

Lineares Modell: Was ändert sich durch Transformation?

Lineare Transf. von Zielvariable und Einflussvariable im einfachen Modell

$$x_i \rightarrow t_i = a_0 + a_1 x_i \quad (a_1 \neq 0) \quad u_i \rightarrow u_i = b_0 + b_1 y_i \quad (b_1 \neq 0)$$

Voraussetzung: $\bar{x}_i, \log(x_i), x_i^2$ etc. sind alles KEINE linearen Transf.

KQ-Schätzer für (linear) transformiertes Modell $y_i = \beta_0 + \beta_1 t_i + \varepsilon_i$

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (t_i - \bar{t})(u_i - \bar{u})}{\frac{1}{n} \sum_{i=1}^n (t_i - \bar{t})^2} = \dots \quad t_i - \bar{t} = a_0 + a_1 x_i - a_0 - a_1 \bar{x} \quad \text{Faktoren rausziehen} = \frac{a_1 b_1}{a_1^2} - \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(u_i - \bar{u})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{b_1}{a_1} \hat{B}_1$$

$$\hat{\beta}_0 = \bar{u} - \hat{\beta}_1 \bar{t} = b_0 + b_1 \bar{y} - \frac{b_1}{a_1} \hat{B}_1 \quad (a_0 + a_1 \bar{x}) = b_0 + b_1 \bar{y} - b_1 \hat{\beta}_1 \bar{x} - \frac{b_1}{a_1} \hat{B}_1 a_0 = b_0 + b_1 (\bar{y} - \hat{\beta}_1 \bar{x}) - \frac{b_1}{a_1} a_0 \hat{B}_1$$

• Parameter-Schätzer und deren Standardfehler ändern sich
Deren Konfidenz (t-Wert) bleibt hingegen identisch Konf.-Intervall nur absolute Unterschiede, in Relationen gleichbleibend

• R^2 ändert sich logischerweise nicht

Beispiel: Logarithmierung von X ($a_0 = -\bar{x}, a_1 = 1$), Y bleibt gleich ($b_0 = 0, b_1 = 1$): $\hat{\beta}_1 = \hat{\beta}_1$, $\hat{\beta}_0 = \hat{\beta}_0 + \frac{1}{n} \bar{x} \hat{\beta}_1$

$$x_i^* = x_i - \bar{x} \Rightarrow \bar{x}^* = 0$$

iii) Standardisierung, also Zentrierung u. Teilen durch Standardabweichung, von beiden Variablen sodass ihre $MW = 0$ und $SD = 1$
 $\tilde{x}_i = \frac{x_i - \bar{x}}{s_x}$ mit $\text{var}(x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ analog $\tilde{u}_i = \frac{u_i - \bar{u}}{s_u}$

$$\tilde{\beta}_1 = \text{corr}(x_i, u_i) \frac{s_u}{s_x} = \text{corr}(x_i, u_i) \quad , \quad \tilde{\beta}_0 = \bar{u} - \tilde{\beta}_1 \bar{x} = 0$$

Standardisierung ist keine lineare Transformation mehr
Teilen durch konstante var(), die von gesamten Datensätzen abhängt

Wozu (non-)lineare Transformationen von Regressoren und/oder Zielvariablie?

Extreme Verteilungen (z.B. weite Range mit vielen Lücken, irreflektivs, superschief) eignen sich nicht für lineare Regression

==> Transformation, klassisch dämpfender Effekt von \log, \sqrt aber auch $X/10000$ kann sinnvoll sein

$\log 10$ bewirkt wegen entfarter Interpretation, $\log 10(x+1)$ bei kleinen Werten, die Wert 0 erhalten

Vorsicht: Interpretation der Parameter von Transformationen, Modell ändert sich!

Beispiel Y und X mit $\log 10$ transformiert: $\log 10(u_i) = \beta_0 + \beta_1 \log 10(x_i) + \varepsilon_i$

β_0 ↓ erwartete Logarithmiepie Y wenn $X=0$ nicht sinnvoll interpretierbar

β_1 ↓ steigt $\log X$ um eine Einheit, so erhöhen sich $\log Y$ entsprechend durchschnittlich um $\beta_1^{(1)}$ bzw. $\log(X+1)$ um $\beta_1^{(2)}$

bzw. steigt $\log X$ um eine Einheit erhöhen sich $E(u_i)$ c.p. um 10^{β_1} bzw. $\log(X+1)$ um 10^{β_1}

bzw. "steigt" X um Faktor a , so steigt $E(u_i)$ c.p. um $a^{\beta_1} = 10^{\beta_1 \cdot \log(a)}$ bzw. $10^{\beta_1 \cdot \log(a+1)}$

Frage] Besteht zum Sigmoidkoeffizienten $\alpha = 0.01$ ein linearer Zusammenhang zwischen Bevölkerungsdichte und Fläche, also ungarthimiert?

$$\log 10(\text{Populi}) = \beta_0 + \beta_1 \log 10(\text{Areal}) \Leftrightarrow \text{Populi} = 10^{\beta_0 + \log 10(\text{Areal})} = 10^{\beta_0} \cdot \text{Areal} \cdot \beta_1$$

Nun auf Bevölkerungsdichte umrechnen: Bevölkerungsdichte $\hat{c} = \frac{\text{Populi}}{\text{Areal}} = \frac{10^{\beta_0} \cdot \text{Areal} \cdot \beta_1}{\text{Areal}} = 10^{\beta_0} \cdot \text{Areal} \cdot \beta_1^{-1} \Rightarrow$ Ja, Fläche hat sign. Einfluss auf Bev. dichte, wenn $\beta_1 < 1$ ist $\beta_1^{-1} = 1$ also ohne Einfluss als Faktor

Was man nicht durch Transformationen fixen kann: kategoriale Variablen "themen" Datensätze in schwer vergleichbare Teile, tauchen aber im Modell nicht auf

Fix i) Eigens Modell für jede Kategorie Fix ii) Kategoriale Var. als Dummy ins Modell aufnehmen

bzw. ein Modell für häufigste Kategorie

Fix iii) Extreme Auslese streichen, wenn

longitudinal sinnvoll (z.B. exstatische Länder)