

Materialien zu

# Lineare Modelle mit Anwendungen

Dr. Christian Kredler

SS 2003

# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
1.1	Das Modell . . . . .	2
1.2	Das Modell mit bedingter Erwartung . . . . .	3
1.3	Notation in der Literatur . . . . .	3
<b>2</b>	<b>Einfache Lineare Regression</b>	<b>4</b>
2.1	Modellannahmen . . . . .	4
2.1.1	Reparametrisierung . . . . .	4
2.1.2	Normalverteilung . . . . .	4
2.1.3	Definitionen und Identitäten . . . . .	5
2.2	Methode der Kleinsten Quadrate, KQ-Schätzer . . . . .	5
2.3	Einfache lineare Regression durch den Ursprung . . . . .	8
2.4	Tests und Konfidenzintervalle . . . . .	9
2.4.1	Testen von Hypothesen . . . . .	9
2.4.2	Analysis of Variance, ANOVA-Tabelle . . . . .	9
2.4.3	t-Test . . . . .	10
2.4.4	Konfidenzintervalle (KI) für Parameter . . . . .	10
2.5	Variablen-Transformationen . . . . .	11
2.6	Ein erster Blick auf die Residuen . . . . .	13
2.6.1	Modellwahl . . . . .	13
2.6.2	Varianzhomogenität . . . . .	15
2.6.3	Ausreißer . . . . .	15
2.6.4	Normalverteilungs-Annahme, QQ-Plots . . . . .	16
2.7	Simultane Inferenz bei Einfacher Linearer Regression . . . . .	17
2.8	Gemeinsame Konfidenzbänder für den <i>mean response</i> nach Working-Hotelling	18
2.9	Bonferroni-Ansatz zu gemeinsamer Interferenz . . . . .	19
<b>3</b>	<b>Multiple Lineare Regression</b>	<b>21</b>
3.1	Modell . . . . .	21
3.2	KQ-Methode (LS estimation) . . . . .	24
3.3	Eigenschaften der KQ-Schätzungen . . . . .	24
3.4	Projektionen und $\chi^2$ -Verteilung . . . . .	25
3.5	Varianzzerlegung, Bestimmtheitsmaß . . . . .	27
3.6	Tests linearer Hypothesen . . . . .	30
3.6.1	Test zur allgemeinen linearen Hypothese . . . . .	31
3.6.2	Tests für Variablen-Subsets . . . . .	32
3.6.3	Partielle F-Tests für einzelne Parameter . . . . .	32
3.7	Konfidenz- und Prognoseintervalle . . . . .	33
3.8	Qualität des Modells, Hat-Matrix . . . . .	34

3.8.1	Fragen . . . . .	34
3.8.2	Ansätze . . . . .	34
3.9	Variablenselektion . . . . .	38
3.9.1	Unterspezifikation . . . . .	39
3.9.2	Überspezifikation . . . . .	39
3.9.3	Kriterien zum Modellvergleich . . . . .	40
3.9.4	Algorithmisches Vorgehen . . . . .	43
3.9.5	Variablenselektion in R und S-Plus . . . . .	43
3.10	Multikollinearität . . . . .	46
3.10.1	Einfache lineare Regression . . . . .	46
3.10.2	Skalierungsinvarianz . . . . .	47
3.10.3	Einfache Berechnung des VIF . . . . .	48
3.11	Korrelierte Fehler, gewichtete Regression . . . . .	49
3.12	Autokorrelierte Fehler . . . . .	51
3.13	Kategoriale Variable; Teil 1: Regressoren . . . . .	55
3.13.1	Einführendes Beispiel mit einer binären Kovariablen . . . . .	55
3.13.2	Beispiel 2: Mehrstufiger Faktor . . . . .	58
3.14	Kategoriale und Indikatorvariablen; Teil 2: ANOVA-Modelle . . . . .	60
3.14.1	Einfaktorielle Varianzanalyse; one way ANOVA . . . . .	60
3.14.2	Zweifaktorielles Design; two way ANOVA . . . . .	62
<b>4</b>	<b>Anhang</b>	<b>70</b>
4.1	Normalgleichungen . . . . .	70
4.2	Rechnen mit Erwartungswerten und Kovarianzmatrizen . . . . .	71
4.2.1	n-dimensionale Normalverteilung . . . . .	71
4.2.2	Erwartungswerte; n-dim. . . . .	71
4.3	Transformationen für Dichten und unabhängig normalverteilte Zufallsvariablen . . . . .	74
4.3.1	Transformationsatz für Dichten . . . . .	74
4.3.2	Lineare Transformationen . . . . .	75
4.4	Zentriertes Modell, orthogonale Designmatrix . . . . .	77
	<b>Literatur</b>	<b>80</b>

# Kapitel 1

## Einführung

**Lineare Modelle mit Anwendungen** schließt direkt an die Einführungsvorlesung Stochastik 1 (zum vorausgesetzten Stoff vgl. [Kredler (1999)]) an und behandelt die wohl häufigste Anwendung der mathematischen Statistik in der Praxis.

Die Vorlesung ist für Diplom-Mathematiker, Techno-, Finanz- und Wirtschaftsmathematiker sowie Studierende des Lehramts Mathematik an Gymnasien konzipiert. Sie kann ab dem vierten Semester gehört werden und bietet sich ebenso für Physiker und Ingenieure an.

Der Praxisbezug wird durch eine intensive Arbeit am Datenmaterial hergestellt. Als Software wurde R bzw. S-Plus gewählt. Das Ausbildungskonzept an der TU München beruht auf einer umfassenden Darstellung im Web (siehe [www-m4.ma.tum.de/nbu/linmod/](http://www-m4.ma.tum.de/nbu/linmod/)). Diese beinhaltet u.a.

- eine Zusammenstellung der notwendigen Vorkenntnisse,
- eine Einführung in R bzw. S-Plus,
- eine ausführliche Diskussion der Vorlesungsbeispiele mit Grafiken und den relevanten S-Plus-Befehlen,
- zum Selbststudium geeignete Übungsaufgaben mit Datenbeispielen,
- eine Sammlung von Beispiel-Datensätzen.

Diese Vorlesungsarbeit erhebt keinen Anspruch auf Vollständigkeit. Sie orientiert sich in Aufbau und Notation am Buch von [Myers (1990)]. Viele Beispiele sind einem Skript von Prof. Claudia Czado entnommen, die eine entsprechende Vorlesung im SS 2000 an der TU München hielt.

## 1.1 Das Modell

*Regressionsanalyse* (mit Hilfe Linearer Modelle) ist ein statistisches Verfahren, das den Zusammenhang bestimmt zwischen

- einer *abhängigen Variablen*  $Y$  (Zielvariable, response)
- und *unabhängigen Variablen*  $X_1, \dots, X_k$  (Kovariablen, Regressoren, erklärende Variablen), wobei eine lineare Beziehung der Form

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + E$$

mit einer Fehlervariablen  $E$  unterstellt wird.

Beispiele mit mehreren Kovariablen werden in Abschnitt 3.1 diskutiert. Der Spezialfall **Einfache Lineare Regression** ( $k = 1$ ) wird in Kapitel 2 behandelt. Schon am Beispiel

$Y$ : Körpergewicht

$X$ : Körpergröße

sieht man, dass Response  $Y$  und Kovariable (hier  $X$ ) nicht vertauschbar sind. In der Regel wird diejenige Variable als Response gewählt, für die Prognosen zu erstellen sind oder die eine natürliche Variabilität aufweist. So liegt die Körpergröße ab einem bestimmten Alter fest, was man vom Körpergewicht meist nicht sagen kann.

Allgemein betrachtet man in der sog. **Multiplen Linearen Regression**

$$(M_k) \quad Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n$$

$x_{ik}$  gegebene, nicht stochastische Daten

$$E(\epsilon_i) = 0, \quad Cov(\epsilon_i, \epsilon_j) = \sigma^2 \delta_{ij}, \quad \delta_{ij} = \begin{cases} 1, & \text{falls } i = j \\ 0, & \text{sonst} \end{cases}, \quad (1.1)$$

d.h. Fehler sind unkorreliert mit homogener Varianz.

Damit sind die  $Y_i$  Zufallsvariablen (ZV) mit homogener Varianz, unkorreliert, aber mit verschiedenen Erwartungswerten.

Das Modell  $(M_k)$  ist linear in den unbekanntem Parametern  $\beta_0, \beta_1, \dots, \beta_k$ .

Die  $Y_i$  und die  $x_{ij}$  können nichtlineare Transformationen anderer Variablen sein.

Spezialfälle:

$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$	Einfache Lineare Regression
$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$	linear in $\beta_0, \beta_1, \beta_2$
$Y_i = \beta_0 e^{\beta_1 x_i} \cdot \epsilon_i$	$\ln(Y_i)$ linear in $\gamma_0 = \ln(\beta_0)$ und $\gamma_1 = \ln(\beta_1)$
$Y_i = \beta_0 + e^{\beta_1 x_i} + \epsilon_i$	$\ln(Y_i)$ nichtlinear in $\beta_0, \beta_1$

Ziele

1. Finde ein "gutes Modell" für  $Y$  (*Modell-Spezifikation*)  $\rightarrow$  "kleiner Fehler"  $\epsilon_i$
2. Finde "gute" Regressoren (*Variablenselektion*)
3. Schätze die unbekannt Parameter (*Parameterschätzung*)
4. Sage  $Y$  für feste  $x$ 's voraus (*Vorhersage, Prediction*)

Einige Probleme

1. Wichtige Regressoren wurden in den Daten nicht miterhoben.
2. Vorhersagende  $Y$ -Werte gehören zu  $x$ -Werten, die außerhalb der bisher beobachteten  $x$ -Werte liegen.
3. Der wahre Zusammenhang zwischen  $Y$  und den  $x$ 's entspricht nicht dem Ansatz von  $(M_k)$ , ist z.B. nicht additiv.

## 1.2 Das Modell mit bedingter Erwartung

In vielen Fällen betrachtet man  $n$  unabhängige Wiederholungen des Zufallsvektors  $(Y, X_1, \dots, X_k)$ . Dabei seien die Realisierungen  $x_{i1}, \dots, x_{ik}$  der ZV  $X_{i1}, \dots, X_{ik}$ ,  $i = 1, \dots, n, j = 1, \dots, k$ , zu einem festgesetzten Status bzw. Zeitpunkt beobachtbar, die Zielvariablen  $Y_i$  aber erst später. Nach Eintritt der  $X_{ij} = x_{ij}$  sind also nur noch die  $Y_i$  stochastisch. In diesem Fall lautet der Modellansatz

$$(M'_k) \quad E(Y_i | X_{i1} = x_{i1}, \dots, X_{ik} = x_{ik}) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}, \quad i = 1, \dots, n.$$

## 1.3 Notation in der Literatur

Als häufigste Schreibweise findet man (siehe auch Kapitel 3)

$$(M''_k) \quad y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n$$

zu gegebenen Kovariablen-Daten  $x_{ik}$ . Diese werden als nicht stochastisch angesehen. Die  $y_i$  sind stochastisch, bezeichnen also ZV. In den meisten Fällen wird nicht zwischen den ZV selbst und deren Realisierungen (Daten) unterschieden.

# Kapitel 2

## Einfache Lineare Regression

### 2.1 Modellannahmen

In der Einfachen Linearen Regression liegen Daten der Form  $(y_i, x_i)_{i=1, \dots, n}$  vor. Für diese wählt man den Ansatz

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad (2.1)$$

(M)  $E(\epsilon_i) = 0$  und  $Cov(\epsilon_i, \epsilon_j) = \delta_{ij} \sigma^2$ ,  $i = 1, \dots, n$ .

$\beta_0$  und  $\beta_1$  sind unbekannte Parameter.

Die Daten  $y_i$  werden als Realisierungen der Zufallsvariablen (ZV)  $Y_i$  angesehen. Die  $x_i$  sind nicht stochastisch, z.B. Messstellen.

#### 2.1.1 Reparametrisierung

Für theoretische Überlegungen erweist sich stets das sog. "zentrierte Modell" als überlegen. Dieses erhält man aus Gleichung 2.1 mit  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$  und

$$y_i = \beta_0^* + \beta_1 (x_i - \bar{x}) + \epsilon_i, \quad (2.2)$$

wobei  $\beta_0 = \beta_0^* - \beta_1 \bar{x}$ .

#### 2.1.2 Normalverteilung

Die Annahme

$$\epsilon_i \text{ iid } N(0, \sigma^2) \quad (2.3)$$

impliziert (1.1) bzw. (2.1), d.h.  $E(\epsilon_i) = 0$  und  $Cov(\epsilon_i, \epsilon_j) = \delta_{ij} \sigma^2$ .

$b_0$  (bzw.  $b_0^*$ ) und  $b_1$  seien statistische Schätzer für die unbekannt Parameter  $\beta_0$  (bzw.  $\beta_0^*$ ) und  $\beta_1$ . Diese werden so gewählt, dass die Regressionsgerade  $y = b_0 + b_1 x$  die Punktwolke  $(y_i, x_i)_{i=1, \dots, n}$  möglichst "gut" anpasst.

### 2.1.3 Definitionen und Identitäten

**Def. 2.1 (Response und Residuen)**

$\hat{Y}_i := b_0 + b_1 x_i = b_0^* + b_1(x_i - \bar{x})$  heißen **Y-Schätzungen (fitted response)**  
 $e_i := Y_i - \hat{Y}_i$  heißen **Residuen**.

$\hat{Y}_i$  und  $e_i$  sind ZV mit denselben Erwartungswerten wie  $Y_i$  und  $\epsilon_i$ . Sie sind aber nicht mehr unkorreliert.

**Def. 2.2 (Sums of Squares)**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$SS_{xy} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) = \sum_{i=1}^n (x_i - \bar{x})Y_i$$

$$SS_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i - \bar{x})x_i$$

$$SS_e = RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - b_0 - b_1 x_i)^2 = SS_{Res}$$

**RSS = residual sum of squares**

$Q(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2$  heißt *Summe der quadratischen Abweichungen*  
für die Parameter  $\beta_0$  und  $\beta_1$ .

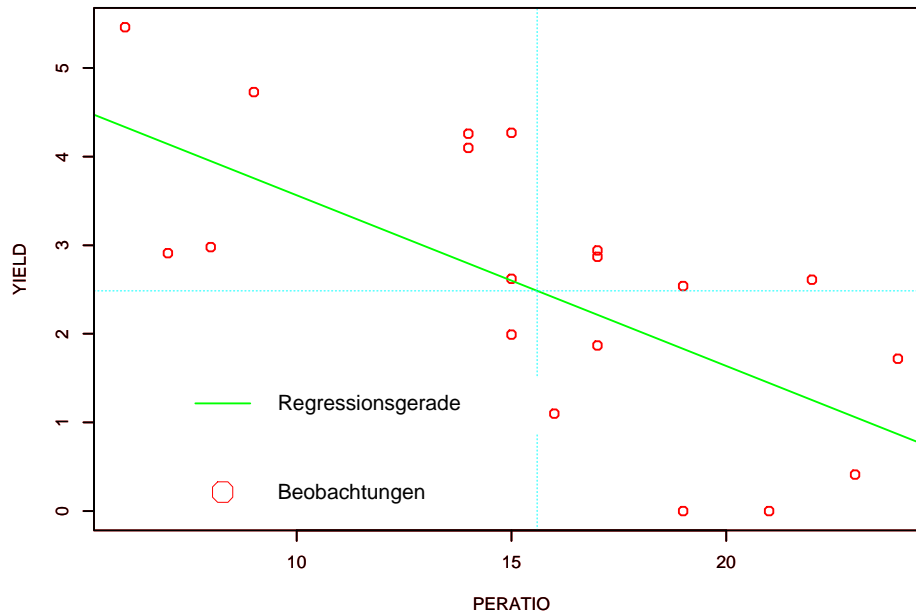
## 2.2 Methode der Kleinsten Quadrate, KQ-Schätzer

Nach der Methode der kleinsten Quadrate (**KQ, Least Squares, LS**) wählt man die Schätzungen  $b_0$  und  $b_1$  so, dass

$$Q(b_0, b_1) \leq Q(\beta_0, \beta_1) \quad \forall \beta_0, \beta_1.$$



## Beispiel 1: Dividende und Price-Expense Ratio

Abbildung 2.1:  $y = YIELD$ ,  $x = PERATIO$ 

$YIELD$  = Jährliche Dividendenrate in %  
bezogen auf Kurswert der Aktie (share price) am 16.03.90

$$PERATIO = \frac{\text{Kurswert (16.03.90)}}{EPS}$$

$$EPS = \frac{\text{Nettogewinn} - \text{Dividenden}}{\text{Anzahl der Aktien}} \text{ bzgl. 1989 earnings}$$

**Satz 2.3 (KQ-Schätzer)**

$$1. b_1 = \frac{SS_{xy}}{SS_{xx}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$b_0^* = \bar{Y}$$

$$b_0 = \bar{Y} - b_1 \bar{x}$$

2. Ein unverzerrter Schätzer für  $\sigma^2$  ist

$$s^2 = \hat{\sigma}^2 = \frac{SS_{Res}}{n-2} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}$$

$n$ : Anzahl der Beobachtungen

2: Anzahl der unbekannt Parameter

**Satz 2.4 (Eigenschaften der KQ-Schätzer)**

$$1. E(b_1) = \beta_1$$

$$2. E(b_0) = \beta_0, E(b_0^*) = \beta_0^*$$

3.  $b_1$  und  $b_0^*$  sind unkorreliert.

4.  $(b_0, b_1)$  und  $\mathbf{e} = (Y_1 - \hat{Y}_1, \dots, Y_n - \hat{Y}_n)^T$  sind unkorreliert.

$$5. \text{Var}(b_1) = \frac{\sigma^2}{SS_{xx}}$$

$$6. \text{Var}(b_0^*) = \frac{\sigma^2}{n}$$

$$7. \text{Var}(b_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}} \right)$$

8. Unter Normalverteilungsannahme gilt:

8a)  $(b_0, b_1)$  sind normalverteilt.

8b)  $(b_0, b_1)$  und  $s^2$  sind unabhängig.

Seien nun  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ ,  $\mathbf{x} = (x_1, \dots, x_n)^T$ ,  $\mathbf{e} = (e_1, \dots, e_n)^T$ ,  $\mathbf{1} = (1, \dots, 1)^T$  und  $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_n)^T$ .

Mit  $\hat{Y}_i = b_0 + b_1 x_i$ ,  $e_i = Y_i - \hat{Y}_i$ ,  $i = 1, \dots, n$  und

$$\hat{\mathbf{Y}} \in \text{span}(\mathbf{1}, \mathbf{x}), \mathbf{e} \in \text{span}(\mathbf{1}, \mathbf{x})^\perp$$

folgen

$$\hat{\mathbf{Y}} \perp \mathbf{e} \text{ und } \hat{\mathbf{Y}} - \bar{Y} \cdot \mathbf{1} \perp \mathbf{e}. \quad (2.4)$$

**Satz 2.5 (Streuungszerlegung)**

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$SS_{Total}$	$SS_{Reg}$	$SS_{Res}$
<i>totale</i>	<i>durch Regression</i>	<i>nicht erklärte</i>
<i>Variabilität</i>	<i>erklärte Variabilität</i>	<i>Variabilität</i>

Das **Bestimmtheitsmaß**  $R^2$  (coefficient of determination) misst die Güte des gewählten Modells.

**Def. 2.6 (Multipler Korrelationskoeffizient R)**

$$R^2 = \frac{SS_{Reg}}{SS_{Total}} = 1 - \frac{SS_{Res}}{SS_{Total}}$$

$R = \text{sign}(b_1) \cdot \sqrt{R^2}$  heißt **multipler Korrelationskoeffizient**.

Es gilt:

$$0 \leq R^2 \leq 1 \tag{2.5}$$

$$R^2 = 1 \Leftrightarrow Y_i = b_0 + b_1 x_i, \quad i = 1, \dots, n$$

$$R^2 = 0 \Leftrightarrow SS_{Res} = SS_{Total} \Leftrightarrow \text{Regression bringt nichts}$$

**2.3 Einfache lineare Regression durch den Ursprung**

Der Vollständigkeit halber werden kurz die Gleichungen für das Modell ohne "intercept"  $\beta_0$  eingeschoben.

$$Y_i = \beta x_i + \epsilon_i \tag{2.6}$$

$$b = \hat{\beta} = \frac{\sum_{i=1}^n x_i Y_i}{\sum_{i=1}^n x_i^2} \tag{2.7}$$

$$\tilde{R}^2 = \frac{\sum_{i=1}^n \hat{Y}_i^2}{\sum_{i=1}^n Y_i^2}, \quad \hat{Y}_i = b x_i \tag{2.8}$$

Ab jetzt gelte wieder das Modell (2.1).

## 2.4 Tests und Konfidenzintervalle

### 2.4.1 Testen von Hypothesen

#### Satz 2.7 (Verteilung von $SS_{Res}$ , $SS_{Total}$ , $SS_{Reg}$ )

Vor.:  $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ ,  $\epsilon_i$  iid  $N(0, \sigma^2)$ ,  $i = 1, \dots, n$

Beh.:

1.  $SS_{Res} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \sim \sigma^2 \chi_{n-2}^2$

2.  $SS_{Res}$  und  $SS_{Reg}$  sind unabhängig.

3. Falls zusätzlich gilt:  $H_0 : \beta_1 = 0$

$$SS_{Reg} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \sim \sigma^2 \chi_1^2 \text{ (unter } H_0)$$

$$SS_{Total} = \sum_{i=1}^n (Y_i - \bar{Y})^2 \sim \sigma^2 \chi_{n-1}^2 \text{ (unter } H_0)$$

$$F = \frac{SS_{Reg}}{SS_{Res}/(n-2)} \sim F_{1, n-2}$$

Testvorschrift:

$$\begin{array}{ll} H_0 : \beta_1 = 0, & H_1 : \beta_1 \neq 0 \\ \text{Nullhypothese} & \text{Alternative} \end{array}$$

$H_0$  ist beim Niveau  $\alpha$  (z.B.  $\alpha = 0.05, 0.01$ ) abzulehnen, falls  $F > F_{1-\alpha; 1, n-2}$

### 2.4.2 Analysis of Variance, ANOVA-Tabelle

df = degrees of freedom, Freiheitsgrade

Source	df	SS	MS	F
Regression	1	$SS_{Reg}$	$SS_{Reg}$	$MS_{Reg}/s^2$
Residual	$n-2$	$SS_{Res}$	$s^2$	
Total	$n-1$	$SS_{Total}$		

MS = mean square

ANOVA-Tabelle zu Beispiel 1

Source	df	SS	MS	F	P-Wert
Regression	1	18.60	18.60	11.36	0.0034
Residual	18	29.48	1.64		
Total	19	48.08			

$H_0$  ist zu verwerfen beim Niveau  $\alpha = 0.05$  und  $\alpha = 0.01$

$H_0$  ist anzunehmen beim Niveau  $\alpha = 0.001$

**Def. 2.8 (P-Wert; p-value)**

*Der P-Wert ist das kleinste  $\alpha$ -Niveau, für das die Hypothese noch verworfen werden kann.*

### 2.4.3 t-Test

Generell gilt für eine mit  $m$  Freiheitsgraden studentverteilte Zufallsvariable

$$T_m^2 = F_{1,m}. \quad (2.9)$$

Für den Signifikanztest zur Hypothese  $H_0 : \beta_1 = 0$  gegen die Alternative  $H_1 : \beta_1 \neq 0$  ist

$$t = \frac{b_1 - \beta_1}{se(b_1)}. \quad (2.10)$$

studentverteilt mit  $n-2$  Freiheitsgraden, wobei  $se(b_1) = \frac{s}{\sqrt{SS_{xx}}}$  den Standardfehler (standard error) von  $b_1$  bezeichnet.

Der Standardfehler eines Schätzers (hier  $b_1$ ) ist die Schätzung für seine (unbekannte) Standardabweichung (hier  $\sigma/\sqrt{SS_{xx}}$ ). Dabei muss die entsprechende Schätzung der Varianz (hier von  $b_1$ ) erwartungstreu sein und sollte selbst möglichst kleine Varianz haben.

### 2.4.4 Konfidenzintervalle (KI) für Parameter

Mit  $a = t_{1-\frac{\alpha}{2};n-2} \cdot se(b_1)$  gilt

$$1 - \alpha = P(b_1 - a \leq \beta_1 \leq b_1 + a) \quad (2.11)$$

allgemein:  $a = \text{Quantil} \cdot \text{Standard-Error Punktschätzer}$

Analog hierzu erhält man ein KI für  $\beta_0$  mit

$$se(b_0) = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_{xx}}}$$

**KI für erwarteten Response (mean response)  $E(Y | x_0)$**

$$\hat{Y}(x_0) = b_0 + b_1 x_0, \quad se(E(Y | x_0)) = s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}$$

$$c = t_{1-\frac{\alpha}{2};n-2} \cdot se(E(Y | x_0))$$

$$1 - \alpha = P(\hat{Y}(x_0) - c \leq \beta_0 + \beta_1 x_0 \leq \hat{Y}(x_0) + c)$$

**Prognose-Intervall (prediction interval) für  $Y_0 = \beta_0 + \beta_1 x_0 + \epsilon_0$**

$$\hat{Y}_0 = b_0 + b_1 x_0$$

$$se(\hat{Y}_0 - Y_0) = s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}}$$

$$d = t_{1-\frac{\alpha}{2}; n-2} \cdot se(\hat{Y}_0 - Y_0)$$

$$1 - \alpha = P(\hat{Y}_0 - d \leq Y_0 \leq \hat{Y}_0 + d)$$

## 2.5 Variablen-Transformationen

Oft deuten Scatterplots auf Varianzheterogenität hin; vgl. Mietspiegel-Daten mit  $x = Wfl$ ,  $Z = NM$ .

In diesem Beispiel führt z.B. die Transformation  $Y := \tilde{Z} := \ln Z$  zu einem zufriedenstellenderen Ergebnis (Vor.  $Z_i > 0$ )

Man hätte auch eine Transformation aus der Klasse

$\tilde{Z} := Z^\gamma$ ,  $0 < \gamma < 1$ , z.B.  $\tilde{Z} = \sqrt{Z}$  wählen können. (Vor.  $Z_i \geq 0$ )

Am praktischen Beispiel sucht man geeignete Transformationen durch Probieren und Vergleich der Scatterplots transformierter Daten.

Dabei verfolgt man das Ziel: Die Residuen sollen keine erkennbaren Muster aufweisen.

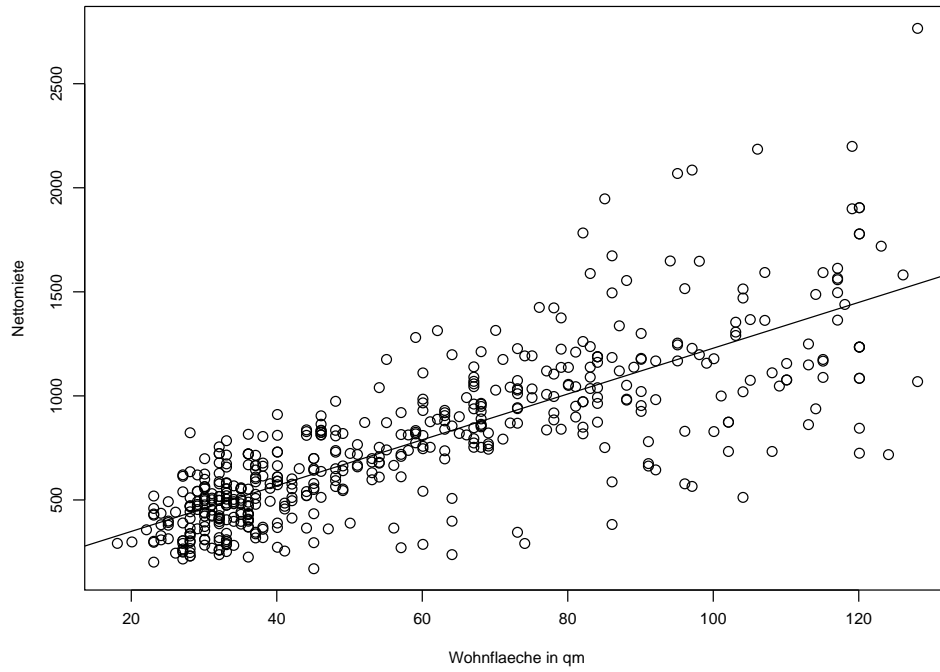


Abbildung 2.2:  $y = \text{Nettomiete}$ ,  $x = \text{Wohnflaeche}$

## 2.6 Ein erster Blick auf die Residuen

Seien

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$E(\epsilon_i) = 0, \text{Cov}(\epsilon_i, \epsilon_j) = \delta_{ij} \sigma^2, \quad i, j = 1, \dots, n \text{ oder } \epsilon_i \text{ iid } N(0, \sigma^2)$$

$b_0, b_1$  KQ-Schätzer für  $\beta_0, \beta_1$

$$\hat{Y}_i = b_0 + b_1 x_i, \quad e_i = Y_i - \hat{Y}_i \text{ (Residuen).}$$

Achtung:

$e_i$	$\neq$	$\epsilon_i$
berechenbar		nicht beobachtbar

Die  $e_i$  sind korrelierte Zufallsvariablen, verwendet zur

1. Verifizierung der Modellwahl (siehe z.B. "Korngrößen")
2. Überprüfung Varianzhomogenität:  
 $Var(\epsilon_i) = \sigma^2$ , konstant; verletzt bei "Mietspiegel"
3. Identifizierung von Ausreißern
4. Überprüfung der Normalverteilungs-Annahme (siehe "Korngrößen")

### 2.6.1 Modellwahl

- a) Die Residuen sollen (aufgetragen gegen die  $x_i$  oder später gegen die  $\hat{Y}_i$ ) *kein Muster* aufweisen (vgl. Abb. 2.3).



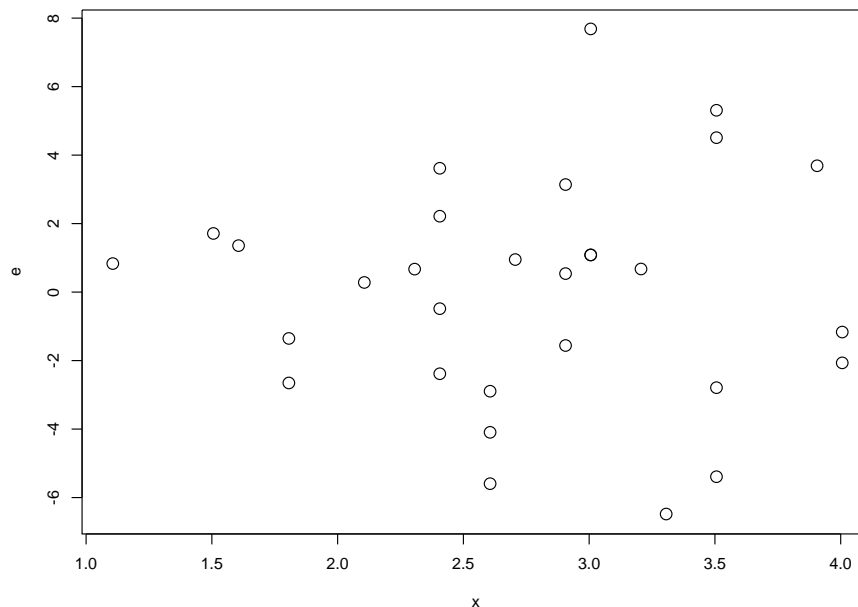


Abbildung 2.3: Kein erkennbares Muster in den Residuen

Die Muster in Abb. 2.4 wären z.B. durch Daten-Transformationen wie  $a \sin(tx_i + c)$  oder  $a(x_i - c)^2$  zu eliminieren:

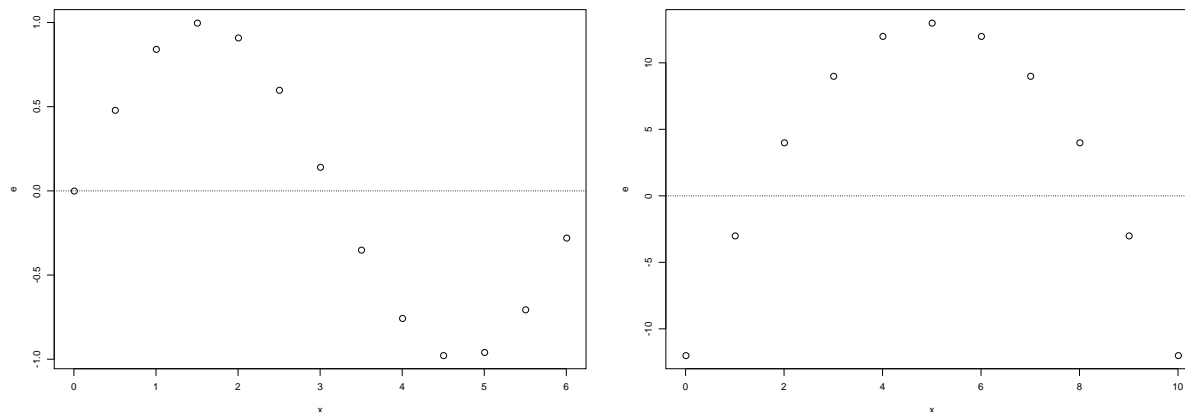


Abbildung 2.4: Funktionale Muster in den Residuen

b) Verletzung der Modellannahmen

z.B. weil die  $x_i$  den  $E(Y_i | x_i)$  nicht gut genug erklären.

Mögliche Abhilfe, z.B.

$$Y_i = \beta_0 + \beta_1 g(x_i) + \epsilon_i,$$

wobei  $g$  nach grafischer Inspektion der Daten durch "Ausprobieren" zu finden ist.

In manchen Fällen hilft für lokale Approximationen ein kubischer Ansatz der Form

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \epsilon_i.$$

Dies ist ein multiples Regressionsmodell; siehe Kap. 3.

## 2.6.2 Varianzhomogenität

In der Regel gilt  $Var(\epsilon_i) = \sigma_i^2$ ,  $i = 1, \dots, n$ .

Beim Beispiel "Mietspiegel" wächst z.B. die Varianz mit der Wohnfläche.

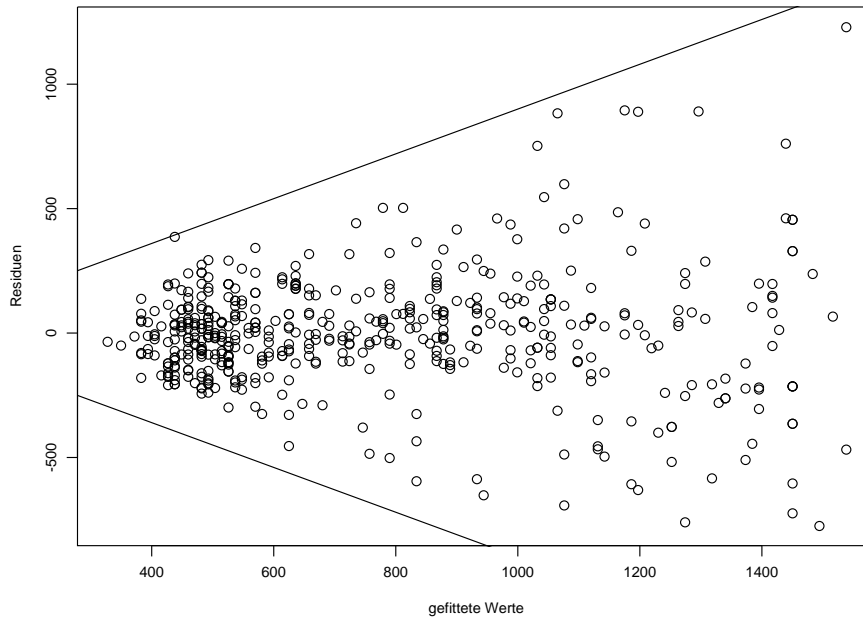


Abbildung 2.5: Inhomogene Varianz

Abhilfe: Transformationen wie  $\tilde{Y}_i = \log(Y_i)$  oder  $\tilde{Y}_i = w_i Y_i$ ,  $w_i > 0$ , bekannte Gewichte.

## 2.6.3 Ausreißer

Seien  $\mathbf{X} = (\mathbf{1}, \mathbf{x} - \mathbf{1} \cdot \bar{x})$  und  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = (h)_{ij}$ , wobei in der Einfachen Linearen Regression gilt  $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_{xx}}$  und

$$\frac{1}{n} \leq h_{ii} \leq 1, \quad \sum_{i=1}^n h_{ii} = 2.$$

$$Cov(Y_i, \hat{Y}_i) = \sigma^2 h_{ii} = Var(\hat{Y}_i),$$

$$Var(e_i) = Var(Y_i - \hat{Y}_i) = \sigma^2(1 - h_{ii}), \quad Cov(e_i, e_j) = \sigma^2(1 - h_{ij}) \neq 0, \quad \text{für } i \neq j.$$

$h_{ii} = 1$  ist möglich. In diesem Fall hat man  $Var(e_i) = 0$  unabhängig von der Fehlerquadratsumme  $SS_{Res}$ .

$$\text{Mit } s^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2, \quad se(e_i) = s\sqrt{1 - h_{ii}}$$

definiert man die "standardisierten Residuen" (standardized residuals)

$$r_i := \frac{e_i}{se(e_i)} = \frac{Y_i - \hat{Y}_i}{s\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n.$$

Achtung: Die  $r_i$  sind weder Student-verteilt noch unkorreliert.

Abhilfe: *Jackknifed Residuals* (manchmal auch: *Studentized Residuals*)

$$r_j^* = \frac{Y_j - \hat{Y}_j^{(j)}}{se(Y_j - \hat{Y}_j^{(j)})} = r_j / \sqrt{\frac{n-2-r_j^2}{n-3}}, \quad j = 1, \dots, n$$

wobei für  $j = 1, \dots, n$ ,  $b_0^{(j)}$ ,  $b_1^{(j)}$  und  $s_{(j)}^2$  aus dem Modell

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n, \quad i \neq j$$

mit  $n-1$  Beobachtungen ohne "  $j$  " berechnet sind; also  $\hat{Y}_j^{(j)} = b_0^{(j)} + b_1^{(j)} x_j$ .

Vorteil: Ein eventuell großes Residuum  $e_j = Y_j - \hat{Y}_j$  geht nicht in die Berechnungen von  $\hat{Y}_j^{(j)}$ ,  $s_{(j)}^2$  und somit von  $r_j^*$  ein.

**Faustregel:** Beobachtung  $j$  mit  $|r_j^*| > 2 \approx t_{0.975; n-3}$  ist Kandidat für *Ausreißer*.

## 2.6.4 Normalverteilungs-Annahme, QQ-Plots

$e_i \sim N(0, \sigma^2(1 - h_{ii}))$ ,  $r_i^* \approx N(0, 1)$

geordnet:  $r_{(1)}^* \leq r_{(2)}^* \leq \dots \leq r_{(n)}^*$ .

Zu  $U_1, \dots, U_n$  iid  $N(0, 1)$  definiert man die Ordnungsstatistiken:

$$U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}.$$

Aus der Statistik 1 ist bekannt:  $E(U_{(k)}) = \Phi^{-1}\left(\frac{k}{n+1}\right)$

Die Punktwolke  $\left(\Phi^{-1}\left(\frac{k}{n+1}\right), U_{(k)}\right)_{k=1, \dots, n}$  sollte also "eng" um die Winkelhalbierende

streuen; im Fall der *jackknifed residuals* also die Punktwolke  $\left(\Phi^{-1}\left(\frac{k}{n+1}\right), r_{(k)}^*\right)$ , vgl. etwa den Datensatz "Korngrößen".

$$n = 39 : \Phi^{-1}(1/40) = \Phi^{-1}(0.025) = -1.96$$

$$n = 99 : \Phi^{-1}(1/100) = \Phi^{-1}(0.01) = -2.33$$

$$n = 999 : \Phi^{-1}(1/1000) = \Phi^{-1}(0.001) = -3.09$$

Die folgenden Grafiken wurden durch Simulationen von unabhängigen normalverteilten Fehlern erzeugt; stellen also den Idealfall dar.

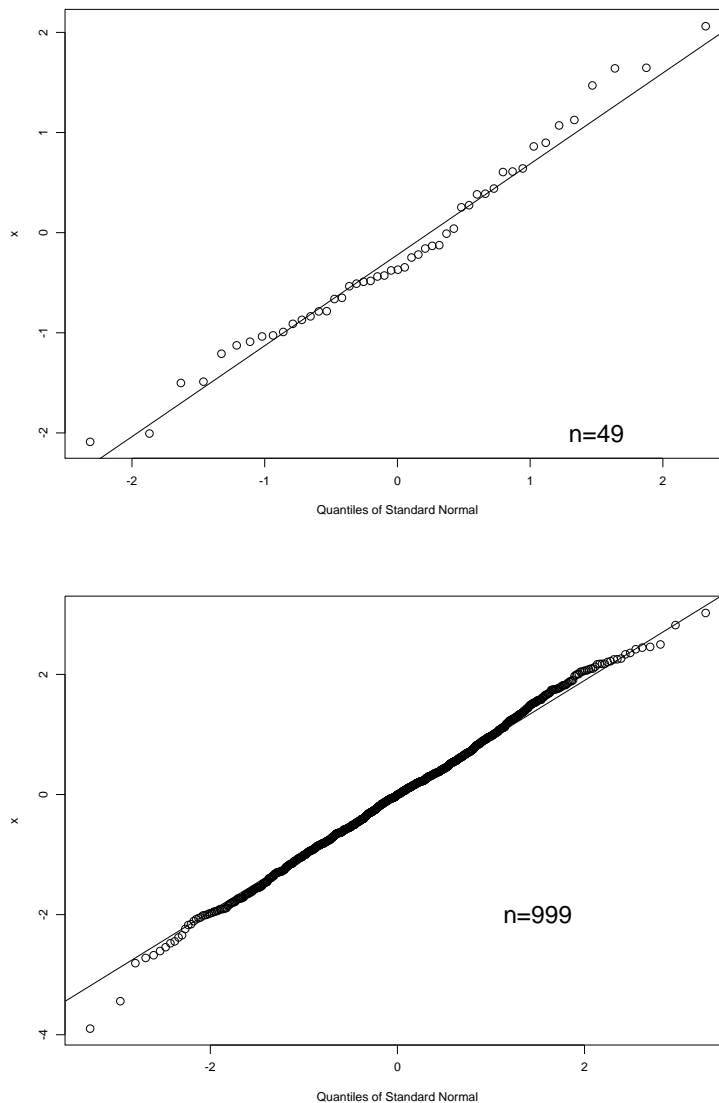


Abbildung 2.6: Ideale QQ-Plots für exakt normalverteilte Fehler

## 2.7 Simultane Inferenz bei Einfacher Linearer Regression

$$Y_i = \beta_0^* + \beta_1(x_i - \bar{x}) + \epsilon_i, \quad \epsilon_i \text{ iid } N(0, \sigma^2), \quad i = 1, \dots, n$$

Die KQ-Schätzer  $b_0^*$  und  $b_1$  sind unabhängig.

$$b_0^* = \bar{Y} \sim N(\beta_0^*, \sigma^2/n), \quad b_1 \sim N(\beta_1, \sigma^2/SS_{xx})$$

Also ist

$$\frac{1}{\sigma^2} [n(\bar{Y} - \beta_0^*)^2 + (b_1 - \beta_1)^2 SS_{xx}] \sim \chi_2^2$$

und unabhängig von  $s^2$  (nach Satz 2.4), womit

$$W := \frac{(n(\bar{Y} - \beta_0^*)^2 + (b_1 - \beta_1)^2 SS_{xx})/2\sigma^2}{((n-2)s^2)/(n-2)\sigma^2} = \frac{n(\bar{Y} - \beta_0^*)^2 + (b_1 - \beta_1)^2 SS_{xx}}{2s^2} \sim F_{2, n-2}$$

also

$$P(n(\bar{Y} - \beta_0^*)^2 + (b_1 - \beta_1)^2 SS_{xx} \leq 2s^2 F_{1-\alpha; 2, n-2}) = 1 - \alpha$$

95 % - Konfidenz-Ellipse

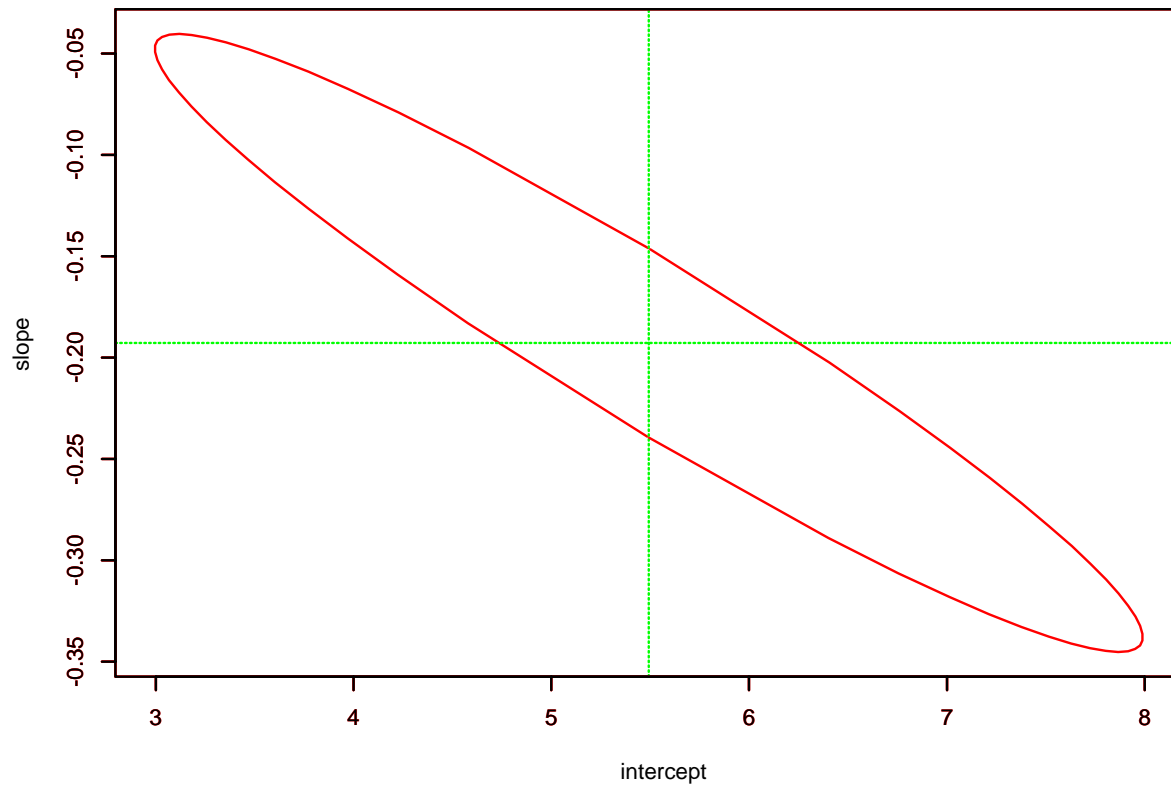


Abbildung 2.7: Simultane Inferenz

**Satz 2.9 (Simultane Konfidenzintervalle)**

Für  $\epsilon_i$  iid  $N(0, \sigma^2)$  erhält man die simultane  $(1 - \alpha)$ -Konfidenzregion für  $(\beta_0^*, \beta_1)$  aus:

$$E := \{(\beta_0^*, \beta_1) \mid n(\bar{Y} - \beta_0^*)^2 + (b_1 - \beta_1)^2 SS_{xx} \leq 2s^2 F_{1-\alpha; 2, n-2}\}$$

$E$  ist eine Ellipse für  $(\beta_0^*, \beta_1)$  mit Mittelpunkt  $(\bar{Y}, b_1)$ .

Die dazu äquivalente Konfidenz-Ellipse für  $(\beta_0, \beta_1)$  um den Mittelpunkt  $(\bar{Y} - b_1 \bar{x}, b_1)$  ergibt sich aus dem Zusammenhang  $\beta_0 = \beta_0^* - \beta_1 \bar{x}$ .

## 2.8 Gemeinsame Konfidenzbänder für den *mean response* nach Working-Hotelling

$E$  sei der  $100(1 - \alpha)\%$  Konfidenzbereich von Satz 2.9. Ein  $100(1 - \alpha)\%$ -Konfidenzband für

$$E(Y \mid x_0) = \beta_0 + \beta_1 x_0$$

ist gegeben durch

$$\begin{aligned} A(x_0) &:= \{\beta_0 + \beta_1 x_0 \mid (\beta_0, \beta_1) \in E\} \\ &= \left\{ \hat{Y}(x_0) \pm \sqrt{2F_{1-\alpha;2,n-2}} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_{xx}}} \right\} \end{aligned}$$

Dies ist eine andere (kleinere) Konstante verglichen mit den individuellen Konfidenzintervallen für  $E(Y \mid x_0)$ .

## 2.9 Bonferroni-Ansatz zu gemeinsamer Interferenz

Allgemein gilt:

Sei  $\theta$  ein Parameter der Verteilung von  $Y$  und sei  $[L(Y), U(Y)]$  ein  $(1-\alpha)$ -Konfidenzintervall für  $\theta$ , dann gilt

$$P(\underbrace{Y : \theta \in [L(Y), U(Y)]}_A) = 1 - \alpha.$$

Seien  $[L_1(Y), U_1(Y)]$   $(1-\alpha)$ -Konfidenzintervalle für  $\theta_1$   
 $[L_2(Y), U_2(Y)]$   $(1-\alpha)$ -Konfidenzintervalle für  $\theta_2$

und  $A_1, A_2$  entsprechend zu  $A$  (oben) definiert.

$\Rightarrow P(A_1) = 1 - \alpha = P(A_2)$  und

$$\begin{aligned} P(A_1 \cap A_2) &= 1 - P(A_1^c \cup A_2^c) \\ &= 1 - \underbrace{[P(A_1^c) + P(A_2^c) - P(A_1^c \cap A_2^c)]}_{\geq 0} \\ &\geq 1 - 2\alpha \end{aligned}$$

Damit ist ein  $100(1 - 2\alpha)\%$ -Konfidenzbereich für  $(\theta_1, \theta_2)$  gegeben durch

$$\{(\theta_1, \theta_2) : \theta_1 \in [L_1(Y), U_1(Y)], \theta_2 \in [L_2(Y), U_2(Y)]\}$$

Anwendung:

$$\{(\beta_0, \beta_1) : \beta_0 \in [b_0 \pm \hat{se}(b_0) t_{1-\frac{\alpha}{4}; n-2}], \beta_1 \in [b_1 \pm \hat{se}(b_1) t_{1-\frac{\alpha}{4}; n-2}]\}$$

ist ein  $100(1 - \alpha)\%$ -Konfidenzbereich für  $(\beta_0, \beta_1)$ .

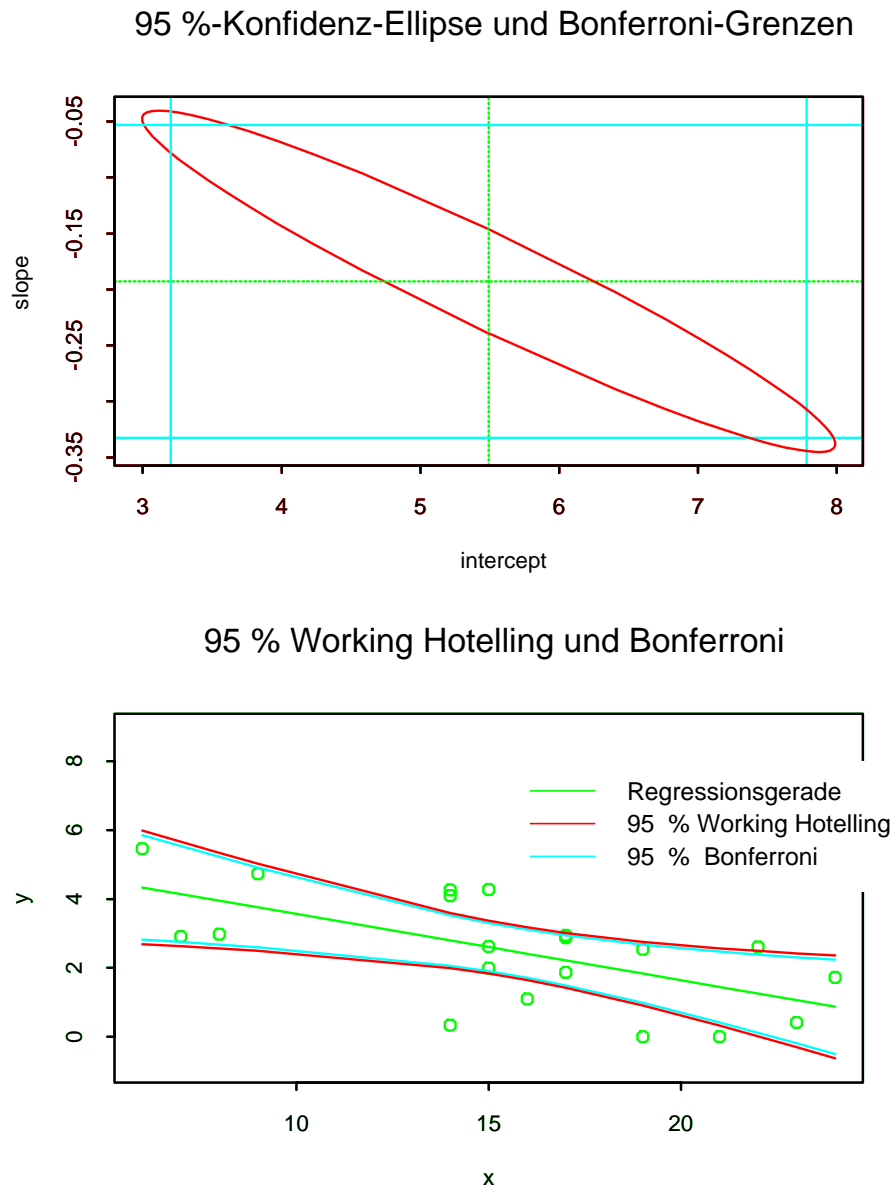


Abbildung 2.8: Vergleich der verschiedenen Methoden

# Kapitel 3

## Multiple Lineare Regression

### 3.1 Modell

Wir betrachten das Modell der Multiplen Linearen Regression für gegebene Daten  $(y_i, x_{i0}, x_{i1}, \dots, x_{ik})_{i=1, \dots, n}$ . In den meisten Fällen wählt man  $x_{i0} = 1$ . Damit erhält man den Ansatz

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad (3.1)$$

wobei  $E(\epsilon_i) = 0$  und  $Cov(\epsilon_i, \epsilon_j) = \delta_{ij} \sigma^2$ ,  $i = 1, \dots, n$ .

$\beta_0, \beta_1, \dots, \beta_k$  sind  $p = k+1$  unbekannte Parameter. Damit sind die  $y_i, i = 1, \dots, n$ , Zufallsvariable (ZV). In der Folge wird nicht mehr zwischen den ZV  $y_i$  und deren Realisierungen (Daten) unterschieden. Die  $x_{ij}, i = 1, \dots, n, j = 1, \dots, k$ , werden als Ausprägungen nicht stochastischer Kovariablen angesehen.

**Tabelle 1: Beispiel: College Spending 1994; vgl. [Dielman et al. (1996)], p. 202**

Variable	Abk.	Erklärung
SAT	S	average SAT score
TOP10	T	freshmen in the top 10% of their high school class (percentage)
ACCRATE	A	acceptance rate (percentage)
PHD	P	faculty with PHD (percentage)
RATIO	R	student faculty ratio
SPEND	S	educational spending per full-time equivalent (FTE) student (in dollars)
GRADRATE	G	Graduation Rate (percentage)
ALUMNI	Al	alumni giving rate (percentage)

SAT = Scholastic Aptitude Test (standardisierter Test in der letzten High School Klasse (= K12 in D) für den Hochschulzugang)

Auf dieses Beispiel werden wir in diesem Abschnitt bis hin zur Variablenselektion zurückgreifen.

Seien nun  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ ,  $\mathbf{1} = (1, \dots, 1)^T$ ,



$$\mathbf{X} = (\vec{x}_0, \vec{x}_1, \dots, \vec{x}_k) = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \cdot \\ \cdot \\ \mathbf{x}_n^T \end{pmatrix}, \quad (3.2)$$

wobei in der Regel, aber nicht unbedingt:  $\vec{x}_0 = \mathbf{1}$ .

Das Modell lässt sich also folgendermaßen schreiben:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad E(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}, \quad (3.3)$$

$\mathbf{y} \in \mathbb{R}^n$ : stochastische Zielvariable (response)

$\boldsymbol{\beta} \in \mathbb{R}^p$ : unbekannter Parametervektor

$\mathbf{X} \in \mathbb{R}^{n \times p}$ : Designmatrix, Regressormatrix (feste, nicht stochastische Daten)  
(fast immer:  $\text{Rang}(\mathbf{X}) = p$ )

$\boldsymbol{\epsilon} \in \mathbb{R}^n$ : un beobachtbare Fehlervariable

Oft gilt zusätzlich:  $\epsilon_i$  normalverteilt,

d.h.  $\epsilon_i$  iid  $N(0, \sigma^2)$ ,  $i = 1, \dots, n$ .

Im College-Spending-Beispiel ist man am Response *SPEND* interessiert.

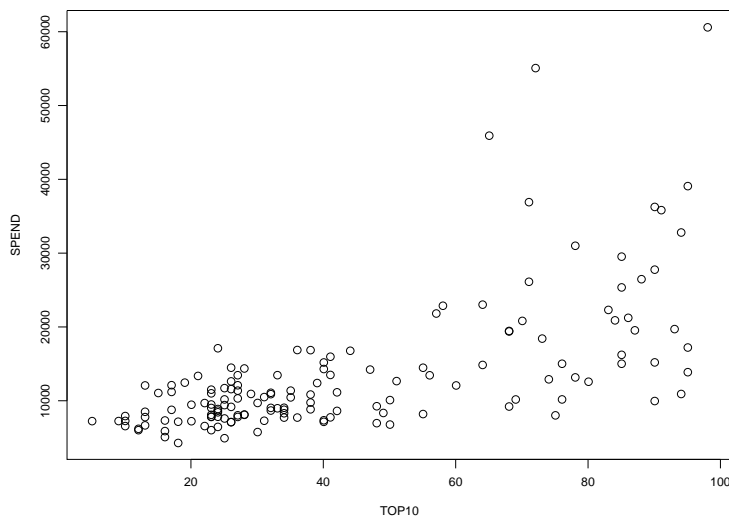
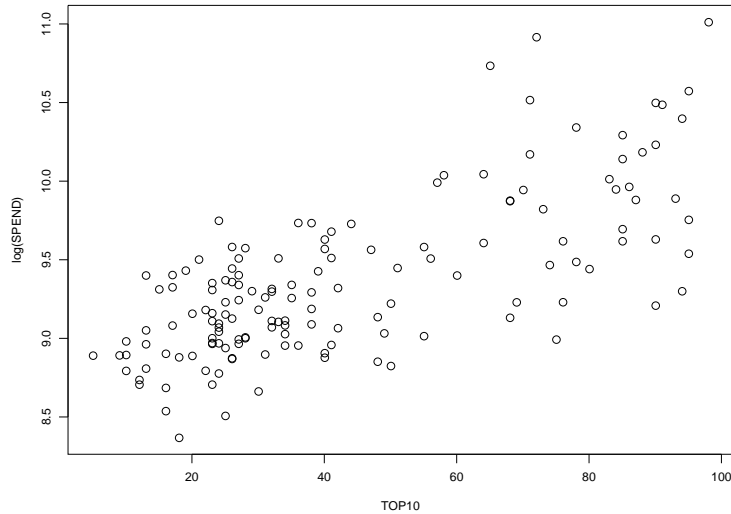
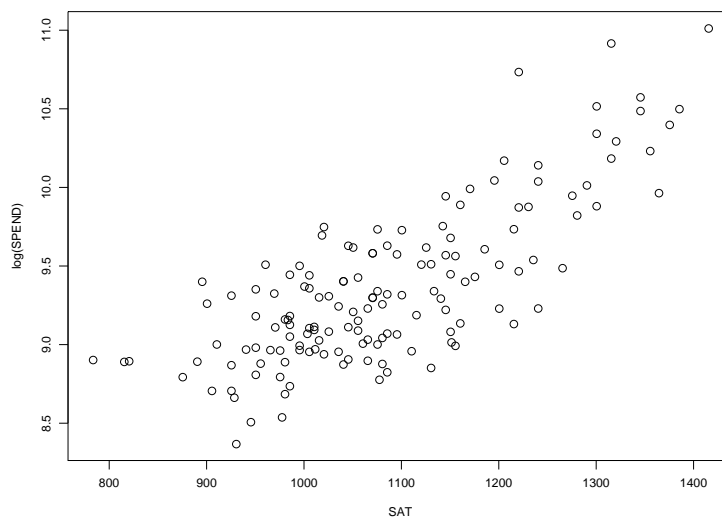


Abbildung 3.1: Scatterplot von SPEND gegen die Kovariable TOP10

An den TOP10-Werten über 60 sieht man, dass die Annahme konstanter Varianz nicht haltbar ist. Dieses Bild verbessert sich wesentlich, wenn man die logarithmierte Zielvariable betrachtet. Wir arbeiten also fortan mit dem Response  $\log(\text{SPEND})$ .

Ähnlich reagieren auch die Plots gegenüber der Kovariablen *SAT*.

Abbildung 3.2: Scatterplot von  $\log(\text{SPENDING})$  gegen die Kovariable TOP10Abbildung 3.3: Scatterplot von  $\log(\text{SPENDING})$  gegen die Kovariable SAT

Diese "Pairs-Plots" gegenüber potenziellen Kovariablen rechtfertigen den linearen (und additiven) Modellansatz noch nicht vollständig. Sie könnten diesen aber widerlegen. Dann müsste man noch weitere Transformationen von Kovariablen bzw. Response suchen.

Bei Prognosen und Vorhersageintervallen ist in diesem Beispiel natürlich zu berücksichtigen, dass nicht mit den Originaldaten gearbeitet wird. Rücktransformationen ergeben sich hier mit der Exponentialfunktion und allgemein aus der entsprechenden Umkehrabbildung.

## 3.2 KQ-Methode (LS estimation)

Wir setzen für das Folgende stets  $\text{Rang}(\mathbf{X}) = p$  voraus.

Damit erhält man die KQ-Schätzungen  $\mathbf{b} = (b_0, b_1, \dots, b_k)^T$  für  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)^T$  aus den Normalgleichungen

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y} \quad (3.4)$$

durch

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.5)$$

Mit der sog. "Hat-Matrix"

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (3.6)$$

lautet die Schätzung  $\hat{\mathbf{y}}$  für  $\mathbf{y}$

$$\hat{\mathbf{y}} = \mathbf{H} \mathbf{y} = \mathbf{X} \mathbf{b}. \quad (3.7)$$

Damit ergibt sich der Residuenvektor zu

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y}, \quad (3.8)$$

wobei

$$SS_{Res} = \mathbf{e}^T \mathbf{e} = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3.9)$$

## 3.3 Eigenschaften der KQ-Schätzungen

### Satz 3.1 (Gauß-Markov-Theorem)

Vor.:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,  $\text{Rang}(\mathbf{X}) = p$ ,  $E(\boldsymbol{\epsilon}) = \mathbf{0}$ ,  $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$

Beh.:

1.  $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  ist BLUE (*best linear unbiased estimator*) für  $\boldsymbol{\beta}$
2.  $\text{Cov}(\mathbf{b}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$
3.  $s^2 = SS_{Res} / (n - p)$  ist unverzerrter Schätzer für  $\sigma^2$

### 3.4 Projektionen und $\chi^2$ -Verteilung

Dieser Abschnitt dient zur Vorbereitung der Verteilungsaussagen für Schätzer.

#### Lemma 3.2 (Rechnen mit Kovarianzen (Wiederholung))

Wir betrachten die Zufallsvariablen  $V, W$  mit existierenden 2. Momenten und

$$E(V) = \mu_V, \quad E(W) = \mu_W, \quad \text{Var}(V) = \sigma_V^2, \quad \text{Var}(W) = \sigma_W^2,$$

$$\text{cov}(V, W) := E[(V - \mu_V)(W - \mu_W)] = \rho\sigma_V\sigma_W.$$

Weiterhin seien

$$\mathbf{Z} := (V, W)^T, \quad \boldsymbol{\mu}_Z = (\mu_V, \mu_W)^T$$

und

$$\text{Cov}(\mathbf{Z}) := E[(\mathbf{Z} - \boldsymbol{\mu}_Z)(\mathbf{Z} - \boldsymbol{\mu}_Z)^T].$$

Dann gilt mit  $\mathbf{a} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2,2}$

1.  $\text{Var}(V) = \text{cov}(V, V)$
2.  $\text{Cov}(\mathbf{Z}) = \begin{pmatrix} \sigma_V^2 & \rho\sigma_V\sigma_W \\ \rho\sigma_V\sigma_W & \sigma_W^2 \end{pmatrix}$
3.  $\text{Var}(\mathbf{a}^T \mathbf{Z}) = \text{Cov}(\mathbf{a}^T \mathbf{Z}) = \mathbf{a}^T \text{Cov}(\mathbf{Z}) \mathbf{a}$
4.  $\text{Cov}(\mathbf{A} \mathbf{Z}) = \mathbf{A} \text{Cov}(\mathbf{Z}) \mathbf{A}^T.$

Die Aussagen von 3. und 4. gelten allgemein für  $\mathbf{Z}, \mathbf{a} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{n,n}$ .

#### Def. 3.3 ( $\chi^2$ -Verteilung)

Seien  $V_i$  iid  $N(0, 1)$ ,  $i = 1, \dots, r$ .

Dann heißt

$$W^2 = \sum_{i=1}^r V_i^2 \sim \chi_r^2,$$

$\chi^2$ -verteilt mit  $r$  Freiheitsgraden.

#### Satz 3.4 (Additionsformel)

Sind  $W_1^2$  bzw.  $W_2^2$  unabhängig  $\chi_{r_1}^2$  bzw.  $\chi_{r_2}^2$ , dann ist

$$W^2 := W_1^2 + W_2^2 \sim \chi_{r_1+r_2}^2$$

**Satz 3.5 (Satz von Cochran)**

Vor.:  $\mathbf{w} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$   
 $\mathbf{P} \in \mathbb{R}^{n,n}$ ,  $\mathbf{P}^T = \mathbf{P}$ ,  $\mathbf{P}^2 = \mathbf{P}$ ,  $\text{Rang}(\mathbf{P}) = r$   
 Beh.:  $W_1^2 := \mathbf{w}^T \mathbf{P} \mathbf{w} / \sigma^2 \sim \chi_r^2$   
 $W_2^2 := \mathbf{w}^T (\mathbf{I} - \mathbf{P}) \mathbf{w} / \sigma^2 \sim \chi_{n-r}^2$   
 $W_1^2$  und  $W_2^2$  sind unabhängig.

Beweis:

Die Komponenten  $w_i$  von  $\mathbf{w}$  sind unabhängig, also auch die Komponenten  $z_i$  von  $\mathbf{z} = \mathbf{A} \mathbf{w}$ , für orthogonale Matrizen  $\mathbf{A}$ . Insbesondere gilt dann  $\mathbf{z} / \sigma \sim N_n(\mathbf{0}, \mathbf{I})$ .

Die Projektionsmatrix  $\mathbf{P}$  ist symmetrisch und hat  $r$  Eigenwerte 1 und  $n - r$  Eigenwerte 0. Insbesondere gibt es eine Orthonormalbasis  $\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_n$ , Matrizen  $\mathbf{U}_r = (\mathbf{u}_1, \dots, \mathbf{u}_r)$ ,  $\mathbf{U}_{n-r} = (\mathbf{u}_{r+1}, \dots, \mathbf{u}_n)$ ,  $\mathbf{U} = (\mathbf{U}_r, \mathbf{U}_{n-r}) \in \mathbb{R}^{n,n}$ ,  $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}$ ,

so dass gilt:

$$\mathbf{P} = \mathbf{U}_r \mathbf{U}_r^T, \quad \mathbf{I} - \mathbf{P} = \mathbf{U}_{n-r} \mathbf{U}_{n-r}^T, \quad \mathbf{U}_r^T \mathbf{U}_{n-r} = \mathbf{0}.$$

Mit  $\mathbf{A} := \mathbf{U}^T$  und  $\mathbf{z} := \mathbf{U}^T \mathbf{w} = \begin{pmatrix} \mathbf{U}_r^T \mathbf{w} \\ \mathbf{U}_{n-r}^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} \mathbf{z}_r \\ \mathbf{z}_{n-r} \end{pmatrix}$

sind alle Komponenten von  $\mathbf{z}$  unabhängig, insbesondere auch  $\mathbf{z}_r$  und  $\mathbf{z}_{n-r}$ .

Außerdem:

$$W_1^2 = \mathbf{w}^T \mathbf{P} \mathbf{w} / \sigma^2 = \mathbf{z}_r^T \mathbf{z}_r \sim \chi_r^2 \text{ ist unabhängig von}$$

$$W_2^2 = \mathbf{w}^T (\mathbf{I} - \mathbf{P}) \mathbf{w} / \sigma^2 = \mathbf{z}_{n-r}^T \mathbf{z}_{n-r} \sim \chi_{n-r}^2 \quad \square$$

**Satz 3.6 (Verteilung von  $\mathbf{b}$  und  $s^2$ )**

Vor.: Seien  $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,  $\mathbf{X} \in \mathbb{R}^{n,p}$ ,  $\text{Rang}(\mathbf{X}) = p$ ,  
 $\epsilon_i$  iid  $N(0, \sigma^2)$ ,  $i = 1, \dots, n$   
 $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ ,  $SS_{Res} = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}$ .

Beh.:

1.  $\mathbf{b} \sim N_p(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$
2.  $Q_{\boldsymbol{\beta}} := (\mathbf{b} - \boldsymbol{\beta})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \boldsymbol{\beta}) / \sigma^2 \sim \chi_p^2$
3.  $\mathbf{b}$  ist unabhängig von  $s^2 = \hat{\sigma}^2 = \frac{SS_{Res}}{n - p}$
4.  $SS_{Res} / \sigma^2 \sim \chi_{n-p}^2$
5.  $Q_{\boldsymbol{\beta}}$  und  $SS_{Res}$  sind unabhängig.

Beweis:

1.  $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{A} \mathbf{y}$  ist als lineare Transformation von  $\mathbf{y}$  normalverteilt. Erwartungswert und Kovarianz folgen aus dem Gauß-Markov-Theorem (Satz 3.1).

2.,4.,5. mit Satz von Cochran:

Mit  $\mathbf{w} := \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$  und  $\mathbf{P} := \mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  sind die Voraussetzungen des Satzes von Cochran (Satz 3.5) erfüllt und es gilt mit  $\mathbf{H}^T = \mathbf{H}$ ,  $\mathbf{H}^2 = \mathbf{H}$  ( $\mathbf{I} - \mathbf{H}$  analog) und  $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ :

$$\begin{aligned} W_1^2 &:= \mathbf{w}^T \mathbf{P} \mathbf{w} / \sigma^2 = \mathbf{w}^T \mathbf{H} \mathbf{H} \mathbf{w} / \sigma^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) / \sigma^2 \\ &= (\mathbf{b}^T - \boldsymbol{\beta}^T) \mathbf{X}^T \mathbf{X} (\mathbf{b} - \boldsymbol{\beta}) / \sigma^2 \\ &= Q_{\boldsymbol{\beta}} \sim \chi_p^2 \end{aligned}$$

ist unabhängig von

$$\begin{aligned} W_2^2 &:= \mathbf{w}^T (\mathbf{I} - \mathbf{P}) \mathbf{w} / \sigma^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) / \sigma^2 \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y} / \sigma^2 \\ &= SS_{Res} / \sigma^2 \sim \chi_{n-p}^2, \end{aligned}$$

wobei  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$  benutzt wurde.

3. Die Unabhängigkeit von  $\mathbf{b}$  und  $SS_{Res}$  zeigt man mit Techniken analog zum Beweis des Satzes von Cochran.

□

### 3.5 Varianzzerlegung, Bestimmtheitsmaß

Betrachte  $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1) \in \mathbb{R}^{n,p}$ ,  $\text{Rang}(\mathbf{X}) = p$ ,  
 $\mathbf{1} = (1, \dots, 1)^T$ ,  $\mathbf{H}_0 = \mathbf{1} \cdot \mathbf{1}^T / n$

#### Lemma 3.7

$$SS_{Total} = \sum_{i=1}^n (y_i - \bar{y})^2 = \mathbf{y}^T (\mathbf{I} - \mathbf{H}_0) \mathbf{y}$$

Beweis: Übung

#### Lemma 3.8

Mit  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  gilt:

1.  $\mathbf{H}\mathbf{H}_0 = \mathbf{H}_0\mathbf{H} = \mathbf{H}_0$
2.  $\mathbf{H} - \mathbf{H}_0$  ist Orthogonalprojektion.

Beweis: Übung

**Satz 3.9**

$$\text{Vor.: } SS_{Reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \quad SS_{Res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{Beh.: } SS_{Total} = SS_{Reg} + SS_{Res}.$$

Beweisskizze:

$$SS_{Res} = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{H} \mathbf{y} \Rightarrow SS_{Reg} = \mathbf{y}^T (\mathbf{H} - \mathbf{H}_0) \mathbf{y}$$

und

$$\mathbf{I} - \mathbf{H} = \mathbf{I} - \mathbf{H}_0 - (\mathbf{H} - \mathbf{H}_0),$$

Rest Übung.

**Def. 3.10 (Bestimmtheitsmaß,  $R^2$ )**

Mit den Bezeichnungen von 3.9

$$SS_{Total} = \sum_{i=1}^n (y_i - \bar{y})^2, \quad SS_{Reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \quad SS_{Res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

und unter der Vor.  $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1)$  heißen

$$R^2 := \frac{SS_{Reg}}{SS_{Total}} = 1 - \frac{SS_{Res}}{SS_{Total}} \quad \text{Bestimmtheitsmaß,}$$

$R$  multipler Korrelationskoeffizient;

und mit  $s_0^2 := SS_{Total}/(n-1)$ ,  $s^2 = SS_{Res}/(n-p)$ :

$$R_{adj}^2 := 1 - \frac{s^2}{s_0^2} = 1 - (1 - R^2) \frac{n-1}{n-p} \quad \text{adjustiertes Bestimmtheitsmaß.}$$

**Bemerkung 3.11**

1.  $0 \leq R^2 \leq 1$
2.  $R^2 = 1 \Leftrightarrow SS_{Res} = 0$
3.  $R^2(\mathbf{X}_0, \mathbf{X}_1) \geq R^2(\mathbf{X}_0)$  für Zerlegungen  $\mathbf{X} = (\mathbf{X}_0, \mathbf{X}_1)$
4.  $R_{adj}^2$  kann für ein Modell mit mehr Variablen fallen (wenn  $SS_{Reg}$  kaum mehr wächst)

**Bemerkung 3.12 (Varianzzerlegung)**

3.9 ist Spezialfall folgenden Prinzips:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbf{X} = (\mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2) \text{ mit } \mathbf{X}_0 = \mathbf{1}, \quad \mathbf{X}_1 \in \mathbb{R}^{n,p_1}, \quad \mathbf{X}_2 \in \mathbb{R}^{n,p_2}$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}, \quad p = 1 + p_1 + p_2, \quad \text{Rang}(\mathbf{X}) = p$$

$$1. \quad \mathbf{y} = \mathbf{X}_0\beta_0 + \boldsymbol{\epsilon}_0 = \mathbf{1} \cdot \beta_0 + \boldsymbol{\epsilon}_0 \text{ (evtl. "underfitting")}$$

$$\mathbf{H}_0 = \mathbf{X}_0(\mathbf{X}_0^T \mathbf{X}_0)^{-1} \mathbf{X}_0^T = \mathbf{1} \cdot \mathbf{1}^T / n$$

$$SS_{Res}(\mathbf{X}_0) = SS_{Total} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_0)\mathbf{y}$$

$$2. \quad \mathbf{y} = \mathbf{X}_0\beta_0 + \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}_1$$

$$\mathbf{X}_{01} = (\mathbf{X}_0, \mathbf{X}_1)$$

$$\mathbf{H}_{01} = \mathbf{X}_{01}(\mathbf{X}_{01}^T \mathbf{X}_{01})^{-1} \mathbf{X}_{01}^T$$

$$SS(\mathbf{X}_1 | \mathbf{X}_0) := SS_{Reg}(\mathbf{X}_{01}) = SS_{Res}(\mathbf{X}_0) - SS_{Res}(\mathbf{X}_{01}) =$$

$$= \mathbf{y}^T(\mathbf{H}_{01} - \mathbf{H}_0)\mathbf{y}$$

$$3. \quad \mathbf{y} = \mathbf{X}_0\beta_0 + \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon} \text{ (evtl. "overfitting")}$$

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$SS(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_0) := SS_{Reg}(\mathbf{X})$$

$$SS(\mathbf{X}_2 | \mathbf{X}_{01}) := SS_{Res}(\mathbf{X}_{01}) - SS_{Res}(\mathbf{X}) = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_{01})\mathbf{y} - \mathbf{y}^T(\mathbf{I} - \mathbf{H})\mathbf{y} =$$

$$\mathbf{y}^T(\mathbf{H} - \mathbf{H}_{01})\mathbf{y}$$

$$SS_{Res}(\mathbf{X}) = \mathbf{y}^T(\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_0 + \mathbf{H}_0 - \mathbf{H}_{01} + \mathbf{H}_{01} - \mathbf{H})\mathbf{y} =$$

$$\mathbf{y}^T(\mathbf{I} - \mathbf{H}_0)\mathbf{y} - \mathbf{y}^T(\mathbf{H} - \mathbf{H}_{01})\mathbf{y} - \mathbf{y}^T(\mathbf{H}_{01} - \mathbf{H}_0)\mathbf{y} =$$

$$SS_{Total} - SS(\mathbf{X}_1 | \mathbf{X}_0) - SS(\mathbf{X}_2 | \mathbf{X}_{01})$$

Damit hat man folgende allgemeinere Varianzzerlegung:

$$SS_{Total} = \quad SS(\mathbf{X}_1 | \mathbf{X}_0) \quad \text{durch } \mathbf{X}_1 \text{ erklärte Varianz}$$

$$+ SS(\mathbf{X}_2 | \mathbf{X}_{01}) \quad \text{durch } \mathbf{X}_2 \text{ zusätzlich erklärte Varianz}$$

$$+ SS_{Res}(\mathbf{X}) \quad \text{von } \mathbf{X} \text{ nicht erklärte Restvarianz}$$

Wegen  $SS_{Total} = SS(\mathbf{X} | \mathbf{X}_0) + SS_{Res}(\mathbf{X})$  folgt durch Vertauschung der Rollen von  $\mathbf{X}_1$  und  $\mathbf{X}_2$

$$SS(\mathbf{X} | \mathbf{X}_0) = SS(\mathbf{X}_1 | \mathbf{X}_0) + SS(\mathbf{X}_2 | \mathbf{X}_{01})$$

$$= SS(\mathbf{X}_2 | \mathbf{X}_0) + SS(\mathbf{X}_1 | \mathbf{X}_{02}).$$

Es kommt also auf die Reihenfolge an, in welcher die Variablen ins Modell kommen.



Das Beispiel College Spending liefert beispielsweise:

```
r1 <- lm(log(SPEND TOP10 + ACCRATE + RATIO + GRADRATE))
anova(r1)
```

Analysis of Variance Table

Response: log(SPEND)

Terms added sequentially (first to last)

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
TOP10	1	18.78383	18.78383	274.3323	0.0000000
ACCRATE	1	1.32889	1.32889	19.4080	0.0000207
RATIO	1	6.79686	6.79686	99.2661	0.0000000
GRADRATE	1	0.15791	0.15791	2.3063	0.1310776
Residuals	142	9.72289	0.06847		

Die Spalte Sum of Sq gibt die sequential Sums of Squares. Sie hängen von der Reihenfolge der Variablen ab, mit der sie in das Modell eingehen. Dies kann man wie folgt sehen:

```
r2 <- lm(log(SPEND) RATIO + GRADRATE + TOP10 + ACCRATE)
anova(r2)
```

Analysis of Variance Table

Response: log(SPEND)

Terms added sequentially (first to last)

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
RATIO	1	15.86151	15.86151	231.6528	0.00000000
GRADRATE	1	6.55064	6.55064	95.6702	0.00000000
TOP10	1	4.46617	4.46617	65.2271	0.00000000
ACCRATE	1	0.18915	0.18915	2.7625	0.09870084
Residuals	142	9.72289	0.06847		

Die so von einem Statistik-Programm-Paket erzeugten Resultate liefern die numerische Grundlage für wichtige Spezialfälle der folgenden Tests.

### 3.6 Tests linearer Hypothesen

Wichtige Hypothesen (z.B. bei der Variablenauswahl) haben etwa die Form

$$\left. \begin{array}{l} H_0 : \beta_2 = 0 \\ \beta_4 = 0 \end{array} \right\} \text{ zugehörige x-Variablen können aus dem Modell entfernt werden}$$

Alternative  $H_1$  : nicht  $H_0$ .

### 3.6.1 Test zur allgemeinen linearen Hypothese

#### Def. 3.13 (Allgemeine lineare Hypothese)

Geg.:  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,  $\text{Rang}(\mathbf{X}) = p$ ,  $\mathbf{X} \in \mathbb{R}^{n,p}$   
 $\mathbf{C} \in \mathbb{R}^{r,p}$ ,  $\text{Rang}(\mathbf{C}) = r$ ,  $\mathbf{d} \in \mathbb{R}^r$

Die Restriktion

$$H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{d} \quad (3.10)$$

heißt allgemeine lineare Hypothese.

Die Alternative lautet stets:  $H_1 : \text{nicht } H_0$ .

#### Satz 3.14 (Restringiertes KQ-Problem)

Unter den Bedingungen von 3.13 gilt mit

$$SS_{Res} = \min_{\boldsymbol{\beta}} \{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2\} = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2$$

$$SS_{H_0} = \min_{\boldsymbol{\beta}} \{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \mid \underbrace{\mathbf{C}\boldsymbol{\beta} = \mathbf{d}}_{\text{sprich: unter der Bed. } H_0}\}$$

$$SS_{H_0} = SS_{Res} + (\mathbf{C}\mathbf{b} - \mathbf{d})^T [\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T]^{-1} (\mathbf{C}\mathbf{b} - \mathbf{d})$$

Beweis: Lagrange-Ansatz und Differenzieren

#### Beispiel 3.15

$$\mathbf{X} = \underbrace{(\mathbf{1}, \vec{\mathbf{x}}_2)}_{\mathbf{X}_1}, \underbrace{(\vec{\mathbf{x}}_1, \vec{\mathbf{x}}_3)}_{\mathbf{X}_2}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \beta_3)^T$$

$$H_0 : \beta_1 = \beta_3 = 0 \Leftrightarrow \underbrace{\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{C} \in \mathbb{R}^{2,4}, r=2} \cdot \boldsymbol{\beta} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \mathbf{d}$$

Falls  $H_0$  richtig ist, gilt das kleinere Modell

$$\mathbf{y} = (\mathbf{1}, \vec{\mathbf{x}}_2) \begin{pmatrix} \beta_0 \\ \beta_2 \end{pmatrix} + \boldsymbol{\epsilon}$$

#### Satz 3.16 (F-Test zur linearen Hypothese)

Vor.:  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,  $\epsilon_i \text{ iid } N(0, \sigma^2)$

$\text{Rang}(\mathbf{X}) = p$ ,  $\mathbf{X} \in \mathbb{R}^{n,p}$ ,  $\mathbf{C} \in \mathbb{R}^{r,p}$ ,  $\text{Rang}(\mathbf{C}) = r$ ,  $\mathbf{d} \in \mathbb{R}^r$

$SS_{H_0}$  aus 3.14

Beh: Falls die Hypothese

$$H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{d}$$

richtig ist, gilt

$$F = \frac{(SS_{H_0} - SS_{Res})/r}{SS_{Res}/(n-p)} \sim F_{r, n-p}.$$

Testvorschrift:

$H_0$  ist beim Niveau  $\alpha$ ,  $0 < \alpha < 1$ , z.B.  $\alpha = 0.05$  oder  $\alpha = 0.01$  abzulehnen, falls

$$F > F_{1-\alpha;r,n-p}$$

Häufigste Anwendung:

### 3.6.2 Tests für Variablen-Subsets

$$\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2), \mathbf{X}_1 \in \mathbb{R}^{n,p-r}, \mathbf{X}_2 \in \mathbb{R}^{n,r}$$

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}, \boldsymbol{\beta}_1 \in \mathbb{R}^{p-r}, \boldsymbol{\beta}_2 \in \mathbb{R}^r, \text{Rang}(\mathbf{X}) = p$$

$$H_0 : \boldsymbol{\beta}_2 = \mathbf{0}, \mathbf{C} = (\mathbf{0}, \mathbf{I}_r)$$

$$SS_{Res} = SS_{Res}(\mathbf{X}) = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}$$

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \mathbf{H}_1 = \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$$

$$SS_{H_0} = SS_{Res}(\mathbf{X}_1) = \mathbf{y}^T (\mathbf{I} - \mathbf{H}_1) \mathbf{y}$$

Damit kann

$$F = \frac{(SS_{Res}(\mathbf{X}_1) - SS_{Res})/r}{SS_{Res}/(n-p)}$$

ohne die umständliche Formel von 3.14 leicht - insbesondere mit jedem gängigen Statistik-Programmpaket bestimmt werden.

Hilfsmittel in S-Plus: `anova(r.red, r.full)`, wobei `r.full <- lm(y ~ x1 + ... + xp)`  
`r.red <- lm(y ~ x1 + ... + x(p-r))`

Noch spezieller und wichtiger:

### 3.6.3 Partielle F-Tests für einzelne Parameter

Hypothese  $H_j : \beta_j = 0$ , für ein  $j \in \{1, \dots, p\}$

$$F_j = \frac{(SS_{H_j} - SS_{Res})/1}{SS_{Res}/(n-p)} \sim F_{1,n-p}$$

Im Zähler: Vergrößerung der Residuenquadratsumme beim Weglassen der Variablen  $j$ .

Bekanntlich:  $F_j = T_j^2$ , wobei

$$T_j = \frac{b_j}{se(b_j)} \sim t_{n-p}$$

analog zur Einfachen Linearen Regression.

$b_j$ , der zugehörige Standardfehler (standard error)  $se(b_j)$  und  $T_j$  bzw.  $F_j$  finden sich im Ausdruck jedes guten Statistik-Programmpakets.

### 3.7 Konfidenz- und Prognoseintervalle

Seien

$$\mathbf{x}_0 = (1, x_{01}, x_{02}, \dots, x_{0,p-1})^T \in \mathbb{R}^p$$

$$y(\mathbf{x}_0) \sim N(\mathbf{x}_0^T \boldsymbol{\beta}, \sigma^2) \text{ unabhängig von } \mathbf{y} = (y_1, \dots, y_n)^T$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n), \mathbf{X} \in \mathbb{R}^{n,p}, \text{Rang}(\mathbf{X}) = p$$

$$\hat{y}(\mathbf{x}_0) = \mathbf{x}_0^T \mathbf{b}, \mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$s^2 = SS_{Res}/(n-p)$$

#### Satz 3.17 (Konfidenz- und Prognoseintervalle)

1.  $E[\hat{y}(\mathbf{x}_0)] = \mathbf{x}_0^T \boldsymbol{\beta} = E[y(\mathbf{x}_0)]$   
 $\hat{y}(\mathbf{x}_0) = \mathbf{x}_0^T \mathbf{b}$ : berechenbar  
 $y(\mathbf{x}_0)$ : tritt erst (in der Zukunft) ein
2.  $\text{Var}[\hat{y}(\mathbf{x}_0)] = \text{Var}[\mathbf{x}_0^T \mathbf{b}] = \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0$
3.  $\text{Var}[y(\mathbf{x}_0) - \hat{y}(\mathbf{x}_0)] = \sigma^2 (1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0)$
4.  $(1 - \alpha)$ -Konfidenz-Intervall für den "mean response"  
 $v := t_{1-\frac{\alpha}{2}; n-p} \cdot s \cdot \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}, 0 < \alpha < 1$   
 $P(\hat{y}(\mathbf{x}_0) - v \leq E[y(\mathbf{x}_0)] = \mathbf{x}_0^T \boldsymbol{\beta} \leq \hat{y}(\mathbf{x}_0) + v) = 1 - \alpha$
5.  $(1 - \alpha)$ -Prognose-Intervall  
 $w := t_{1-\frac{\alpha}{2}; n-p} \cdot s \cdot \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}, 0 < \alpha < 1$   
 $P(\hat{y}(\mathbf{x}_0) - w \leq y(\mathbf{x}_0) \leq \hat{y}(\mathbf{x}_0) + w) = 1 - \alpha$
6. Für  $\mathbf{1} = (1, \dots, 1)^T$ ,  $\mathbf{x} = (x_1, \dots, x_n)^T$ ,  $\mathbf{x}_0 = (1, x_0)^T$  und  $\mathbf{X} = (\mathbf{1}, \mathbf{x})$  erhält man die Aussagen von Abschnitt 2.4.4

## 3.8 Qualität des Modells, Hat-Matrix

### 3.8.1 Fragen

1. Ist das gewählte Modell richtig?  
Wurden die Daten so *transformiert* und die möglichen Variablen so ausgewählt, dass der lineare Modellansatz gerechtfertigt ist? Gibt es Ausreißer unter den Beobachtungen?
2. Wie lautet ein optimales Modell hinsichtlich einer möglichst kleinen Varianz des Vorhersagefehlers  $Var[y(\mathbf{x}_0) - \hat{y}(\mathbf{x}_0)]$  ?

### 3.8.2 Ansätze

Mit  $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$  gilt für den Vorhersagefehler

$$d_i^2 = Var(y(\mathbf{x}_i) - \hat{y}(\mathbf{x}_i)) = \sigma^2(1 - \underbrace{\mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i}_{h_{ii}})$$

$$\hat{d}_i^2 = s^2(1 - h_{ii}), \quad h_{ii} \text{ Diagonalelemente der Hat-Matrix } \mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

Die Qualität der Beobachtung  $i$  hängt ab von  $h_{ii}$  und  $s^2 = \frac{1}{n-p} \sum_{i=1}^n [y_i - \hat{y}(\mathbf{x}_i)]^2$ .

#### Satz 3.18 (Eigenschaften von $\mathbf{H}$ )

1.  $\mathbf{H}$  ist Projektionsmatrix, d.h.  $\mathbf{H}^T = \mathbf{H}$ ,  $\mathbf{H}^2 = \mathbf{H}$  (Orthogonalprojektion)
2.  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$
3.  $\mathbf{e} = (\mathbf{I} - \mathbf{H})\mathbf{y}$ ,  $SS_{Res} = \mathbf{e}^T \mathbf{e} = \mathbf{y}^T (\mathbf{I} - \mathbf{H})\mathbf{y}$
4.  $0 \leq h_{ii} \leq 1$
5.  $\frac{1}{n} \leq h_{ii} \leq 1$  (für Modelle mit intercept)

Konsequenzen:

Wegen  $Rang(\mathbf{H}) = tr(\mathbf{H}) = \sum_{i=1}^n h_{ii} = p$  sollten die  $h_{ii}$  idealerweise  $\approx p/n$  sein.

#### Def. 3.19 (High Leverage Points)

Beobachtungen  $i$  mit  $h_{ii} > 2p/n$  heißen **Hebelpunkte** oder **high leverage points**.

**Def. 3.20 (Jackknife Methode)**

Sei  $\mathbf{b}$  der KQ-Schätzer zum vollen Modell

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad i = 1, \dots, n$$

und  $\mathbf{b}_{(j)}$  der KQ-Schätzer zum reduzierten Modell

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad i = 1, \dots, n, \quad i \neq j$$

mit  $n-1$  Beobachtungen.

**Def. 3.21 (Cook's Distance)**

Mit den Bezeichnungen von 3.20 heißt

$$D_j := \frac{(\mathbf{b} - \mathbf{b}_{(j)})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \mathbf{b}_{(j)})}{ps^2}$$

Cook's Distance der Beobachtung  $j$ .

**Bemerkung 3.22**

1.  $D_j$  ist Standardisierung von  $\mathbf{b} - \mathbf{b}_{(j)}$ , vgl. Satz 3.6

2.  $D_j$  ist einfach zu berechnen. Es gilt:

$$D_j = \frac{e_j^2}{(1 - h_{jj})^2} \cdot \frac{h_{jj}}{s^2 \cdot p}$$

3.  $D_j$  ist groß, falls  $h_{jj} \approx 1$  oder das  $j$ -te Residuum  $e_j = y_j - \hat{y}(\mathbf{x}_j)$  betragsmäßig groß ist.

4.  $D_j$  misst die Verschiebung der Konfidenz-Region für  $\boldsymbol{\beta}$ , wenn die  $j$ -te Beobachtung weggelassen wird.

5.  $D_j \sim F_{p, n-p}$

6. Mit  $s^2$  anstelle von  $\sigma^2$  in Satz 3.6 erhält man durch

$$\left\{ \boldsymbol{\beta} : \frac{(\mathbf{b} - \boldsymbol{\beta})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \boldsymbol{\beta})}{ps^2} \leq F_{1-\alpha; p, n-p} \right\}$$

eine  $(1 - \alpha)$ -Konfidenzregion für  $\boldsymbol{\beta}$ .

**Def. 3.23 (Influential observation)**

Beobachtungen mit

$$D_j > F_{0.5; p, n-p}$$

heißen **einflussreiche Beobachtungen (influential observations)**.

Achtung, nicht immer gilt: "influential observation" = "Ausreißer"!

Ein Beispiel zur Demonstration: Die Beobachtungen mit den Nummern 1, 2 und 10 fallen durch ihre extreme Lage auf.

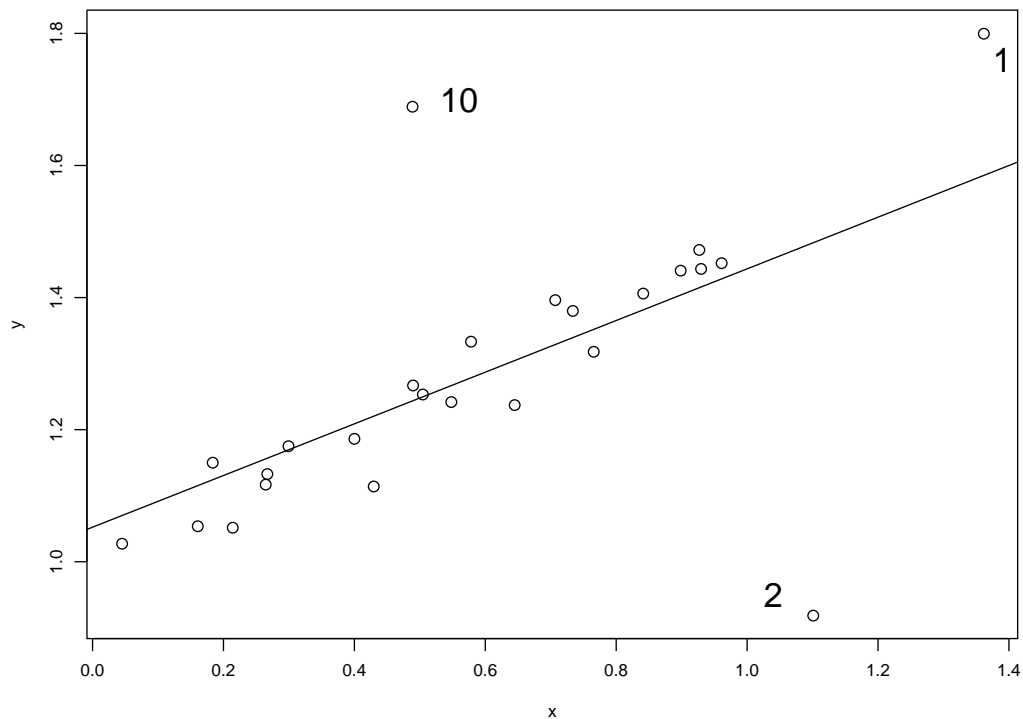


Abbildung 3.4: Plot mit auffälligen Beobachtungen

Beobachtung 1 ist wegen seiner extremen Lage am Rand der konvexen Hülle der Kovariablen-Daten ein Hebelpunkt (high leverage point):  $h_{11} = 0.27 > 2p/n = 4/25 = 0.16$ . Weitere Hebelpunkte existieren in diesem Datensatz nicht.

Nach dem Cook-Kriterium ist nur Beobachtung 2 eine "influential observation". Obwohl das Auge auch Beobachtung 10 als Kandidat für einen Ausreißer ansieht, schlägt das Cook-Kriterium nicht an, da die entsprechende x-Koordinate im Zentrum der anderen Beobachtungen liegt.

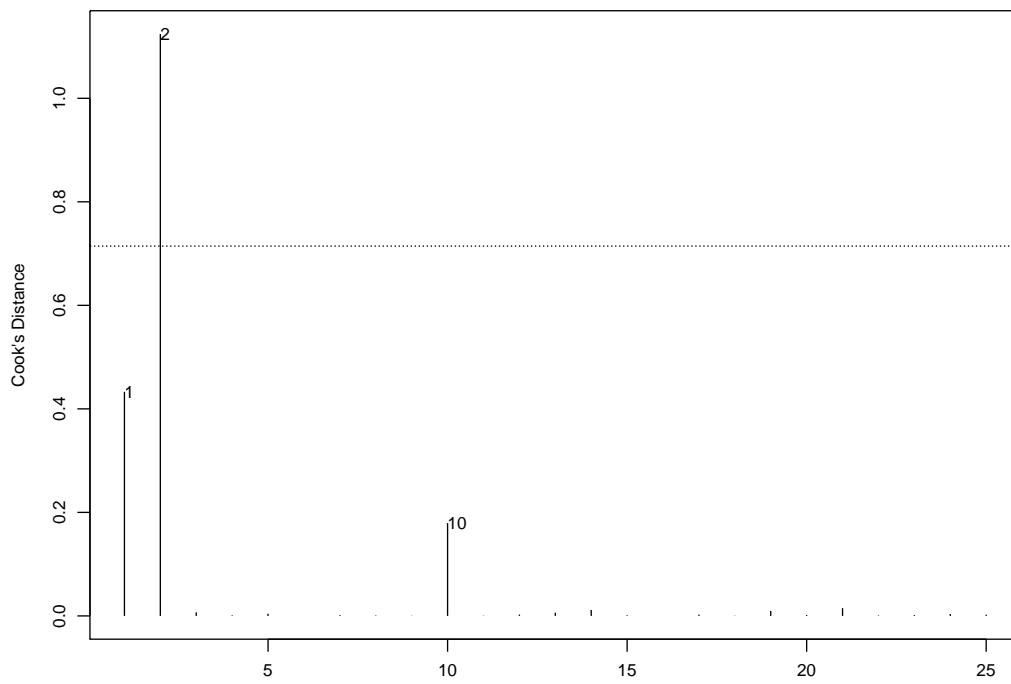


Abbildung 3.5: Cook'sche Distanz



### 3.9 Variablenselektion

Gegeben:

1. Daten  $(y_i, x_{i1}, \dots, x_{ik})_{i=1, \dots, n}$
2. Zielvariable  $y$  und potenzielle Kovariablen  $x_1, \dots, x_k$
3. Modell-Terme wie  $\log(x_j)$ ,  $x_j^2$ ,  $\sqrt{y}$ ,  $\dots$
4. Faktoren und Interaktionen  $x_j * x_l$

Ziele:

1. Welche Transformationen sind notwendig?  $\rightarrow$  2-dim. Scatter-Plots
2. Welche Variablen sollen ins Modell, welche nicht?  $\rightarrow$  siehe unten

Siehe auch bsp2 in Handout 4.

Hinsichtlich Variablenauswahl gibt es zwei Fehlerarten:

- Unterspezifikation (zu wenige Variablen; "underfitting")
- Überspezifikation (zu viele Variablen; "overfitting")

Faustregel: "Lieber zu wenige als zu viele Variablen!"

Vereinbarung

Aus  $m = k + 1$  Einflussgrößen (k Kovariablen, Regressoren + intercept) werden  $p \leq m$  (Achtung: Hier ist  $p < m$  ausdrücklich zugelassen!) ausgewählt. Seien

$$J := \{j_1, j_2, \dots, j_p\} \subseteq \{0, 1, \dots, k\}$$

und  $\vec{x}_{j_1}, \dots, \vec{x}_{j_p}$  die dem Variablensatz J entsprechenden Spalten von  $\mathbf{X}$  mit

$$\mathbf{X} = (\vec{x}_0 \equiv \mathbf{1}, \dots, \vec{x}_k)$$

und  $\bar{J} := \{0, 1, \dots, k\} \setminus J$ .

Wir ordnen die Größen  $RSS_J$  und  $s_J^2$  dem Modell mit Variablensatz J zu. Insgesamt o.B.d.A.

$$\mathbf{X} = (\mathbf{X}_J, \mathbf{X}_{\bar{J}}) \in \mathbb{R}^{n,m}, \quad \mathbf{X}_J \in \mathbb{R}^{n,p}.$$

In S-Plus und R werden oft Abkürzungen

wie  $RSS_p = RSS(p) = RSS_J$

bzw.  $s_p^2 = s^2(p) = s_J^2 = \frac{RSS_J}{n-p}$

verwendet, da  $|J| = p$ .

### 3.9.1 Unterspezifikation

Wahres Modell:

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}$$

$$\mathbf{X} = (\mathbf{X}_J, \mathbf{X}_{\bar{J}}) = (\mathbf{X}_1, \mathbf{X}_2), \quad \mathbf{X}_1 \in \mathbb{R}^{n \times p_1}, \quad \mathbf{X}_2 \in \mathbb{R}^{n \times p_2}, \quad p = p_1 + p_2, \quad p_2 > 0$$

Gewähltes Modell:

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}_1.$$

Nachteile für KQ-Schätzung  $\mathbf{b}_1$  für  $\boldsymbol{\beta}_1$ :

- $\mathbf{b}_1$  ist verzerrt, d.h.  $E(\mathbf{b}_1) \neq \boldsymbol{\beta}_1$ ,  $E(\mathbf{y}) \neq \mathbf{X}_1\boldsymbol{\beta}_1$
- $E(s_{p_1}^2) > \sigma^2$
- Prognosen falsch

### 3.9.2 Überspezifikation

$$\mathbf{X} = (\mathbf{X}_J, \mathbf{X}_{\bar{J}}) = (\mathbf{X}_1, \mathbf{X}_2), \quad \mathbf{X}_1 \in \mathbb{R}^{n \times p}, \quad \mathbf{X}_2 \in \mathbb{R}^{n \times n-p}$$

Wahres Modell:  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}_1$

Angepasstes Modell:  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}_2$

$$\mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} \text{ KQ-Schätzung für } \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} = \boldsymbol{\beta}.$$

$$\begin{aligned} \mathbf{H}_J = \mathbf{H}_1 &:= \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \\ \mathbf{E} &:= \mathbf{X}_2^T (\mathbf{I} - \mathbf{H}_1) \mathbf{X}_2 \\ \mathbf{F} &:= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \end{aligned} \quad \text{Orthogonalprojektion von } \mathbf{X}_2 \text{ auf } \text{span}(\mathbf{X}_1)^\perp$$

Es gilt:

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{pmatrix} (\mathbf{X}_1^T \mathbf{X}_1)^{-1} + \mathbf{F} \mathbf{E}^{-1} \mathbf{F}^T & -\mathbf{F} \mathbf{E}^{-1} \\ -\mathbf{E}^{-1} \mathbf{F}^T & \mathbf{E}^{-1} \end{pmatrix}$$

Spezialfall:  $\text{span}(\mathbf{X}_2) \perp \text{span}(\mathbf{X}_1)$

$$\Rightarrow \mathbf{F} = \mathbf{0} \text{ und } \mathbf{E} = \mathbf{X}_2^T \mathbf{X}_2$$

Kovarianzmatrizen von  $\mathbf{b}_{(1)}$  und  $\mathbf{b} = (\mathbf{b}_1^T, \mathbf{b}_2^T)^T$ :

$\text{Cov}(\mathbf{b}_{(1)}) = \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1}$  (kleines Modell) und

$\mathbf{C} := \mathbf{F} \mathbf{E}^{-1} \mathbf{F}^T$  ist positiv semidefinit,

$$\mathbf{C} \neq \mathbf{0}, \text{ falls } \mathbf{X}_1^T \mathbf{X}_2 \neq \mathbf{0}$$

$\text{Cov}(\mathbf{b}_1) = (\mathbf{X}^T \mathbf{X})^{-1} + \mathbf{C}$  (aus großem Modell) und

$$\text{Cov}(\mathbf{b}_1) - \text{Cov}(\mathbf{b}_{(1)}) \geq \mathbf{0},$$

d.h. die (Ko-)Varianz von  $\mathbf{b}_1$  wächst mit jeder hinzukommenden Variablen.

### 3.9.3 Kriterien zum Modellvergleich

#### 1. Bestimmtheitsmaß $R^2$

$R_J^2 = 1 - \frac{RSS_J}{RSS_{\{0\}}}$  ist ungeeignet, da

$R_{J \cup \{k\}}^2 \geq R_J^2$  mit der Anzahl der Variablen wächst.

#### 2. $s_J^2 = RSS_J / (n - |J|)$

Da der Nenner mit wachsender Variablenzahl kleiner wird, kann es  $J_1$  und  $J_2$  geben:  $s_{J_1 \cup J_2}^2 > s_{J_1}^2$

#### 3. Adjustiertes Bestimmtheitsmaß $R_{adj}^2$

$$R_{adj}^2(J) = 1 - [1 - R^2(J)] \frac{n-1}{n-p} = 1 - \frac{s_J^2}{s_{\{0\}}^2},$$

$$\text{wobei } s_{\{0\}}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2.$$

Falls  $s_J^2$  kaum mehr wächst, kann  $R_{adj}^2$  auch für steigende Variablenzahl fallen.

#### 4. Mallow's $C_p$ -Statistik

$$\hat{\sigma}^2 := s_m^2 = \frac{1}{n-m} \sum_{i=1}^n (y_i - \bar{y})^2 \text{ (volles Modell)}$$

$$C_p = p + (n-p) \frac{s_J^2 - \hat{\sigma}^2}{\hat{\sigma}^2} = C_J = \frac{RSS_J}{\hat{\sigma}^2} - n + 2p, \text{ wobei } |J| = p.$$

Ab jetzt sei  $RSS_p := RSS_J$ .

#### 5. Das AIC-Kriterium

Das Akaike information criterion wurde 1973 bzw. 1974 von Akaike eingeführt.

In S-Plus:

$$AIC := RSS_p + scale \cdot p \cdot k$$

Eingabeparameter: *scale* und *k*

Voreinstellung:  $k = 2$  (und meist:  $scale = \hat{\sigma}^2$ )

Bemerkung:  $scale = 0$  liefert bei Vorwärtsselektion stets alle Variablen

In R (anders als in S-Plus):

a)  $scale = 0$  :

$$AIC = n \ln(RSS_p/n + 2p)$$

(siehe Log-Likelihoodfunktion)

b)  $AIC = RSS_p - n + scale \cdot p \cdot k$ , sonst

Andere Programmpakete verwenden noch:

$$AIC = \hat{\sigma}^2(C_p + n),$$

wobei wieder  $\hat{\sigma}^2 = \frac{RSS_m}{n - m}$  aus dem vollen Modell verwendet wird.

Hier gilt dann die Identität

$$AIC = RSS_p + 2p\hat{\sigma}^2.$$

Dies ist übrigens die Voreinstellung in S-Plus bei "backward" ausgehend vom vollen Modell.

### Satz 3.24 (Eigenschaften von Mallows's $C_p$ )

1.  $C_m = m$
2.  $C_p$  schätzt den mittleren  $MSE(\mathbf{x}_i^T \mathbf{b})$  für alle  $i = 1, \dots, n$ .
3. Für festgehaltenes  $p$  ist  $C_p$  streng monoton fallend in  $s_p^2$ .
4. Die optimale Anzahl  $p$  an Einflussgrößen sollte so gewählt werden, dass  $C_p$  minimal wird.

5. Darstellungen:

$$C_p = C_J = p + (n - p) \frac{s_J^2 - \hat{\sigma}^2}{\hat{\sigma}^2} = \frac{RSS_J}{\hat{\sigma}^2} - n + 2p$$

$$|J| = p$$

6. Vergleich von  $C_p$  mit Gerade  $C_p = p$

Gemäß 3.24, 2 schätzt  $C_p$  den sog. *Mean Square Error (MSE)*.  
Ein kleinerer MSE sichert bessere Vorhersagen.

Sei etwa  $y = y(\mathbf{x}_0) \sim N(\underbrace{\mathbf{x}_0^T \boldsymbol{\beta}}_{\Theta}, \sigma^2)$

Wir suchen eine optimale Vorhersage  $\hat{y}$  für  $y$ , mit  $\hat{y} = \hat{\Theta} = \mathbf{x}_0^T \mathbf{b}$ ,  $\mathbf{b} = \hat{\boldsymbol{\beta}}$ .

### Def. 3.25 (MSE, Bias)

Sei  $\hat{\Theta}$  mit  $E(\hat{\Theta}) = \theta$  ein Schätzer für einen unbekanntem Parameter  $\Theta$ . Dann heißen

$$MSE(\hat{\Theta}) := E[(\hat{\Theta} - \Theta)^2] \text{ mean square error}$$

und

$$Bias(\hat{\Theta}) := E[\hat{\Theta} - \Theta] \text{ die Verzerrung (Bias) von } \hat{\Theta}.$$

**Satz 3.26 (Zerlegung)**

$$MSE(\hat{\Theta}) = Var(\hat{\Theta}) + Bias^2(\hat{\Theta}) \quad (3.11)$$

Beweis: Sei  $\theta := E(\hat{\Theta})$  :

$$\begin{aligned} MSE(\hat{\Theta}) &= E[(\hat{\Theta} - \Theta)^2] = E[(\hat{\Theta} - \theta + \theta - \Theta)^2] = \\ &= E[(\hat{\Theta} - \theta)^2 + 2(\hat{\Theta} - \theta)(\theta - \Theta) + (\theta - \Theta)^2] = \\ &\stackrel{\theta = E(\hat{\Theta})}{=} \underbrace{E(\hat{\Theta} - \theta)^2}_{Var(\hat{\Theta})} + 2(\theta - \Theta) \underbrace{E(\hat{\Theta} - \theta)}_0 + \underbrace{(\theta - \Theta)^2}_{Bias^2(\hat{\Theta})}. \end{aligned}$$

Theoretische Grundlage des AIC ist die Log-Likelihoodfunktion für normalverteilte Fehler:

$$y_i \sim N(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2), \quad f(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}\right)$$

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2) &:= \ln\left(\prod_{i=1}^n f(y_i)\right) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 - \frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) \end{aligned}$$

$\hat{\boldsymbol{\beta}}_{ML}$  und  $\hat{\sigma}_{ML}^2$  maximieren  $L(\boldsymbol{\beta}, \sigma^2)$ , d.h.

$$\hat{\boldsymbol{\beta}}_{ML} = \hat{\boldsymbol{\beta}}_{KQ} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (\text{KQ-Schätzer})$$

$$\hat{\sigma}_{ML}^2 = \frac{RSS_p}{n} \neq s_p^2 = \frac{RSS_p}{n-p}$$

$$L(\hat{\boldsymbol{\beta}}_{ML}, \hat{\sigma}_{ML}^2) = -\frac{n}{2}(1 - \ln(2\pi)) - \frac{n}{2} \ln(RSS_p/n)$$

Weglassen von Konstanten unabhängig von p:

$$AIC = -2L(\hat{\boldsymbol{\beta}}_{ML}, \hat{\sigma}_{ML}^2) + \underbrace{2p}_{\text{Strafterm für zu viele Variablen}}$$

(vgl.  $scale = 0$  in R)

Optimiert man die Log-Likelihoodfunktion nur nach  $\boldsymbol{\beta}$  und setzt  $\sigma^2$  als Parameter ein, so ergibt sich:

$$L(\boldsymbol{\beta}_{ML}, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln\left(\frac{\sigma^2}{n}\right) - \frac{RSS_p}{2\sigma^2}$$

Weglassen der beiden ersten Konstanten liefert mit  $\hat{\boldsymbol{\beta}}_{ML} = \hat{\boldsymbol{\beta}}_{KQ}$  und dem *Strafterm*  $2p$ :

$$-2L(\hat{\boldsymbol{\beta}}_{ML}, \sigma^2) \hat{=} \frac{RSS_p}{\sigma^2} + 2p.$$

Dieser letzte Ausdruck erklärt den Gebrauch von

$$AIC = RSS_p + 2 \cdot p \cdot scale$$

in S-Plus, wobei in der Regel

$$scale = \hat{\sigma}^2 = \frac{RSS_m}{n-m}$$

aus dem vollen Modell gewählt wird.

### 3.9.4 Algorithmisches Vorgehen

1. All subset regression mit der Prozedur "leaps" (nur in S-Plus verfügbar)  
Ausdruck der jeweils r besten Variablensätze bzgl. eines oder mehrerer der Kriterien  $C_p$ ,  $AIC$ ,  $R_{adj}^2$ ,  $s_p^2$
2. Schrittweise Regression mit der Prozedur "step" (auch in R)  
Varianten mit AIC:
  - (a) Forward selection  
ausgehend vom Modell zum Kommando  $\text{lm}(y \sim 1)$
  - (b) Backward elimination  
ausgehend vom vollen Modell  $\text{lm}(y \sim x_1 + x_2 + \dots + x_k)$   
 $m = k+1$  Variablen mit intercept

### 3.9.5 Variablenselektion in R und S-Plus

in R verfügbar: "leaps" für  $R_{adj}^2$  und

"step" für AIC und schrittweise Regression

nur in Splus: "leaps" für  $C_p$  und Auswahl der besten Variablensätze

Splus-Befehle für Mallows's  $C_p$  am Beispiel von "College Spending"

Zusammenfügen der möglichen Variablen

```
x.cbind(SAT, TOP10, ACCRATE, PHD, RATIO, GRADRATE, ALUMNI)
```

die zu langen Variablen-Namen werden durch die "names"-Option verkürzt

```
name_c("S", "T", "A", "P", "R", "G", "A")
```

in der Regel will man nur die besten zwei oder drei Variablensätze sehen; z.B.  $nbest=3$

```
r.cp.leaps(x, log(SPEND), method="Cp", nbest=3, names=name)
```

Zusammenfassen der nach Variablensätzen geordneten Listenelemente

```
list.cp.cbind(round(r.cp$Cp, 3), r.cp$size, r.cp$label)
```

Umformen in ein Data-Frame

```
list.cp.as.data.frame(list.cp)
```

Grafik erstellen

nur kleine Werte von  $C_p$  wegen Skalierung der Grafik

```
cp.r.cp$Cp[r.cp$Cp<15]
```

```
p.r.cp$size[r.cp$Cp<15]
```

```
plot(p,cp,ylab="Mallows Cp",xlab="Parameter im Modell")
```

```
title("Cp-Plot für log(SPEND)-Modell")
```

```
abline(0,1)
```

Zunächst zeigt sich, dass verschiedene Kriterien zu unterschiedlichen besten Variablensätzen führen.

<b>Tabelle 2: Beste Variablensätze für "College Spending"</b>					
p	Variable	$R_{adj}^2$	Rang( $R_{adj}^2$ )	$C_p$	Rang( $C_p$ )
2	S	0.6005		75.554	
2	T	0.5072		126.606	
2	G	0.4657		149.300	
3	TR	0.7199		11.201	
3	SR	0.6848		30.265	
3	ST	0.6197		65.632	
4	STR	0.7339		4.589	3
4	TPR	0.7272		8.173	
4	SPR	0.6984		23.749	
5	STPR	0.7374	2	3.690	1
5	STAR	0.7342		5.430	
5	STRG	0.7325		6.297	
6	STAPR	0.7385	1	4.131	2
6	STPRG	0.7357		5.626	
6	STPRAI	0.7356		5.644	
7	STAPRAI	0.7367	3	6.072	
7	STAPRG	0.7366	4	6.098	
7	STPRGAI	0.7340		7.530	
8	STAPRGAI	0.7350	5	8.000	

Zum Schluss noch eine grafische Darstellung von Mallows'  $C_p$ . Die besten Variablensätze weisen  $C_p$ -Werte unterhalb der Geraden  $C_p = p$  auf; siehe unter anderem *STPR* für  $p = 5$ .

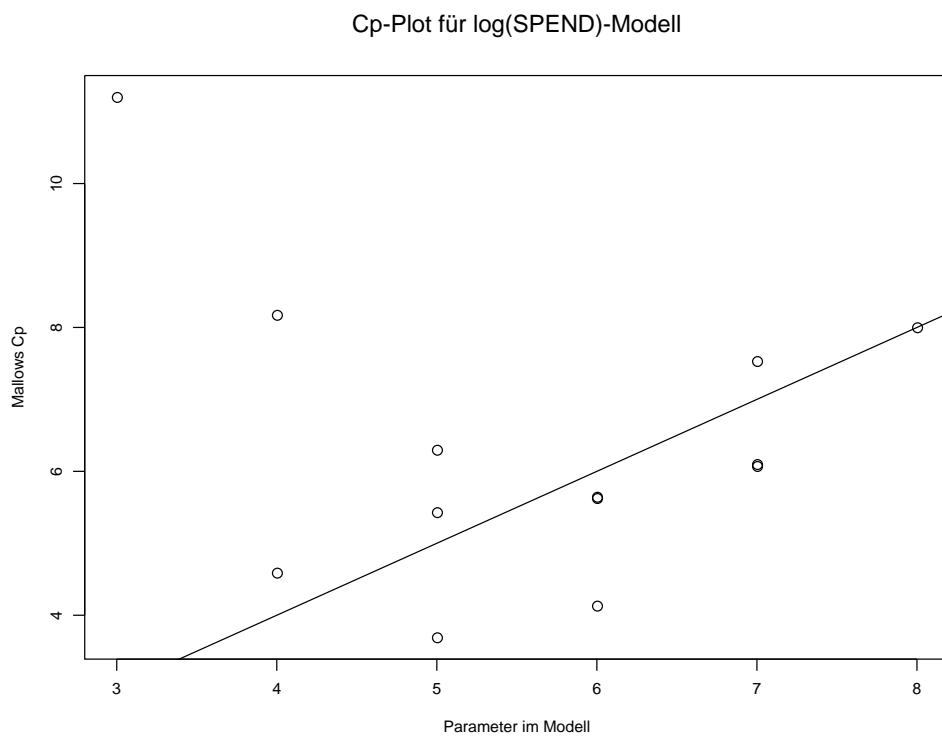


Abbildung 3.6: Mallows'  $C_p$  im Beispiel College Spending



## 3.10 Multikollinearität

Eines der Ziele statistischer Analysen ist die Konstruktion erwartungstreuer Schätzer mit minimaler Varianz. Wir haben bereits am Beispiel des "Overfitting" gesehen, dass man lieber mit zu wenigen als mit zu vielen Variablen arbeiten sollte.

Ein anderes Phänomen, das zur sprunghaften Vergrößerung der Varianzen der KQ-Schätzungen führen kann, ist die sog. "Multikollinearität" (M-K). Um M-K zu kontrollieren, versucht man, Größen wie

$$\sum_j \text{Var}(b_j)$$

klein zu halten. Natürlich müssen erst einmal einige Regressoren ins Modell aufgenommen werden, damit das Design den Erwartungswert von  $\mathbf{y}$  hinreichend gut modelliert. M-K verhindert präzise statistische Aussagen. Diese Tatsache ist vielen Nutzern statistischer Software bekannt, ohne dass die meisten genau wissen, welche geometrischen und numerischen Zusammenhänge sich dahinter verbergen.

Hier können nicht alle Phänomene der M-K im Detail diskutiert werden. Die für die Praxis wichtigsten Aspekte werden jedoch angesprochen. Diese sind:

1. M-K verursacht eine "Explosion" der Varianzen der Schätzer (siehe VIF = variance inflation factor), und dementsprechend ungenau werden Vorhersagen.
2. M-K beruht nur auf Eigenschaften der Kovariablen (Regressoren) und hängt wie die leverage points in keiner Weise vom Response  $\mathbf{y}$  ab.
3. M-K entsteht, wenn die Daten zu zwei oder mehreren Kovariablen nahezu linear abhängig sind. Wären die Kovariablen Realisierungen von Zufallsvariablen, so würde man sagen: Zwei oder mehrere Kovariablen sind hochkorreliert.
4. M-K ist der Gegenpol zur Orthogonalität. In sog. orthogonalen Designs (Spalten der Design-Matrix  $\mathbf{X}$  stehen aufeinander senkrecht) tritt keine M-K auf.
5. Im Falle von M-K vergrößert sich die Varianz aufgrund einer geometrischen "Schiefstellung" der Kovariablen-Daten, ohne einen Gewinn bei der Modellierung von  $E(\mathbf{y})$  zu liefern. Bei orthogonalen Designs wird dagegen die zusätzliche Erklärung von  $E(\mathbf{y})$  durch die minimal mögliche Varianzvergrößerung erreicht.

Im Unterschied zu den übrigen Abschnitten ist es hier günstiger, den intercept nicht als erste Spalte der Design-Matrix  $\mathbf{X}$  zu modellieren.  $\mathbf{X}$  enthalte also in diesem Abschnitt ausnahmsweise nur echte Kovariablenspalten. Also

$$\mathbf{X} = (\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_k), \text{Rang}(\mathbf{1}, \mathbf{X}) = k + 1 = p$$

$$\mathbf{y} = \mathbf{1} \cdot \beta_0 + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, E(\boldsymbol{\epsilon}) = \mathbf{0}, \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}.$$

Ausnahmsweise:  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)^T$

### 3.10.1 Einfache lineare Regression

$$k = 1, \mathbf{x} = (x_1, \dots, x_n)^T \neq \mathbf{1} \cdot \bar{x}.$$

Damit stehen im zentrierten Modell

$$y_i = \beta_0 + \beta_1 (x_i - \bar{x}) + \epsilon_i, i = 1, \dots, n$$

die Spalten  $\mathbf{1}$  und  $\mathbf{x} - \mathbf{1} \cdot \bar{x}$  aufeinander senkrecht. Man erkennt aber z.B. an der Formel für die Varianz von  $b_1$ , dass es beispielsweise im Fall frei wählbarer Designpunkte  $x_1, \dots, x_n$  günstiger ist,  $\sum (x_i - \bar{x})^2$  möglichst groß zu halten.

### 3.10.2 Skalierungsinvarianz

Einige Größen wie die KQ-Schätzungen  $b_j$  und deren Varianzen hängen von der Wahl des Koordinatensystems für die Kovariablen ab. Andere wie RSS, t-Werte der  $b_j$ ,  $R^2$ ,  $R_{adj}^2$  oder Mallows  $C_p$  sind "skalierungsinvariant" bezüglich folgenden Transformationen der Kovariablen

$$\vec{x}_j \rightarrow \gamma \cdot \vec{x}_j + \delta \cdot \mathbf{1}, \quad \gamma \neq 0, \quad \delta \in \mathbf{R},$$

wobei  $\vec{x}_j = (x_{1j}, \dots, x_{nj})^T$  die Daten zur j-ten Kovariablen sind.

t-Werte,  $R^2$  und  $R_{adj}^2$  sind sogar auch noch invariant gegen Skalierungen im Response  $\mathbf{y}$ .

Transformationen der genannten Art erhält man etwa durch den Übergang von \$ zu Euro oder von °C zu °F. Man überzeugt sich leicht, dass z.B. die Varianzen der  $b_j$  von der Wahl der Skalierung der Kovariablen abhängen. Damit ist auch unser ins Auge gefasstes Zielkriterium  $\sum Var(b_j)$  skalierungsabhängig und damit so in der Praxis wertlos.

Es gibt jedoch eine eindeutige kanonische Form  $\mathbf{X}^*$ , zu der man ausgehend von beliebigen Skalierungen für die Kovariablen gelangen kann. Bezüglich dieses eindeutig bestimmten Designs ist dann  $\sum Var(b_j)$  klein zu halten.

#### Satz 3.27 (Zentriertes, standardisiertes Modell)

Vor.:

$$\begin{aligned} \mathbf{X} &= (\vec{x}_1, \dots, \vec{x}_k) \in \mathbf{R}^{n,k}, \quad Rang(\mathbf{1}, \mathbf{X}) = p = k + 1 \\ \mathbf{C} &:= \text{diag}\{\gamma_1, \dots, \gamma_k\} \in \mathbf{R}^{k,k}, \\ \mathbf{D} &:= \text{diag}\{\delta_1, \dots, \delta_k\} \in \mathbf{R}^{k,k}, \\ \mathbf{Z} &:= [\mathbf{X} - (\mathbf{1}, \dots, \mathbf{1}) \cdot \mathbf{D}] \cdot \mathbf{C} = (\vec{z}_1, \dots, \vec{z}_k). \end{aligned}$$

Beh.:

Für beliebiges  $\mathbf{D}$  und nichtsinguläres  $\mathbf{C}$  führt die Zentrierung mit

$$\bar{z}_j := \frac{1}{n} \sum_{i=1}^n z_{ij}, \quad \vec{z}_j = (z_{1j}, \dots, z_{nj})^T$$

und Standardisierung mit

$$S_j := \sqrt{\sum_{i=1}^n (z_{ij} - \bar{z}_j)^2}$$

und

$$\vec{x}_j^* := (\vec{z}_j - \mathbf{1} \cdot \bar{z}_j) / S_j$$

zu einem eindeutigen

$$\mathbf{X}^* = (\vec{x}_1^*, \dots, \vec{x}_k^*).$$

Gemäß diesem Satz führt also jede Ausgangsskalierung der Kovariablen zu demselben  $\mathbf{X}^*$ .

Man versucht nun im Modell

$$\mathbf{y} = \mathbf{1} \cdot \beta_0^* + \mathbf{X}^* \boldsymbol{\beta}^* + \boldsymbol{\epsilon}, \boldsymbol{\beta}^* \in \mathbf{R}^k,$$

das Kriterium

$$S^2(\mathbf{b}^*) = \sum_{i=1}^k \text{Var}(b_i^*) = E(\mathbf{b}^* - \boldsymbol{\beta}^*)^T (\mathbf{b}^* - \boldsymbol{\beta}^*)$$

zu minimieren.

$S^2(\mathbf{b}^*)$  ist, wie gesehen, nicht mehr von der Ausgangsskalierung abhängig.

### 3.10.3 Einfache Berechnung des VIF

Wegen der Zentrierung stehen alle Spalten von  $\mathbf{X}^*$  senkrecht auf  $\mathbf{1}$ . Damit gilt

$$\text{Cov} \begin{pmatrix} b_0^* \\ \mathbf{b}^* \end{pmatrix} = \sigma^2 \begin{pmatrix} n & \mathbf{0}^T \\ \mathbf{0} & \mathbf{X}^{*T} \mathbf{X}^* \end{pmatrix}^{-1},$$

also

$$S^2(\mathbf{b}^*)/\sigma^2 = \text{tr}[(\mathbf{X}^{*T} \mathbf{X}^*)^{-1}] =: w_{11} + \dots + w_{kk}.$$

$VIF_j := w_{jj}$  bezeichnet den *variance inflation factor*, d.h. den Anteil der  $j$ -ten Kovariablen am Zielkriterium  $S^2(\mathbf{b}^*)$ .

Dabei ist die explizite Berechnung der Inversen von  $\mathbf{X}^{*T} \mathbf{X}^*$  gar nicht notwendig, denn es gilt

$$VIF_j = \frac{1}{1 - R_j^2},$$

wobei mit  $\mathbf{X}_{(j)} := (\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_{j-1}, \vec{\mathbf{x}}_{j+1}, \dots, \vec{\mathbf{x}}_k) \in \mathbf{R}^{n, k-1}$  die Design-Matrix ohne Daten zur Kovariablen  $j$  bezeichnet wird und  $R_j^2$  das Bestimmtheitsmaß zu  $k$  fiktiven Modellen

$$\vec{\mathbf{x}}_j = \mathbf{1}\beta_{0j} + \mathbf{X}_{(j)}\boldsymbol{\beta}_{(j)} + \boldsymbol{\epsilon}_j, \quad j = 1, \dots, k$$

bezeichnet.

Wegen der Skalierungsinvarianz des Bestimmtheitsmaßes können die  $R_j^2$  mit den ursprünglichen Daten berechnet werden. Eine Zentrierung und Skalierung ist also in der Praxis gar nicht unbedingt notwendig.

### 3.11 Korrelierte Fehler, gewichtete Regression

Das verallgemeinerte lineare Modell lautet:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad E(\boldsymbol{\epsilon}) = \mathbf{0}, \quad Cov(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}, \quad \mathbf{V} \text{ positiv definit und bekannt.}$$

Wir diskutieren den wichtigsten Spezialfall:

$$\mathbf{V} = \text{diag}\{v_{11}, \dots, v_{nn}\}, \quad v_{ii} > 0, \quad i = 1, \dots, n.$$

Da die  $v_{ii}$  a priori bekannt sind (in der Regel Transformationen einer oder mehrerer Kovariablen), ist eine Rücktransformation zum Standardmodell möglich:

Mit

$$\mathbf{V}^{-1/2} = \text{diag}\{1/\sqrt{v_{11}}, \dots, 1/\sqrt{v_{nn}}\}$$

und

$$\mathbf{y}_v := \mathbf{V}^{-1/2} \mathbf{y}, \quad \mathbf{X}_v := \mathbf{V}^{-1/2} \mathbf{X}, \quad \boldsymbol{\epsilon}_v := \mathbf{V}^{-1/2} \boldsymbol{\epsilon}$$

gilt:

$$\mathbf{y}_v = \mathbf{X}_v \boldsymbol{\beta}_v + \boldsymbol{\epsilon}_v, \quad E(\boldsymbol{\epsilon}_v) = \mathbf{0}, \quad Cov(\boldsymbol{\epsilon}_v) = \sigma^2 \mathbf{I}. \quad (3.12)$$

Der KQ-Schätzer zum transformierten Modell 3.12 lautet:

$$\mathbf{b}_v = (\mathbf{X}_v^T \mathbf{X}_v)^{-1} \mathbf{X}_v^T \mathbf{y}_v = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

Die gewichtete Residuenquadratsumme lässt sich demgemäß darstellen als

$$\|\mathbf{y}_v - \mathbf{X}_v \mathbf{b}_v\|^2 := \|\mathbf{y} - \mathbf{X} \mathbf{b}_v\|_v^2 := (\mathbf{y} - \mathbf{X} \mathbf{b}_v)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b}_v).$$

In R oder S-Plus ist die Gewichtung in einfacher Weise darstellbar (siehe Beispiel bsp5).

`attach(data5)`      Bereitstellen der Daten

1. Ohne Gewichtung (Auslassen der Beobachtungen mit fehlenden Daten)

```
r <- lm(Salary ~ Price + Work, na.action = na.omit)
```

```
summary(r)
```

	Value	Std. Error	t value	Pr(>  t )
(Intercept)	10.0942	31.3257	0.3222	0.7488
Price	0.8688	0.1159	7.4974	0.0000
Work	-0.0167	0.0142	-1.1768	0.2458

Multiple R-Squared: 0.6572

2. Mit Gewichtung:  $1/Price$

```
r.x <- lm(Salary ~ Price + Work, weights = 1/Price)
```

```
summary(r.x)
```

	Value	Std. Error	t value	Pr(>  t )
(Intercept)	3.1874	28.0665	0.1136	0.9101
Price	0.9373	0.1096	8.5528	0.0000
Work	-0.0156	0.0126	1.2381	0.2224

Multiple R-Squared: 0.7163

3. Zu  $r.x$  entsprechendes transformiertes Regressionsmodell:

```
r.t <- -lm(Salary/sqrt(Price) ~ (I(1/sqrt(Price))
+ sqrt(Price) + I(Work/sqrt(Price)) - 1)
summary(r.t)
```

	Value	Std. Error	t value	Pr(>  t )
I(1/sqrt(Price))	3.1874	28.0665	0.1136	0.9101
sqrt(Price)	0.9373	0.1096	8.5528	0.0000
I(Work/sqrt(Price))	-0.0156	0.0126	1.2381	0.2224

\*\* Wert falsch \*\* Multiple R-Squared: 0.9012 \*\* wegen "-1" in lm \*\*

4. Suche nach "optimalen" Weights:

```
r.x2 <- -lm(Salary ~ Price + Work, weights = 1/(Price^2))
Multiple R-Squared: 0.7459
```

```
r.x3 <- -lm(Salary ~ Price + Work, weights = 1/(Price^3))
Multiple R-Squared: 0.7471
```

```
r.x4 <- -lm(Salary ~ Price + Work, weights = 1/(Price^4))
Multiple R-Squared: 0.7235
```

5. Fazit Nach dem  $R^2$ -Kriterium scheint das Modell  $r.x3$  am besten geeignet. Das zeigt sich u. a. auch beim Vergleich der entsprechenden QQ-Plots der Residuen. Abbildung 3.7 zeigt das Ergebnis von `plot(r.x3)`.

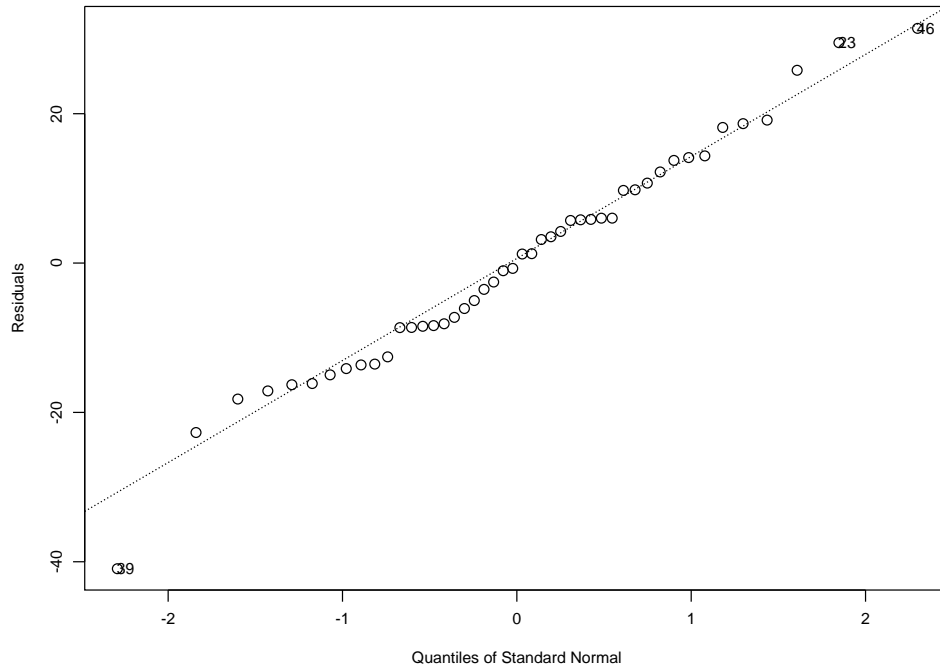


Abbildung 3.7: QQ-Plot

### 3.12 Autokorrelierte Fehler

Besonders bei Zeitreihen sind die Fehler oft korreliert; d.h.  $E(\epsilon_i, \epsilon_j) \neq 0$ . Somit ist die Unabhängigkeitsannahme verletzt. In der einfachsten Abhängigkeitsstruktur sind die Fehler *autokorreliert*.

Modell (autoregressive Fehlerstruktur 1. Ordnung):

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \\ \epsilon_t &= \rho\epsilon_{t-1} + U_t, \end{aligned}$$

wobei  $U_t \text{ iid } N(0, \sigma_U^2)$ ;  $-1 < \rho < 1$ .

Indizien für Fehlerkorrelation

Bei positiver Korrelation ändern sich  $\epsilon_t, \epsilon_{t+1}$  kaum (vgl. Grafik 3.8)

Dagegen alternieren  $\epsilon_t$  und  $\epsilon_{t+1}$  bei negativer Korrelation wegen  $\epsilon_t = -|\rho|\epsilon_{t-1} + U_t$ , vielfach.

Formaler Test: Durbin-Watson-Test

Der formalen Einfachheit wegen betrachten wir unendlich viele  $t$ .

$$\epsilon_t = \rho\epsilon_{t-1} + U_t = \rho(\rho\epsilon_{t-2} + U_{t-1}) + U_t = \rho^2\epsilon_{t-2} + \rho U_{t-1} + U_t$$

usw. liefert

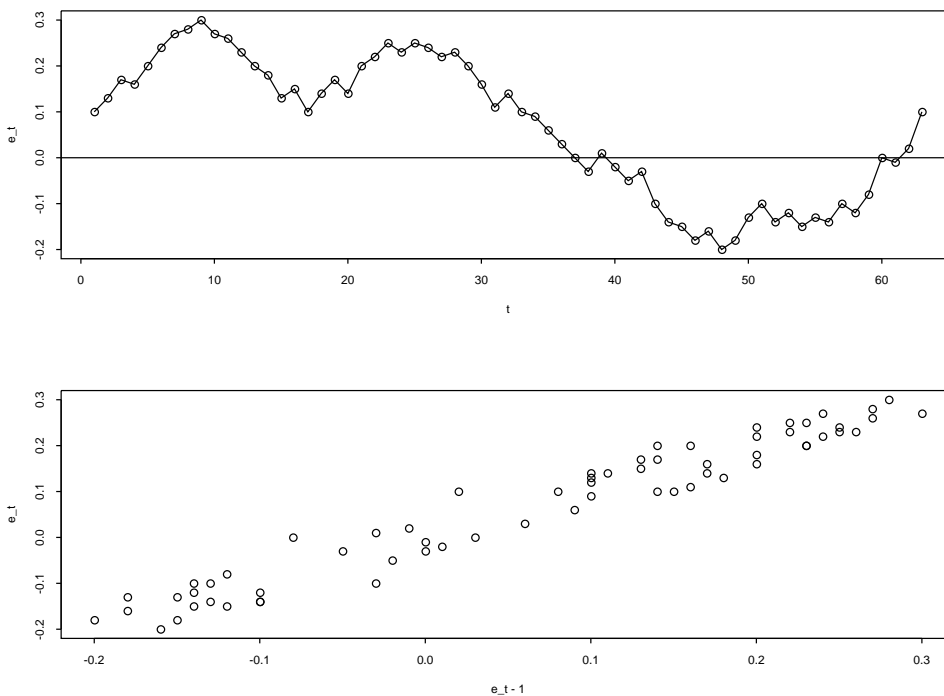


Abbildung 3.8: Autokorrelierte Fehler

$$\epsilon_t = \sum_{j=0}^{\infty} \rho^j U_{t-j}$$

woraus

$$E(\epsilon_t) = \sum_{j=0}^{\infty} \rho^j \underbrace{E(U_{t-j})}_{=0} = 0$$

und

$$\text{Var}(\epsilon_t) = \sigma_U^2 \sum_{j=0}^{\infty} (\rho^2)^j \stackrel{\text{geom. Reihe}}{=} \sigma_U^2 \frac{1}{1 - \rho^2}$$

folgen. In ähnlicher Weise kann man leicht zeigen, dass

$$\text{Cov}(\epsilon_t, \epsilon_{t+j}) = \rho^{|j|} \sigma_U^2 \frac{1}{1 - \rho^2}, \quad |\rho| < 1.$$

Zu testen ist

$$H_0 : \rho = 0 \quad \text{gegen} \quad H_1 : \rho > 0.$$

Mit  $e_t := Y_t - \hat{Y}_t$  schlagen Durbin und Watson folgende Teststatistik vor

$$d := \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}.$$

Bemerkung:

Wegen  $e_t - e_{t-1} = \rho e_{t-1} + U_t - e_{t-1} = (\rho - 1)e_{t-1} + U_t$  ist  $d$  für  $\rho \approx 1$  klein und  $H_0$  in

diesem Fall zu verwerfen. Andererseits erwartet man für  $\rho = 0$  einen eher großen Wert, denn unabhängige Fehler liefern

$$E((\epsilon_t - \epsilon_{t-1})^2) = 2\sigma_U^2 \text{ und } E(\epsilon_t^2) = \sigma_U^2.$$

### Durbin–Watson–Test

Verwerfe  $H_0$  (d.h.  $\rho > 0$ ), falls  $d < d_L$ .

Akzeptiere  $H_0$  ( $\rho = 0$ ), falls  $d > d_U$ .

Test ist nicht zu entscheiden, falls  $d_L < d < d_U$ .

$d_L$  und  $d_U$  entnehme man geeigneten Tabellen. In der Praxis ist die Bedeutung des Durbin–Watson–Tests umstritten. Grafische Analysemethoden gewinnen immer mehr an Bedeutung. Natürlich möchte man der Einfachheit halber am liebsten die Hypothese  $\rho = 0$  bestätigt haben.

### Maßnahmen bei Autokorrelation

Einfachstes Modell:

$$\begin{aligned} Y_t &= \beta_0 + \beta_1 x_t + \epsilon_t \\ \epsilon_t &= \rho \epsilon_{t-1} + U_t \end{aligned}$$

Wegen  $\epsilon_t - \rho \epsilon_{t-1} = U_t$  betrachtet man

$$\begin{aligned} Y_t^* := Y_t - \rho Y_{t-1} &= \beta_0 + \beta_1 x_t + \epsilon_t - \rho(\beta_0 + \beta_1 x_{t-1} + \epsilon_{t-1}) \\ &= \underbrace{\beta_0(1-\rho)}_{\beta_0^*} + \underbrace{\beta_1}_{\beta_1^*} \underbrace{(x_t - \rho x_{t-1})}_{x_t^*} + \underbrace{\epsilon_t - \rho \epsilon_{t-1}}_{=U_t \sim N(0, \sigma_U^2)}. \end{aligned}$$

Das ist eine Einfache Lineare Regression.

### Schätzung von $\rho$

Ausgehend von  $\epsilon_t = \rho \epsilon_{t-1} + U_t$  liefert eine Einfache Lineare Regression durch den Ursprung eine KQ-Schätzung

$$\hat{\rho} = \frac{\sum_{t=2}^n e_t e_{t-1}}{\sum_{t=1}^{n-1} e_t^2}$$

für  $\rho$ .

Dies ergibt folgendes iterative Schema:

1. Berechne  $\hat{\rho}$ .
2. Bilde  $Y_t^* = \beta_0^* + \beta_1^* X_t^* + U_t$  und berechne die KQ-Schätzer  $b_0^*$  und  $b_1^*$ .



3. Falls die Residuen  $e_t^* := Y_t^* - b_0^* - b_1^* X_t$  unkorreliert sind, schätze  $\beta_0$  und  $\beta_1$  durch  $b_0 = \frac{b_0^*}{1 - \rho}$  und  $b_1 = b_1^*$ .
4. Ansonsten setze  $e_t = e_t^*$  und gehe zurück zu Schritt 1.

Die Gewichtungen aus dem vorangegangenen Abschnitt sind völlig unproblematisch und tragen oft zu erheblichen Modellverbesserungen bei.

Die Behandlung abhängiger Strukturen führt aber schon im einfachsten Fall autokorrelierter Fehler nicht nur zu einem erheblichen Zusatzaufwand. Schlimmer ist noch, dass eine Interpretation der berechneten Schätzungen immer schwieriger wird. Für die Behandlung autoregressiver Fehler ist eine ausreichende Erfahrung (bei Datentransformationen und effizientem Umgang mit Programmpaketen) notwendig, die in einem Einführungskurs wie diesem nicht von allen Teilnehmern erwartet werden kann.

### 3.13 Kategoriale Variable; Teil 1: Regressoren

Die Behandlung qualitativer Regressoren wird in zwei Abschnitten behandelt. Dieser Abschnitt befasst sich mit linearen Modellen mit genau einer kategorialen Kovariablen.

#### 3.13.1 Einführendes Beispiel mit einer binären Kovariablen

Beispiel:  $y$ : Geburtsgewicht

$x_1$ : Alter (metrisch)

$\tilde{x}_2$ : Geschlecht (kategorial)

mit  $n_1$  "männlichen" und  $n - n_1$  "weiblichen" Datensätzen.

Dummy-Kodierung:

$$x_{i2} = \begin{cases} 1, & \text{falls } \tilde{x}_{i2} = m \text{ (männlich)} \\ 0, & \text{sonst} \end{cases}$$

Dies entspricht der Voreinstellung in R.

In S-Plus muss vor dem `lm`-Kommando noch

```
options(contrasts = c("contr.treatment", "contr.poly"))
```

gesetzt werden.

Man erhält folgende Design-Matrix

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 & 1 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & x_{n_1} & 1 \\ 1 & x_{n_1+1} & 0 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & x_n & 0 \end{pmatrix}$$

Effekt-Kodierung:

Die Variante

```
options(contrasts = c("contr.sum", "contr.poly"))
```

liefert die sog. Effekt-Kodierung

$$x_{i2} = \begin{cases} 1, & \text{falls } \tilde{x}_{i2} = m \\ -1, & \text{sonst} \end{cases}$$

In beiden Fällen gilt  $p = 3$  und  $\text{Rang}(\mathbf{X}) = p$ .

Ein einfacher Ansatz für das Beispiel auf dem Übungsblatt liefert

```
out.add <- lm(y ~ age + sex)
```

$$y = -1773 + 121 \cdot \text{age} + 163 \cdot \text{sex}$$

Dieses rein additive Modell geht von identischen Steigungen in den beiden Gruppen  $m$  bzw.  $w$  aus. Man erhält die einfachen Regressionen

$$y_m = -1610 + 121 \cdot \text{age}$$

$$y_w = -1773 + 121 \cdot age.$$

Den allgemeinen Fall modelliert man mit einem zusätzlichen **Interaktionsterm** `age : sex`

$$\begin{aligned} \text{out.ia} &<- \text{lm}(y \sim \text{age} * \text{sex}) \\ &\hat{=} \text{lm}(y \sim \text{age} + \text{sex} + \text{age} : \text{sex}) \end{aligned}$$

und der Design-Matrix

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 & 1 & x_1 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 1 & x_{n_1} & 1 & x_{n_1} \\ 1 & x_{n_1+1} & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 1 & x_n & 0 & 0 \end{pmatrix} \in \mathbb{R}^{n,4}.$$

Die diesbezügliche KQ-Schätzung

$$y = -2142 + 873 \cdot \text{sex} + 130 \cdot \text{age} - 18 \cdot \text{sex} \cdot \text{age}$$

Führt zu den nun völlig unabhängigen KQ-Geraden

$$y_w = -2142 + 130 \cdot \text{age}$$

$$y_m = -1269 + 112 \cdot \text{age} = (-2142 + 873) + (130 - 18) \cdot \text{age}$$

in den verschiedenen Geschlechts-Gruppen.

Die zwei KQ-Geraden erhält man direkt über

$$\text{lm}(y \sim \text{sex} / \text{age} - 1)$$

Übrigens sind die Steigungen nicht signifikant voneinander verschieden, Dies zeigt sich aus einer sog. "Varianzanalyse":

```
out.add <- lm(y ~ age + sex)
out.ia <- lm(y ~ age * sex)
aov(out.ia, out.add)
```

mit dem Resultat:

	Terms	Resid. Df	RSS	Test	Df	Sum of Sq	F Value	Pr(F)
1	sex * age	20	652424.5					
2	sex + age	21	658770.7	-sex:age	-1	-6346.225	0.1945428	0.6638934

Varianzanalyse (analysis of variance, ANOVA) vergleicht ineinandergeschachtelte lineare Modelle anhand spezieller F-Tests (wir hatten die Prozedur `aov()` schon im Zusammenhang der Variablenselektion kennengelernt). In den nächsten Abschnitten wird noch näher auf diese Methode eingegangen.

Man nennt kategoriale Variablen wie `sex` auch *Faktoren*. Wegen der alphanumerischen Kodierung  $\text{sex} \in \{m, w\}$  ist für R bzw. S-Plus klar, dass es sich um eine kategoriale

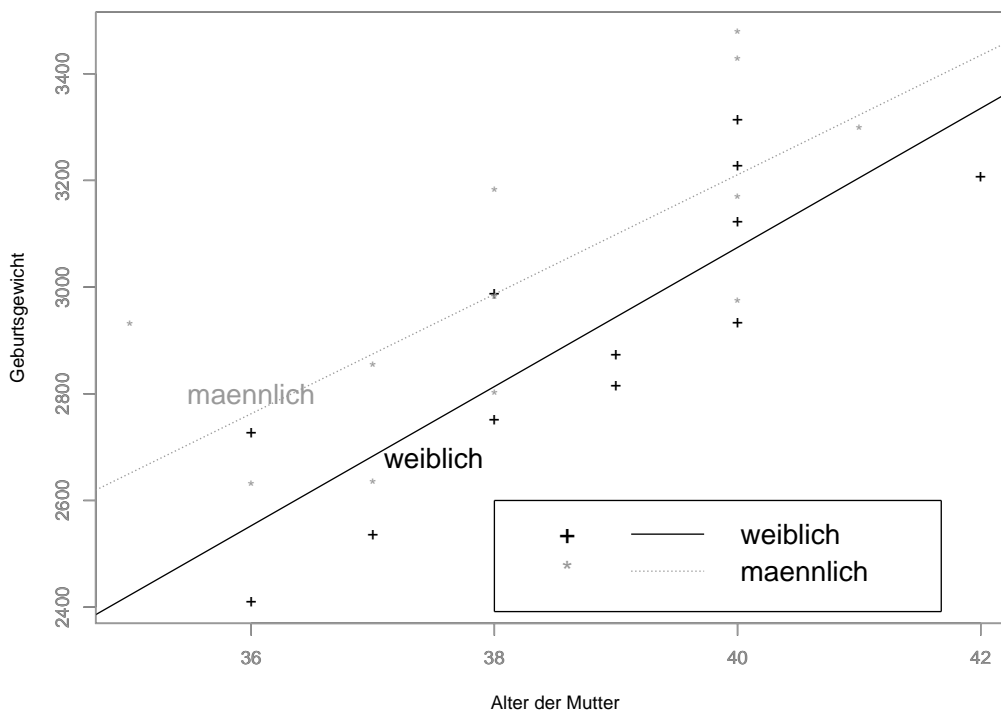


Abbildung 3.9: KQ-Geraden in den 2 Gruppen

Variable handelt.

Dies trifft nicht auf numerische Kodierungen zu wie

$$\tilde{x}_2 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 3 \\ 3 \\ 3 \\ 4 \end{pmatrix} \begin{array}{l} 1 \hat{=} \text{„ledig“} \\ 2 \hat{=} \text{„verheiratet“} \\ 3 \hat{=} \text{„geschieden“} \\ 4 \hat{=} \text{„verwitwet“} \end{array}$$

Eine Behandlung der zunächst numerischen Variablen  $\tilde{x}_2 = \text{famstand}$  als Faktor muss durch die Kommandos

```
f.famstand <- factor(famstand)
options(contrasts = c("contr.treatment", "contr.poly"))
out.add <- lm(y ~ age + f.famstand)
```

z.B. für y (Einkommen) erzwungen werden.

Die zugehörige Design-Matrix lautet:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & 1 & 0 & 0 \\ 1 & . & 1 & 0 & 0 \\ 1 & . & 1 & 0 & 0 \\ 1 & x_{41} & 1 & 0 & 0 \\ 1 & x_{51} & 0 & 1 & 0 \\ 1 & x_{61} & 0 & 1 & 0 \\ 1 & x_{71} & 0 & 1 & 0 \\ 1 & x_{81} & 0 & 0 & 1 \\ 1 & x_{91} & 0 & 0 & 1 \\ 1 & x_{10,1} & 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{n,5}, \text{Rang}(X) = 5$$

Zur exakten Modellierung mit beliebigen Achsenabschnitten und Steigungen in den möglichen Gruppen müssen wiederum Interaktionsterme berücksichtigt werden.

$$\begin{aligned} \text{out.ia} &<- \text{lm} ( y \sim \text{age} * \text{f.famstand}, x = T) \\ &\hat{=} \text{lm} ( y \sim \text{age} + \text{f.famstand} + \text{age} : \text{f.famstand}, x = T) \end{aligned}$$

Der Befehl `out.ia $ x` liefert die resultierende Design-Matrix (nur möglich, falls im vorangegangenen `lm`-Statement "`..., x = T`" gesetzt wurde).

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & 1 & 0 & 0 & x_{11} & 0 & 0 \\ 1 & . & 1 & 0 & 0 & . & 0 & 0 \\ 1 & . & 1 & 0 & 0 & . & 0 & 0 \\ 1 & x_{41} & 1 & 0 & 0 & x_{41} & 0 & 0 \\ 1 & x_{51} & 0 & 1 & 0 & 0 & x_{51} & 0 \\ 1 & x_{61} & 0 & 1 & 0 & 0 & x_{61} & 0 \\ 1 & x_{71} & 0 & 1 & 0 & 0 & x_{71} & 0 \\ 1 & x_{81} & 0 & 0 & 1 & 0 & 0 & x_{81} \\ 1 & x_{91} & 0 & 0 & 1 & 0 & 0 & x_{91} \\ 1 & x_{10,1} & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Die konkrete Behandlung mehrstufiger kategorialer Kovariablen wird nochmals an einem Beispiel erläutert.

### 3.13.2 Beispiel 2: Mehrstufiger Faktor

Beispiel:  $Y$  = Benzinverbrauch auf 100 km

$X_1$  = Autogewicht

$X_2$  = Automarke (VW, BMW, Mazda)

Wir beginnen wieder mit einer numerischen Kodierung, die aber kein befriedigendes Ergebnis liefert.

Erster Versuch: 1-2-3-Kodierung

$$X_2 = \begin{cases} 1 & \text{VW} \\ 2 & \text{BMW} \\ 3 & \text{Mazda} \end{cases}$$

$$\text{Modell: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

Interpretation: VW  $Y = \beta_0 + \beta_1 X_1 + \beta_2 + \epsilon = (\beta_0 + \beta_2) + \beta_1 X_1 + \epsilon$   
 BMW  $Y = \beta_0 + \beta_1 X_1 + \beta_2 \cdot 2 + \epsilon = (\beta_0 + 2\beta_2) + \beta_1 X_1 + \epsilon$   
 Mazda  $Y = \beta_0 + \beta_1 X_1 + \beta_2 \cdot 3 + \epsilon = (\beta_0 + 3\beta_2) + \beta_1 X_1 + \epsilon$

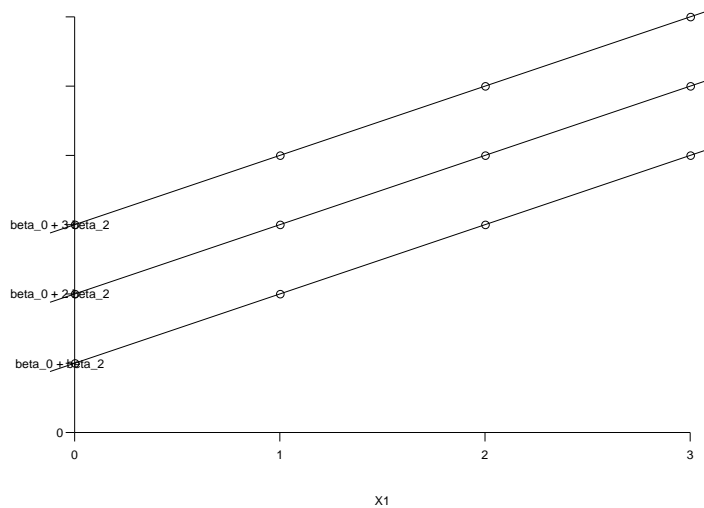


Abbildung 3.10: Unzulässige Festlegung des Intercepts; parallele KQ-Geraden

Die Grafik 3.10 zeigt, dass der Abstand der Intercepts durch die Festlegung (1,2,3) schon festgelegt ist. Außerdem impliziert der linear additive Ansatz die Parallelität der drei KQ-Geraden, die natürlich a priori nicht klar ist.

Man möchte verschiedene Intercepts und Steigungen zulassen. Dies lässt sich – wie gesehen – mit Dummy- Variablen modellieren.

$$z_{i1} = \begin{cases} 1 & \text{falls Auto } i \text{ ein VW} \\ 0 & \text{sonst} \end{cases}$$

$$z_{i2} = \begin{cases} 1 & \text{falls Auto } i \text{ ein BMW} \\ 0 & \text{sonst} \end{cases}$$

$$z_{i3} = \begin{cases} 1 & \text{falls Auto } i \text{ ein Mazda} \\ 0 & \text{sonst} \end{cases}$$

Wegen  $z_{i1} + z_{i2} + z_{i3} = 1 \forall i = 1, \dots, n$ , hätte man eine perfekte Multikollinearität. Die Entfernung der dritten Dummy-Variablen stellt keinen Informationsverlust dar, denn

$$z_{i1} = z_{i2} = 0 \implies z_{i3} = 1.$$

Zweiter Versuch:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 Z_1 + \beta_3 Z_2 + \epsilon$$

Intercepts für VW  $\beta_0 + \beta_2$ ,  
 BMW  $\beta_0 + \beta_3$  und  
 Mazda  $\beta_0$ .

Damit ist erreicht, dass im Allgemeinen alle Intercepts verschieden sind, die Steigungen jedoch noch nicht. Hierzu müssten analog zum einführenden Beispiel noch Interaktionen berücksichtigt werden.

## 3.14 Kategoriale und Indikatorvariablen; Teil 2: ANOVA-Modelle

### 3.14.1 Einfaktorielle Varianzanalyse; one way ANOVA

#### Modell

Einfache lineare Regression mit einem qualitativen Regressor  $Z$ .

$Z$  habe  $k$  Kategorien. Definiere

$$z_{il} := \begin{cases} 1 & \text{Beobachtung } i \text{ fällt in Kategorie } l \\ 0 & \text{sonst} \end{cases} \quad l = 1, \dots, k, \quad i = 1, \dots, n$$

Modell:  $y_i = \beta_0 + \alpha_1 z_{i1} + \dots + \alpha_{k-1} z_{i,k-1} + \epsilon_i$

Bemerkungen:

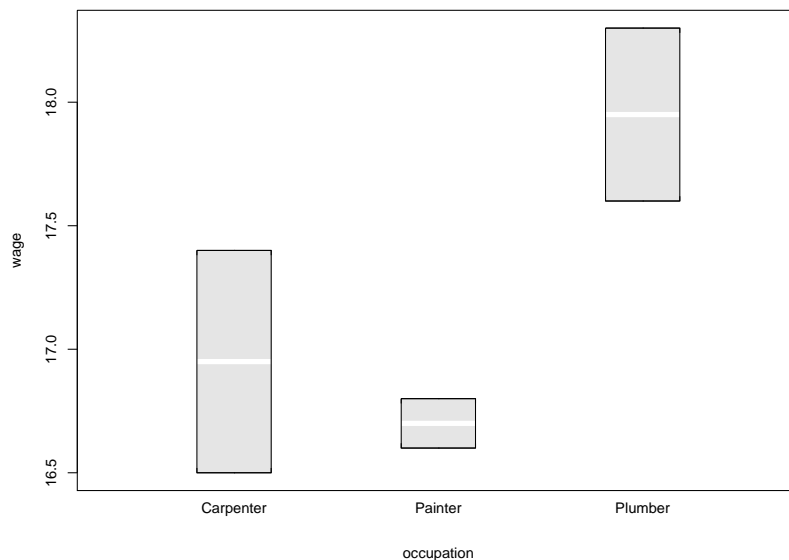
- Erwarteter Response in Kategorie  $l = \beta_0 + \alpha_l$ ,  $l = 1, \dots, k - 1$   
Erwarteter Response in Kategorie  $k = \beta_0$
- $H_0$ : keine Unterschiede zwischen den Kategorien  $\Leftrightarrow H_0 : \alpha_1 = \dots = \alpha_{k-1} = 0$
- Dies ist ein Spezialfall des allgemeinen F-Tests:

$$\alpha_1 = \dots = \alpha_{k-1} = 0 \Leftrightarrow \mathbf{C}\boldsymbol{\beta} = \mathbf{0}, \quad \mathbf{C} = \begin{pmatrix} 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{pmatrix} \in \mathbb{R}^{(k-1),k}$$

$$\boldsymbol{\beta}^T = (\beta_0, \alpha_1, \dots, \alpha_{k-1})$$

Gemischte Modelle mit quantitativen Kovariablen und genau einem qualitativen Regressor werden analog zu obigen Beispielen behandelt.

## Beispiel: one way ANOVA

Abbildung 3.11: Boxplot: Response *wage* innerhalb der Faktorstufen

## Die Daten

```
wage  occupation
17.6   Plumber
18.3   Plumber
16.8   Painter
16.7   Painter
16.6   Painter
17.4   Carpenter
16.5   Carpenter
```

geben die Einkünfte verschiedener Handwerker wieder. Die Variabilität der Zielvariablen *wage* innerhalb der Faktorstufen wird in der Regel durch Boxplots dargestellt; siehe Grafik 3.11.

Man erkennt, dass die Installateure mehr verdienen als die anderen genannten Berufe. Aber ist der im Plot sichtbare Unterschied auch statistisch signifikant?

Die S-Plus-Kommandos

```
r <- aov(wage ~ occupation)
summary(r)
```

liefern

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
occupation	2	1.964286	0.9821429	5.863539	0.06468802
Residuals	4	0.670000	0.1675000		

Der P-Wert über 5% besagt, dass die Unterschiede innerhalb der Faktorstufen **nicht signifikant** sind. Der Grund hierfür ist die geringe Fallzahl von nur zwei Beobachtungen pro Zelle.



### 3.14.2 Zweifaktorielles Design; two way ANOVA

#### Balanced Design; K Beobachtungen pro Zelle

In der Literatur ist es üblich, mit folgendem Spezialfall zu beginnen:

Daten:  $y_{ijk}$  = k-ter Response des Faktors A in Stufe (Level) i und Faktor B in Stufe j,  
 $i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K$ .

Es gibt also in jeder Zelle bzgl. der Faktorstufen genau  $K$  Beobachtungen.

Beispiel: Aktivitätslevel von Kindern

		Drug	
		Placebo	Ritalin
Group	normal	50	67
		45	60
		55	58
		52	65
hyperaktiv	70	51	
	72	57	
	68	48	
	75	55	

In Faktor A = Group und in Faktor B = Drug gibt es jeweils 2 Stufen (levels).  
 $K = 4, I = 2 = J$ , z.B.  $y_{121} = 67, y_{124} = 65$ .

Explorative Graphische Analyse:

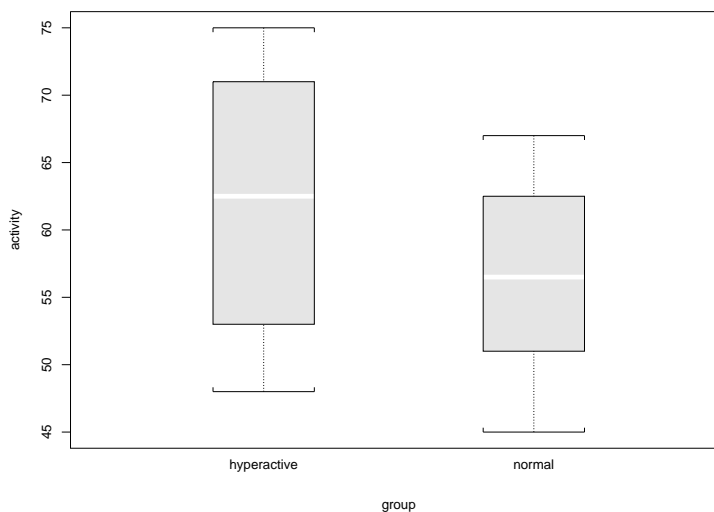


Abbildung 3.12: Boxplot: Response *activity* innerhalb der Faktorstufen von *group*

Die Bilder 3.12 und 3.13 zeigen gewisse Unterschiede in der Streuung des Response innerhalb der Faktorstufen, während die Mediane nur "wenig" voneinander abweichen. Dementsprechend liefert auch

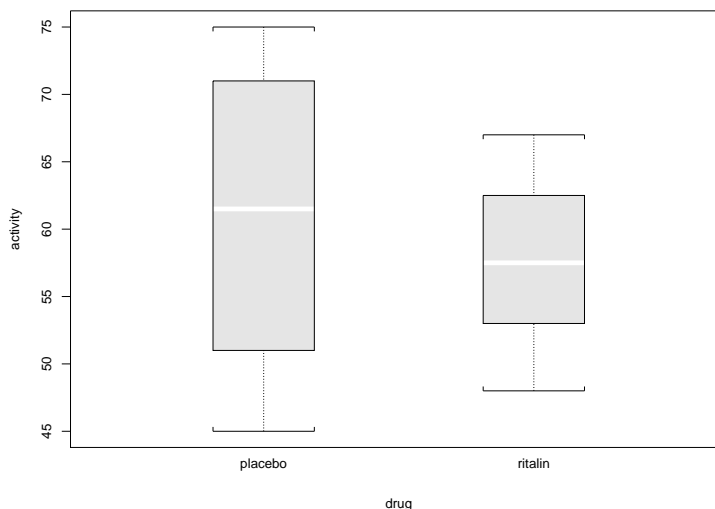


Abbildung 3.13: Boxplot: Response *activity* innerhalb der Faktorstufen von *drug*

```
r <- aov(activity ~ group + drug)
```

```
summary(r)
```

das Ergebnis

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
group	1	121.00	121.0000	1.414886	0.2555190
drug	1	42.25	42.2500	0.494041	0.4945256
Residuals	13	1111.75	85.5192		

Die P-Werte von 25% bzw. 49% scheinen darauf hinzudeuten, dass kein Zusammenhang zwischen dem Response und den Faktoren besteht. Diese Analyse ist jedoch falsch.

Da  $K$  Wiederholungen für jede Faktorkombination vorhanden sind, kann man diese benutzen, um eine Schätzung des erwarteten Response für jede Faktorkombination zu bekommen, d.h.

$$\bar{y}_{ij\cdot} := \frac{1}{K} \sum_{k=1}^K y_{ijk}$$

schätzt den Erwartungswert des Response in der  $ij$ -ten Zelle. Davon ausgehend dokumentiert die Grafik 3.14 einen lupenreinen Interaktionseffekt zwischen den Faktoren *group* und *drug*. Während das Placebo die hohe Aktivität bei den Hyperaktiven erwartungsgemäß nicht reduziert und bei der Kontrollgruppe die auch sonst übliche Aktivität auslöst, wirkt Ritalin umgekehrt: Die hohe Aktivität in der hyperaktiven Gruppe wird stark reduziert, dagegen steigt seltsamerweise die Aktivität in der Kontrollgruppe über das Mittelmaß hinaus, wenn Ritalin verabreicht wird.

Dies wird auch in S-Plus mit der allgemeineren Interaktions-Modellierung deutlich:

```
r <- aov(activity ~ group * drug)
```

```
summary(r)
```

liefert

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
group	1	121.00	121.000	8.00000	0.0152201
drug	1	42.25	42.250	2.79339	0.1205062
group:drug	1	930.25	930.250	61.50413	0.0000046
Residuals	12	181.50	15.125		

und damit einen hochsignifikanten Interaktionseffekt. In diesem größeren Modell verändern sich auch die Signifikanzen der Faktoren selbst im Vergleich zum rein additiven Modell.

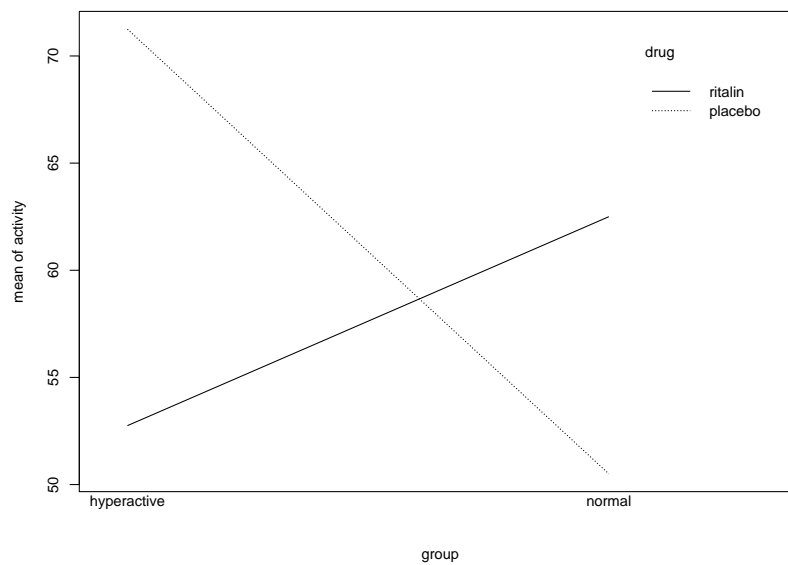


Abbildung 3.14: Interaktions-Plot: Response *activity* bzgl. der Faktoren *group* und *drug*

### Modellansätze für die Parameter in ANOVA

#### 1. ohne Interaktion

$$y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk} \quad \begin{array}{l} i = 1, \dots, I - 1 \\ j = 1, \dots, J - 1 \\ k = 1, \dots, K \end{array}$$

$$\alpha_I = 0, \quad \beta_J = 0$$

$\mu$  = overall mean

$\alpha_i$  = Effekt von Faktor A in Stufe (level)  $i$

$\beta_j$  = Effekt von Faktor B in Stufe (level)  $j$

$$H_0 : \text{kein Faktor A Effekt} \Rightarrow H_0 : \alpha_1 = \dots = \alpha_{I-1} = 0$$

$$H_0 : \text{kein Faktor B Effekt} \Rightarrow H_0 : \beta_1 = \dots = \beta_{J-1} = 0$$

Für ein ANOVA-Modell ohne Interaktion gilt:

$$E(y_{ijk}) = \mu + \alpha_i + \beta_j.$$

Also müssen  $\bar{y}_{ij}$  parallel verlaufen, d.h. der Plot sollte im Idealfall qualitativ wie die Grafik 3.15 aussehen. Falls dies nicht hinreichend gilt, so muss man Interaktionen zwischen den Faktoren erlauben.

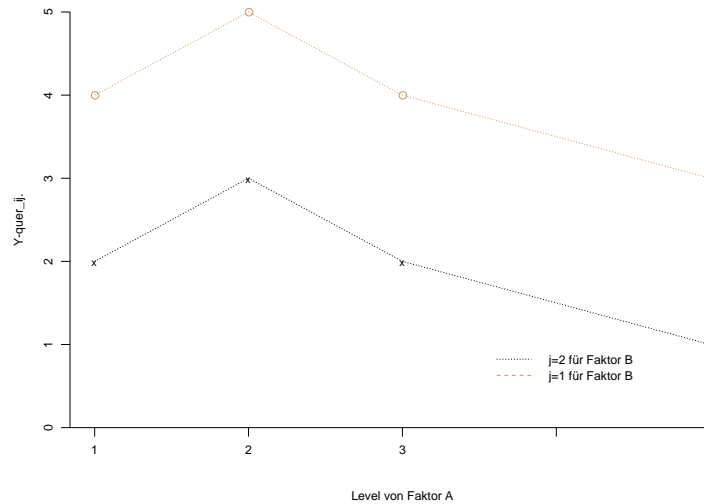


Abbildung 3.15: Two-Way-Anova ohne Interaktion

## 2. mit Interaktion

$$\begin{aligned}
 y_{ijk} &= \mu + \alpha_i + \beta_j + \theta_{ij} + \epsilon_{ijk}, & i &= 1, \dots, I \\
 & & j &= 1, \dots, J \\
 & & k &= 1, \dots, K
 \end{aligned}$$

Identifizierbarkeitsbedingungen:

$$\begin{aligned}
 \alpha_1 + \dots + \alpha_I &= 0 \\
 \beta_1 + \dots + \beta_J &= 0 \\
 \sum_{i=1}^I \theta_{ij} &= \sum_{j=1}^J \theta_{ij} = 0
 \end{aligned}$$

Auch hier werden in der Regel die Gleichungsnebenbedingungen durch Einsetzen eliminiert, so dass man für obiges Beispiel (diesmal mit Effekt-Kodierung) folgende Designmatrix erhält:

$$\begin{array}{l}
 i = 1, j = 1 \\
 \\
 i = 2, j = 1 \\
 \\
 i = 1, j = 2 \\
 \\
 i = 2, j = 2
 \end{array}
 \begin{pmatrix}
 50 \\
 45 \\
 55 \\
 52 \\
 \\
 70 \\
 72 \\
 68 \\
 75 \\
 \\
 67 \\
 60 \\
 58 \\
 65 \\
 \\
 51 \\
 57 \\
 48 \\
 55
 \end{pmatrix}
 =
 \begin{pmatrix}
 1 & 1 & 1 & 1 \\
 1 & 1 & 1 & 1 \\
 1 & 1 & 1 & 1 \\
 1 & 1 & 1 & 1 \\
 \\
 1 & -1 & 1 & -1 \\
 1 & -1 & 1 & -1 \\
 1 & -1 & 1 & -1 \\
 1 & -1 & 1 & -1 \\
 \\
 1 & 1 & -1 & -1 \\
 1 & 1 & -1 & -1 \\
 1 & 1 & -1 & -1 \\
 1 & 1 & -1 & -1 \\
 \\
 1 & -1 & -1 & 1 \\
 1 & -1 & -1 & 1 \\
 1 & -1 & -1 & 1 \\
 1 & -1 & -1 & 1
 \end{pmatrix}
 \cdot
 \begin{pmatrix}
 \mu \\
 \alpha_1 \\
 \beta_1 \\
 \theta_{11}
 \end{pmatrix}
 + \epsilon$$

Einige denkbare Hypothesen sind:

$$\begin{array}{ll}
 H_0 : \text{keine Interaktion} & \Leftrightarrow H_0 : \theta_{11} = 0 \\
 H_0 : \text{keine Interaktion} & \Leftrightarrow H_0 : \theta_{11} = 0 \\
 \quad \text{kein Faktor A Effekt} & \alpha_1 = 0 \\
 H_0 : \text{keine Interaktion} & \Leftrightarrow H_0 : \theta_{11} = 0 \\
 \quad \text{kein Faktor B Effekt} & \beta_1 = 0 \\
 H_0 : \text{keine Interaktion} & \Leftrightarrow H_0 : \theta_{11} = 0 \\
 \quad \text{keine Faktor A und B Effekte} & \alpha_1 = \beta_1 = 0
 \end{array}$$

Identifizierbarkeitsbedingungen:

$$\begin{array}{l}
 \alpha_1 + \alpha_2 = 0 \\
 \beta_1 + \beta_2 = 0 \\
 \theta_{11} + \theta_{12} = 0, \quad \theta_{11} + \theta_{21} = 0 \\
 \theta_{21} + \theta_{22} = 0, \quad \theta_{12} + \theta_{22} = 0
 \end{array}$$

## Unbalanced Design; beliebig viele Beobachtungen pro Zelle

Beispiel: Genotypus mit zwei Faktoren

In einem Beispiel aus Scheffé (1959, pp. 139-140) wird das 28-Tage-Gewicht/Wurf von Tieren untersucht, die von Pflegemüttern aufgezogen wurden. Man unterscheidet in der betreffenden Art 4 Genotypen (Faktor A: Genotypus des Wurfes, Faktor B: Genotypus der Pflegemutter).

Man möchte anhand des Experimentes feststellen, ob die Pflegemütter Kinder des gleichen Typs besser betreuen. In diesem Fall müsste das Gewicht signifikant vom Mittelwert abweichen.

Einige Daten:

weight	A	B
61.5	1	1
68.2	1	1
64.0	1	1
65.0	1	1
59.7	1	1
55.0	1	2
42.0	1	2
60.2	1	2
52.5	1	3
61.8	1	3
49.5	1	3
52.7	1	3
42.0	1	4
...	..	..

Wir haben 5 Beobachtungen in der 1-1-Zelle, 3 Beobachtungen in der 1-2-Zelle, 4 Beobachtungen in der 1-3-Zelle usw., also keine "balanced Design" wie oben. Die Analyse kann aber mit S-Plus formal in gleicher Weise erfolgen.

Explorative Graphische Analyse:

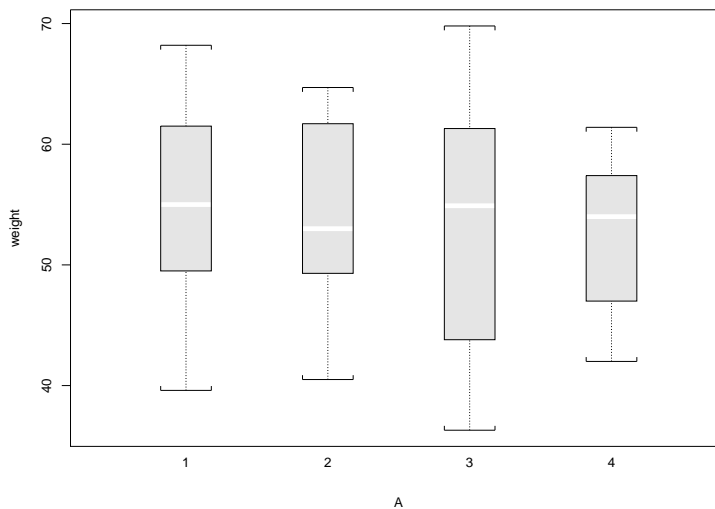
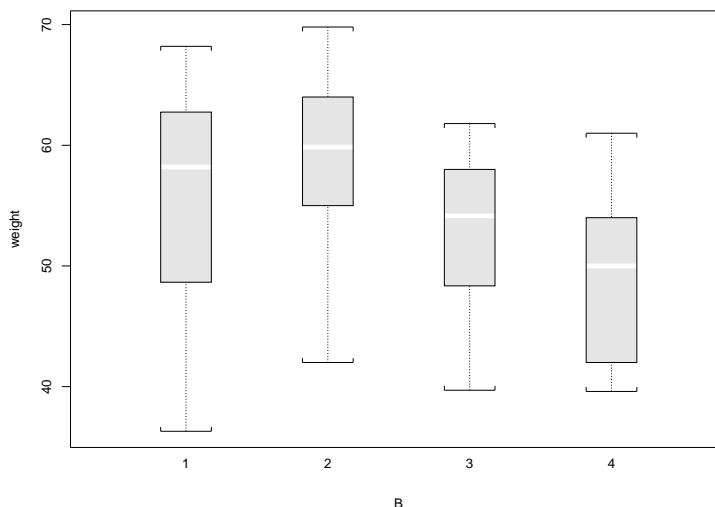
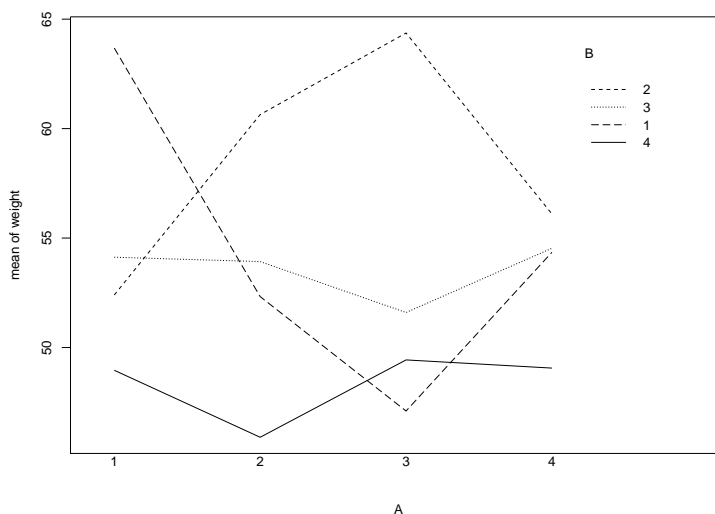


Abbildung 3.16: Response *weight* innerhalb des Faktors *A*: Genotypus Wurf

Abbildung 3.17: Response *weight* innerhalb des Faktors *B*: Genotypus PflegemutterAbbildung 3.18: Interaktions-Plot: Response *weight* bzgl. der Faktoren *A* und *B*

Natürlich werden konkrete Interaktions-Plots von den Idealdarstellungen der Bilder 3.13 und 3.14 abweichen. Trotzdem deutet in der Abbildung 3.18 nur wenig auf mögliche Interaktionen hin. Dies bestätigen auch die folgenden Rechnungen in S-Plus.

```
r <- aov(weight ~ A * B)
```

```
summary(r)
```

liefert mit Interaktionen

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
A	3	60.157	20.0524	0.369696	0.775221
B	3	775.081	258.3602	4.763246	0.005736
A:B	9	824.073	91.5636	1.688108	0.120053
Residuals	45	2440.816	54.2404		

Der Interaktionsterm ist mit einem P-Wert von 12% nicht signifikant. Also kann das rein additive Modell gewählt werden:

```
r <- aov(weight ~ A + B)
summary(r)
```

mit dem Resultat:

	Df	Sum of Sq	Mean Sq	F Value	Pr(F)
A	3	60.157	20.0524	0.331659	0.8024695
B	3	775.081	258.3602	4.273178	0.0088605
Residuals	54	3264.889	60.4609		

Dieses weist nur auf einen signifikanten Einfluss des Faktors B hin.

Fazit:

1. Die Ausgangsfragestellung des Experiments kann nicht bestätigt werden:  
Würden nämlich Pflegemütter Kinder des gleichen Typs besser betreuen, so wären signifikante Interaktionen aufgetreten. Das war nicht der Fall.
2. Der signifikante Einfluss des Faktors B besagt, dass gewisse Genotypen von Müttern besser zur Aufzucht geeignet sind. Der Boxplot 3.17 weist den Genotypus 2 als besonders gute, und den Genotypus 4 als weniger geeignete Pflegemutter aus.



# Kapitel 4

## Anhang

### 4.1 Normalgleichungen

Seien  $\mathbf{X} \in \mathbb{R}^{n,p}$ ,  $n \geq p$ ,  $\mathbf{y} \in \mathbb{R}^n$  gegebene Daten und  $\boldsymbol{\beta} \in \mathbb{R}^p$  ein unbekannter Parametervektor. Dann ist jede Lösung der Normalgleichungen

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y} \quad (4.1)$$

Lösung des KQ-Problems

$$\min_{\boldsymbol{\beta}} \{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2\}. \quad (4.2)$$

Die Normalgleichungen (4.1) haben in jedem Fall mindestens eine Lösung. Falls  $\text{Rang}(\mathbf{X}) = p$ , so lautet die eindeutige Lösung von (4.2)

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (4.3)$$

Besonders einfach stellt sich die Lösung von (4.2) dar für

$$\mathbf{X} = \mathbf{Q} = (\mathbf{q}_1, \dots, \mathbf{q}_p) \quad (4.4)$$

und

$$\mathbf{q}_i^T \mathbf{q}_j = \delta_{ij},$$

wenn also  $\mathbf{Q}$  orthonormierte Spalten hat. In diesem Fall vereinfacht sich (4.3) zu:

$$\mathbf{b} = \mathbf{Q}^T \mathbf{y}.$$

#### Lemma 4.1

*Lösungen des KQ-Problems*

$$\min_{\boldsymbol{\beta}} \{Q(\boldsymbol{\beta})\}, \text{ wobei } Q(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2,$$

*sind durch 0 nach unten beschränkt. Nach einem Satz aus der Theorie quadratischer Funktionen muss die Funktion  $Q$  ihr Minimum in mindestens einem Punkt  $\mathbf{b} \in \mathbb{R}^p$  annehmen. Hier gilt notwendigerweise*

$$\nabla Q(\mathbf{b}) = \mathbf{0} \quad \text{und das ist äquivalent zu} \quad (4.1).$$

*Damit haben die Normalgleichungen immer eine Lösung.*

## 4.2 Rechnen mit Erwartungswerten und Kovarianzmatrizen

### 4.2.1 n-dimensionale Normalverteilung

#### Def. 4.2 (n-dimensionale Normalverteilung $N(\boldsymbol{\mu}, \mathbf{C})$ )

Seien  $\boldsymbol{\mu} \in \mathbb{R}^n$  und  $\mathbf{C} = (c_{ij}) \in \mathbb{R}^{n,n}$  eine positiv definite Matrix (damit ist  $\mathbf{C}$  auch symmetrisch).  $\mathbf{X} = (X_1, \dots, X_n)^T$  heißt n-dimensional normalverteilt, wenn es eine Dichte der Form

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det(\mathbf{C})}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (4.5)$$

besitzt. Dabei bezeichnet  $\det(\mathbf{C})$  die Determinante von  $\mathbf{C}$ .

#### Def. 4.3 (2-dimensionale Normalverteilung)

Mit  $(\mu_Z, \mu_Y) \in \mathbb{R}^2$ ,  $\sigma > 0, \tau > 0$  und  $-1 < \rho < 1$  ist die Kovarianzmatrix (vgl. Def. 4.4)

$$\mathbf{C} := \begin{pmatrix} \sigma^2 & \rho \sigma \tau \\ \rho \sigma \tau & \tau^2 \end{pmatrix}$$

positiv definit und  $\det(\mathbf{C}) = \sigma^2 \tau^2 (1 - \rho^2)$ . Die Inverse  $\mathbf{C}^{-1}$  ist explizit berechenbar und es gilt gemäß (4.5)

$$f_{Z,Y}(z, y) = \frac{1}{2\pi\sigma\tau\sqrt{1-\rho^2}} \exp \left[ -\frac{1}{2(1-\rho^2)} \left( \frac{(z-\mu_Z)^2}{\sigma^2} - \frac{2\rho(z-\mu_Z)(y-\mu_Y)}{\sigma\tau} + \frac{(y-\mu_Y)^2}{\tau^2} \right) \right] \quad (4.6)$$

### 4.2.2 Erwartungswerte; n-dim.

Es sei hier nochmals an die Vereinbarung erinnert, daß  $\mathbf{X} = (X_1, \dots, X_n)^T \in \mathbb{R}^n$  und  $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n$  mit der Indizierung  $x_k, k = 1, \dots, n$ .

**Def. 4.4 (Kovarianzen, Kovarianzmatrix)**

Für  $\mathbf{X} = (X_1, \dots, X_n)^T$  mögen alle zweiten Momente

$$\tau_{kl} := E(X_k X_l), \quad 1 \leq k, l \leq n,$$

existieren. Mit  $\mu_k := E(X_k)$  nennt man

$$\sigma_{kl} := \text{Cov}(X_k, X_l) := E[(X_k - \mu_k)(X_l - \mu_l)], \quad 1 \leq k, l \leq n,$$

die **Kovarianzen** von  $X_k$  und  $X_l$ .

$$\mathbf{C} := \mathbf{Cov}(\mathbf{X}) := (\sigma_{kl})_{1 \leq k, l \leq n} \in \mathbb{R}^{n, n}$$

heißt **Kovarianzmatrix** von  $\mathbf{X}$ .

$X_k$  und  $X_l$  heißen **unkorreliert**, falls  $\text{Cov}(X_k, X_l) = 0$ .

Die Existenz aller zweiten Momente von  $\mathbf{X}$  sichert also, dass alle Erwartungswerte  $\mu_k$ , alle Varianzen  $\sigma_{kk}$  und alle Kovarianzen  $\sigma_{kl}$  der Komponenten  $X_k, X_l$ ,  $1 \leq k, l \leq n$  von  $\mathbf{X}$  wohldefiniert sind.

**Folgerung 4.5 ((Ko-)Varianzen, Verschiebungsregel)**

1. Für  $k = l$  gilt

$$\begin{aligned} \sigma_k^2 := \sigma_{kk} = \text{Cov}(X_k, X_k) &= E[(X_k - \mu_k)(X_k - \mu_k)] \\ &= E[(X_k - \mu_k)^2] = \text{Var}(X_k). \end{aligned} \quad (4.7)$$

2. Auch für Kovarianzen gibt es eine **Verschiebungsregel**, nämlich

$$\sigma_{kl} = \text{Cov}(X_k, X_l) = E[(X_k - \mu_k)(X_l - \mu_l)] = E(X_k X_l) - \mu_k \mu_l. \quad (4.8)$$

**Def. 4.6 (Korrelationskoeffizient)**

Für zwei ZV, z.B.  $(Z, Y) = (X_k, X_l)$ , sollen alle zweiten Momente existieren. Ferner gelte  $\sigma_Z^2 := \text{Var}(Z) > 0$  und  $\sigma_Y^2 := \text{Var}(Y) > 0$ . Dann heißt

$$\rho(Z, Y) := \frac{\text{Cov}(Z, Y)}{\sigma_Z \sigma_Y} = \frac{\text{Cov}(Z, Y)}{D(Z) D(Y)} = \frac{\text{Cov}(Z, Y)}{\sqrt{\text{Var}(Z) \text{Var}(Y)}}$$

der **Korrelationskoeffizient** von  $Z$  und  $Y$ .

**Satz 4.7 (Korrelationskoeffizient)**

Falls alle zweiten Momente der ZV  $(Z, Y)$  existieren und falls  $\text{Var}(Z) > 0$  und  $\text{Var}(Y) > 0$ , so gilt:

$$-1 \leq \rho(Z, Y) \leq 1$$

und  $\rho(Z, Y) = \pm 1$  genau dann, wenn

$$P(Y = a + bZ) = P(\{\omega \mid Y(\omega) = a + bZ(\omega)\}) = 1$$

für geeignete  $a \in \mathbb{R}$  und  $b \neq 0$ . Ferner gilt für  $|\rho(Z, Y)| = 1$ :

$$\rho(Z, Y) = 1 \iff b > 0 \quad \text{und} \quad \rho(Z, Y) = -1 \iff b < 0.$$

**Beispiel 4.8 (Zweidimensionale Normalverteilung)**

von  $(Z, Y)$  bzgl. Def. 4.3 mit den Parametern  $\boldsymbol{\mu} := (\mu_z, \mu_y)^T \in \mathbb{R}^2$  sowie  $\sigma_z > 0, \sigma_y > 0$  und  $-1 < \rho < 1$  und der positiv definiten Kovarianzmatrix

$$\mathbf{C} := \mathbf{Cov}(Z, Y) = \begin{pmatrix} \sigma_z^2 & \rho \sigma_z \sigma_y \\ \rho \sigma_z \sigma_y & \sigma_y^2 \end{pmatrix}.$$

Dort gilt  $E(Z) = \mu_z, E(Y) = \mu_y, \text{Var}(Z) = \sigma_z^2, \text{Var}(Y) = \sigma_y^2, \rho(Z, Y) = \rho$ .

**Def. 4.9 (EW von Zufallsvektoren, -matrizen)**

Vorausgesetzt für  $k = 1, \dots, n$  existieren die Erwartungswerte  $\mu_k := E(X_k)$ , so setzt man für  $\mathbf{X} = (X_1, \dots, X_n)^T$

$$E(\mathbf{X}) := \boldsymbol{\mu} := (\mu_1, \dots, \mu_n)^T.$$

Analog (also auch elementweise) wird der Erwartungswert von Matrizen definiert, deren Elemente ZV sind.

**Satz 4.10 (EW von Zufallsvektoren)**

Für  $\mathbf{X} = (X_1, \dots, X_n)^T$  existiere  $\boldsymbol{\mu} = E(\mathbf{X})$ , dann gilt für (nichtstochastisches)  $\mathbf{A} \in \mathbb{R}^{m,n}, m \in \mathbb{N}$

$$E(\mathbf{A} \mathbf{X}) = \mathbf{A} E(\mathbf{X}) = \mathbf{A} \boldsymbol{\mu}. \quad (4.9)$$

**Satz 4.11 (Kovarianzen von Zufallsvektoren; Varianz von Summen)**

Vorausgesetzt alle zweiten Momente von  $\mathbf{X} = (X_1, \dots, X_n)^T$  existieren, dann gilt mit  $\boldsymbol{\mu} = E(\mathbf{X})$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{a} = (a_1, \dots, a_n)^T \in \mathbb{R}^n$ :

1. Die Kovarianzmatrix  $\mathbf{C}$  ist positiv semidefinit.

2.

$$\begin{aligned} \mathbf{C} = (\sigma_{kl}) = \mathbf{Cov}(\mathbf{X}) &= E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] \\ &= E[\mathbf{X} \mathbf{X}^T] - E[\boldsymbol{\mu} \boldsymbol{\mu}^T] \end{aligned}$$

3.

$$\mathbf{Cov}(\mathbf{A} \mathbf{X}) = \mathbf{A} \mathbf{C} \mathbf{A}^T \quad (4.10)$$

4.

$$\begin{aligned} \text{Var}(\mathbf{a}^T \mathbf{X}) = \text{Var}\left(\sum_{i=1}^n a_i X_i\right) &= \text{Cov}(\mathbf{a}^T \mathbf{X}) = \mathbf{a}^T \mathbf{C} \mathbf{a} = \sum_{k,l} a_k a_l \sigma_{kl} \\ &= \sum_{k=1}^n \sum_{l=1}^n a_k a_l \text{Cov}(X_k, X_l) \quad (4.11) \\ &= \sum_{k=1}^n a_k^2 \text{Var}(X_k) + 2 \sum_{k < l} a_k a_l \text{Cov}(X_k, X_l) \end{aligned}$$

5.

$$\begin{aligned} \text{Var}(\mathbf{a}^T \mathbf{X}) = \text{Var}\left(\sum_{i=1}^n a_i X_i\right) &= \sum_{k=1}^n a_k^2 \text{Var}(X_k), \quad (4.12) \\ &\text{falls } (X_k, X_l) \text{ unkorreliert.} \end{aligned}$$

## 4.3 Transformationen für Dichten und unabhängige normalverteilte Zufallsvariablen

### 4.3.1 Transformationssatz für Dichten

Die  $n$ -dimensionale Zufallsvariable  $\mathbf{X} = (X_1, \dots, X_n)^T$  habe eine Dichte  $f(\mathbf{x})$ ,

$\mathbf{x} = (x_1, \dots, x_n)^T$ . Weiterhin seien  $B := \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) > 0\}$  und  $\mathbf{h} : B \rightarrow D$ ,  $D \subset \mathbb{R}^n$ , differenzierbar und bijektiv. Dann ist  $\mathbf{Y} := \mathbf{h}(\mathbf{X})$  wieder eine  $n$ -dim. ZV. Mit  $\mathbf{y} = (y_1, \dots, y_n)^T$  sei die Umkehrabbildung

$$\mathbf{x} := \mathbf{h}^{-1}(\mathbf{y}) = (x_1(\mathbf{y}), \dots, x_n(\mathbf{y}))^T$$

auf  $D$  (abgesehen von endlich vielen glatten Hyperflächen) stetig differenzierbar.

$$\frac{\partial(x_1, \dots, x_n)}{\partial(y_1, \dots, y_n)} = \frac{\partial(x_1, \dots, x_n)}{\partial(y_1, \dots, y_n)}(\mathbf{y}) := \det \begin{pmatrix} \frac{\partial x_1}{\partial y_1}(\mathbf{y}) & \cdots & \frac{\partial x_1}{\partial y_n}(\mathbf{y}) \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial y_1}(\mathbf{y}) & \cdots & \frac{\partial x_n}{\partial y_n}(\mathbf{y}) \end{pmatrix}$$

sei die Funktionaldeterminante der Umkehrabbildung  $\mathbf{h}^{-1}$ . Dann hat die ZV  $\mathbf{Y} = (Y_1, \dots, Y_n)^T := \mathbf{h}(\mathbf{X})$  wieder eine Dichte  $g(\mathbf{y})$ , die gegeben ist durch

$$g(\mathbf{y}) = g(y_1, \dots, y_n) = f(x_1(\mathbf{y}), \dots, x_n(\mathbf{y})) \left| \frac{\partial(x_1, \dots, x_n)}{\partial(y_1, \dots, y_n)} \right|, \quad (4.13)$$

wobei  $|\cdot|$  hier den Betrag einer Determinante bezeichnet.

### 4.3.2 Lineare Transformationen

Wir betrachten nun eine weitere wichtige Klasse linearer Transformationen und folgen in der Darstellung [Stirzaker (1994)], pp. 287.

#### Satz 4.12 (Lineare Transformation)

Zur Matrix  $\mathbf{A} = (a_{ij})$  existiere die Inverse  $\mathbf{A}^{-1} = (b_{ij}) = \mathbf{B}$ . Weiterhin betrachten wir  $\mathbf{X} = (X_1, \dots, X_n)^T$ ,  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  und

$$\mathbf{Y} = \mathbf{A} \mathbf{X}, \quad \mathbf{X} = \mathbf{B} \mathbf{Y} \quad \text{bzw. für } i = 1, \dots, n : \quad Y_i = \sum_{j=1}^n a_{ij} X_j, \quad X_i = \sum_{j=1}^n b_{ij} Y_j.$$

Aus der gemeinsamen Dichte  $f_X(x_1, \dots, x_n)$  der  $(X_1, \dots, X_n)$  kann dann wegen  $\det(\mathbf{A}) = 1/\det(\mathbf{B}) \neq 0$  die gemeinsame Dichte  $f_Y(y_1, \dots, y_n)$  von  $(Y_1, \dots, Y_n)$  berechnet werden. Es gilt gemäß Gleichung (4.13) mit  $\mathbf{y} = (y_1, \dots, y_n)^T$

$$\begin{aligned} f_Y(y_1, \dots, y_n) &= \frac{1}{|\det(\mathbf{A})|} f_X(x_1(y_1, \dots, y_n), \dots, x_n(y_1, \dots, y_n)) \\ &= |\det(\mathbf{B})| f_X(x_1, \dots, x_n) = |\det(\mathbf{A}^{-1})| f_X(\mathbf{A}^{-1} \mathbf{y}). \end{aligned}$$

Für affin lineare Transformationen  $\mathbf{Y} = \mathbf{A} \mathbf{X} + \mathbf{b}$  mit  $\det(\mathbf{A}) \neq 0$  gilt entsprechend

$$f_Y(\mathbf{y}) = |\det(\mathbf{A}^{-1})| f_X(\mathbf{A}^{-1}(\mathbf{y} - \mathbf{b})). \quad (4.14)$$

**Satz 4.13 (Orthogonale Transformation bei Normalverteilung)**

Seien  $(X_1, \dots, X_n)$   $N(0, 1)$ -iid Zufallsvariablen und  $\mathbf{A} = (a_{ij})$  eine orthogonale Matrix mit  $\det(\mathbf{A}) = \pm 1$  sowie  $\mathbf{A}^{-1} = \mathbf{A}^T$ . Weiterhin gelte mit  $\mathbf{X} = (X_1, \dots, X_n)^T$  und  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  die Beziehung  $\mathbf{Y} = \mathbf{A} \mathbf{X}$ , d. h.

$$Y_i = \sum_{j=1}^n a_{ij} X_j, \quad 1 \leq i \leq n. \quad (4.15)$$

Dann erhält man:

1.  $(Y_1, \dots, Y_n)$  sind unabhängige  $N(0, 1)$ -verteilte Zufallsvariablen.

2. Das Stichprobenmittel  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  und die

Stichprobenvarianz  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  sind unabhängig.

3. Für  $N(\mu, \sigma^2)$ -iid ZV  $(X_1, \dots, X_n)$  gilt  $E(S^2) = \sigma^2$ .

Beweis:

1. Es gilt  $\mathbf{X} = \mathbf{A}^T \mathbf{Y}$  und

$$\sum_{i=1}^n X_i^2 = \mathbf{X}^T \mathbf{X} = \mathbf{Y}^T \mathbf{A} \mathbf{A}^T \mathbf{Y} = \mathbf{Y}^T \mathbf{Y} = \sum_{i=1}^n Y_i^2.$$

Die unabhängigen  $N(0, 1)$ -verteilten  $X_i$  haben die Dichte  $(2\pi)^{-n/2} \exp(-\frac{1}{2} \sum x_i^2)$ . Somit gilt für die Dichte von  $(Y_1, \dots, Y_n)$  nach Satz 4.12

$$f_Y(y_1, \dots, y_n) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n y_i^2\right).$$

Deshalb sind  $(Y_1, \dots, Y_n)$  unabhängig  $N(0, 1)$ -verteilt.

2. Nun sei  $\mathbf{A} = (a_{ij})$  orthogonal mit spezieller erster Zeile der Form  $a_{1j} = 1/\sqrt{n}$ , was

$$Y_1 = \sum_{j=1}^n \frac{1}{\sqrt{n}} X_j = \sqrt{n} \bar{X}$$

zur Folge hat. Weiterhin gilt

$$\begin{aligned} (n-1)S^2 &= \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 = \sum_{i=1}^n X_i^2 - n\bar{X}^2 \\ &= \sum_{i=1}^n Y_i^2 - Y_1^2 = \sum_{i=2}^n Y_i^2. \end{aligned} \quad (4.16)$$

$S^2$  ist unabhängig von  $\bar{X}$ , da  $Y_1$  gemäß 1. unabhängig von  $(Y_2, \dots, Y_n)$  ist.

3. Ausgehend von  $N(\mu, \sigma^2)$ -iid  $\tilde{X}_i$  führt man zunächst eine Translation der Form  $\tilde{X}_i - \mu \rightarrow X_i, i = 1, \dots, n$  durch. Dann sind die  $X_i$  iid und  $N(0, \sigma^2)$ . Durch Einsetzen gemäß 1. verifiziert man sofort, daß die orthogonale Transformation (4.15)  $N(0, \sigma^2)$  iid  $Y_i$  liefert. Damit folgt gemäß 2. und (4.16)

$$E[(n-1)S^2] = E\left[\sum_{i=2}^n Y_i^2\right] = (n-1)\sigma^2. \quad \square$$

## Lineare Transformationen normalverteilter ZV

Aus Satz 4.12 folgt noch, dass allgemeine lineare Transformationen normalverteilter ZV wieder normalverteilt sind:

Seien  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{C})$  und  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\text{Rang}(\mathbf{A}) = m > 0$  gegeben. Dann gilt

$$\mathbf{Y} := \mathbf{A}\mathbf{X} \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\mathbf{C}\mathbf{A}^T). \quad (4.17)$$

Zum Beweis ergänze man  $\mathbf{A}$  zu einer nichtsingulären Matrix.

## 4.4 Zentriertes Modell, orthogonale Designmatrix

Wir betrachten das zentrierte Modell der Einfachen Linearen Regresstion mit normalverteilten Fehlern. Für gegebene Daten  $(x_1, y_1), \dots, (x_n, y_n)$ ,  $i = 1, \dots, n$  ist dies ein auf  $\bar{x}$  zentrierter Ansatz der Form

$$Y_i = \beta_0^* + \beta_1(x_i - \bar{x}) + \epsilon_i, \quad \epsilon_i \text{ iid } N(0, \sigma^2), \quad i = 1, \dots, n.$$

Seien nun  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ ,  $\mathbf{x} = (x_1, \dots, x_n)^T$ ,  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ ,  $\mathbf{1} = (1, \dots, 1)^T$ ,  
 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  sowie  $SS_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$  und  $SS_{xy} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})$ .

Weiterhin seien  $b_0^*$  und  $b_1$  die KQ-Schätzungen für  $\beta_0^*$  und  $\beta_1$ .  $\hat{Y}_i = b_0^* + b_1(x_i - \bar{x})$ ,  
 $e_i = Y_i - \hat{Y}_i$ ,  $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_n)^T$ ,  $\mathbf{e} = (e_1, \dots, e_n)^T$  und  $SS_e = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ .

### Lemma 4.14

$$\mathbf{X} = (\mathbf{1}, \mathbf{x} - \bar{x} \cdot \mathbf{1})$$

hat orthogonale Spalten und

$$\mathbf{b} = \begin{pmatrix} b_0^* \\ b_1 \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{pmatrix} \bar{Y} \\ SS_{xy}/SS_{xx} \end{pmatrix}.$$

Beweis: Einsetzen und Anwendung von Abschnitt 4.1.



**Satz 4.15** *Es gilt:*

1.  $E(\mathbf{b}) = E \begin{pmatrix} b_0^* \\ b_1 \end{pmatrix} = \begin{pmatrix} \beta_0^* \\ \beta_1 \end{pmatrix}$
2.  $E(SS_e) = (n - 2)\sigma^2$
3.  $b_0^*, b_1$  und  $SS_e$  sind unabhängig.

Beweis:

Die Lösung besteht darin, den Zufallsvektor  $\mathbf{Y} \in \mathbb{R}^n$  durch eine orthogonale Transformation  $\mathbf{T} \in \mathbb{R}^{n,n}$  so zu "drehen", dass die beiden ersten Komponenten  $Z_1, Z_2$  von  $\mathbf{Z} := \mathbf{T}^T \mathbf{Y}$  in  $\text{span}(\mathbf{X})$  liegen und die übrigen  $Z_3, \dots, Z_n$  senkrecht dazu stehen, wobei gemäß Lemma 4.14

$$\mathbf{X} = (\mathbf{1}, \mathbf{x} - \bar{x} \cdot \mathbf{1}) \in \mathbb{R}^{n,2}$$

$$\mathbf{x} = (x_1, \dots, x_n)^T, \mathbf{Y} = (Y_1, \dots, Y_n)^T.$$

Dazu reparametrisiert man das zentrierte Modell

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \boldsymbol{\beta} = (\beta_0^*, \beta_1)^T$$

mit

$$\mathbf{q}_1 := \mathbf{1}/\sqrt{n}, \mathbf{q}_2 := (\mathbf{x} - \bar{x} \cdot \mathbf{1})/\sqrt{SS_{xx}}$$

$$\mathbf{Q} := (\mathbf{q}_1, \mathbf{q}_2) \in \mathbb{R}^{n,2}, \mathbf{D} = \begin{pmatrix} 1/\sqrt{n} & 0 \\ 0 & 1/\sqrt{SS_{xx}} \end{pmatrix}.$$

Nach Definition hat  $\mathbf{Q}$  orthogonale Spalten. Weiterhin gilt

$$\mathbf{Q} = \mathbf{X}\mathbf{D}$$

und

$$\mathbf{X}\boldsymbol{\beta} = \underbrace{\mathbf{X}\mathbf{D}}_{\mathbf{Q}} \underbrace{\mathbf{D}^{-1}\boldsymbol{\beta}}_{\boldsymbol{\gamma}} = \mathbf{Q}\boldsymbol{\gamma},$$

wobei

$$\boldsymbol{\gamma} = \begin{pmatrix} \gamma_1 \\ \gamma_2 \end{pmatrix} = \begin{pmatrix} \sqrt{n}\beta_0^* \\ \sqrt{SS_{xx}}\beta_1 \end{pmatrix}$$

und  $\text{span}(\mathbf{Q}) = \text{span}(\mathbf{X})$ . Also

$$\mathbf{Y} = \mathbf{Q}\boldsymbol{\gamma} + \boldsymbol{\epsilon} \iff \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Nach Lemma 4.14 gilt für die KQ-Schätzungen

$$\mathbf{g} = \hat{\boldsymbol{\gamma}} = (\mathbf{Q}^T \mathbf{Q})^{-1} \mathbf{Q}^T \mathbf{Y} = \mathbf{Q}^T \mathbf{Y} = \mathbf{D} \mathbf{X}^T \mathbf{Y} =$$

$$= \mathbf{D} \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}}_{\mathbf{D}^{-1} \mathbf{Q}^T \mathbf{Q} \mathbf{D}^{-1} \mathbf{b}} = \mathbf{D} \mathbf{D}^{-2} \mathbf{b} = \mathbf{D}^{-1} \mathbf{b}$$

Nun ergänzen wir  $\mathbf{Q}$  durch

$$\mathbf{R} := (\mathbf{q}_3, \dots, \mathbf{q}_n) \in \mathbb{R}^{n, n-2}$$

zu einer orthogonalen Matrix  $\mathbf{T} = (\mathbf{Q}, \mathbf{R})$ .

Wegen  $\mathbf{Y} \sim N(\mathbf{Q}\boldsymbol{\gamma}, \sigma^2 \mathbf{I})$  und Satz 4.13 ist

$$\mathbf{Z} = \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_n \end{pmatrix} := \mathbf{T}^T \mathbf{Y} = \begin{pmatrix} \mathbf{Q}^T \mathbf{Y} \\ \mathbf{R}^T \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{g} \\ \mathbf{Z}_e \end{pmatrix} \sim N(\boldsymbol{\delta}, \sigma^2 \mathbf{I});$$

und alle Komponenten von  $\mathbf{Z}$  sind unabhängig.

$$1. E(\mathbf{g}) = E(\mathbf{Q}^T \mathbf{Y}) = \mathbf{Q}^T E(\mathbf{Y}) = \mathbf{Q}^T \mathbf{Q} \boldsymbol{\gamma} = \boldsymbol{\gamma} = \begin{pmatrix} \sqrt{n} \beta_0^* \\ \sqrt{SS_{xx}} \beta_1 \end{pmatrix}.$$

Insbesondere

$$\begin{aligned} Z_1 = g_0 &\sim N(\sqrt{n} \beta_0^*, \sigma^2), \\ Z_2 = g_1 &\sim N(\sqrt{SS_{xx}} \beta_1, \sigma^2) \text{ und} \\ SS_Z &= \sum_{i=3}^n Z_i^2 = \mathbf{Z}_e^T \mathbf{Z}_e. \end{aligned}$$

Es gilt  $SS_Z/\sigma^2 \sim \chi_{n-2}^2$ , da  $E(\mathbf{Z}_e) = \mathbf{R}^T \mathbf{Q} \boldsymbol{\gamma} = \mathbf{0}$ .

Wegen  $\mathbf{b} = \mathbf{D} \mathbf{g}$  folgt  $b_0^*$  und  $b_1$  sind unabhängig mit

$$\begin{aligned} b_0^* &= g_0/\sqrt{n} \sim N(\beta_0^*, \sigma^2/\sqrt{n}) \text{ und} \\ b_1 &= g_1/\sqrt{SS_{xx}} \sim N(\beta_1, \sigma^2/\sqrt{SS_{xx}}). \end{aligned}$$

Es bleibt noch zu zeigen:

$$SS_Z = SS_e = \mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y},$$

wobei

$$\begin{aligned} \mathbf{H} &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{D}^{-1} \underbrace{(\mathbf{D}^{-1} \mathbf{Q}^T \mathbf{Q} \mathbf{D}^{-1})^{-1}}_{\mathbf{D}^{-2}} \mathbf{D}^{-1} \mathbf{Q}^T = \\ &= \mathbf{Q} \mathbf{D}^{-1} \mathbf{D}^2 \mathbf{D}^{-1} \mathbf{Q}^T = \mathbf{Q} \mathbf{Q}^T. \end{aligned}$$

Mit

$$\mathbf{I} = \mathbf{T} \mathbf{T}^T = (\mathbf{Q}, \mathbf{R}) \begin{pmatrix} \mathbf{Q}^T \\ \mathbf{R}^T \end{pmatrix} = \mathbf{Q} \mathbf{Q}^T + \mathbf{R} \mathbf{R}^T$$

folgt schließlich mit  $\mathbf{e} = \mathbf{Y} - \mathbf{X} \mathbf{b}$ :

$$\begin{aligned} SS_e &= \mathbf{e}^T \mathbf{e} = \mathbf{Y}^T \mathbf{Y} - 2\mathbf{Y}^T \mathbf{X} \mathbf{b} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b} = \\ &= \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \\ &= \mathbf{Y}^T (\mathbf{Q} \mathbf{Q}^T + \mathbf{R} \mathbf{R}^T - \mathbf{Q} \mathbf{Q}^T) \mathbf{Y} = \\ &= \mathbf{Y}^T \mathbf{R} \mathbf{R}^T \mathbf{Y} = \mathbf{Z}_e^T \mathbf{Z}_e = SS_Z. \end{aligned}$$

Insbesondere sind  $b_0^*$ ,  $b_1$  und  $SS_e$  unabhängig. Damit ist 3. gezeigt.

$$2. E(SS_e) = E(SS_Z) = (n-2)\sigma^2, \text{ da } SS_Z/\sigma^2 \sim \chi_{n-2}^2.$$

# Literaturverzeichnis

- [Dielman et al. (1996)] T.E.: Applied Regression Analysis for Business and Economics (2nd ed). Wadsworth Publishing Company
- [Falk et al. (1995)] Falk M., Becker R. und Mahrohn F.: Angewandte Statistik mit SAS. Springer, Berlin.
- [Fahrmeir et al. (1996)] Fahrmeir. L, Hamerle A. und Tutz G. (Hrsg.): Multivariate statistische Verfahren, 2., erweiterte Auflage. De Gruyter, Berlin.
- [Kredler (1999)] Ch.: Einführung in die Wahrscheinlichkeitsrechnung und Statistik. Ausgearbeitetes Vorlesungsskript zur Stochastik 1. Zentrum Mathematik, TU München.
- [Seber (1977)] G.A.F.: Linear Regression Analysis. Wiley, New York.
- [Myers (1990)] R.H. (1990): Classical and Modern Regression with Applications. Duxbury Press, Belmont, CA.
- [Draper/Smith (1998)] Draper N.R. and Smith H. (1998): Applied Regression Analysis (3rd ed), Wiley, N.Y.
- [Cook/Weisberg (1999)] Cook, R.D. and Weisberg S.D. (1999): Applied Regression Including Computing and Graphics, Wiley, N.Y.
- [Kleinbaum/Kupper (1998)] Kleinbaum D.G. and Kupper L.L.: Applied Regression Analysis and Other Multivariate Methods, Duxbury Press, Belmont, CA.
- [Toutenburg (1992)] H. : Lineare Modelle, Physica-Verlag, Heidelberg
- [S-PLUS (2000)] S-PLUS 2000 - Guide to Statistics, Vol.1, MathSoft, Seattle