

# Zum linearen Modell, klassischer Kleinste-Quadrate- und Maximum-Likelihood-Schätzung

1. Vortrag des Seminars Erweiterungen des linearen  
Regressionsmodells und genomische Anwendungen in der  
Biomedizin, WS 2014/2015

Matthias Birkner

3.11.2014



JOHANNES GUTENBERG  
UNIVERSITÄT MAINZ

- 1 Einführungsbeispiel: Einfache lineare Regression
- 2 Allgemeines zum linearen Modell
- 3 Die Welt des normalverteilten Rauschens
- 4 Beispiel

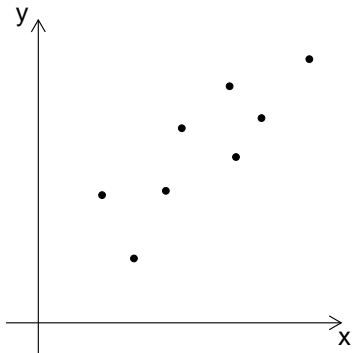
# Inhalt

- 1 Einführungsbeispiel: Einfache lineare Regression
- 2 Allgemeines zum linearen Modell
- 3 Die Welt des normalverteilten Rauschens
- 4 Beispiel

# Erinnerung: Lineare Regression

Daten:

$n$  Wertepaare  $(x_i, y_i)$



# Erinnerung: Lineare Regression

Daten:

$n$  Wertepaare  $(x_i, y_i)$

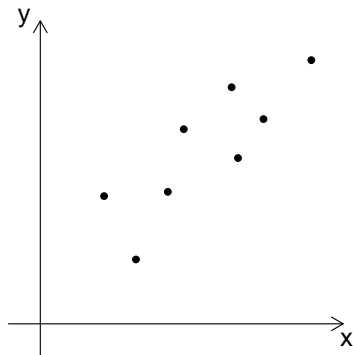
Modell: (approximativ) linearer  
Zusammenhang

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

mit  $\varepsilon_i$  u.a.,  $\mathbb{E}[\varepsilon_i] = 0$ ,

$$\text{Var}[\varepsilon_i] = \sigma^2$$

(hier:  $p = 2$  Parameter  $\beta_0, \beta_1$ )



# Erinnerung: Lineare Regression

Daten:

$n$  Wertepaare  $(x_i, y_i)$

Modell: (approximativ) linearer  
Zusammenhang

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

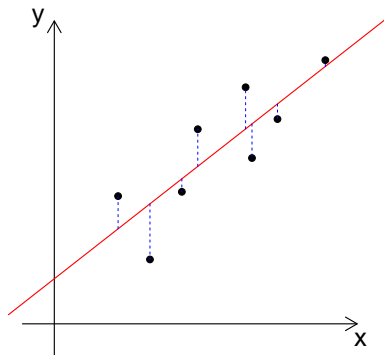
mit  $\varepsilon_i$  u.a.,  $\mathbb{E}[\varepsilon_i] = 0$ ,

$$\text{Var}[\varepsilon_i] = \sigma^2$$

(hier:  $p = 2$  Parameter  $\beta_0, \beta_1$ )

Kleinste-Quadrate-Schätzer:  $\hat{\beta}_0$  und  $\hat{\beta}_1$  derart, dass

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \stackrel{!}{=} \text{minimal}$$



# Erinnerung: Lineare Regression

Daten:

$n$  Wertepaare  $(x_i, y_i)$

Modell: (approximativ) linearer  
Zusammenhang

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

mit  $\varepsilon_i$  u.a.,  $\mathbb{E}[\varepsilon_i] = 0$ ,

$$\text{Var}[\varepsilon_i] = \sigma^2$$

(hier:  $p = 2$  Parameter  $\beta_0, \beta_1$ )

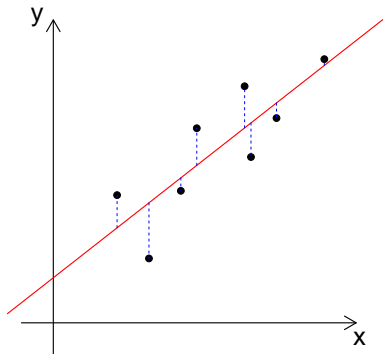
Kleinste-Quadrate-Schätzer:  $\hat{\beta}_0$  und  $\hat{\beta}_1$  derart, dass

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \stackrel{!}{=} \text{minimal}$$

(mit Lösung  $\hat{\beta}_1 = \text{cov}(x, y) / \sigma_x^2$ ,  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ ,

wo  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ ,

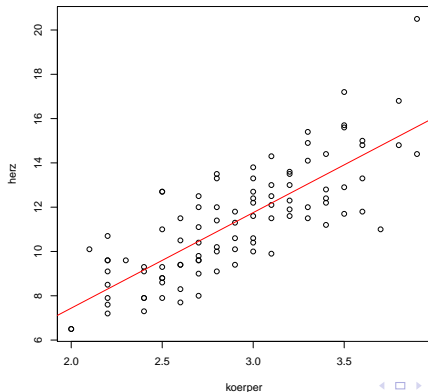
$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}))$$



# Lineare Regression mit R

RS cats-Datensatz, aus Fisher, Biometrics, 3, 65-68 (1947)  
(Wir betrachten nur die 97 Kater):

```
> data(cats, package="MASS"); attach(cats)
> koerper <- Bwt[Sex=="M"]; herz <- Hwt[Sex=="M"]
```





# Lineare Regression mit R

```
> modell <- lm(herz ~ koerper)
> summary(modell)
```

Call:

```
lm(formula = herz ~ koerper)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.7728	-1.0478	-0.2976	0.9835	4.8646

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-1.1841	0.9983	-1.186	0.239
koerper	4.3127	0.3399	12.688	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1

# Inhalt

- 1 Einführungsbeispiel: Einfache lineare Regression
- 2 Allgemeines zum linearen Modell**
- 3 Die Welt des normalverteilten Rauschens
- 4 Beispiel

# Modellvorstellung

In einer Population besitzt jedes Individuum/Objekt einen (reellen)  $y$ -Wert („Zielgröße“) und  $p$  (reelle)  $x$ -Werte  $x_1, \dots, x_p$  ( $p$  „erklärende Variablen“).

Seien  $((X_1, \dots, X_p), Y)$  die Merkmale eines (zufällig) aus der Population gezogenen Individuums.

# Modellvorstellung

In einer Population besitzt jedes Individuum/Objekt einen (reellen)  $y$ -Wert („Zielgröße“) und  $p$  (reelle)  $x$ -Werte  $x_1, \dots, x_p$  ( $p$  „erklärende Variablen“).

Seien  $((X_1, \dots, X_p), Y)$  die Merkmale eines (zufällig) aus der Population gezogenen Individuums.

Wir nehmen an, dass für die Verteilung von  $Y$ , gegeben beobachtete (oder je nach Situation auch von uns vorgegebene) Werte  $(X_1, \dots, X_p) = (x_1, \dots, x_p)$  gilt

$$\mathbb{E}_{\beta}[Y \mid (X_1, \dots, X_p) = (x_1, \dots, x_p)] = \beta_1 x_1 + \dots + \beta_p x_p$$

und

$$\text{Var}_{\beta}[Y \mid (X_1, \dots, X_p) = (x_1, \dots, x_p)] = \sigma^2$$

für einen Parametervektor  $\beta = (\beta_1, \dots, \beta_p)^T$  und ein  $\sigma^2 > 0$ .

Bemerkung:

$$\mathbb{E}_\beta[Y | (X_1, \dots, X_p) = (x_1, \dots, x_p)] = \beta_1 x_1 + \dots + \beta_p x_p$$

Der Einfachheit der Formulierung halber integrieren wir hier den Parameter für einen konstanten Wert („Achsenabschnitt“) in das Modell, indem wir annehmen, dass für eine der erklärenden Variablen  $\equiv 1$  erfüllt.

# Modellvorstellung, 2

Gegeben seien  $n$  Beobachtungen, zur  $i$ -ten Beobachtung gehört  $y_i$  (beobachteter Wert der Zielgröße) und  $x_{i,1}, x_{i,2}, \dots, x_{i,p}$  (Werte der  $p$  erklärenden Variablen),

wir interpretieren diese als  $n$  unabhängige Realisierungen von  $((X_1, \dots, X_p), Y)$  und schreiben

$$y_i = \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \varepsilon_i$$

mit  $\varepsilon_i$  (Realisierungen von) unabhängige(n) Zufallsvariablen mit  $\mathbb{E}[\varepsilon_i] = 0$ ,  $\text{Var}[\varepsilon_i] = \sigma^2$ ,

wobei wir das „wahre“  $\beta$  und das „wahre“  $\sigma^2$  nicht kennen.

(Je nach Design des Experiments/der Studie stellen wir uns ggfs. vor, dass die Objekte/Individuen bedingt auf gewisse  $x$ -Werte ausgewählt wurden.)

# Lineares Modell in Matrix-Schreibweise

$n$  Beobachtungen, zur  $i$ -ten Beobachtung gehört  $y_i$  und  $x_{i,1}, x_{i,2}, \dots, x_{i,p}$  (und wir denken im Moment an  $n > p$ ),

$$y_i = \sum_{j=1}^p x_{i,j} \beta_j + \varepsilon_i, \quad i = 1, \dots, n,$$

zusammengefasst

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

mit

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix},$$

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$$

$\mathbf{X}$  ist die „Designmatrix“ ( $i$ -te Zeile gehört zur  $i$ -ten Beobachtung)

# Kleinste-Quadrate-Ansatz

$n \geq p$  (wir denken an  $n \gg p$ )

Gegeben beobachtetes  $\mathbf{y}$  und bekannte/beobachtete Matrix  $\mathbf{X}$   
mit  $n \geq p$  (wir denken an  $n \gg p$ )

Typischerweise besitzt das lineare Gleichungssystem  
(aufgefasst als Bestimmungsgleichungen für die  $\beta_i$ )  $\mathbf{y} = \mathbf{X}\beta$   
keine Lösung.

Kleinste-Quadrate-Schätzer (“least squares”)

$$\hat{\beta}_{\text{LS}} := \operatorname{argmin} \{ \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \}$$



# Berechnung (wenigstens prinzipiell)

$\mathbf{X}$  habe (maximalen) Rang  $p$ , d.h. die  $p$  Spalten von  $\mathbf{X}$  sind linear unabhängig, dann gilt:

Die  $p \times p$ -Matrix  $\mathbf{X}^T \mathbf{X}$  hat vollen Rang  $p$ , ist insbesondere invertierbar, es ist

$$\hat{\beta}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

[Falls  $\mathbf{X}^T \mathbf{X} \mathbf{c} = \mathbf{0}$  für ein  $\mathbf{c} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$ , so wäre  $\|\mathbf{X} \mathbf{c}\|_2^2 = \mathbf{c}^T \mathbf{X}^T \mathbf{X} \mathbf{c} = 0$  im Widerspruch zur Ann., dass  $\mathbf{X}$  Rang  $p$  hat.]

## Argument (analytisch)

Die „Zielfunktion“

$$\beta \mapsto \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{i,j} \beta_j \right)^2$$

ist quadratisch, das Gleichungssystem

$$\frac{\partial}{\partial \beta_k} \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{i,j} \beta_j \right)^2 = -2 \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{i,j} \beta_j \right) x_{i,k} \stackrel{!}{=} 0$$

für  $k = 1, \dots, p$  ist in Matrixform geschrieben äquivalent zu

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \beta \quad \text{„Normalgleichungen“}$$

[Bem.: Diese bestimmen die Lösungsmenge  $\{\hat{\beta}_{LS}\}$  auch in dem Fall, dass  $\mathbf{X}^T \mathbf{X}$  nicht invertierbar ist.]

## Argument (geometrisch)

Sei  $L = \{\mathbf{X}\boldsymbol{\eta} : \boldsymbol{\eta} \in \mathbb{R}^p\} \subset \mathbb{R}^n$  der Spaltenraum von  $\mathbf{X}$ ,

$\Pi_L : \mathbb{R}^n \rightarrow L$  die orthogonale Projektion auf  $L$

Für  $\mathbf{y} \in \mathbb{R}^n$ :

$$(i) \quad \Pi_L \mathbf{y} \in L \text{ und } \|\mathbf{y} - \Pi_L \mathbf{y}\|_2^2 \leq \min_{\mathbf{u} \in L} \|\mathbf{y} - \mathbf{u}\|_2^2$$

$$(ii) \quad \Pi_L \mathbf{y} \in L \text{ und } \mathbf{y} - \Pi_L \mathbf{y} \perp L$$

(und  $\Pi_L \mathbf{y}$  ist durch (i) und durch (ii) charakterisiert)

$$\Leftrightarrow (ii)' \quad \Pi_L \mathbf{y} \in L \text{ und } \mathbf{u}^T (\mathbf{y} - \Pi_L \mathbf{y}) = 0 \text{ f\"ur } \mathbf{u} \text{ aus einer Basis von } L$$

$$\iff \mathbf{X}^T (\mathbf{y} - \Pi_L \mathbf{y}) = \mathbf{0}$$

Es ist  $\Pi_L \mathbf{y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  für  $\mathbf{y} \in \mathbb{R}^n$ :

$$\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \in L \text{ und } \mathbf{X}^T (\mathbf{y} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) = \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

[Details ggfs. an der Tafel]

$\hat{\beta}_{\text{LS}}$  der kleinste-Quadrate-Schätzer,

$$\hat{\mathbf{y}} := \mathbf{X}\hat{\beta}_{\text{LS}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

die angepassten Werte der  $y$ -Beobachtungen („gefittete“ Werte),  $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  die „Hut-Matrix“,

$$\mathbf{r} := \mathbf{y} - \hat{\mathbf{y}} \quad ( = \mathbf{y} - \mathbf{X}\hat{\beta}_{\text{LS}} )$$

die „Residuen“

(die Abweichungen der Beobachtungen vom Modell, sozusagen der „unerklärte Rest“)

Bem.:

Die Kovarianzmatrix von  $\hat{\beta}_{\text{LS}}$  ist  $\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$ .

# Satz von Gauß-Markov

Im linearen Modell (mit unkorrelierten, zentrierten  $\varepsilon_i$  mit derselben Varianz  $\sigma^2$ ) gilt

- 1 Der kleinste-Quadrate-Schätzer  $\hat{\beta}_{\text{LS}}$  ist erwartungstreu für  $\beta^*$ , d.h.  $\mathbb{E}_{\beta^*}[\hat{\beta}_{\text{LS}}] = \beta^*$  (für jede Wahl des „wahren“  $\beta^*$ ).
- 2  $\tau : \mathbb{R}^p \rightarrow \mathbb{R}$  eine lineare Kenngröße von  $\beta$  (d.h.  $\tau(\beta) = \mathbf{c}^T \beta$  für ein  $\mathbf{c} \in \mathbb{R}^p$ , z.B.  $\tau(\beta) = \beta_1$ ), dann ist  $\hat{\tau} := \mathbf{c}^T \hat{\beta}_{\text{LS}}$  erwartungstreuer Schätzer für  $\tau(\beta^*)$  und hat unter allen erwartungstreuen linearen Schätzern für  $\tau(\beta^*)$  die kleinste Varianz.

Weiterhin ist ( $\Pi_L = \text{orth. Proj. auf den Spaltenraum von } \mathbf{X}$ )

$$\hat{\sigma}^2 := \frac{\|\mathbf{r}\|_2^2}{n-p} = \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}_{\text{LS}}\|_2^2}{n-p} = \frac{\|\mathbf{y} - \Pi_L \mathbf{y}\|_2^2}{n-p} = \frac{\|\mathbf{y}\|_2^2 - \|\Pi_L \mathbf{y}\|_2^2}{n-p}$$

ein erwartungstreuer Schätzer für die Varianz  $\sigma^2$  (der  $\varepsilon_i$ ).

[Details ggfs. an der Tafel]

# Inhalt

- 1 Einführungsbeispiel: Einfache lineare Regression
- 2 Allgemeines zum linearen Modell
- 3 Die Welt des normalverteilten Rauschens**
- 4 Beispiel

# Exkurs/Erinnerung: Rund um die multivariate Normalverteilung

$$\mathbf{Z} = (Z_1, \dots, Z_n)^T \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

hat Dichte

$$\begin{aligned} \varphi_{\mathbf{Z}}(z_1, \dots, z_n) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} e^{-z_i^2/2\sigma^2} \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z}\|_2^2\right) \end{aligned}$$

Die Verteilung von  $\mathbf{Z}$  ist invariant unter orthogonalen Transformationen

(insbesondere: Projektionen auf verschiedene, zueinander orthogonale Teilräume sind unabhängig)

# Rund um die multivariate Normalverteilung, 2

Seien  $Z_1, \dots, Z_m, Z'_1, \dots, Z'_n$  u.a. Standard-normalverteilt.

- $Z_1^2 + \dots + Z_m^2$  ist  $\chi_m^2$ -verteilt

(Dichte  $\frac{(1/2)^{m/2}}{\Gamma(m/2)} x^{m/2-1} e^{-(1/2)x}$  auf  $(0, \infty)$  [ $\chi_m^2 = \Gamma(m/2, 1/2)$ ])

- $F_{m,n} := \frac{\frac{1}{m} \sum_{i=1}^m Z_i^2}{\frac{1}{n} \sum_{j=1}^n Z_j'^2}$  ist Fisher $_{m,n}$ -verteilt

(Dichte  $f_{m,n}(x) = \frac{\Gamma(\frac{n+m}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} m^{\frac{m}{2}} n^{\frac{n}{2}} \frac{x^{\frac{m}{2}-1}}{(n+mx)^{\frac{(m+n)}{2}}}$  auf  $(0, \infty)$ )

- $T := \frac{Z_1}{\sqrt{\frac{1}{n} \sum_{j=1}^n Z_j'^2}}$  ist Student- $t_n$ -verteilt

(Dichte  $\frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{1}{2})} \frac{1}{\sqrt{n}} (1 + \frac{x^2}{n})^{-\frac{n+1}{2}}$  auf  $(-\infty, \infty)$ )



# ML-Schätzung im normalen Modell

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Wir nehmen nun zusätzlich an, dass die  $\varepsilon_i$  u.i.v.,  $\sim \mathcal{N}(0, 1)$  sind.

Gegeben  $\mathbf{X}$  (und festes  $\boldsymbol{\beta}$ ) hat  $\mathbf{y} = (y_1, \dots, y_n)^T$  dann die Dichte

$$\begin{aligned} & \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left(- (y_i - (\mathbf{X}\boldsymbol{\beta})_i)^2 / (2\sigma^2)\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(- \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 / (2\sigma^2)\right) \end{aligned}$$

Demnach ist  $\hat{\boldsymbol{\beta}}_{\text{LS}}$  auch der Maximum-Likelihood-Schätzer für  $\boldsymbol{\beta}$  (und es ist  $\hat{\sigma}_{\text{ML}}^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 / n$ ).

# Verteilungseigenschaften (im normalen Modell)

- 1  $\hat{\beta}_{\text{LS}} \sim \mathcal{N}(\beta^*, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$
- 2  $\frac{n-p}{\sigma^2} \hat{\sigma}^2$  ist  $\chi_{n-p}^2$ -verteilt und unabhängig von  $\hat{\beta}_{\text{LS}}$ .
- 3  $\frac{1}{\sigma^2} \|\mathbf{X}(\hat{\beta}_{\text{LS}} - \beta^*)\|_2^2 = \frac{1}{\sigma^2} \|\Pi_L \mathbf{y} - \mathbb{E}_{\beta^*}[\mathbf{y}]\|_2^2$  ist  $\chi_p^2$ -verteilt und unabhängig von  $\hat{\sigma}^2$ , insbesondere

$$\frac{\|\mathbf{X}(\hat{\beta}_{\text{LS}} - \beta^*)\|_2^2}{p \hat{\sigma}^2} \sim \text{Fisher}_{p, n-p}$$

- 4  $H \subset L$  linearer Teilraum von  $L$ ,  $\dim H = r < p$  und  $\mathbf{X}\beta^* \in H$ , dann ist  $\frac{1}{\sigma^2} \|\Pi_L \mathbf{y} - \Pi_H \mathbf{y}\|_2^2 \sim \chi_{p-r}^2$  und unabhängig von  $\hat{\sigma}^2$ , insbesondere ist

$$\frac{n-p}{p-r} \frac{\|\Pi_L \mathbf{y} - \Pi_H \mathbf{y}\|_2^2}{\|\mathbf{y} - \Pi_L \mathbf{y}\|_2^2} = \frac{\|\mathbf{X}\hat{\beta}_{\text{LS}} - \Pi_H \mathbf{y}\|_2^2}{(p-r)\hat{\sigma}^2} \sim \text{Fisher}_{p-r, n-p}.$$

# Korollar: Konfidenzbereiche

Sei  $\alpha \in (0, 1)$ .

- 1  $C_{\hat{\beta}_{LS}} := \{\tilde{\beta} \in \mathbb{R}^p : \|\mathbf{X}(\tilde{\beta} - \hat{\beta}_{LS})\|_2^2 < p\hat{\sigma}^2 f_{p,n-p;1-\alpha}\}$  ist ein Konfidenzbereich für  $\beta^*$  zum Sicherheitsniveau  $1 - \alpha$ , wo  $f_{p,n-p;1-\alpha} = (1 - \alpha)$ -Quantil der Fisher $_{p,n-p}$ -Verteilung.
- 2  $\tau(\beta) = \mathbf{c}^T\beta$  (mit einem  $\mathbf{c} \in \mathbb{R}^p$ ) ein lineares Parametermerkmal.

$$C_{\tau(\hat{\beta}_{LS})} = (\mathbf{c}^T\hat{\beta}_{LS} - \delta\sqrt{\hat{\sigma}^2}, \mathbf{c}^T\hat{\beta}_{LS} + \delta\sqrt{\hat{\sigma}^2}),$$

mit  $\delta = \sqrt{\mathbf{c}^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{c}} t_{n-p,1-\frac{\alpha}{2}}$  und  $t_{n-p,1-\frac{\alpha}{2}} = (1 - \frac{\alpha}{2})$ -Quantil der Student- $T$ -Verteilung mit  $n - p$ -Freiheitsgraden, ist ein Konfidenzintervall für  $\tau(\beta^*)$  zum Sicherheitsniveau  $1 - \alpha$ .

Bem.:  $C_{\hat{\sigma}^2} := (\frac{n-p}{\chi_{n-p,1-\frac{\alpha}{2}}^2} \hat{\sigma}^2, \frac{n-p}{\chi_{n-p,\frac{\alpha}{2}}^2} \hat{\sigma}^2)$  ist ein Konfidenzintervall für  $\sigma^2$  zum Sicherheitsniveau  $1 - \alpha$ .

# Inhalt

- 1 Einführungsbeispiel: Einfache lineare Regression
- 2 Allgemeines zum linearen Modell
- 3 Die Welt des normalverteilten Rauschens
- 4 Beispiel**

# Prostata-Krebs-Daten von Stamey et al

Datensatz beschrieben in T. Hastie, R. Tibshirani, J. Friedman, *The elements of statistical learning*, Example 3.2.1, <http://statweb.stanford.edu/~tibs/ElemStatLearn/data.html>

$n = 97$  Beobachtungen (Prostata-Krebs-Patienten),  
 $p = 8$  erklärende Variablen (plus “intercept”)

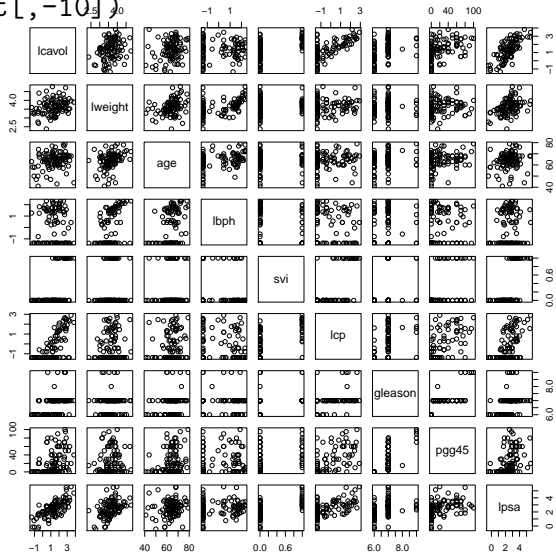
*log cancer volume* ( $lcavol$ ), *log prostate weight* ( $lweight$ ), *age*, *log of the amount of benign prostatic hyperplasia* ( $lbph$ ), *seminal vesicle invasion* ( $svi$ ), *log of capsular penetration* ( $lcp$ ), *Gleason score* ( $gleason$ ), *percent of Gleason scores 4 or 5* ( $pgg45$ )

Zielgröße (“response”) *log prostate specific antigen*  $lpsa$

Stamey, T., Kabalin, J., McNeal, J., Johnstone, I., Freiha, F., Redwine, E. and Yang, N. (1989). Prostate specific antigen in the diagnosis and treatment of adenocarcinoma of the prostate II radical prostatectomy treated patients, *Journal of Urology* 16: 1076–1083

```
> dat <- read.table("prostate.data", header=T)
```

```
> pairs(dat[, -10])
```



# Zentrierung/Standardisierung

Wir folgen Hastie, Tibshirani & Friedman (und „üblicher Praxis“):

Zentrierung der Spalten von  $\mathbf{X}$ :

ersetze  $x_{i,j}$  durch  $x_{i,j} - \bar{x}_j$  mit  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{i,j}$

Standardisierung der Spalten von  $\mathbf{X}$ : Skaliere so, dass

$\frac{1}{n} \sum_{i=1}^n x_{i,j}^2 = 1$  für  $j = 1, \dots, p$

(d.h. ersetze  $x_{i,j}$  durch  $x_{i,j} / (\frac{1}{n} \sum_{i=1}^n x_{i,j}^2)^{1/2}$ , dann hat  $\mathbf{X}^T \mathbf{X}$  auf der Diagonale 1en stehen)

```
> dat <- read.table("prostate.data", header=T)
> x <- dat[,1:8]
> xp <- scale(x,TRUE,TRUE)
> sdat <- data.frame(xp, lpsa=dat$lpsa)
```

Es gibt nicht-triviale Korrelationen zwischen den Spalten von  $\mathbf{X}$ :

```
> round(cor(sdat[1:8]), digits=3)
```

	lcavol	lweight	age	lbph	svi	lcp	gleason	pgg45
lcavol	1.000	0.281	0.225	0.027	0.539	0.675	0.432	0.434
lweight	0.281	1.000	0.348	0.442	0.155	0.165	0.057	0.107
age	0.225	0.348	1.000	0.350	0.118	0.128	0.269	0.276
lbph	0.027	0.442	0.350	1.000	-0.086	-0.007	0.078	0.078
svi	0.539	0.155	0.118	-0.086	1.000	0.673	0.320	0.458
lcp	0.675	0.165	0.128	-0.007	0.673	1.000	0.515	0.632
gleason	0.432	0.057	0.269	0.078	0.320	0.515	1.000	0.752
pgg45	0.434	0.107	0.276	0.078	0.458	0.632	0.752	1.000



```
> summary(modell)
```

```
Call:
```

```
lm(formula = lpsa ~ lcavol + lweight + age + lbph + svi + lcp +
    gleason + pgg45, data = sdat)
```

```
Residuals:
```

	Min	1Q	Median	3Q	Max
	-1.76644	-0.35510	-0.00328	0.38087	1.55770

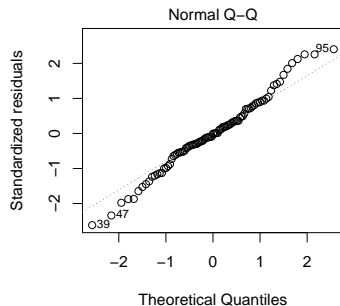
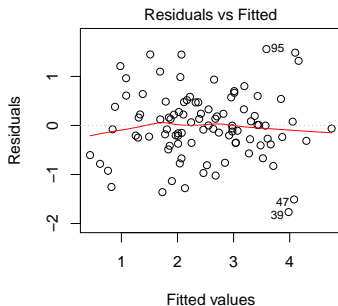
```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.47839	0.07102	34.895	< 2e-16	***
lcavol	0.66515	0.10352	6.425	6.55e-09	***
lweight	0.26648	0.08607	3.096	0.00263	**
age	-0.15820	0.08252	-1.917	0.05848	.
lbph	0.14031	0.08402	1.670	0.09848	.
svi	0.31533	0.09985	3.158	0.00218	**
lcp	-0.14829	0.12566	-1.180	0.24115	
gleason	0.03555	0.11218	0.317	0.75207	
pgg45	0.12572	0.12312	1.021	0.31000	

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Betrachten wir die Anpassung graphisch sowie einen Q-Q-Plot der Residuen

```
> plot(modell)
```



sa ~ lcaivol + lweight + age + lbph + svi + lcp + gleasorsa ~ lcaivol + lweight + age + lbph + svi + lcp + gleasor

Wird der Fit schlechter, wenn wir `age`, `lcp`, `gleason` und `pgg45` weglassen (z.B. weil wir glauben, dass sie in Wirklichkeit keinen Einfluss haben, d.h. die entsprechenden  $\beta_i = 0$  sind)?

```
> summary(modell2)
```

Call:

```
lm(formula = lpsa ~ lcavol + lweight + lbph + svi, data = sdatt)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-1.85483	-0.35990	-0.01492	0.44467	1.53051

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.47839	0.07148	34.673	< 2e-16	***
lcavol	0.62289	0.08781	7.094	2.63e-10	***
lweight	0.22964	0.08415	2.729	0.00761	**
lbph	0.11403	0.08139	1.401	0.16454	
svi	0.29207	0.08610	3.392	0.00102	**

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Wird der Fit schlechter, wenn wir age, lcp, gleason und pgg45 weglassen (z.B. weil wir glauben, dass sie in Wirklichkeit keinen Einfluss haben, d.h. die entsprechenden  $\beta_i = 0$  sind)?

```
> anova(modell2, modell)
```

```
Analysis of Variance Table
```

```
Model 1: lpsa ~ lcavol + lweight + lbph + svi
```

```
Model 2: lpsa ~ lcavol + lweight + age + lbph + svi + lcp
         pgg45
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	92	45.595				
2	88	43.058	4	2.537	1.2963	0.2777

# Bemerkung: $\ell_2$ -Regularisierung, Ridge regression

$\lambda \geq 0$  ein „Regularisierungsparameter,“

$$\hat{\beta}_{\text{ridge}}(\lambda) := \operatorname{argmin} \{ \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 \}$$

Für  $\lambda > 0$  ist  $\hat{\beta}_{\text{ridge}}(\lambda)$  (im Gegensatz zu  $\hat{\beta}_{\text{LS}} = \hat{\beta}_{\text{ridge}}(0)$ ) nicht erwartungstreu,

kann aber ggfs. kleinere (quadratische) „Verlustfunktion“

$$\mathbb{E}_{\beta^*} [\|\hat{\beta}_{\text{ridge}}(\lambda) - \beta^*\|_2^2]$$

haben (“bias-variance trade-off”), speziell wenn die Designmatrix  $\mathbf{X}$  (einige) stark korrelierte Spalten hat.

Beachte: Man läßt typischerweise den “intercept” weg, da der Regularisierungsterm  $\lambda \|\beta\|_2^2$  die „Verschiebungsinvarianz“ des Schätzers bricht.

Die Berechnung von  $\hat{\beta}_{\text{ridge}}(\lambda)$  ist eine quadratische Minimierungsaufgabe, die wieder (im Prinzip) explizit gelöst werden kann:

$$\hat{\beta}_{\text{ridge}}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

(mit  $\mathbf{I} = p \times p$ -Einheitsmatrix)

# Literaturhinweise

- Überblick z.B.: T. Hastie, R. Tibshirani, J. Friedman, *The elements of statistical learning*, Springer, 2009, speziell Ch. 3
- Für lineares Modell mit normalverteiltem Rauschen z.B.: H.O. Georgii, *Stochastik*, de Gruyter, 2002, G. Kersting, A. Wakolbinger, *Elementare Stochastik*, Birkhäuser, 2008
- Monographie: J.S. Stapleton, *Linear statistical models*, Wiley, 1995