

# Blatt 11

Für eine lineare Regression ( $\hat{y} = ax + b$ ) ist der Pearsonsche

Korrelationskoeffizient definiert als  $\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$ . Das

Bestimmtheitsmaß ist definiert als  $R^2 = \rho^2 = \frac{\text{Cov}(X, Y)^2}{\text{Var}(X) \text{Var}(Y)}$

Alternativ kann man mithilfe der Variablen

$$SQE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n ((Ax_i + B) - \bar{y})^2 \rightarrow \text{erklärte Streuung}$$

$$SQR = SSR = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (Ax_i + B))^2 \rightarrow \text{residuale Streuung}$$

je nach Literatur  
versch. Bezeichnungen

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

rechnen:

$$R^2 = \frac{SQE}{SQT} \Rightarrow \text{Das Bestimmtheitsmaß ist also der Anteil}$$

der durch die Regression beschriebenen Streuung an der

Gesamtstreuung  $\Rightarrow 0 \leq R^2 \leq 1$ ; je größer  $R^2$  desto besser

erklärt die Regressionsgerade den Zusammenhang von  $X$  und  $Y$

Um das Konfidenzintervall von  $\rho$  um einen Schätzer  $r$

zu berechnen, muss man die Fisher-Transformation  $z = \frac{1}{2} \ln\left(\frac{1+r}{1-r}\right)$

nutzen, da bekannt ist, dass  $z \sim \mathcal{N}\left(\frac{1}{\sqrt{n-3}}, \frac{1}{2} \ln\left(\frac{1+\rho}{1-\rho}\right)\right)$ .

Es müssen also folgende Schritte durchgeführt werden

a) Transformation:  $z = \frac{1}{2} \ln\left(\frac{1+r}{1-r}\right)$

b) Berechnung  $(1-\alpha)$ -Konfidenzintervall (in Abhängigkeit

Von  $z$ ):  $\left[ \underbrace{z - \frac{z_{1-\alpha/2}}{\sqrt{n-3}}}_{z_1}; \underbrace{z + \frac{z_{1-\alpha/2}}{\sqrt{n-3}}}_{z_2} \right]$   $z$ -Wert der Standardnormalverteilung

c) Retransformation:  $r_1 = \frac{e^{2z_1} - 1}{e^{2z_1} + 1}$ ,  $r_2 = \frac{e^{2z_2} - 1}{e^{2z_2} + 1}$

⇒ Das Konfidenzintervall von  $\rho$  zum Niveau  $1-\alpha$  ist dann  $[r_1; r_2]$

Achtung: das Konfidenzintervall von  $\rho$  liegt normalerweise nicht symmetrisch bzgl dem Schätzer  $r$   
[d.h.  $r-r_1 \neq r_2-r$ ]

• Test:  $H_0: \rho=0$   $H_1: \rho \neq 0$

Teststatistik:  $T = \frac{r \sqrt{n-2}}{\sqrt{1-r^2}}$

$H_0$  ablehnen, falls  $|T| > t_{1-\alpha/2}(n-2)$

• F-Test Testet, ob das Regressionsmodell einen Erklärungswert für  $Y$  hat

Teststatistik:  $T = \frac{R^2}{1-R^2} \frac{n-p-1}{p}$

$p =$  'Anzahl Parameter'  
[d.h. bei  $Y = AX + B \Rightarrow 1$   
bei  $Y = A_1 X_1 + A_2 X_2 + B \Rightarrow 2$   
...]

Es hat Erklärungswert, falls  $T > F_{1-\alpha}(p, n-p-1)$   
F-Wert der Fisher-Verteilung

•  $(1-\alpha)$ -Konfidenzintervall für  $Y$  in  $x_0$

Im Gegensatz zum Prognoseintervall, ist das Konfidenzintervall

$$\left[ (Ax_0 + B) \pm \left( \sqrt{s^2 \cdot \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\text{Var}(X) \cdot n} \right)} + t_{n-2}(1-\alpha/2) \right) \right]$$

In Sage werden die meisten Regressionsberechnungen über  $R$  getätigt, daher ist es am einfachsten mittels '%r' auf den  $R$ -Interpreter zu wechseln und dort

mittels der  $\text{lm}()$ -Funktion die notwendigen  
Werte berechnen zu lassen (z.B. die Schätzer  $A, B$   
der Regression  $Y = AX + B$  mit `'coef(lm(y~x))'`)