

# Kapitel 1 - Das einfache lineare Regressionsmodell

## Einfaches lineares Regressionsmodell

Das einfache lineare Regressionsmodell hat die Form

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

für ein festes numerisches  $x_i$  und  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . Beachte, dass per Definition gilt  $Y_i | x_i \sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$

## Kleinste Quadrate (KQ) Schätzer

Wir schätzen die Parameter  $(\beta_0, \beta_1)$  durch

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(\beta_0, \beta_1)} \sum_{i=1}^n \underbrace{(Y_i - (\beta_0 + \beta_1 x_i))^2}_{\varepsilon_i^2} \quad (1)$$

und nennen  $(\hat{\beta}_0, \hat{\beta}_1)$  den KQ-Schätzer von  $(\beta_0, \beta_1)$  und  $\hat{\varepsilon}_i := Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$  die Residuen.

## Existenz und Berechnung vom KQ Schätzer

Der KQ-Schätzer existiert und ist eindeutig, falls  $\sum_{i=1}^n (x_i - \bar{x})^2 \neq 0$ . Dieser lässt sich berechnen als

$$\begin{aligned} \hat{\beta}_1 &= \frac{S_{xy}}{S_x^2} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \hat{\varepsilon}_i \sim \mathcal{N}(0, \sigma^2) \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{x}. \end{aligned}$$

Durch Differenzieren von der Gleichung (1) erhält man  $(\hat{\beta}_0, \hat{\beta}_1)$  als Lösung der Normalengleichungen

$$\begin{aligned} \sum_{i=1}^n \hat{\varepsilon}_i &= 0 \\ \sum_{i=1}^n \hat{\varepsilon}_i x_i &= 0 \quad \left( \begin{array}{c|c} x_1 & \hat{\varepsilon}_1 \\ x_2 & \hat{\varepsilon}_2 \\ \vdots & \vdots \\ x_n & \hat{\varepsilon}_n \end{array} \right) \perp \left( \begin{array}{c|c} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{array} \right) \end{aligned}$$

## Interpretation der Modellparameter

Für  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i = 1, \dots, n$  mit  $E(Y_i | x_i) = \beta_0 + \beta_1 x_i$  gilt,

- wenn  $x$  um eine Einheit steigt, dann steigt  $Y$  im Erwartungswert um  $\beta_1$  Einheiten.
- Es gilt  $\beta_0 = E(Y | X = 0)$ .
- Der Parameter  $\sigma$  die erwartete Abweichung der

## Eigenschaften des KQ-Schätzers

Gegeben dem einfachen linearen Modell, gilt für den KQ-Schätzer  $(\hat{\beta}_0, \hat{\beta}_1)$

- Erwartungstreue:  $E(\hat{\beta}_0, \hat{\beta}_1) = (\beta_0, \beta_1)$ .
- $V(\hat{\beta}_1) = \frac{\sigma^2}{n S_x^2}$  und  $V(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{n S_x^2} \right)$ .
- $(\hat{\beta}_0, \hat{\beta}_1)$  ist der maximum-likelihood Schätzer.

## Schätzer für $\sigma^2$

Gegeben dem einfachen linearen Modell mit  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , gilt

$$\sigma^2 = \frac{RSS}{n-2}$$

$$\hat{\sigma}^2 := \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad \left( Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2$$

ist ein erwartungstreuer Schätzer von  $\sigma^2$  und

$$\frac{RSS}{\sigma^2} = \frac{n-2}{\sigma^2} \hat{\sigma}^2 \sim \chi^2_{n-2}.$$

Der KQ-Schätzer  $(\hat{\beta}_0, \hat{\beta}_1)$  und der Schätzer  $\hat{\sigma}^2$  sind stoch.unabhängig.

## Konfidenzintervalle für $\beta_0$ und $\beta_1$

Gegeben dem einfachen linearen Modell mit  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , gilt für  $\hat{\beta}_1$  und  $\hat{\beta}_0$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

$$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}_{\hat{\beta}_1}} \sim t_{n-2} \text{ mit } \hat{\sigma}_{\hat{\beta}_1} := \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad \begin{matrix} \text{FG} \\ \text{verbraucht für} \\ \text{Schätzung von } \beta_1, \beta_0 \end{matrix}$$

$$\frac{\hat{\beta}_0 - \beta_0}{\hat{\sigma}_{\hat{\beta}_0}} \sim t_{n-2} \text{ mit } \hat{\sigma}_{\hat{\beta}_0} := \sqrt{\hat{\sigma}^2 \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}}$$

Damit können wir Konfidenzintervalle zum Niveau  $1 - \alpha$  für  $\beta_1$  und  $\beta_0$  erzeugen:

$$[\hat{\beta}_1 - \hat{\sigma}_{\hat{\beta}_1} t_{1-\alpha/2}(n-2); \hat{\beta}_1 + \hat{\sigma}_{\hat{\beta}_1} t_{1-\alpha/2}(n-2)]$$

$$[\hat{\beta}_0 - \hat{\sigma}_{\hat{\beta}_0} t_{1-\alpha/2}(n-2); \hat{\beta}_0 + \hat{\sigma}_{\hat{\beta}_0} t_{1-\alpha/2}(n-2)]$$

Variante von Prognosefehler =  $\text{Var}((\hat{\beta}_0 + \hat{\beta}_1 x_{n+1}) - (\beta_0 + \beta_1 x_{n+1})) + \text{Var}(\varepsilon_{n+1})$

$$= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_{n+1}) + \sigma^2 = \text{Var}(\hat{\beta}_0) + \text{Var}(\hat{\beta}_1 x_{n+1}) + \text{Var}(\varepsilon_{n+1})$$

$$= \text{Var}(\hat{\beta}_0) + \text{Var}(\hat{\beta}_1 x_{n+1}) + \text{Var}(\varepsilon_{n+1})$$

## Quadratsummenzerlegung

Gegeben sei ein einfaches lineares Modell mit  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  und  $\hat{Y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i$ . Dann gilt

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{\text{SST}} + \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{\text{SSE}} + \underbrace{\sum_{i=1}^n (\hat{Y}_i - Y_i)^2}_{\text{SSM}}$$

SST(otal):

Gesamtstreuung von  $Y$

SSError):

Streuung der Residuen

SSModel):

Streuung, die das Modell erklärt

## Bestimmtheitsmaß

Unter Verwendung der obigen Notation definieren wir das Bestimmtheitsmaß als

$$R^2 = \frac{\text{SSM}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}.$$

Es gilt

$$R^2 = r_{xy}^2 = \frac{S_{xy}}{S_x S_y},$$

wobei  $r_{xy}$  der Bravais-Pearson Korrel.koeffizient ist.

## Interpretation von $R^2$

- $R^2$  beschreibt den Anteil der Varianz von  $Y$ , die durch  $x$  erklärt wird.
- $R$  ist invariant gegenüber linearen Transformationen von  $x$  und  $Y$  (weil  $r_{xy}$  invariant)
- $R$  ist symmetrisch bzgl.  $x$  und  $Y$ .
- !  $R^2$  hängt auch von der Streuung von  $x$  in der Stichprobe ab.

## Prognosewert

Gegeben sei ein einfaches lineares Modell mit  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  und  $\hat{Y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i, i = 1, \dots, n$ . Sei nun eine weitere Beobachtung  $x_{n+1}$  mit zugehörigem  $Y_{n+1} = \beta_0 + \beta_1 x_{n+1} + \varepsilon_{n+1}$  gegeben. Der Prognosewert von  $Y_{n+1}$  ist definiert als  $\hat{Y}_{n+1} = \hat{\beta}_0 + \hat{\beta}_1 x_{n+1}$

$$\text{Prognosefehler: } \hat{\varepsilon}_{n+1} = \hat{Y}_{n+1} - Y_{n+1} = [\hat{\beta}_0 + \hat{\beta}_1 x_{n+1} - (\beta_0 + \beta_1 x_{n+1})] - \varepsilon_{n+1}$$

#### Prognosefehler

Gegeben sei ein einfaches linearen Modell, sowie eine weitere Beobachtung  $x_{n+1}$  mit zugehörigem  $Y_{n+1}$  sowie der Prognosewert  $\hat{Y}_{n+1}$ . Dann gilt

$$\begin{aligned} \text{sowie der Prognosewert } \hat{Y}_{n+1}. \text{ Dann gilt } \\ E(\hat{Y}_{n+1} - Y_{n+1}) = 0 \quad E(B_0 + B_1 X_{n+1} - B_0 - B_1 X_{n+1} - e_{n+1}) = 0 \\ = V(Y_{n+1}) + V(e_{n+1}) \\ V(\hat{Y}_{n+1} - Y_{n+1}) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ \text{Einschr. von Schätzung } (B_0, B_1) \quad V(Y_{n+1}) = \sigma^2 \left[ 1 + (x_{n+1} - \bar{x})^2 \right], \quad V(e_{n+1}) = \sigma^2 \end{aligned}$$

### Prognoseintervall

Gegeben sei ein einfaches linearen Modell, sowie eine weitere Beobachtung  $x_{n+1}$  mit zugehörigem  $Y_{n+1}$  sowie der Prognosewert  $\hat{Y}_{n+1}$ . Dann können wir für  $Y_{n+1}$  ein Konfidenzintervall zum Niveau  $1 - \alpha$  konstruieren:

$$[\hat{Y}_{n+1} - \hat{\sigma}_{\hat{Y}_{n+1}} t_{1-\alpha/2}(n-2); \hat{Y}_{n+1} + \hat{\sigma}_{\hat{Y}_{n+1}} t_{1-\alpha/2}(n-2)]$$

mit

$$\text{eigtl. } \hat{\sigma}_{\text{pred}, n+1} = \sqrt{\text{Var}(Y_{n+1} - \bar{y}_{n+1})}$$

$$\hat{\sigma}_{Y_{n+1}} = \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

unbekanntes  $\sigma^2$  schätzen

## R-Code

```

# simuliere aus einfachem lin. Modell
beta0 <- 3
beta1 <- 1
sigma <- 2
x <- seq(from = 0, to = 10, by = 0.5)
e <- rnorm(length(x), sd = sigma)
y <- beta0 + beta1 * x + e
dat <- data.frame(x, y)

# Lineares Modell erzeugen
reg = lm(y ~ x, data = dat)
summary(reg)

# Konfidenzintervalle
confint(reg, level = 0.95)

```

Interpretation von transformierten Modellen

- Log-Log-Modell:  $\log(10x) = \log(10) + \log(x)$  um eine Einheit  
 $\Rightarrow \beta_1$  steigt im EW um  $e^{\beta_1} = 10^{\beta_1}$

$\log(Y_i) = \beta_0 + \beta_1 \log(x_i) + \varepsilon_i$

$E(Y_i | x_i) = e^{\beta_0 + \beta_1 \log(x_i)}$   $= e^{\beta_0} \cdot e^{\beta_1 \log(x_i)} = e^{\beta_0} \cdot x_i^{\beta_1} = e^{\beta_0} \cdot x^{\beta_1} = e^{\beta_0 + \beta_1 \log(x)}$   $E(Y_i | x_i)$

Wenn  $x_i$  um den Faktor  $a$  steigt, dann steigt  $Y_i$  im Erwartungswert um den Faktor  $a^{\beta_1} = e^{\beta_1 \log(a)}$ .  $e^{\beta_1 \log(x_i)} = e^{\log(x_i^{\beta_1})} = x_i^{\beta_1}$  genau bei  $\log(10)$ !

Alternativ: Wenn  $x_i$  um 1% steigt, dann steigt  $Y_i$  im Erwartungswert um  $100 \cdot (e^{\beta_1 \log(1.01)} - 1)\%$ .

- Linear-Log-Modell:

$$Y_i = \beta_0 + \beta_1 \log(x_i) + \varepsilon_i$$

Wenn  $x_i$  um  $p\%$  steigt, dann steigt  $Y_i$  im Erwartungswert um  $\beta_1 \cdot \log(1 + p\%)$ .

Alternativ: Wenn  $x_i$  um 1% steigt, dann steigt  $Y_i$  im Erwartungswert um approximativ  $\beta_1\%$  Einheiten.

- Log-Linear-Modell:

$$\log(Y_i) = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Wenn  $x_i$  um eine Einheit steigt, dann steigt  $Y_i$  im Erwartungswert um den Faktor  $e^{\beta_1}$ .

Anmerkungen aus der Vorlesung

$R^2$  ist abhangig von  $X$ . Das heit uber mehrere Studien hinweg, die das gleiche messen, ist  $R^2$  nur vergleichbar, wenn auch  $X$  vergleichbar ist. Je sicherer wir mit unserem Schatzer sein wollen, desto hohrer sollten wir die Varianz von  $X$  einstellen. Gegeben, dass der Zusammenhang tatsachlich linear ist, wurde eine hohere Varianz von  $X$  zu einer geringeren Varianz von  $\hat{\beta}_1$  fuhren (false  $\text{Var}(X) \uparrow$  doppelt gut)

Im multiplen Reg.modell ist es KEINE Annahme, dass  $x_i, x_j$  unabhängig voneinander sind. Es wäre nur praktisch für die Interpretation der Effekte. Das "magische" am multiplen Reg.modell ist, dass ich für verschiedene Größen kontrollieren/korrigieren kann.

## Annahmen des linearen Regressionsmodells

Gegeben sei das einfache (oder multiple) lineare Regressionsmodell mit

$$Y = X\beta + \varepsilon \quad (1)$$

$$\mathbb{E}(\varepsilon) = \mathbf{0} \quad (2)$$

$$\mathbb{V}(\varepsilon_i) = \sigma^2, \text{ für alle } i = 1, \dots, n \quad (3)$$

$\varepsilon_i$  sind paarweise unabhängig voneinander (4)

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2), \text{ für alle } i = 1, \dots, n \quad (5)$$

Die folgende Tabelle gibt an, welche Annahmen für die jeweiligen Schätzer, Eigenschaften, Größen etc. benötigt werden.

	(1)	(2)	(3)	(4)	(5)
KQ-Schätzer	✓				
ML-Schätzer	✓	✓	✓	✓	✓
$E(\hat{\beta}) = \beta$	✓	✓			
$\hat{\beta} \sim \mathcal{N}(\beta, V(\hat{\beta}))$	✓	✓	✓	✓	(✓)
Konfidenzintervalle	✓	✓	✓	✓	(✓)
Prädiktionsintervalle	✓	✓	✓	✓	✓

(✓) bedeutet, dass die Annahme nicht benötigt wird, wenn der Stichprobenumfang  $n$  groß genug ist.

**Beachte:** Grundsätzlich gelten alle Aussagen nur unter der zentralen Modellannahme des Linearen Zusammenhangs von  $E(Y)$  und  $X$ !