

# Stochastik für Informatik 2023

## Schätzer im linearen Modell

Hanno Gottschalk

July 13, 2023

<b>Transformationseigenschaften der Normalverteilung</b>	<b>3</b>
Multivariate Standardnormalverteilung . . . . .	4
Lineare Transformationen von Gauss-Zufallsvariablen . . . . .	5
Unkorreliert impliziert Unabhängig . . . . .	6
Unkorreliert impliziert Unabhängig II. . . . .	7
<b>Verteilung von <math>\hat{\beta}</math></b>	<b>8</b>
Wiederholung lin. Modell . . . . .	9
Erwartungstreue im linearen Modell . . . . .	10
Kovarianz und Verteilung des Schätzers im Linearen Modell . . . . .	11
<b>Versuchsplanung beim linearen Modell</b>	<b>12</b>
Verteilung der Parameterschätzer . . . . .	13
$(\underline{\underline{M}}' \underline{\underline{M}})^{-1}$ und die Versuchsplanung . . . . .	14
Grobe Orientierungshilfen zur Versuchsplanung . . . . .	15
Beispiel Blockdesign . . . . .	16
<b>Die BLUP Eigenschaft</b>	<b>17</b>
Effizienz von Schätzern . . . . .	18
Wer ist BLUP? . . . . .	19
Kleinste Quadrate Schätzer ist BLUP . . . . .	20
Beweis BLUP-Eigenschaft . . . . .	21
<b>Erwartungstreuer Schätzer für die Residuenvarianz</b>	<b>22</b>
Schätzer der Residuenvarianz . . . . .	23
Beweis erwartungstreue Varianzschätzer . . . . .	24
Darstellung des Projektors . . . . .	25
Beweis Projektordarstellung . . . . .	26
$\chi^2$ -Verteilung . . . . .	27
Verteilung des Varianzschätzers . . . . .	28
<b>Explorative Validierung der Modellannahmen</b>	<b>29</b>
Vorbemerkungen . . . . .	30
Explorative Validierung von Modellhypothesen im linearen Modell . . . . .	31

QQ-Plot. . . . .	32
Residuals over Fitted Plot. . . . .	33
TimeSeries Plot. . . . .	34

## Inhaltsverzeichnis der Vorlesung

- Transformationseigenschaften der Normalverteilung
- Verteilung von  $\hat{\beta}$
- Die BLUP Eigenschaft
- Erwartungstreuer Schätzer für die Residuenvarianz
- Explorative Validierung der Modellannahmen
- Konfidenzintervalle für Parameter im Linearen Modell
- Schätzung und Vorhersage - verschiedene Fragen
- Simulierte Konfidenzintervalle bei linearer Regression
- Konfidenzbereich für die Modellvorhersage
- Konfidenzbereich für den nächsten Wert

Hanno Gottschalk

Stochastik für Informatik – 2 / 34

## Transformationseigenschaften der Normalverteilung 3 / 34

### Multivariate Standardnormalverteilung

**Def:** Die Multivariate Standardnormalverteilung ist die Normalverteilung auf  $\mathbb{R}^q$  mit  $\Sigma = \mathbf{1}$  und  $\mu = 0$ , wobei  $\mathbf{1}$  die Einheitsmatrix in  $\mathbb{R}^q$  ist.

Für  $\underline{x} = (x_1, \dots, x_q)$  erhalten wir

$$f_{\underline{X}}(x) = \prod_{j=1}^q \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_j^2} = \frac{1}{(2\pi)^{q/2}} \exp \left\{ -\frac{1}{2} \|\underline{x}\|^2 \right\}$$

**Satz:** Es sei  $\underline{X} \sim N(0, \mathbf{1})$  und  $\underline{U}$  eine orthogonale  $q \times q$ -Matrix,  $\underline{U}\underline{U}' = \mathbf{1}$  ( $\underline{U}' = \underline{U}^{-1}$ ). Dann  $\underline{U}\underline{X} \sim N(0, \mathbf{1})$ .

**Beweis:** Nach der Transformationsformel für Dichten:

$$f_{\underline{U}\underline{X}}(\underline{x}) = \frac{1}{|\underline{U}|} f_{\underline{X}}(\underline{U}'\underline{x}) = \frac{1}{(2\pi)^{q/2}} \exp \left\{ -\frac{1}{2} \underbrace{\|\underline{U}'\underline{x}\|^2}_{=\|\underline{x}\|^2} \right\} \checkmark$$

Hanno Gottschalk

Stochastik für Informatik – 4 / 34

## Lineare Transformationen von Gauss-Zufallsvariablen

**Satz:** Es sei  $\underline{X} \sim N(\underline{\mu}, \sigma^2 \text{Id})$  und  $\underline{A} : \mathbb{R}^d \rightarrow \mathbb{R}^q$  eine Matrix mit vollem Rang. Dann ist  $\underline{Y} = \underline{A}\underline{X} \sim N(\underline{A}\underline{\mu}, \sigma^2 \underline{A}\underline{A}')$ .

**Beweis:** OBdA:  $\underline{\mu} = 0$ . Setze  $M = \underline{A}'\mathbb{R}^q$  und wähle Orthonormalbasis so dass  $\underline{u}_1, \dots, \underline{u}_d$  so dass  $\underline{u}_1, \dots, \underline{u}_q$  eine ONB von  $M$  ist.  $\underline{U}$  zug. Matrix  $\Rightarrow \underline{Y} = \underline{U}\underline{X} \sim N(0, \sigma^2 \text{Id})$  (Transformationssatz + Invarianz Gaußdichte)

Da  $\underline{A}\underline{u}_j = 0$  für  $j = q+1, \dots, n \Rightarrow$

Da  $\underline{A}\underline{U}_q : \mathbb{R}^q \rightarrow \mathbb{R}^q$  bijektiv folgt nun die Aussage nach dem Transformationssatz und  $\underline{U}_q$  orthogonal + Invarianz  $\Rightarrow$

$$\begin{aligned} \underline{A}\underline{U}_q(Y_1, \dots, Y_q) &\sim N(0, \sigma^2 \underline{A}\underline{U}_q \text{Id}_q \underline{U}_q' \underline{A}') \\ &= N(0, \sigma^2 \underline{A}\underline{U} \text{Id}_d \underline{U}' \underline{A}') = N(0, \sigma^2 \underline{A}\underline{A}') \end{aligned}$$

Hanno Gottschalk

Stochastik für Informatik – 5 / 34

## Unkorreliert impliziert Unahängig

**Satz:** Zwei normalverteilte Zufallsvektoren  $\underline{X}_1, \underline{X}_2$  seien unkorreliert  $\Rightarrow \underline{X}_1, \underline{X}_2$  sind *unabhängig*.

Setze  $\underline{X} = (\underline{X}_1, \underline{X}_2)$ , dann

$$\begin{aligned} \Sigma_{\underline{X}} &= \begin{pmatrix} \Sigma_{\underline{X}_1} & 0 \\ 0 & \Sigma_{\underline{X}_2} \end{pmatrix} \\ \langle (\underline{x} - \underline{\mu}), \Sigma_{\underline{X}}^{-1}(\underline{x} - \underline{\mu}) \rangle &= \langle (\underline{x}_1 - \underline{\mu}_1), \Sigma_{\underline{X}_1}^{-1}(\underline{x}_1 - \underline{\mu}_1) \rangle \\ &\quad + \langle (\underline{x}_2 - \underline{\mu}_2), \Sigma_{\underline{X}_2}^{-1}(\underline{x}_2 - \underline{\mu}_2) \rangle \end{aligned}$$

Hanno Gottschalk

Stochastik für Informatik – 6 / 34

## Unkorreliert impliziert Unabhängig II

$$|\Sigma_{\underline{X}}| = |\Sigma_{\underline{X}_1}| |\Sigma_{\underline{X}_2}|$$

$$\begin{aligned} & \frac{1}{(2\pi)^{q/2} |\Sigma_{\underline{X}}|^{1/2}} \exp \left\{ -\frac{1}{2} \langle \underline{x} - \underline{\mu}, \Sigma_{\underline{X}}^{-1} (\underline{x} - \underline{\mu}) \rangle \right\} \\ = & \frac{1}{(2\pi)^{q_1/2} |\Sigma_{\underline{X}_1}|^{1/2}} \exp \left\{ -\frac{1}{2} \langle \underline{x}_1 - \underline{\mu}_1, \Sigma_{\underline{X}_1}^{-1} (\underline{x}_1 - \underline{\mu}_1) \rangle \right\} \\ \times & \frac{1}{(2\pi)^{q_2/2} |\Sigma_{\underline{X}_2}|^{1/2}} \exp \left\{ -\frac{1}{2} \langle \underline{x}_2 - \underline{\mu}_2, \Sigma_{\underline{X}_2}^{-1} (\underline{x}_2 - \underline{\mu}_2) \rangle \right\} \end{aligned}$$

$$f_{\underline{X}}(\underline{x}) = f_{\underline{X}_1}(\underline{x}_1) f_{\underline{X}_2}(\underline{x}_2)$$

Faktorisierung der Dichten ist äquivalent zur Unabhängigkeit. **qed**

Hanno Gottschalk

Stochastik für Informatik – 7 / 34

## Verteilung von $\hat{\underline{\beta}}$

8 / 34

### Wiederholung lin. Modell

$$Y_j = \beta_0 + \underline{x}'_j \underline{\beta} + \epsilon_j, \quad \epsilon_j \sim N(0, \sigma^2) \text{ i.i.d.}$$

Vektorwertige Schreibweise:

$$\underline{Y} = \underline{M} \underline{\beta} + \underline{\epsilon} \quad \underline{M} = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,q-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,q-1} \end{pmatrix}, \quad \underline{\epsilon} \sim N(0, \sigma^2 \mathbf{1})$$

Kleinste-Quadrate Schätzer:

$$\hat{\underline{\beta}} = (\underline{M}' \underline{M})^{-1} \underline{M}' \underline{Y}$$

**Frage:** Da Verteilung von  $\underline{Y} \sim N(\underline{M} \underline{\beta}, \sigma^2 \mathbf{1})$  bekannt ist, kann man dann die Verteilung von  $\hat{\underline{\beta}}$  explizit berechnen?

Hanno Gottschalk

Stochastik für Informatik – 9 / 34

## Erwartungstreue im linearen Modell

**Satz:**  $\hat{\underline{\beta}}$  ist im linearen Modell *erwartungstreu*, also

$$\mathbb{E}_{(\underline{\beta}, \sigma^2)}[\hat{\underline{\beta}}] = \underline{\beta} \quad (1)$$

(Erwartungswert eines Zufallsvektors wird komponentenweise genommen).

**Beweis:** Dank Linearität von  $\mathbb{E}[\cdot]$  und  $\mathbb{E}_{(\underline{\beta}, \sigma^2)}[\underline{Y}] = \underline{M} \underline{\beta}$

$$\begin{aligned} \mathbb{E}_{(\underline{\beta}, \sigma^2)}[\hat{\underline{\beta}}] &= (\underline{M}' \underline{M})^{-1} \underline{M}' \mathbb{E}_{(\underline{\beta}, \sigma^2)}[\underline{Y}] \\ &= (\underline{M}' \underline{M})^{-1} \underline{M}' \underline{M} \underline{\beta} = \underline{\beta} \end{aligned}$$

qed.

Hanno Gottschalk

Stochastik für Informatik – 10 / 34

## Kovarianz und Verteilung des Schätzers im Linearen Modell

**Satz:** Die *Kovarianz* von  $\hat{\underline{\beta}}$  ist

$$\Sigma_{\hat{\underline{\beta}}} = \sigma^2 (\underline{M}' \underline{M})^{-1} \quad (2)$$

**Beweis:** Nach dem Kov-Transformationsgesetz

$$\begin{aligned} \Sigma_{\hat{\underline{\beta}}} &= ((\underline{M}' \underline{M})^{-1} \underline{M}') \Sigma_{\underline{Y}} ((\underline{M}' \underline{M})^{-1} \underline{M}')' \\ &= \sigma^2 ((\underline{M}' \underline{M})^{-1} \underline{M}') \text{Id} ((\underline{M}' \underline{M})^{-1} \underline{M}')' \\ &= \sigma^2 ((\underline{M}' \underline{M})^{-1} \underline{M}' \underline{M} (\underline{M}' \underline{M})^{-1}) \\ &= \sigma^2 (\underline{M}' \underline{M})^{-1} \end{aligned}$$

qed.

**Satz (Folgerung)**  $\hat{\underline{\beta}} \sim N(\underline{\beta}, \sigma^2 (\underline{M}' \underline{M})^{-1})$

Hanno Gottschalk

Stochastik für Informatik – 11 / 34

## Verteilung der Parameterschätzer

Haben in folgenden Satz abgeleitet:

$$\underline{\hat{\beta}} \sim N(\underline{\beta}, \sigma^2(\underline{M}'\underline{M})^{-1})$$

Wende dies an auf den Abnutzungsversuch

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	254.7500	6.1872	41.17	0.0000
run2	-2.2500	5.5340	-0.41	0.6984
run3	12.5000	5.5340	2.26	0.0647
run4	-9.2500	5.5340	-1.67	0.1457
position2	26.2500	5.5340	4.74	0.0032
position3	8.5000	5.5340	1.54	0.1755
position4	8.2500	5.5340	1.49	0.1866
materialB	-45.7500	5.5340	-8.27	0.0002
materialC	-24.0000	5.5340	-4.34	0.0049
materialD	-35.2500	5.5340	-6.37	0.0007

In diesem Falle ist der *Versuchsplan* (latin hypercube) so konstruiert, dass die Standardabweichung (hier: Std. Error) aller Effekte gleich ist.

Hanno Gottschalk

Stochastik für Informatik – 13 / 34

## $(\underline{M}'\underline{M})^{-1}$ und die Versuchsplanung

**Wiederholung:** Interpretierbarkeit der Effekte hängt eng mit *Orthogonalität* der Spaltenvektoren in  $\underline{M}$  zusammen

**Beobachtung:** Bis auf eine (unbekannte) Zahl  $\sigma^2$ , die auf alle Faktoren gleich wirkt, ist die Standardabweichung des Schätzers  $\hat{\beta}_i$  allein durch die Modellmatrix bestimmt

$$\text{Var}(\hat{\beta}_i) = \sigma^2(\underline{M}'\underline{M})_{i,i}^{-1}$$

**Folgerung:** Je nach Fragestellung kann der Versuchsplan schon im Vorhinein optimiert werden, da  $\underline{M}$  nur vom Versuchsplan abhängt!!

Oft existiert sogar eine grobe Vorstellung von  $\beta_i$  und  $\sigma^2$ , so dass man oft *im Vorhinein* das zu erwartende *effect to error*  $\hat{\beta}_i/\text{sd}(\hat{\beta}_i)$  abschätzen kann

Hanno Gottschalk

Stochastik für Informatik – 14 / 34

## Grobe Orientierungshilfen zur Versuchsplanung

### Bauernregeln für Nachweis eines Effektes

- Effect to Std.Error ratio  $\gtrsim 3$
- Mehr Versuche  $\Rightarrow$  weniger Std. Error ("noise")
- Störeffekte berücksichtigen (etwa Blockdesign)  $\Rightarrow \sigma^2$  verringert, aber mehr Parameter
- Bei lin. Reg., möglichst große Parametervariation

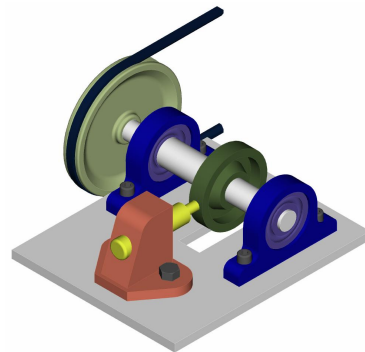
Hanno Gottschalk

Stochastik für Informatik – 15 / 34

## Beispiel Blockdesign

### Blockdesign für Abnutzungstest

- Zielgröße: Abnutzung
- Effekt: Materialeinfluss
- Störeffekt: Testmaschine
- Blockdesign: Jedes Material auf jeder Maschine



**Std. Error für Materialeffekte:** 10.84 ohne Berücksichtigung von 'position' und 8.67 mit! Das Verhältnis 1.25 kann one empirische Daten Vorhergesagt werden!

Hanno Gottschalk

Stochastik für Informatik – 16 / 34



## Effizienz von Schätzern

**Def.:** Zwei Schätzer  $S$  und  $T$  sollen denselben Parameter  $\theta_i$  eines statistischen Modelles schätzen.  $S$  heißt *effizient* verglichen mit  $T$ , falls (für gegebenes  $n$ )

$$\text{MSE}_\theta(S) \leq \text{MSE}_\theta(T) \forall \theta \in \Theta \quad (3)$$

**Def.:** Ein erwartungstreuer Schätzer heißt bester Schätzer, wenn er gegenüber jedem anderen erwartungstreuen Schätzer effizient ist.

Ein Schätzer  $\hat{\theta}$  heißt asymptotisch effizient, wenn für  $n \rightarrow \infty$

$$\limsup \frac{\text{MSE}[\hat{\theta}_n](\theta)}{\text{MSE}[\hat{\theta}'_n](\theta)} \leq 1 \quad \forall \theta \in \Theta$$

für alle erwartungstreuen Schätzer  $\hat{\theta}'$ .

Hanno Gottschalk

Stochastik für Informatik – 18 / 34

## Wer ist BLUP?

Wir schränken die Vergleiche der Effizienz nun auf eine Unterklasse von Schätzern ein:

**Def.:** Ein Schätzer  $S$  für einen Parameter  $\beta_i$  in einem linearen Modell  $\underline{Y} = \underline{M}\beta + \underline{\epsilon}$  heißt *linear*, falls  $S = \langle b, Y \rangle$  für  $b \in \mathbb{R}^n$ .

Insbesondere ist  $\hat{\beta}_i = \langle e_i, (\underline{M}'\underline{M})^{-1}\underline{M}'\underline{Y} \rangle$  linear,  $e_i = (0, 0, \dots, 1, \dots, 0)$  ( $i$ -te Stelle),  $i = 1, \dots, q$

**Def.:** Ein linearer, erwartungstreuer Schätzer  $S$  ist BLUP (Best Linear Unbiased Predictor), wenn  $S$  gegenüber jedem anderen linearen, erwartungstreuen Schätzer effizient ist.

Hanno Gottschalk

Stochastik für Informatik – 19 / 34

## Kleinste Quadrate Schätzer ist BLUP

### Satz: (Markov)

$\hat{\underline{\beta}}$  ist der einzige BLUP im linearen Modell

D.h.  $\forall \underline{c} \in \mathbb{R}^q$  ist  $\langle \underline{c}, \hat{\underline{\beta}} \rangle$  BLUP für  $\langle \underline{c}, \underline{\beta} \rangle$

Insbesondere ist  $\hat{\beta}_i = \langle e_i, \hat{\underline{\beta}} \rangle$  BLUP für  $\beta_i$

**Beweis:** Erwartungstreue ✓

**zu zeigen:**  $\langle \underline{c}, \hat{\underline{\beta}} \rangle$  hat unter allen erwartungstreuen, linearen Schätzern  $\langle \underline{b}, \underline{Y} \rangle$  für  $\langle \underline{c}, \underline{\beta} \rangle$  die kleinste Varianz.

$$\underline{a} := \underline{M} (\underline{M}' \underline{M})^{-1} \underline{c} \Rightarrow \underline{M}' \underline{a} = \underline{c}.$$

$$\begin{aligned} \langle \underline{b}, \underline{M} \underline{\beta} \rangle &= \mathbb{E}_{(\underline{\beta}, \sigma^2)} [\langle \underline{b}, \underline{Y} \rangle] \quad (\langle \underline{b}, \underline{Y} \rangle \text{ erwartungstreu}) \\ &= \langle \underline{c}, \underline{\beta} \rangle = \langle \underline{a}, \underline{M} \underline{\beta} \rangle \quad \forall \underline{\beta} \in \mathbb{R}^q \end{aligned}$$

$$\Rightarrow \underline{b} - \underline{a} \in M^\perp \Rightarrow p_M \underline{b} = \underline{a} \Rightarrow |\underline{b}| \geq |\underline{a}| \text{ mit Gleichheit genau dann wenn } \underline{a} = \underline{b}.$$

Hanno Gottschalk

Stochastik für Informatik – 20 / 34

## Beweis BLUP-Eigenschaft

Zeige nun  $\langle \underline{a}, \underline{Y} \rangle$  ist BLUP:

$$\mathbb{E}_{(\underline{\beta}, \sigma^2)} [\langle \underline{a}, \underline{Y} \rangle] = \langle \underline{a}, \underline{M} \underline{\beta} \rangle = \langle \underline{M}' \underline{a}, \underline{\beta} \rangle = \langle \underline{c}, \underline{\beta} \rangle$$

$$\begin{aligned} \text{Var}_{(\underline{\beta}, \sigma^2)} (\langle \underline{b}, \underline{Y} \rangle) - \text{Var}_{(\underline{\beta}, \sigma^2)} (\langle \underline{a}, \underline{Y} \rangle) &= \sigma^2 (\underline{b}' \text{Id} \underline{b} - \underline{a}' \text{Id} \underline{a}) \\ &= \sigma^2 (|\underline{b}|^2 - |\underline{a}|^2) \geq 0 \end{aligned}$$

Zeige  $\langle \underline{a}, \underline{Y} \rangle = \langle \underline{c}, \hat{\underline{\beta}} \rangle$ :

$$\begin{aligned} \langle \underline{a}, \underline{Y} \rangle &= \langle \underline{M} (\underline{M}' \underline{M})^{-1} \underline{c}, \underline{Y} \rangle \\ &= \langle \underline{c}, (\underline{M}' \underline{M})^{-1} \underline{M}' \underline{Y} \rangle = \langle \underline{c}, \hat{\underline{\beta}} \rangle \end{aligned}$$

qed.

Hanno Gottschalk

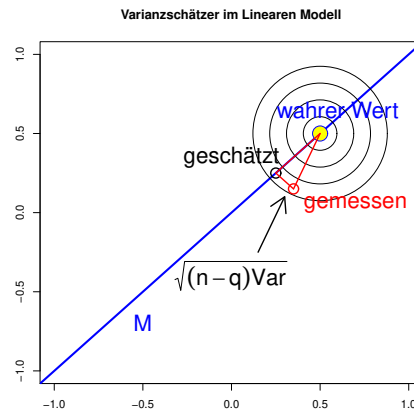
Stochastik für Informatik – 21 / 34

## Schätzer der Residuenvarianz

**Satz:** Im Linearen Modell ist ein erwartungstreuer Schätzer der Varianz gegeben durch

$$\hat{\sigma}^2 = \frac{|\underline{Y} - p_M \underline{Y}|^2}{n - q} = \frac{|p_M^\perp \underline{Y}|^2}{n - q} \quad (4)$$

**Motivation:**



Hanno Gottschalk

Stochastik für Informatik – 23 / 34

## Beweis erwartungstreue Varianzschätzer

OBdA:  $\beta = 0$ , sonst betrachte  $\underline{Y} - \underline{M}\beta$ .

Es sei  $\underline{u}_1, \dots, \underline{u}_n$  ONB so dass  $\underline{u}_1, \dots, \underline{u}_q$  Basis von  $M$ .

$$|p_M^\perp \underline{Y}|^2 = \sum_{j=q+1}^n |\langle \underline{u}_j, \underline{Y} \rangle|^2$$

Aufgrund der Invarianz der Verteilung von  $\underline{Y}$  unter der orthogonalen Transformationen  $\underline{u}$  gilt

$$(\langle \underline{u}_{q+1}, \underline{Y} \rangle, \dots, \langle \underline{u}_n, \underline{Y} \rangle)' \sim N(\underline{0}, \sigma^2 \text{Id}_{(n-q)})$$

$\Rightarrow$

$$\mathbb{E}_{(\underline{0}, \sigma^2)} [|\langle \underline{u}_j, \underline{Y} \rangle|^2] = \sigma^2, \quad j = q + 1, \dots, n$$

**qed.**

Hanno Gottschalk

Stochastik für Informatik – 24 / 34

## Darstellung des Projektors

**Satz:** Der Projektor  $p_M$  auf  $M = \text{Bild}(\underline{\underline{M}})$  hat folgende Darstellung

$$p_M = \underline{\underline{M}} (\underline{\underline{M}}' \underline{\underline{M}})^{-1} \underline{\underline{M}}' \quad (5)$$

**Beweis:** Zu zeigen  $p_M \underline{x} \in M$  ✓,  $p_M \underline{x} = \underline{x}$  für  $\underline{x} \in M$ , und  $\underline{x} - p_M \underline{x} \in M^\perp$  mit  $p_M$  gleich rechte Seite von (1).

Hanno Gottschalk

Stochastik für Informatik – 25 / 34

## Beweis Projektordarstellung

**zu zeigen:**  $p_M \underline{x} = \underline{x}$  für  $\underline{x} \in M$ :

$$\underline{x} \in M \rightarrow \exists \underline{\gamma} \in \mathbb{R}^q, \underline{x} = \underline{\underline{M}} \underline{\gamma}$$

$\Rightarrow$

$$\begin{aligned} p_M \underline{x} &= \underline{\underline{M}} (\underline{\underline{M}}' \underline{\underline{M}})^{-1} \underline{\underline{M}}' \underline{x} \\ &= \underline{\underline{M}} (\underline{\underline{M}}' \underline{\underline{M}})^{-1} \underline{\underline{M}}' \underline{\underline{M}} \underline{\gamma} = \underline{\underline{M}} \underline{\gamma} = \underline{x} \end{aligned}$$

**zu zeigen:**  $\underline{x} - p_M \underline{x} \in M^\perp = \text{Kern}(\underline{\underline{M}}')$

$$\begin{aligned} \underline{\underline{M}}' (\underline{x} - p_M \underline{x}) &= \underline{\underline{M}}' \underline{x} - \underline{\underline{M}} \underline{\underline{M}}' (\underline{\underline{M}}' \underline{\underline{M}})^{-1} \underline{\underline{M}}' \underline{x} \\ &= \underline{\underline{M}}' \underline{x} - \underline{\underline{M}}' \underline{x} = 0 \end{aligned}$$

**qed.**

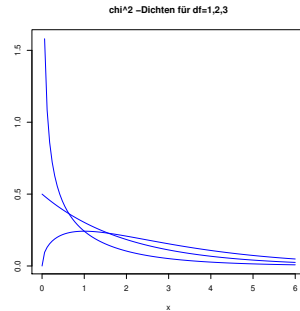
Hanno Gottschalk

Stochastik für Informatik – 26 / 34

## $\chi^2$ -Verteilung

**Def.:** Es seien  $X_1, \dots, X_n$  standardnormalverteilte Zufallsvariablen. Dann ist die  $\chi^2(n)$ -Verteilung –  $\chi^2$ -Verteilung mit  $n$  Freiheitsgraden – definiert durch

$$\sum_{j=1}^n X_j^2 \sim \chi^2(n), \quad f_{\chi^2(n)}(x) = \frac{x^{n/2-1}}{\Gamma(n/2)\sqrt{2}^n} e^{-x/2}, \quad x > 0 \quad (6)$$



Hanno Gottschalk

Stochastik für Informatik – 27 / 34

## Verteilung des Varianzschätzers

**Satz:** (i)  $\frac{(n-q)}{\sigma^2} \hat{\sigma}^2 \sim \chi^2(n-q)$

(ii)  $\hat{\sigma}^2$  ist unabhängig von  $\hat{\beta}$

**Beweis:** (i) ✓

(ii)

$$\underline{\underline{M}} \hat{\beta} = p_M Y = p_M (\underline{\underline{M}} \beta + \epsilon) = \underline{\underline{M}} \beta + p_M \epsilon$$

Da  $\underline{\underline{M}}$  auf  $M$  invertierbar  $\Rightarrow \hat{\beta}$  hängt nur von  $p_M \epsilon$  ab.

$\hat{\sigma}^2$  hängt nur von  $p_M^\perp \epsilon$  ab

Da  $p_M \epsilon$  und  $p_M^\perp \epsilon$  normalvert. + unkorreliert  $\Rightarrow$  unabh.

**qed.**

Hanno Gottschalk

Stochastik für Informatik – 28 / 34

## Vorbemerkungen

**Frage:** Existieren lineare Modelle in der wirklichen Welt?

Wenn etwas zu *schön* ist um wahr zu sein. . . dann *ist* es nicht wahr!

Das lineare Modell *ist* zu schön um wahr zu sein. . .

**Neue Frage:** Kann mich jemand dran kriegen, wenn ich es trotzdem mache?

Mache die Tests, die andere machen würden, und wenn ich mich selbst nicht drankriege, dann können mich auch andere nicht drankriegen

Zeige: Mein Modell ist '**state of the art**' (aber nicht besser)!

Hanno Gottschalk

Stochastik für Informatik – 30 / 34

## Explorative Validierung von Modellhypothesen im linearen Modell

Dem Statistischen Modell 'lineares Modell' liegen drei Hypothesen zugrunde

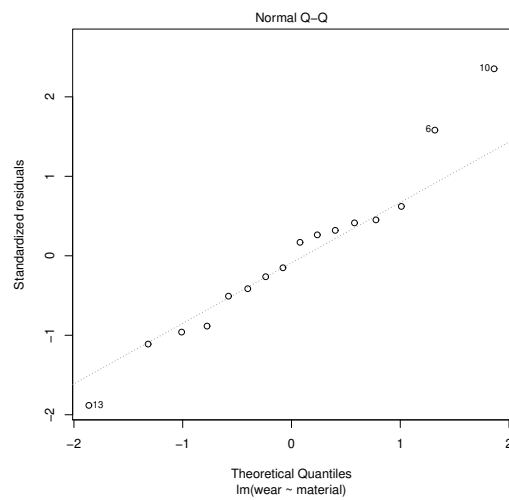
- Die Residuen  $\epsilon_j$  sind *normalverteilt*  $\rightarrow$  Überprüfe mit QQ-plot (s.u.)
- Die Residuen haben alle dieselbe Varianz  $\sigma^2$   $\rightarrow$  Überprüfe mit 'Residuals over Fitted' plot
- Die Residuen sind unabhängig (unkorreliert)  $\rightarrow$  Überprüfe mit time-series-plot (s.u.)

**Def.:** QQ-Plot: *Falls* die Residuen einer Normalverteilung  $N(0, \sigma^2)$  entstammen, dann sollte für die *empirischen Quantile*  $q_\alpha$  der empirischen Residuen  $\epsilon_j$  gelten  $q_\alpha \approx \sigma z_\alpha$  mit  $z_\alpha$  dem  $\alpha$ -Quantil der Standardnormalverteilung.

Hanno Gottschalk

Stochastik für Informatik – 31 / 34

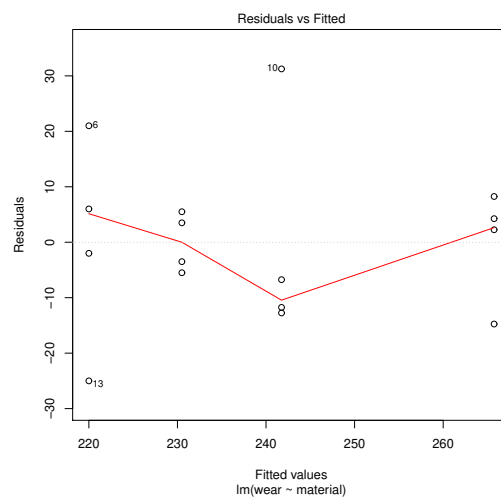
## QQ-Plot



Hanno Gottschalk

Stochastik für Informatik – 32 / 34

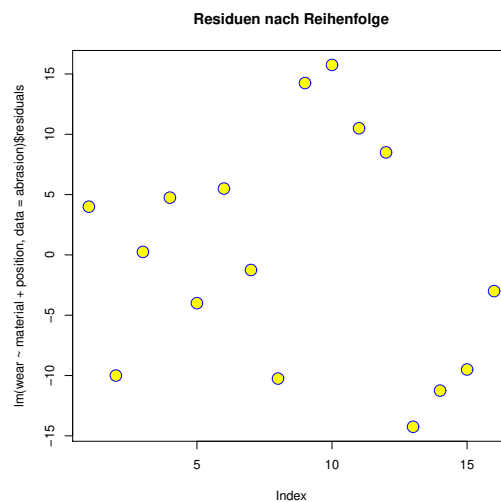
## Residuals over Fitted Plot



Hanno Gottschalk

Stochastik für Informatik – 33 / 34

## TimeSeries Plot



Hanno Gottschalk

Stochastik für Informatik – 34 / 34

## Konfidenzbereiche für Parameter im linearen Modell 35 / 34

### Verteilung der standardisierten Parameter

Haben gesehen

$$\underline{\hat{\beta}} \sim N(\underline{\beta}, \sigma^2 (\underline{\underline{M}}' \underline{\underline{M}})^{-1})$$

Außerdem

$$(n - q) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n - q)$$

**Satz:** Es gilt

$$\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma} \sqrt{(\underline{\underline{M}}' \underline{\underline{M}})^{-1}_{i,i}}} \sim t(n - q) \quad (7)$$

**Denn:**

$$\frac{\frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{(\underline{\underline{M}}' \underline{\underline{M}})^{-1}_{i,i}}}}{\sqrt{n - q} \frac{\hat{\sigma}}{\sigma}} \sim t(n - q)$$

Hanno Gottschalk

Stochastik für Informatik – 36 / 34



## Konfidenzintervalle für die Parameter im lin. Modell

**Folgerung:** Aus dem vorangehenden Satz ergeben sich folgende Konfidenzintervalle für die Parameter im lin. Modell:

Beidseitiges Konfidenzintervall für  $\hat{\beta}_i$

$$\hat{\beta}_i \pm t_{1-\alpha/2}(n-q)\hat{\sigma}\sqrt{(\underline{\underline{M}}'\underline{\underline{M}})_{i,i}^{-1}} \quad (8)$$

Einseitiges Konfidenzintervall: Linksoffen

$$\left[ -\infty, \hat{\beta}_i + t_{1-\alpha}(n-q)\hat{\sigma}\sqrt{(\underline{\underline{M}}'\underline{\underline{M}})_{i,i}^{-1}} \right) \quad (9)$$

Einseitiges Konfidenzintervall: rechtsoffen

$$\left[ \hat{\beta}_i - t_{1-\alpha}(n-q)\hat{\sigma}\sqrt{(\underline{\underline{M}}'\underline{\underline{M}})_{i,i}^{-1}}, \infty \right) \quad (10)$$

Hanno Gottschalk

Stochastik für Informatik – 37 / 34

## Die 'summary' Tabelle und Sternchen-Code im lin. Modell

In der Summary Tabelle finden wir Schätzwert, Wert der T-Statistik ( $p$ -Wert) und Sternchen-Code.

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	255.0000	8.1129	31.43	0.0000	***
materialB	-45.7500	8.6731	-5.27	0.0005	***
materialC	-24.0000	8.6731	-2.77	0.0219	*
materialD	-35.2500	8.6731	-4.06	0.0028	**
position2	26.2500	8.6731	3.03	0.0143	*
position3	8.5000	8.6731	0.98	0.3527	
position4	8.2500	8.6731	0.95	0.3663	

- Kein Stern: 0 ist im 2 seitigen 90% Konfiintervall
- '.': 0 ist nicht im 2 seitigen 90% Konfiintervall
- '\*\*': 0 ist nicht im 2 seitigen 95% Konfiintervall
- '\*\*\*': 0 ist nicht im 2 seitigen 99% Konfiintervall
- '\*\*\*\*': 0 ist nicht im 2 seitigen 99.9% Konfiintervall

Hanno Gottschalk

Stochastik für Informatik – 38 / 34

## Interpretation des Sternchencodes

$\beta_i = 0 \Rightarrow$  den  $\beta_i$ -Effekt gibt es nicht!

Wenn 0 *nicht* im 2-seitigen Konfibereich zu  $1 - \alpha$  Konfidenz liegt, dann können wir uns zu  $(1 - \alpha) \times 100\%$  sicher sein, dass es  $\beta_i$  doch gibt!

Diese Art zu schließen wird später in der Testtheorie noch formalisiert. . .

Hanno Gottschalk

Stochastik für Informatik – 39 / 34

## Schätzung und Vorhersage - verschiedene Fragen 40 / 34

### Fragen an lineare Modelle

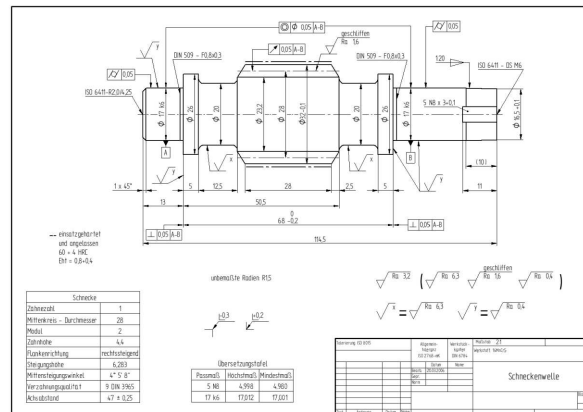
Lineare Modelle können verschiedenen Zwecken dienen

- Herausfinden, welche Effekte wichtig sind (Screening)
- Quantifizierung der Effekte (Schätzen)
- Bei einem noch nicht beobachteten Fall mit bekannten Einflußgrößen  $x_1^{\text{neu}}, \dots, x_d^{\text{neu}}$  die response  $y^{\text{neu}}$  *vorhersagen*
- Reststreuung quantifizieren (Streuung, die in dem Modell nicht reduziert werden kann)

Hanno Gottschalk

Stochastik für Informatik – 41 / 34

Die Kommunikation zwischen Designer/innen und Produzent/innen erfolgt über technische Zeichnungen.



Hanno Gottschalk

Stochastik für Informatik – 42 / 34

## Toleranzen nach ISO 2768

Toleranzen werden entweder explizit, oder durch Verweis auf ISO 2768 vorgegeben

Grenzmaße für Längenmaße entsprechend DIN ISO 2768-1				
Nennmaßbereich in mm	Toleranzklassen			
	f (fein)	m (mittel)	c (grob)	v (sehr grob)
	Toleranzen in mm			
0,5 bis 3	± 0,05	± 0,10	± 0,15	-
über 3 bis 6	± 0,05	± 0,10	± 0,20	± 0,50
über 6 bis 30	± 0,10	± 0,20	± 0,50	± 1,00
über 30 bis 120	± 0,15	± 0,30	± 0,80	± 1,50
über 120 bis 400	± 0,20	± 0,50	± 1,20	± 2,50
über 400 bis 1000	± 0,30	± 0,80	± 2,00	± 4,00
über 1000 bis 2000	± 0,50	± 1,20	± 3,00	± 6,00
über 2000 bis 4000	-	± 2,00	± 4,00	± 8,00

Hanno Gottschalk

Stochastik für Informatik – 43 / 34

## Qualitätskontrolle

Produktion von Schneckenwellen

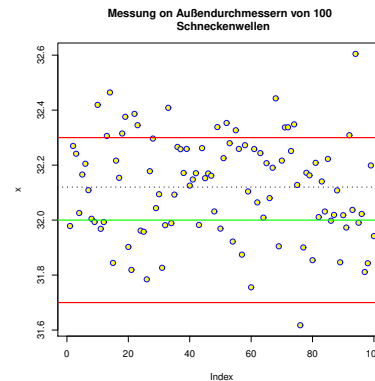
- 32mm Außendurchmesser
- Toleranz (Stufe m)  $\pm 0.3\text{mm}$

**Frage:** Wo liegt Prozessmitte?

→ Konfidenzintervall Mittelwert

**Frage:** 90% Sicherheit, dass neues Teil unter Toleranzobergrenze liegt?

→ Konfidenzintervall Vorhersage neues Teil



Hanno Gottschalk

Stochastik für Informatik – 44 / 34

## Simulierte Konfidenzintervalle bei linearer Regression 45 / 34

### Ein 'Gedankenexperiment' bei lin. Reg.

Gedankenexperiment: Gegeben sei ein 'wahres' einfaches lineares Modell mit  $\alpha = \beta = \sigma^2 = 1$

$$y_i = 1 + x_i + \epsilon_i$$

Messe  $Y$  jeweils 11 mal für  $x_1 = 0, x_2 = 0.1, \dots, x_{11} = 1$

Erzeuge hierfür jeweils 11 Standardnormalverteilte Residuenwerte  $\epsilon_i$

Fitte gerade und wiederhole diese Simulation 100 mal!

Hanno Gottschalk

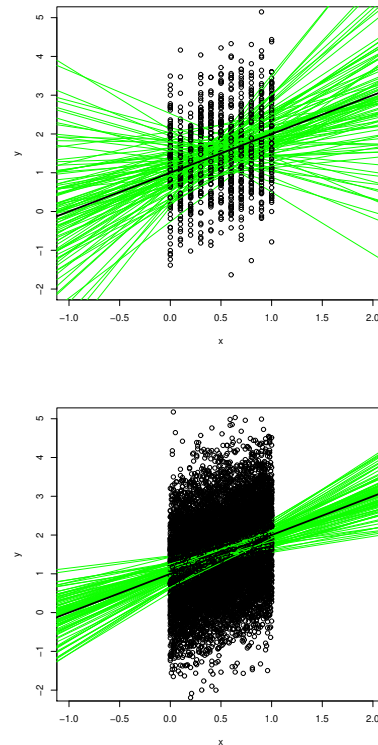
Stochastik für Informatik – 46 / 34

## Ergebnis Simulation

Grüne Zone: Unsicherheit Vorhersage  
lin.Mod

Schwarze Zone: Residuenstreuung

Gesamtunsicherheit Vorhersage neuer  
Punkt: 'Grün+Schwarz'



Die grüne Zone lässt sich mit mehr Messpunkten verkleinern, die schwarze nicht!

Hanno Gottschalk

Stochastik für Informatik – 47 / 34

## Konfidenzbereich für die Modellvorhersage

48 / 34

### Fragestellung Modellvorhersage

Gegeben ein Lineares Modell  $\underline{Y} = \underline{\underline{M}}\underline{\beta} + \underline{\epsilon}$ , wie groß ist die Unsicherheit für die Vorhersage des *Mittelwertes* der Response  $Y$  für neue Daten  $\underline{x}^{\text{new}} = (x_1^{\text{new}}, \dots, x_d^{\text{new}})$ ?

Aus den Daten erhalte ich  $\underline{m}' = (g_1(\underline{x}^{\text{new}}), \dots, g_q(\underline{x}^{\text{new}}))$

**Def:** Schätzung des Mittelwertes

$$\hat{y} = \underline{m}'\hat{\underline{\beta}} \quad (11)$$

**Satz:** Nach dem Transformationsgesetz für normalvert. Z.V. und dem Satz über die Verteilung von  $\hat{\underline{\beta}}$  gilt

$$\hat{y} \sim N(\underline{m}'\underline{\beta}, \sigma^2 \underline{m}'(\underline{\underline{M}}'\underline{\underline{M}})^{-1}\underline{m}) \quad (12)$$

Hanno Gottschalk

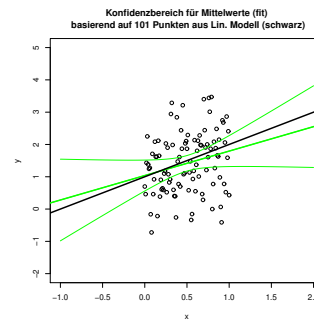
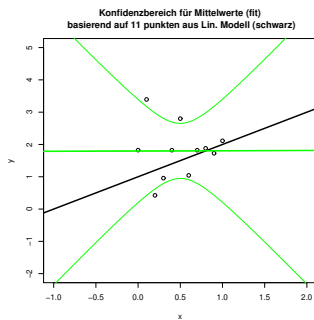
Stochastik für Informatik – 49 / 34

## Konfidenzintervalle Vorhersage Mittelwerte

**Satz:** Ein 2-seitiges Konfidenzbereich zum Konfidenzniveau  $1 - \alpha$  für die Modellvorhersage des Mittelwerts bei  $\underline{x}^{\text{new}}$  ist

$$\underline{m}'\hat{\underline{\beta}} \pm t_{1-\alpha/2}(n-q)\hat{\sigma} \sqrt{\underline{m}'(\underline{M}'\underline{M})^{-1}\underline{m}} \quad (13)$$

Die Aussage  $\underline{m}'\hat{\underline{\beta}}$  liegt in diesem Intervall hat Konfidenzniveau  $1 - \alpha$  (Analog für einseitige Kofi-Intervalle)



Hanno Gottschalk

Stochastik für Informatik – 50 / 34

## Konfidenzbereich für Vorhersage neuer Wert

51 / 34

### Modell für Vorhersage neuer Wert

Man stelle sich vor, dass eine Zeile 'new' aus dem linearen Modell gelöscht worden wäre. . .

$$Y^{\text{new}} = \underline{m}'\underline{\beta} + \epsilon^{\text{new}} \quad (14)$$

Die Unsicherheit der Vorhersage  $\underline{m}'\hat{\underline{\beta}}$  von  $Y^{\text{new}}$  ergibt sich

- 1) Aus der Unsicherheit von  $\hat{\underline{\beta}}$
- 2) Aus der Streuung von  $\epsilon^{\text{new}}$

$\epsilon^{\text{new}} \sim N(0, \sigma^2)$  gleichverteilt zu und unabhängig von  $\epsilon_i$

Hanno Gottschalk

Stochastik für Informatik – 52 / 34

## Konfidenzbereich für die Vorhersage neuer Werte

**Satz:** Die Verteilung von  $Y^{\text{new}}$  ist im linearen Modell

$$Y^{\text{new}} \sim N(\underline{m}'\underline{\beta}, \sigma^2(\underline{m}'(\underline{M}'\underline{M})^{-1}\underline{m} + 1)) \quad (15)$$

**Satz:** Im lin. Modell ist die Vorhersage mit W-keit  $1 - \alpha$  wahr, dass der neue, unbeobachtete Wert zu  $\underline{x}^{\text{new}}$  mit zugehöriger Modell-Matrix-Zeile  $\underline{m}'$  im folgenden Bereich liegt

$$\underline{m}'\hat{\underline{\beta}} \pm t_{1-\alpha/2}(n-q)\hat{\sigma}\sqrt{\underline{m}'(\underline{M}'\underline{M})^{-1}\underline{m} + 1} \quad (16)$$

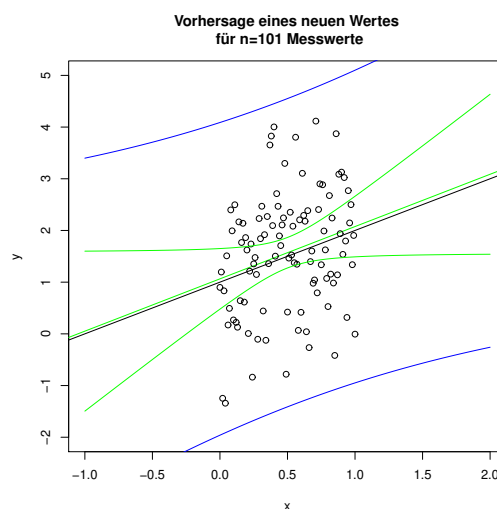
Analog für halboffene Konfi-Intervalle

Hanno Gottschalk

Stochastik für Informatik – 53 / 34

## Konfidenz Vorhersage lin. Reg. neuer Wert

In Blau die Konfidenzbereiche für neue Werte, in grün die für den Mittelwert



Hanno Gottschalk

Stochastik für Informatik – 54 / 34