

# 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

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

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

$$\hat{\beta}_1 = \frac{S_{xY}}{S_x^2} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$
$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}.$$

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

$$\sum_{i=1}^n \hat{\varepsilon}_i = 0$$
$$\sum_{i=1}^n \hat{\varepsilon}_i x_i = 0$$

## 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  $Y_i$ -Werte von der Regressionsgerade an.

## 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

$$\hat{\sigma}^2 := \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

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

$$\frac{n-2}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-2}^2.$$

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  $\hat{\beta}_0$

$$\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}_{\hat{\beta}_1}} \sim t_{n-2} \text{ mit } \hat{\sigma}_{\hat{\beta}_1} := \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

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

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

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

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

## 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

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

SST(otal): Gesamtstreuung von Y  
SSE(rror): Streuung der Residuen  
SSM(odel): 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 linearen Transformationen von  $x$  und  $Y$ .
- $R$  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  $\hat{Y}_{n+1} = \hat{\beta}_0 + \hat{\beta}_1 x_{n+1}$

### Prognosefehler

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

$$E(\hat{Y}_{n+1} - Y_{n+1}) = 0$$

$$V(\hat{Y}_{n+1} - Y_{n+1}) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

### Prognoseintervall

Gegeben sei ein einfaches lineares 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

$$\hat{\sigma}_{\hat{Y}_{n+1}} = \hat{\sigma}^2 \left[ 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

### 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:

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

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

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

- 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  $e^{\beta_1}$  Einheiten.

### Vorlesung

$R^2$  ist abhängig von  $X$ . Das heißt über mehrere Studien hinweg, die das gleiche messen, ist  $R^2$  nur vergleichbar, wenn auch  $X$  vergleichbar ist. Je sicherer wir mit unserem Schätzer sein wollen, desto höher sollten wir die Varianz von  $X$  einstellen. Gegeben, dass der Zusammenhang tatsächlich linear ist, würde eine höhere Varianz von  $X$  zu einer geringeren Varianz von  $\hat{\beta}_1$  führen.

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.

Erwartungstreue gilt auch bei Abhängigkeit und normalverteilt ist nicht nötig. Varianzformel benötigt Unabhängigkeit.

# Kapitel 2 - Das multiple lineare Regressionsmodell

## Multiple lineares Regressionsmodell

Das **multiple lineare Regressionsmodell** hat die Form

$$Y_i = \beta_0 + \underbrace{\beta_1 x_{i1} + \dots + \beta_p x_{ip}}_{\mathbf{x}_i^\top = (1, x_{i1}, \dots, x_{ip})} + \varepsilon_i; i = 1, \dots, n$$

oder in Matrix-Vektor Notation:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  mit

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix},$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}, \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

Wir nehmen dabei an, dass  $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$  eine feste Design-Matrix ist und dass  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . Wir definieren außerdem  $p' := p + 1$ .

## Kleinste Quadrate (KQ) Schätzer

Wir schätzen den Parameter(vektor)  $\boldsymbol{\beta}$  durch

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (2)$$

und nennen  $\hat{\boldsymbol{\beta}}$  den **KQ-Schätzer von  $\boldsymbol{\beta}$**  und  $\hat{\varepsilon}_i := Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$  die **Residuen**.

## Existenz und Berechnung vom KQ Schätzer

Der KQ-Schätzer existiert und ist eindeutig, falls  $\mathbf{X}^\top \mathbf{X}$  invertierbar ist. Dieser lässt sich berechnen als

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Durch differenzieren von der Gleichung (2) erhält man  $\hat{\boldsymbol{\beta}}$  als Lösung der **Normalengleichung**

$$\mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = 0$$

## Interpretation der Modellparameter

- $Y_i$  hängt linear von  $x_{i1}, \dots, x_{in}$  ab.
- Steigt  $x_k$  um eine Einheit, so steigt  $Y$  (ceteris paribus) im Erwartungswert um  $\beta_k$  Einheiten, **wenn** alle anderen  $x$ -Variablen festgehalten werden.
- !**  $\beta_k$  charakterisiert den Einfluss von  $x_k$  unter Berücksichtigung der übrigen Variablen (Confounder-Korrektur). Das heißt, dass in einem einfachen linearen Regressionsmodell mit  $Y_i = \beta_0 + \beta'_k x_{ik} + \varepsilon_i$  wäre im Allgemeinen  $\beta'_k \neq \beta_k$ .

## Eigenschaften des KQ-Schätzers

Gegeben dem multiplen linearen Modell, gilt für den KQ-Schätzer  $\hat{\boldsymbol{\beta}}$

- Erwartungstreue:  $\mathbb{E}(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ .  
**!** Gilt auch ohne die Annahme  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , solange  $\mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$
- $\mathbb{V}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ .  
**!** Gilt auch ohne die Annahme  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , solange  $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$
- $\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$

## Hat-Matrix und Residualmatrix

Gegeben dem multiplen linearen Modell mit  $\text{rang}(\mathbf{X}) = p'$  gilt

$$\hat{\mathbf{Y}} := \mathbf{X} \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}}_{\hat{\boldsymbol{\beta}}}$$

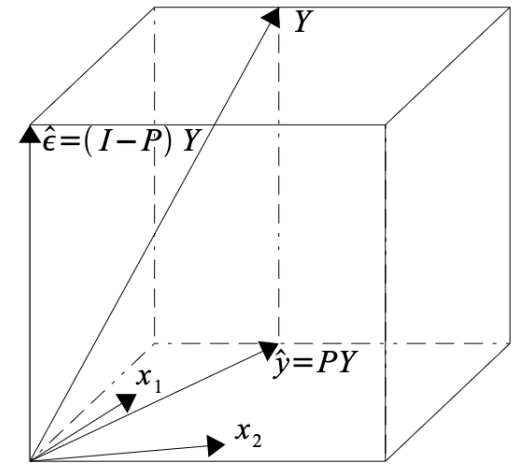
$$\mathbf{P} := \underbrace{\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{n \times n}$$

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{P})\mathbf{Y}$$

$$\mathbf{Q} := \mathbf{I} - \mathbf{P}$$

$\mathbf{P}$  heißt **Hat-Matrix** und  $\mathbf{Q}$  heißt **Residualmatrix**.

## Geometrische Interpretation



Die KQ-Schätzung ist eine orthogonale Projektion von  $\mathbf{Y}$  auf den von den  $\mathbf{x}$ -Vektoren aufgespannten Unterraum.

## Eigenschaften von $\mathbf{P}$ und $\mathbf{Q}$

Die Hat-Matrix  $\mathbf{P}$  und die Residualmatrix  $\mathbf{Q}$  sind Projektionsmatrizen und zueinander orthogonal:

$$\mathbf{P}^\top = \mathbf{P} \text{ und } \mathbf{P}^2 = \mathbf{P}$$

$$\mathbf{Q}^\top = \mathbf{Q} \text{ und } \mathbf{Q}^2 = \mathbf{Q}$$

$$\mathbf{P}\mathbf{Q} = \mathbf{Q}\mathbf{P} = \mathbf{0}.$$

Daraus folgt

$$\mathbb{V}(\hat{\mathbf{Y}}) = \sigma^2 \mathbf{P}$$

$$\mathbb{V}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2 \mathbf{Q}, \text{ da } \hat{\boldsymbol{\varepsilon}} = \mathbf{Q}\boldsymbol{\varepsilon}$$

## Schätzer für $\sigma^2$

Gegeben dem multiplen linearen Modell, gilt

$$\hat{\sigma}^2 := \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n - (p + 1)} = \frac{1}{n - (p + 1)} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

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

**!** Gilt auch ohne die Annahme  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , solange  $\mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$  und  $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$

# Kapitel 3 - Quadratsummenzerlegung und statistische Inferenz im multiplen linearen Regressionsmodell

## Quadratsummenzerlegung

Gegeben sei das multiple lineare Regressionsmodell mit  $\text{rang}(\mathbf{X}) = p'$ . Dann gilt

$$\underbrace{(\mathbf{Y} - \bar{\mathbf{Y}})^\top (\mathbf{Y} - \bar{\mathbf{Y}})}_{SST} = \underbrace{(\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})}_{SSE} + \underbrace{(\hat{\mathbf{Y}} - \bar{\mathbf{Y}})^\top (\hat{\mathbf{Y}} - \bar{\mathbf{Y}})}_{SSM}.$$

SST(otal):	Gesamt-Quadratsumme (korrigiert)
SSE(rror):	Fehler-Quadratsumme
SSM(odell):	Modell-Quadratsumme

## Quadratsummenzerlegung ohne $\beta_0$

Gegeben sei das multiple lineare Regressionsmodell mit, aber ohne Absolutglied  $\beta_0$ . Dann gilt

$$\underbrace{\mathbf{Y}^\top \mathbf{Y}}_{SST^*} = \underbrace{(\mathbf{Y} - \hat{\mathbf{Y}})^\top (\mathbf{Y} - \hat{\mathbf{Y}})}_{SSE} + \underbrace{\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}}_{SSM^*}.$$

SST*:	Gesamt-Quadratsumme (nicht korrigiert)
SSE:	Fehler-Quadratsumme (wie zuvor)
SSM*:	Modell-Quadratsumme (nicht korrigiert)

## Erwartungswerte der Quadratsummen

Gegeben sei das multiple lineare Regressionsmodell mit den üblichen Annahmen. Wir definieren

$$\mathbf{e} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \text{ und } \mathbf{P}_e = \mathbf{e}(\mathbf{e}^\top \mathbf{e})^{-1} \mathbf{e}^\top \text{ und } \mathbf{Q}_e = \mathbf{I} - \mathbf{P}_e.$$

Dann gilt

$$\begin{aligned} \mathbb{E}(SST^*) &= \sigma^2 n + \beta^\top \mathbf{X}^\top \mathbf{X} \beta \\ \mathbb{E}(SST) &= \sigma^2 (n - 1) + \beta^\top (\mathbf{Q}_e \mathbf{X})^\top (\mathbf{Q}_e \mathbf{X}) \beta \\ \mathbb{E}(SSE) &= \sigma^2 (n - p') \\ \mathbb{E}(SSM^*) &= \sigma^2 p' + \beta^\top \mathbf{X}^\top \mathbf{X} \beta \\ \mathbb{E}(SSM) &= \sigma^2 (p' - 1) + \beta^\top (\mathbf{Q}_e \mathbf{X})^\top (\mathbf{Q}_e \mathbf{X}) \beta \end{aligned}$$

Wir können diese Eigenschaften zur Konstruktion von Tests verwenden. Es gilt nämlich unter anderem

$$\begin{aligned} \beta &= 0 \implies \mathbb{E}(SST^*) = \sigma^2 n \\ \beta_1 = \dots = \beta_p &= 0 \implies \mathbb{E}(SSM) = \sigma^2 (p' - 1) \end{aligned}$$

## Chi-Quadrat Verteilung

Sei  $\mathbf{Z} \sim \mathcal{N}_n(\boldsymbol{\mu}, \mathbf{I})$ , so heißt  $\mathbf{W} = \mathbf{Z}^\top \mathbf{Z} = \sum_{i=1}^n Z_i^2$  (nicht-zentral) **Chi-Quadrat-verteilt** und wir schreiben

$$W \sim \chi^2(n, \delta).$$

Wir nennen  $n$  die **Zahl der Freiheitsgrade** und  $\delta = \boldsymbol{\mu}^\top \boldsymbol{\mu}$  den **Nicht-Zentralitätsparameter**. Es gilt

$$\begin{aligned} \mathbb{E}(W) &= n + \delta \\ \mathbb{V}(W) &= 2n + 4\delta \end{aligned}$$

## t-Verteilung

Seien  $Z \sim \mathcal{N}(\delta, 1)$  und  $W \sim \chi^2(n, 0)$  unabhängig. Dann heißt  $T = \frac{Z}{\sqrt{\frac{W}{n}}}$  (nicht-zentral) **t-verteilt** mit  $n$  **Freiheitsgraden** und **Nicht-Zentralitätsparameter**  $\delta$  und wir schreiben

$$T \sim t(n, \delta).$$

Es gilt

$$\mathbb{E}(T) = \delta \sqrt{\frac{n}{2}} \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n}{2})} \text{ für } n > 1$$

## F-Verteilung

Sei  $W_1 \sim \chi^2(n_1, \delta)$  und  $W_2 \sim \chi^2(n_2, 0)$  unabhängig. Dann heißt  $X = \frac{W_1/n_1}{W_2/n_2}$  (nicht-zentral) **F-verteilt** mit  $n_1$  und  $n_2$  **Freiheitsgraden** und **Nicht-Zentralitätsparameter**  $\delta$  und wir schreiben

$$X \sim F(n_1, n_2, \delta).$$

Es gilt

$$\mathbb{E}(X) = \frac{n_2 + \frac{n_2 \delta}{n_1}}{n_2 - 2} \text{ für } n_2 > 2$$

# Kapitel 4 - Diskrete Einflußgrößen

## Kodierung

Sei  $C$  eine nominale Variable mit  $K$  Ausprägungen.

### Dummy-Kodierung:

Wir definieren  $K$  neue Variablen  $Z_1, \dots, Z_K$  als

$$Z_k(C) = \begin{cases} 1, & \text{falls } C = k \\ 0, & \text{sonst} \end{cases}$$

$Z_1, \dots, Z_K$  sind abhängig, da  $Z_K = 1 - \sum_{k=1}^{K-1} Z_k$

**Effekt-Kodierung:** Wir definieren  $K - 1$  neue Variablen  $Z_1^e, \dots, Z_{K-1}^e$  als

$$Z_k^e(C) = \begin{cases} 1, & \text{falls } C = k \\ -1, & \text{falls } C = K \\ 0, & \text{sonst} \end{cases}$$

Note:  $Z_k(C) = \begin{pmatrix} Z_k(C_1) \\ \vdots \\ Z_k(C_n) \end{pmatrix}$  und  $Z_k^e(C) = \begin{pmatrix} Z_k^e(C_1) \\ \vdots \\ Z_k^e(C_n) \end{pmatrix}$

## Setup einfache Varianzanalyse

Im folgenden betrachten wir die einfache Varianzanalyse mit nur einer diskreten Einflußgröße  $C = \begin{pmatrix} C_1 \\ \vdots \\ C_n \end{pmatrix}$  mit  $K$  Ausprägungen. Sei  $n_k$  dabei die Anzahl der Beobachtungen mit  $C_i = k$ .

## Mittelwertsmodell

Das **Mittelwertsmodell** ist gegeben durch

$$Y_{kl} = \mu_k + \epsilon_{kl} \quad l = 1, \dots, n_k \quad k = 1, \dots, K$$

oder in Matrix-Vektor Notation:

$$\mathbf{Y} = (Z_1(C) \cdots Z_K(C)) \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_K \end{pmatrix} + \epsilon$$

Bei dem Mittelwertsmodell gibt es keinen Intercept und die  $\mu_k$  sind die Mittelwerte der  $k$ -ten Gruppe. Der Effekt der  $k$ -ten Gruppe ist also  $\mu_k$ .

## Mittelwertsmodell Beispiel

Für  $K = 3$  Ausprägungen und  $n_k = 2$  für alle  $k = 1, 2, 3$  erhalten wir als Mittelwertsmodell:

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{31} \\ \epsilon_{32} \end{pmatrix}$$

## Modell mit Effekt-Kodierung

Das **Modell mit Effekt-Kodierung** ist gegeben durch

$$Y_{kl} = \mu + \tau_k + \epsilon_{kl}; \quad \tau_K = - \sum_{k=1}^{K-1} \tau_k$$

für  $l = 1, \dots, n_k \quad k = 1, \dots, K$  oder in Matrix-Vektor Notation:

$$\mathbf{Y} = (e \ Z_1^e(C) \cdots Z_{K-1}^e(C)) \begin{pmatrix} \mu \\ \tau_1 \\ \vdots \\ \tau_{K-1} \end{pmatrix} + \epsilon$$

Bei dem Modell mit Effekt-Kodierung gibt es einen Intercept  $\mu$  und die  $\tau_k$  sind die Abweichungen der  $k$ -ten Gruppe vom Gesamtmittelwert bzw. vom Intercept  $\mu$ . Der Effekt der  $k$ -ten Gruppe ist also  $\mu + \tau_k$ .

## Modell mit Effekt-Kodierung Beispiel

Für  $K = 3$  Ausprägungen und  $n_k = 2$  für alle  $k = 1, 2, 3$  erhalten wir als Modell mit Effekt-Kodierung:

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \end{pmatrix} \begin{pmatrix} \mu \\ \tau_1 \\ \tau_2 \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{31} \\ \epsilon_{32} \end{pmatrix}$$

## Modell mit Referenz-Kodierung

Das **Modell mit Referenz-Kodierung** ist gegeben durch

$$Y_{kl} = \mu_K + \tau_k + \epsilon_{kl}; \quad \tau_K = 0$$

für  $l = 1, \dots, n_k \quad k = 1, \dots, K$  oder in Matrix-Vektor Notation:

$$\mathbf{Y} = (e \ Z_1(C) \cdots Z_{K-1}(C)) \begin{pmatrix} \mu_K \\ \tau_1 \\ \vdots \\ \tau_{K-1} \end{pmatrix} + \epsilon$$

Beim Modell mit Referenz-Kodierung gibt es einen Intercept  $\mu_K$  der den Mittelwert der  $K$ -ten Gruppe angibt und die  $\tau_k$  sind die Abweichungen der  $k$ -ten Gruppe vom Mittelwert der  $K$ -ten Referenz-Gruppe. Der Effekt der  $k$ -ten Gruppe ist also  $\mu_K + \tau_k$  für  $k = 1, \dots, K - 1$  und  $\mu_K$  für  $k = K$ .

## Modell mit Referenz-Kodierung Beispiel

Für  $K = 3$  Ausprägungen und  $n_k = 2$  für alle  $k = 1, 2, 3$  erhalten wir als Modell mit Referenz-Kodierung:

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} \mu_3 \\ \tau_1 \\ \tau_2 \end{pmatrix} + \begin{pmatrix} \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{21} \\ \epsilon_{22} \\ \epsilon_{31} \\ \epsilon_{32} \end{pmatrix}$$

## Bemerkungen-Kodierung

Alle Modellvarianten führen zur gleichen Modellanpassung ( $R^2$ ). Die Parameter haben aber unterschiedliche Interpretationen. Parameter und deren Schätzer sind aber ineinander umrechenbar.

### Setup zweifaktorielle Varianzanalyse

Im folgenden betrachten wir zwei diskrete Einflußgrößen  $\mathbf{C} = \begin{pmatrix} C_1 \\ \vdots \\ C_n \end{pmatrix}$  und  $\mathbf{D} = \begin{pmatrix} D_1 \\ \vdots \\ D_n \end{pmatrix}$  mit  $K_C$  bzw.  $K_D$  Ausprägungen. Sei  $n_{k,l}$  dabei die Anzahl der Beobachtungen mit  $C_i = k$  und  $D_j = l$ .

! Hier ist die Mittelwertsdarstellung bzw. das Mittelwertsmodell nicht möglich, da dieses davon abhängig ist, welche Variable zuerst kodiert wird.

### Modell mit Effekt-Kodierung (mehrfaktoriell)

Das **Modell mit Effekt-Kodierung** ist gegeben durch

$$\mathbf{Y} = (e \ Z_1^e(\mathbf{C}) \cdots Z_{K-1}^e(\mathbf{C})) \begin{pmatrix} \mu \\ \tau_1 \\ \vdots \\ \tau_{K-1} \end{pmatrix} + \epsilon$$

Bei dem Modell mit Effekt-Kodierung gibt es einen Intercept  $\mu$  und die  $\tau_k$  sind die Abweichungen der  $k$ -ten Gruppe vom Gesamtmittelwert bzw. vom Intercept  $\mu$ . Der Effekt der  $k$ -ten Gruppe ist also  $\mu + \tau_k$ .