{ii}\sigma_{jj}} \right)_{1 \leq i, j \leq n} = \begin{pmatrix} 1 & \frac{\sigma_{12}^2}{\sigma_{11}\sigma_{22}} & \cdots & \frac{\sigma_{1n}^2}{\sigma_{11}\sigma_{nn}} \\ \frac{\sigma_{21}^2}{\sigma_{22}\sigma_{11}} & 1 & \cdots & \frac{\sigma_{2n}^2}{\sigma_{22}\sigma_{nn}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\sigma_{n1}^2}{\sigma_{nn}\sigma_{11}} & \frac{\sigma_{n2}^2}{\sigma_{nn}\sigma_{22}} & \cdots & 1 \end{pmatrix}. \quad (45)$$

Bemerkungen

- Das Theorem folgt direkt mit dem Theorem zu Korrelationsmatrix und Kovarianzmatrix

Definition und Eigenschaften

Bemerkungen (fortgeführt)

- Umgekehrt gilt bei Definition eines Korrelationsmatrixparameters

$$R := (\rho_{ij})_{1 \leq i, j \leq n} = \begin{pmatrix} \rho_{11} & \rho_{12} & \cdots & \rho_{1n} \\ \rho_{21} & \rho_{22} & \cdots & \rho_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{n1} & \rho_{n2} & \cdots & \rho_{nn} \end{pmatrix} \quad (46)$$

mit

$$\rho_{ij} := 1 \text{ für } i = j \text{ und } \rho_{ij} = \rho_{ji} \in]-1, 1[\text{ für } i \neq j \text{ für } 1 \leq i, j \leq n \quad (47)$$

und zusätzlicher Definition von

$$S_y := \text{diag}(\sigma_{11}, \dots, \sigma_{nn}) \text{ mit } \sigma_{11}, \dots, \sigma_{nn} > 0, \quad (48)$$

dass

$$\Sigma := S_y R S_y = \begin{pmatrix} \sigma_{11}\sigma_{11} & \rho_{12}\sigma_{11}\sigma_{22} & \cdots & \rho_{1n}\sigma_{11}\sigma_{nn} \\ \rho_{21}\sigma_{22}\sigma_{11} & \sigma_{22}\sigma_{22} & \cdots & \rho_{2n}\sigma_{22}\sigma_{nn} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{n1}\sigma_{nn}\sigma_{11} & \rho_{n2}\sigma_{nn}\sigma_{22} & \cdots & \sigma_{nn}\sigma_{nn} \end{pmatrix} \quad (49)$$

einen Kovarianzmatrixparameter eines multivariaten normalverteilten Zufallsvektors definiert.

Definition (Regressionsgerade zweier Zufallsvariablen)

X und Y seien zwei Zufallsvariablen. Dann heißt

$$Y = \beta_0 + \beta_1 X \text{ mit} \quad (50)$$

mit

$$\beta_1 := \frac{\mathbb{C}(X, Y)}{\mathbb{V}(X)} \text{ und } \beta_0 := \mathbb{E}(Y) - \beta_1 \mathbb{E}(X) \quad (51)$$

die *Regressionsgerade der Zufallsvariablen X auf Y* , β_0 und β_1 heißen die zugehörigen *Regressionskoeffizienten*, und die Zufallsvariable

$$E := Y - \beta_0 - \beta_1 X \quad (52)$$

heißt die *Residualvariable*.

Bemerkungen

- X und Y sind Zufallsvariablen, β_0 und β_1 sind keine Zufallsvariablen.

Theorem (Optimalität der Regressionsgerade zweier Zufallsvariablen)

Unter allen Geraden der Form

$$Y = \beta_0 + \beta_1 X \quad (53)$$

ist die Gerade mit

$$\beta_1 := \frac{\mathbb{C}(X, Y)}{\mathbb{V}(X)} \text{ und } \beta_0 := \mathbb{E}(Y) - \beta_1 \mathbb{E}(X) \quad (54)$$

diejenige, für die

$$\tilde{q} : \mathbb{R}^2 \rightarrow \mathbb{R}, (\beta_0, \beta_1) \mapsto \tilde{q}(\beta_0, \beta_1) := \mathbb{E} \left((Y - (\beta_0 + \beta_1 X))^2 \right) \quad (55)$$

ein Minimum hat.

Korrelation und Regression

Beweis

Wir halten zunächst fest, dass

$$\begin{aligned}\tilde{q}(\beta_0, \beta_1) &= \mathbb{E}(Y - \beta_0 - \beta_1 X) \\ &= \mathbb{E}(Y - \beta_1 X - \beta_0 + \beta_1 \mathbb{E}(X) - \beta_1 \mathbb{E}(X) + \mathbb{E}(Y) - \mathbb{E}(Y)) \\ &= \mathbb{E}((Y - \mathbb{E}(Y)) - \beta_1 (X - \mathbb{E}(X)) + (\mathbb{E}(Y) - \beta_1 \mathbb{E}(X) - \beta_0))\end{aligned}\quad (56)$$

Ausmultiplizieren und Anwendung des Theorems zu den Eigenschaften des Erwartungswerts ergibt dann

$$\tilde{q}(\beta_0, \beta_1) = \mathbb{V}(Y) + \beta_1^2 \mathbb{V}(X) - 2\beta_1 \mathbb{C}(X, Y) + (\mathbb{E}(Y) - \beta_1 \mathbb{E}(X) - \beta_0)^2 \quad (57)$$

Berechnen der partiellen Ableitungen von \tilde{q} hinsichtlich von β_0 und β_1 ergibt dann

$$\frac{\partial}{\partial \beta_0} \tilde{q}(\beta_0, \beta_1) = -2(\mathbb{E}(Y) - \beta_1 \mathbb{E}(X) - \beta_0) \quad (58)$$

und

$$\frac{\partial}{\partial \beta_1} \tilde{q}(\beta_0, \beta_1) = 2\beta_1 \mathbb{V}(X) - 2\mathbb{C}(X, Y) - 2\mathbb{E}(X)(\mathbb{E}(Y) - \beta_1 \mathbb{E}(X) - \beta_0) \quad (59)$$

Nullsetzen von (58) ergibt dann als notwendige Bedingungen für ein Minimum von \tilde{q}

$$\begin{aligned}\frac{\partial}{\partial \beta_0} \tilde{q}(\beta_0^*, \beta_1^*) &= 0 \Leftrightarrow \mathbb{E}(Y) - \beta_1^* \mathbb{E}(X) - \beta_0^* = 0 \\ \frac{\partial}{\partial \beta_1} \tilde{q}(\beta_0^*, \beta_1^*) &= 0 \Leftrightarrow 2\beta_1^* \mathbb{V}(X) - 2\mathbb{C}(X, Y) - 2\mathbb{E}(X)(\mathbb{E}(Y) - \beta_1^* \mathbb{E}(X) - \beta_0^*) = 0\end{aligned}\quad (60)$$

Die erste Gleichung impliziert dann für die zweite Gleichung, dass

$$2\beta_1^* \mathbb{V}(X) - 2\mathbb{C}(X, Y) = 0 \Leftrightarrow \beta_1^* = \frac{\mathbb{C}(X, Y)}{\mathbb{V}(X)} \quad (61)$$

Einsetzen in die erste Gleichung ergibt dann

$$\mathbb{E}(Y) - \beta_1^* \mathbb{E}(X) - \beta_0^* = 0 \Leftrightarrow \beta_0^* = \mathbb{E}(Y) - \beta_1^* \mathbb{E}(X) \quad (62)$$

Theorem (Zusammenhang von Korrelation und Regression)

X und Y seien zwei Zufallsvariablen,

$$Y = \beta_0 + \beta_1 X \text{ mit } \beta_1 := \frac{\mathbb{C}(X, Y)}{\mathbb{V}(X)} \text{ und } \beta_0 := \mathbb{E}(Y) - \tilde{\beta}_1 \mathbb{E}(X) \quad (63)$$

sei die Regressionsgerade der Zufallsvariablen Y bezüglich der Zufallsvariablen X mit den Regressionskoeffizienten β_0 und β_1 und

$$X = \tilde{\beta}_0 + \tilde{\beta}_1 Y \text{ mit } \tilde{\beta}_1 := \frac{\mathbb{C}(X, Y)}{\mathbb{V}(Y)} \text{ und } \tilde{\beta}_0 := \mathbb{E}(X) - \tilde{\beta}_1 \mathbb{E}(Y) \quad (64)$$

sei die Regressionsgerade der Zufallsvariablen X bezüglich der Zufallsvariablen Y mit den Regressionskoeffizienten $\tilde{\beta}_0$ und $\tilde{\beta}_1$. Dann gilt

$$\beta_1 \tilde{\beta}_1 = \frac{\mathbb{C}(X, Y)}{\mathbb{V}(X)} \frac{\mathbb{C}(X, Y)}{\mathbb{V}(Y)} = \frac{\mathbb{C}(X, Y)^2}{\mathbb{V}(X)\mathbb{V}(Y)} = \rho(X, Y)^2. \quad (65)$$

Bemerkungen

- $\rho(X, Y)$ kann aus den Regressionskoeffizienten von X auf Y und von Y auf X errechnet werden.

Definition (Stichprobenregressionsgerade)

$(x, Y_1), \dots, (X_n, Y_n)$ sei eine Stichprobe von zweidimensionalen Zufallsvektoren mit identischen unabhängigen Verteilungen. Weiterhin sei für $i = 1, \dots, n$

$$Y_1 = \beta_0 + \beta_1 x \text{ mit } \beta_1 := \frac{\mathbb{C}(x, Y_1)}{\mathbb{V}(x)} \text{ und } \beta_0 := \beta_1 \mathbb{E}(x) + \mathbb{E}(Y_1) \quad (66)$$

die Regressionsgerade der Zufallsvariablen Y_1 bezüglich der Zufallsvariablen x mit den Regressionskoeffizienten β_0 und β_1 . Schließlich seien

- \bar{x} und \bar{y} die Stichprobenmittel von Realisierungen der Komponenten der Stichprobe,
- s_X^2 und s_Y^2 die Stichprobenvarianzen von Realisierungen der Komponenten der Stichprobe und
- $c_{X,Y}$ die Stichprobenkovarianz von Realisierungen der Stichprobe.

Dann heißt für $x \in \mathbb{R}$

$$y = b_0 + b_1 x \text{ mit } b_1 := \frac{c_{X,Y}}{s_X^2} \text{ und } b_0 := \bar{y} - b_1 \bar{x} \quad (67)$$

die Regressionsgerade der y -Werte bezüglich der x_i Werte in der Stichprobe.

Korrelation und Regression

Simulation einer Regressionsgerade

```
library(MASS)                                # multivariate Normalverteilungen

# Modellformulierung
n      = 1e2                                # Anzahl an Stichprobenvektoren
C_XY   = 1                                  # Kovarianz von X und Y
EX      = 2                                  # Erwartungswert von X
EY      = 1                                  # Erwartungswert von Y
VX      = 2                                  # Varianz von X
VY      = 2                                  # Varianz von Y
beta_1  = C_XY/VX                            # Regressionskoeffizient
beta_0  = -beta_1*EX + EY                    # Regressionskoeffizient

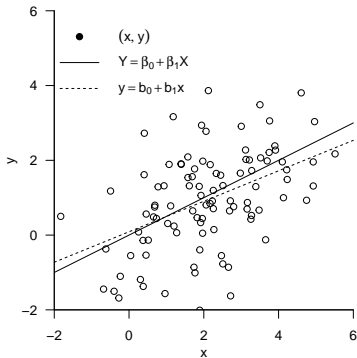
# Realisierungsgeneration
mu      = c(EX, EY)                          # Erwartungswertparameter
Sigma   = matrix(c(VX, C_XY, C_XY, VY), nrow = 2) # Kovarianzmatrixparameter
xy      = mvrnorm(n, mu, Sigma)

# Stichprobenstatistiken
x_bar   = mean(xy[,1])                       # Stichprobenmittel der x,...,x_n
y_bar   = mean(xy[,2])                       # Stichprobenmittel der y_1,...,y_n
s2X     = var(xy[,1])                         # Stichprobenvarianz der x,...,x_n
s2Y     = var(xy[,2])                         # Stichprobenvarianz der y_1,...,y_n
c_xy    = cov(xy[,1], xy[,2])                # Stichprobenkovarianz

# Stichprobenregressionsgeradenparameter
b_1     = c_xy/s2X                           # Stichprobenregressionskoeffizient
b_0     = -b_1*x_bar + y_bar                 # Stichprobenregressionskoeffizient

> beta_0 : 0
> beta_1 : 0.5
> b_0    : 0.136
> b_1    : 0.514
```

Simulation einer Regressionsgerade



Grundlagen

Korrelation und lineare Abhängigkeit

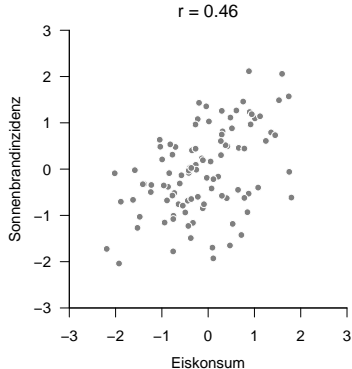
Korrelation und Regression

Korrelation und Bestimmtheitsmaß

Bedingte Korrelation und Partielle Korrelation

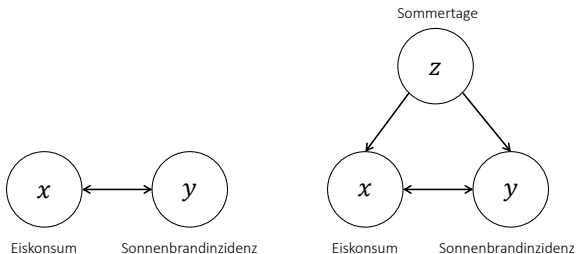
Selbstkontrollfragen

Jährlicher Eiskonsum und jährliche Sonnenbrandinzidenz



- Korrelation impliziert keine Kausalität.
- Kausalität wird zumeist als Koinzidenz mit zeitlicher Rangefolge modelliert.
- Einstiege in die kausale Inferenz geben z.B. Pearl (2000) und Imbens and Rubin (2015).

Jährlicher Eiskonsum und jährliche Sonnenbrandinzidenz



- Korrelation von Eiskonsum und Sonnenbrandinzidenz nach Korrektur für Sommertage?
- "Herausrechnen" des Einflusses von z auf die Kovariation von x und y ?

⇒ Bedingte Korrelation und Partielle Korrelation im Falle dreier Zufallsvariablen.

Definition (Bedingte Kovarianz und bedingte Korrelation)

Gegeben seien drei Zufallsvariablen x, y, z einer gemeinsamen Verteilung $\mathbb{P}_{x,y,z}(x, y, z)$. Weiterhin sei $\mathbb{P}_{x,y|z}(x, y)$ die bedingte Verteilung von x und y gegeben z . Dann heißt die Kovarianz von x und y in der Verteilung $\mathbb{P}_{x,y|z}(x, y)$ die *bedingte Kovarianz von x und y gegeben z* und wird mit $\mathbb{C}(x, y|z)$ bezeichnet. Weiterhin seien $\mathbb{P}_{x,y|z}(y)$ und $\mathbb{P}_{x,y|z}(x)$ die marginalen Verteilungen von x und y gegeben z , respektive, und $\mathbb{S}(x|z)$, $\mathbb{S}(y|z)$ die Standardabweichungen von x und y hinsichtlich $\mathbb{P}_{x,y|z}(y)$ und $\mathbb{P}_{x,y|z}(x)$, respektive. Dann heißt die Korrelation von x und y in der Verteilung $\mathbb{P}_{x,y|z}(x, y)$,

$$\rho(x, y|z) := \frac{\mathbb{C}(x, y|z)}{\mathbb{S}(y|z)\mathbb{S}(x|z)} \quad (68)$$

die *bedingte Korrelation von x und y gegeben z*

Bemerkungen

- Die bedingte Kovarianz zweier ZVen ist die Kovarianz zweier ZVen in einer bedingten Verteilung
- Die bedingte Korrelation zweier ZVen ist die Korrelation zweier ZVen in einer bedingten Verteilung.
- Durch Vertauschen der Variablennamen kann man analog $\rho(y, z|x)$ und $\rho(x, z|y)$ definieren.

Bedingte Korrelation und Partielle Korrelation

Beispiel

Die Zufallsvariablen x, y, z seien multivariat normalverteilt. Wir wollen die bedingte Korrelation von x und y gegeben z bestimmen. Für $v := (x, y, z)^T$ gelte also, dass

$$v \sim N(\mu, \Sigma) \quad (69)$$

mit

$$\mu := \begin{pmatrix} \mu_y \\ \mu_x \\ \mu_z \end{pmatrix} \text{ und } \Sigma := \begin{pmatrix} \sigma_x^2 & \sigma_{x,y}^2 & \sigma_{x,z}^2 \\ \sigma_{y,x}^2 & \sigma_y^2 & \sigma_{y,z}^2 \\ \sigma_{z,x}^2 & \sigma_{z,y}^2 & \sigma_z^2 \end{pmatrix} \quad (70)$$

Um die Kovarianzmatrix der bedingten Verteilung von x und y gegeben z zu bestimmen definieren wir zunächst

$$\Sigma_{x,y} := \begin{pmatrix} \sigma_x^2 & \sigma_{x,y}^2 \\ \sigma_{y,x}^2 & \sigma_y^2 \end{pmatrix}, \Sigma_z := (\sigma_z^2) \text{ und } \Sigma_{(x,y),z} := \Sigma_{z,(x,y)}^T := \begin{pmatrix} \sigma_{x,z}^2 \\ \sigma_{y,z}^2 \end{pmatrix}, \quad (71)$$

so dass

$$\Sigma = \begin{pmatrix} \Sigma_{x,y} & \Sigma_{(x,y),z} \\ \Sigma_{z,(x,y)} & \Sigma_z \end{pmatrix} \quad (72)$$

Mit dem Theorem zu bedingten Normalverteilungen (vgl. (4) Normalverteilungen) ist dann die Kovarianzmatrix des Zufallsvektors (x, y) gegeben durch

$$\Sigma_{x,y|z} = \Sigma_{x,y} - \Sigma_{(x,y),z} \Sigma_z^{-1} \Sigma_{z,(x,y)}. \quad (73)$$

Beispiel (fortgeführt)

Mit den Eigenschaften der multivariaten Normalverteilung gilt dann, dass die Diagonaleinträge von $\Sigma_{x,y|z}$ den bedingten Varianzen von x und y gegeben z entsprechen und dass der Nichtdiagonaleintrag die bedingte Kovarianz von x und y gegeben z ist. In anderen Worten gilt

$$\Sigma_{x,y|z} = \begin{pmatrix} \mathbb{C}(x, x|z) & \mathbb{C}(x, y|z) \\ \mathbb{C}(y, x|z) & \mathbb{C}(y, y|z) \end{pmatrix}. \quad (74)$$

Die bedingte Korrelation $\rho(x, y|z)$ von x und y gegeben z ergibt sich dann aus den Einträgen von $\Sigma_{x,y|z}$ gemäß

$$\rho(x, y|z) = \frac{\mathbb{C}(x, y|z)}{\sqrt{\mathbb{C}(x, x|z)} \sqrt{\mathbb{C}(y, y|z)}} \quad (75)$$

Für

$$\Sigma := \begin{pmatrix} 1.0 & 0.5 & 0.9 \\ 0.5 & 1.0 & 0.5 \\ 0.9 & 0.5 & 1.0 \end{pmatrix} \quad (76)$$

ergibt sich beispielsweise

$$\rho(x, y) = 0.50 \text{ und } \rho(x, y|z) \approx 0.13. \quad (77)$$

Beispiel (fortgeführt)

```
# Bedingte Korrelation bei Normalverteilung
S      = matrix(c( 1,.5,.9,           # \Sigma
                  .5, 1,.5,
                  .9,.5, 1), nrow = 3, byrow = TRUE)
rho_xy  = S[1,2]/(sqrt(S[1,1])*sqrt(S[2,2]))      # \rho(x,y)
S_xy_z  = S[1:2,1:2] - S[1:2,3] %*% solve(S[2,2]) %*% S[3,1:2] # \Sigma_{x,y|z}
rho_xy_z = S_xy_z[1,2]/(sqrt(S_xy_z[1,1])*sqrt(S_xy_z[2,2]))    # \rho(x,y|z)

# Ausgabe
cat("rho(x,y)   :", rho_xy,
    "\nrho(x,y|z) :", rho_xy_z)

> rho(x,y)   : 0.5
> rho(x,y|z) : 0.132
```

Definition (Partielle Korrelation)

x, y, z seien Zufallsvariablen mit linear-affinen Abhängigkeiten zwischen x und z sowie zwischen y und z ,

$$\begin{aligned}x &= \beta_0^{x,z} + \beta_1^{x,z} z \\ y &= \beta_0^{y,z} + \beta_1^{y,z} z\end{aligned}\tag{78}$$

mit Residualvariablen

$$\begin{aligned}e^{x,z} &= x - \beta_0^{x,z} - \beta_1^{x,z} z \\ e^{y,z} &= y - \beta_0^{y,z} - \beta_1^{y,z} z\end{aligned}\tag{79}$$

Dann ist die *partielle Korrelation von x und y mit auspartialisiertem z* definiert als

$$\rho(x, y \setminus z) := \rho(e^{x,z}, e^{y,z}).\tag{80}$$

Bemerkungen

- $e^{x,z}$ ist die Zufallsvariable x , aus der der Einfluss von z "herausgerechnet" wurde.
- $e^{y,z}$ ist die Zufallsvariable y , aus der der Einfluss von z "herausgerechnet" wurde.
- $\rho(x, y \setminus z)$ ist also die Korrelation von x und y , aus denen jeweils der Einfluss von z "herausgerechnet" wurde

Definition (Partielle Stichprobenkorrelation)

x, y, z seien Zufallsvariablen mit linear-affinen Abhängigkeiten zwischen y und z sowie zwischen x und z wie in der Definition der partiellen Korrelation. Weiterhin seien

- $\{(x_i, y_i, z_i)\}_{i=1, \dots, n}$ eine Menge von Realisierungen des Zufallsvektors $(x, y, z)^T$,
- $\hat{\beta}_0^{x,z}, \hat{\beta}_1^{x,z}$ die Ausgleichsgeradenparameter für $\{(x_i, z_i)\}_{i=1, \dots, n}$,
- $\hat{\beta}_0^{y,z}, \hat{\beta}_1^{y,z}$ die Ausgleichsgeradenparameter für $\{(y_i, z_i)\}_{i=1, \dots, n}$.

Schließlich seien für $i = 1, \dots, n$

- $e_i^{x,z} := x_i - \hat{\beta}_0^{x,z} + \hat{\beta}_1^{x,z} z_i$
- $e_i^{y,z} := y_i - \hat{\beta}_0^{y,z} + \hat{\beta}_1^{y,z} z_i$

die Residualwerte der jeweiligen Ausgleichsgeraden. Dann heißt die Stichprobenkorrelation der Wertemenge $\{(e_i^{y,z}, e_i^{x,z})\}_{i=1, \dots, n}$ *partielle Stichprobenkorrelation der x_i und y_i mit auspartialisierten z_i* .

Bemerkungen

- Die partielle Stichprobenkorrelation wird als Schätzer der partiellen Korrelation genutzt.

Theorem (Bedingte und Partielle Korrelation bei Normalverteilung)

x, y, z seien drei gemeinsam multivariat normalverteilte Zufallsvariablen. Dann gilt

$$\rho(x, y|z) = \rho(x, y \setminus z) \quad (81)$$

Bemerkungen

- Wir verzichten auf einen Beweis.
- Generell sind bedingte und partielle Korrelationen nicht identisch.
- Für Details, siehe zum Beispiel Lawrance (1976) und Baba, Shibata, and Sibuya (2004).

Theorem (Bedingte Korrelation und Korrelationen bei Normalverteilung)

x, y, z seien drei gemeinsam multivariat normalverteilte Zufallsvariablen. Dann gilt

$$\rho(x, y|z) = \frac{\rho(x, y) - \rho(x, z)\rho(y, z)}{\sqrt{(1 - \rho(x, z)^2)}\sqrt{(1 - \rho(y, z)^2)}} \quad (82)$$

Bemerkungen

- $\rho(x, y|z)$ kann bei Normalverteilung aus den Korrelationen $\rho(x, y)$, $\rho(x, z)$, $\rho(y, z)$ berechnet werden.
- Ein entsprechender Schätzer für $\rho(x, y|z)$ ergibt sich mit den Stichprobenkorrelationen $r_{x,y}$, $r_{x,z}$, $r_{y,z}$ als

$$r_{x,y|z} = \frac{r_{x,y} - r_{x,z}r_{y,z}}{\sqrt{(1 - r_{x,z}^2)}\sqrt{(1 - r_{y,z}^2)}} \quad (83)$$

- Mit $\rho(x, y|z) = \rho(x, y \setminus z)$ bei Normalverteilung die Formel auch für $\rho(x, y \setminus z)$.

Bedingte Korrelation und Partielle Korrelation

Beweis

Ohne Beschränkung der Allgemeinheit betrachten wir den Fall eines standardisierten multivariaten normalverteilten Zufallsvektors $v := (x, y, z)^T$ mit Kovarianzmatrixparameter

$$\Sigma := \begin{pmatrix} 1 & \rho(x, y) & \rho(x, z) \\ \rho(y, x) & 1 & \rho(y, z) \\ \rho(z, x) & \rho(z, y) & 1 \end{pmatrix}. \quad (84)$$

Wir definieren nun zunächst

$$\Sigma_{x,y} := \begin{pmatrix} 1 & \rho(x, y) \\ \rho(y, x) & 1 \end{pmatrix}, \Sigma_z := (1) \text{ und } \Sigma_{(x,y),z} := \Sigma_{z,(x,y)}^T := \begin{pmatrix} \rho(x, z) \\ \rho(y, z) \end{pmatrix}, \quad (85)$$

so dass

$$\Sigma = \begin{pmatrix} \Sigma_{x,y} & \Sigma_{(x,y),z} \\ \Sigma_{z,(x,y)} & \Sigma_z \end{pmatrix}. \quad (86)$$

Mit dem Theorem zu bedingten Normalverteilungen (vgl. (4) Normalverteilungen) ist dann die Kovarianzmatrix des Zufallsvektors (x, y) gegeben durch

$$\Sigma_{x,y|z} = \Sigma_{x,y} - \Sigma_{(x,y),z} \Sigma_z^{-1} \Sigma_{z,(x,y)}. \quad (87)$$

Beweis (fortgeführt)

Es ergibt sich also

$$\begin{aligned} \begin{pmatrix} \sigma_{x,x|z}^2 & \sigma_{x,y|z}^2 \\ \sigma_{y,x|z}^2 & \sigma_{y,y|z}^2 \end{pmatrix} &= \begin{pmatrix} 1 & \rho(x,y) \\ \rho(y,x) & 1 \end{pmatrix} - \begin{pmatrix} \rho(x,z) \\ \rho(y,z) \end{pmatrix} \begin{pmatrix} 1 \end{pmatrix}^{-1} \begin{pmatrix} \rho(x,z) & \rho(y,z) \end{pmatrix} \\ &= \begin{pmatrix} 1 & \rho(x,y) \\ \rho(y,x) & 1 \end{pmatrix} - \begin{pmatrix} \rho(x,z)\rho(x,z) & \rho(x,z)\rho(y,z) \\ \rho(y,z)\rho(x,z) & \rho(y,z)\rho(y,z) \end{pmatrix} \\ &= \begin{pmatrix} 1 - \rho(x,z)^2 & \rho(x,y) - \rho(x,z)\rho(y,z) \\ \rho(y,x) - \rho(y,z)\rho(x,z) & 1 - \rho(y,z)^2 \end{pmatrix}. \end{aligned} \quad (88)$$

Es ergibt sich also

$$\rho(x,y|z) = \frac{\sigma_{x,y|z}^2}{\sqrt{\sigma_{x,x|z}^2} \sqrt{\sigma_{y,y|z}^2}} = \frac{\rho(x,y) - \rho(x,z)\rho(y,z)}{\sqrt{1 - \rho(x,z)^2} \sqrt{1 - \rho(y,z)^2}}. \quad (89)$$

□

Bedingte Korrelation und Partielle Korrelation

Beispiel

```
# Modellformulierung und Datenrealisierung
library(MASS)
set.seed(1)
S = matrix(c( 1,.5,.9,
             .5, 1,.5,
             .9,.5, 1),nrow=3,byrow=TRUE)

n = 1e6
xyz = mvrnorm(n,rep(0,3),S)

# Multivariate Normalverteilung
# reproduzierbare Daten
# Kovarianzmatrixparameter \Sigma

# Anzahl Realisierungen von  $v := (x,y,z)^T$ 
# Realisierungen von  $v := (x,y,z)^T$ 

# Partielle Stichprobenkorrelation als Residualstichprobenkorrelation
bars = apply(xyz, 2, mean) # Stichprobenmittel
s = apply(xyz, 2, sd) # Stichprobenstandardabweichungen
c = cov(xyz) # Stichprobenkovarianzen
b_xz1 = c[1,3]/c[3,3] #  $\beta_1(x,z)$ 
b_xz0 = bars[1] - b_xz1*bars[3] #  $\beta_0(x,z)$ 
b_yz1 = c[2,3]/c[3,3] #  $\beta_1(y,z)$ 
b_yz0 = bars[2] - b_yz1*bars[3] #  $\beta_0(y,z)$ 
e_xz = xyz[,1] - b_xz1*xyz[,3] - b_xz0 # Residualwerte  $e^{\{x,z\}}$ 
e_yz = xyz[,2] - b_yz1*xyz[,3] - b_yz0 # Residualwerte  $e^{\{y,z\}}$ 
pr_e = cor(e_xz,e_yz) #  $\rho(x,y|z)$ 

# Partielle Stichprobenkorrelation aus Stichprobenkorrelationen
r = cor(xyz) # Stichprobenkorrelationsmatrix
pr_r_n = r[1,2]-r[1,3]*r[2,3] #  $\rho(x,y|z)$  Formel Zähler
pr_r_d = sqrt((1-r[1,3]^2)*(1-r[2,3]^2)) #  $\rho(x,y|z)$  Formel Nenner
pr_r = pr_r_n/pr_r_d #  $\rho(x,y|z)$ 

# partielle Stichprobenkorrelation aus Toolbox
library(ppcor) # Laden der Toolbox

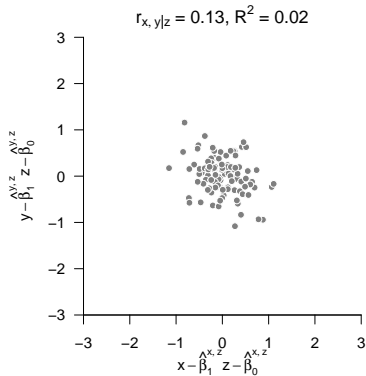
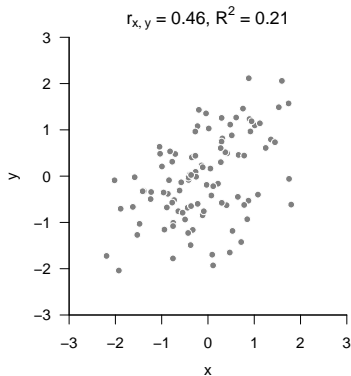
> Warning: Paket 'ppcor' wurde unter R Version 4.1.3 erstellt

pr_t = pcor(xyz) #  $\rho(x,y|z), \rho(x,z|y), \rho(y,z|x)$ 

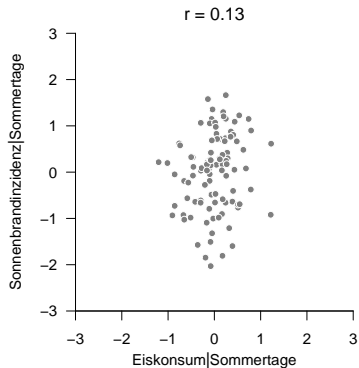
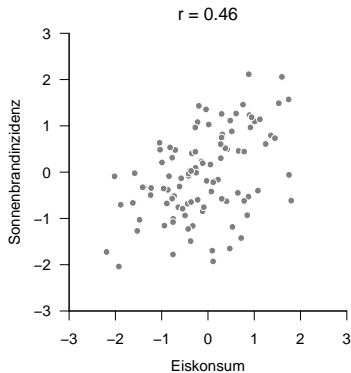
# Ausgabe
cat("r(x,y) :", r[1,2],
    "\nr(x,y/z) aus Residuenkorrelation :", pr_e,
    "\nr(x,y/z) aus Korrelationen :", pr_r,
    "\nr(x,y/z) aus Toolbox :", pr_t$estimate[1,2])

> r(x,y) : 0.5
> r(x,y/z) aus Residuenkorrelation : 0.133
> r(x,y/z) aus Korrelationen : 0.133
> r(x,y/z) aus Toolbox : 0.133
```

Bedingte Korrelation und Partielle Korrelation



Bedingte Korrelation und Partielle Korrelation



Grundlagen

Korrelation und Bestimmtheitsmaß

Korrelation und lineare Abhängigkeit

Korrelation und Regression

Bedingte Korrelation und Partielle Korrelation

Selbstkontrollfragen

Selbstkontrollfragen

1. Geben Sie die Definition der Korrelation zweier Zufallsvariablen wieder.
2. Geben Sie die Definitionen von Stichprobenmittel, -standardabweichung, -kovarianz und -korrelation wieder.
3. Erläutern Sie anhand der Mechanik der Kovariationsterme, wann eine Stichprobenkorrelation einen hohen absoluten Wert annimmt, einen hohen positiven Wert annimmt, einen hohen negativen Wert annimmt und einen niedrigen Wert annimmt.
4. Berechnen Sie die Korrelation von Anzahl der Therapiestunden und Symptomreduktion anhand der Daten in Beispieldatensatz.csv.
5. Geben Sie das Theorem zur Stichprobenkorrelation bei linear-affinen Transformationen wieder.
6. Erläutern Sie das Theorem zur Stichprobenkorrelation bei linear-affinen Transformationen.
7. Geben Sie die Definitionen von erklärten Werten und Residuen einer Ausgleichsgerade wieder.
8. Geben Sie das Theorem zur Quadratsummenzerlegung bei einer Ausgleichsgerade wieder.
9. Erläutern Sie die intuitiven Bedeutungen von SQT , SQE und SQR .
10. Geben Sie die Definition des Bestimmtheitsmaßes R^2 wieder.
11. Geben Sie das Theorem zum Zusammenhang von Stichprobenkorrelation und Bestimmtheitsmaß wieder.
12. Erläutern Sie die Bedeutung von hohen und niedrigen R^2 Werten im Lichte der Ausgleichsgerade.
13. Berechnen Sie in einem R-Skript R^2 für die Daten in der Datei Beispieldatensatz.csv anhand der Definition von R^2 . Überprüfen Sie Ihr Ergebnis anhand des Theorems zum Zusammenhang von Stichprobenkorrelation und Bestimmtheitsmaß.
14. Geben Sie das Theorem zum Zusammenhang von Korrelation und linear-affiner Abhängigkeit wieder.
15. Geben Sie die Definition der Regressionsgerade zweier Zufallsvariablen wieder.
16. Geben Sie das Theorem zur Optimalität der Regressionsgerade zweier Zufallsvariablen wieder.
17. Geben Sie das Theorem zum Zusammenhang von Korrelation und Regression an.
18. Erläutern Sie, wie aus den Ergebnissen einer Regressionsanalyse das Ergebnis einer Korrelationsanalyse errechnet werden kann.

References

- Baba, Kunihiro, Ritei Shibata, and Masaaki Sibuya. 2004. "Partial Correlation and Conditional Correlation as Measures of Conditional Independence." *Australian & New Zealand Journal of Statistics* 46 (4): 657–64. <https://doi.org/10.1111/j.1467-842X.2004.00360.x>.
- Imbens, Guido, and Donald B. Rubin. 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Academic Press.
- Lawrance, A. J. 1976. "On Conditional and Partial Correlation." *The American Statistician* 30 (3): 146. <https://doi.org/10.2307/2683864>.
- Pearl, Judea. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge, U.K. ; New York: Cambridge University Press.