

Gemischte Modelle

Oktober 2016

Sonja Greven (V) und David Rügamer (Ü)

Institut für Statistik

Ludwig-Maximilians-Universität München

<https://www.elab.moodle.elearning.lmu.de/course/view.php?id=1074>



Mit Dank an Susanne Konrath und Fabian Scheipl für Material vergangener Jahre.

Literatur Gemischte Modelle

- Demidenko, E. (2004). *Mixed Models: Theory and Applications*. Wiley.
- Diggle, P. J.; Heagerty, P.; Liang, K. L.; Zeger, S. L. (2002). *Analysis of Longitudinal Data*. Oxford University Press, Oxford.
- **Fahrmeir, L., Kneib, T. und Lang, S. (2009)**. *Regression: Modelle, Methoden und Anwendungen (2. Auflage)*. Springer. - Begleitend zur Vorlesung. Erhältlich als Ebook bei der Universitätsbibliothek: <https://opac.ub.uni-muenchen.de/>
- McCulloch, C. E.; Searle, S. R. (2001). *Generalized, Linear, and Mixed Models*. John Wiley.
- Pinheiro, J. C.; Bates, D. M. (2000). *Mixed-Effects Models in S and S-PLUS*. Springer, New York. - Praxisorientierte Einführung in die Analyse gemischter Modelle und ausführliche Beschreibung des R-Pakets `nlme` für LMMs.
- Ruppert, D.; Wand, M. P.; Carroll, R. J. (2003). *Semiparametric Regression*. Cambridge University Press. - Verbindung gemischte Modelle und Penalisierung.
- Verbeke, G.; Molenberghs, G. (2000). *Linear Mixed Models for Longitudinal Data*. Springer, New York.
- Wood, S. (2006). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC. - Verbindung gemischte Modelle und Penalisierung, R-Paket `mgcv`.

Inhalt der Vorlesung Gemischte Modelle

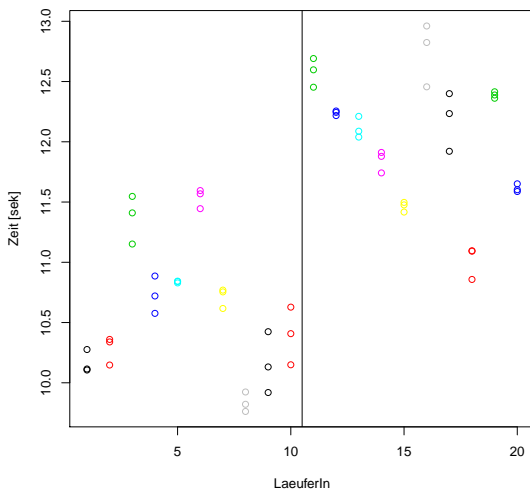
- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Inhalt

- 1 Das lineare gemischte Modell
 - Motivation
 - Das allgemeine lineare gemischte Modell (LMM)
 - Spezialfälle
 - Die Kovarianzstruktur
 - Konditionale und Marginale Perspektive
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Motivation: 100-Meter-Lauf

Drei 100-Meter-Läufe mit 10 Männern und 10 Frauen. Die Zeiten sehen so aus:



Motivation: 100-Meter-Lauf

Wir wollen ein Modell, das uns ermöglicht:

- die Schätzung des Geschlechtseffekts (Populationsparameter)
- die Schätzung der LäuferInneneffekte (individuelle Effekte)
- die Schätzung der Korrelationsstruktur
- valide Inferenz.

100-Meter-Lauf: Ein lineares gemischtes Modell

Überlegung: Hätten wir pro Person ihre Durchschnittszeit μ_i (d.h. ein Wert pro Person), wäre ein sinnvolles Modell (iid = unabhängig identisch verteilt)

$$\mu_i = \beta_{g_i} + b_i, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad (1)$$

da die Durchschnittszeiten einzelner Personen um das Geschlechtsmittel β_{g_i} , $g_i \in \{1, 2\}$, variieren.

Die beobachteten Zeiten pro Lauf streuen um die individuelle Durchschnittszeit:

$$y_{ij} = \mu_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2). \quad (2)$$

Zusammen genommen ergibt sich damit das Modell (\perp = unabhängig)

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2), \quad b_i \perp \varepsilon_{ij}. \quad (3)$$

Ein Blick auf das Modell

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2). \quad (3)$$

Die Effekte b_i sind

- die „Fehler“terme in (1)
- Teil des Erwartungswertes in (2)
- sogenannte **zufällige Effekte** in Modell (3). Sie spiegeln hier wieder, dass die LäuferInnen aus einer Population kommen (in der wir die individuellen Durchschnittszeiten als normalverteilt um das Geschlechtsmittel annehmen).

Ein Modell mit zufälligen und festen Effekten nennt man **gemischtes Modell**.

Interpretation der Parameter

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2). \quad (3)$$

- β_1, β_2 die Durchschnittszeiten für Männer / Frauen (**Populationsparameter**)
- b_i die Abweichung der Durchschnittszeit von Person i vom Mittel der Geschlechtsgruppe (**individuelle Effekte**)
- τ^2 die Varianz der Durchschnittszeiten pro Geschlechtsgruppe (in den zwei Gruppen als gleich **angenommen**)
- ε_{ij} die Abweichung der j -ten Laufzeit von der Durchschnittszeit für Person i
- σ^2 die Varianz der persönlichen Laufzeiten (für alle i gleich **angenommen**)

Bedingte Sicht

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2). \quad (3)$$

Betrachte die Verteilung von y_{ij} bedingt auf b_i (\rightarrow **überlegen**):

$$y_{ij}|b_i \stackrel{iid}{\sim} N(\beta_{g_i} + b_i, \sigma^2)$$

b_i modelliert individuelle Effekte im **Erwartungswert** (EW) analog zum linearen Modell mit festen Effekten

$$y_{ij} = \beta_{g_i} + \beta_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2).$$

Aber:

- Die Geschlechtseffekte sind identifizierbar (nicht kollinear mit den β_i).
- Das Modell ist bei 1-2 Messungen für einige i schätzbar (mehr später).
- Da wir im Wesentlichen τ^2 schätzen (zur Vorhersage der b_i ; mehr später), wächst die Zahl der Parameter nicht mit der Anzahl der LäuferInnen.

Marginale Sicht

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2). \quad (3)$$

Betrachte die marginale Verteilung der y_{ij} (\rightarrow **überlegen**):

$$\mathbf{y}_i = \begin{pmatrix} y_{i1} \\ y_{i2} \\ y_{i3} \end{pmatrix} \stackrel{iid}{\sim} N\left(\begin{pmatrix} \beta_{g_i} \\ \beta_{g_i} \\ \beta_{g_i} \end{pmatrix}, \begin{pmatrix} \tau^2 + \sigma^2 & \tau^2 & \tau^2 \\ \tau^2 & \tau^2 + \sigma^2 & \tau^2 \\ \tau^2 & \tau^2 & \tau^2 + \sigma^2 \end{pmatrix} \right) \quad (4)$$

b_i induziert eine **Kovarianzstruktur** analog zum allgemeinen linearen Modell

$$y_{ij} = \beta_{g_i} + \varepsilon_{ij}, \quad \boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \varepsilon_{i2}, \varepsilon_{i3}) \sim N(\mathbf{0}, \mathbf{V}_i). \quad (5)$$

Hier:

- Begründung für eine mögliche Kovarianzstruktur \mathbf{V}_i (auch unbalanziert).
- Vorhersage individueller Effekte möglich (mehr später).

(3) impliziert das allgemeine lineare Modell (5), aber nicht umgekehrt!

Wozu gemischte Modelle?

Anhand des Beispiels zeigt sich bereits: Gemischte Modelle werden gerne verwendet für die Analyse korrelierter Daten. Z.B.

- **Longitudinaldaten:** Wiederholte Beobachtungen in zeitlicher Abfolge an denselben Subjekten/Beobachtungseinheiten. (z.B. Patienten über die Zeit)
- **Clusterdaten / gruppierte Daten:** Gruppen (Cluster) mit mehreren Beobachtungen pro Gruppe. (z.B. Läuferdaten, Daten für Familien)
- **Hierarchische Daten:** Gruppierte Daten mit mehreren geschachtelten Ebenen. (z.B. Schüler in Klassen in Schulen, Patienten in Ärzten in Krankenhäusern)
- **Gekreuzte Designs:** Gruppierte Daten mit mehreren gekreuzten Ebenen. (z.B. alle Personen bearbeiten die gleichen Aufgaben)

Daten vom gleichen Subjekt/Cluster/Beobachtungseinheit sind sich tendenziell ähnlicher als Daten verschiedener Subjekte/Cluster/Beobachtungseinheiten.

Das allgemeine lineare gemischte Modell (LMM)

Definition: In allgemeiner Form ist das **lineare gemischte Modell** gegeben durch

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times r}{\mathbf{Z}} \underset{r \times 1}{\mathbf{b}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}. \quad (6)$$

Mit $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R})$, $\mathbf{b} \sim N(\mathbf{0}, \mathbf{G})$ und der Annahme, dass die zufälligen Effekte \mathbf{b} und die Fehler $\boldsymbol{\varepsilon}$ **unabhängig** sind, ist die **Verteilungsannahme** gegeben durch

$$\begin{pmatrix} \mathbf{b} \\ \boldsymbol{\varepsilon} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix} \right). \quad (7)$$

Die Kovarianzmatrizen \mathbf{G} für \mathbf{b} und \mathbf{R} für $\boldsymbol{\varepsilon}$ werden als **positiv semi-definit** bzw. **positiv definit** angenommen.

Zu den Verteilungsannahmen

- Die **Normalverteilungsannahme der Fehler** ε ist nicht für alle Aussagen zur Inferenz in LMMs notwendig. Da Likelihood-basierte Schätzung üblich ist, nehmen wir sie jedoch in die Definition mit auf.
- Für die **zufälligen Effekte** \mathbf{b} ist die **Normalverteilungsannahme** nicht zwingend. Alternative Verteilungen, z.B. Mischungsverteilungen, möglich. I.d.R. werden dann die Algorithmen zur Berechnung der Schätzer komplexer.
- ε und \mathbf{b} sind als unabhängig angenommen.

Vorteile der Analyse mit gemischten Modellen

Zufällige Effekte können als Platzhalter für die Effekte von unbeobachteten oder unzureichend gemessenen Kovariablen dienen, die **Korrelation zwischen Beobachtungen an den gleichen Beobachtungseinheiten** verursachen.

Im Gegensatz zum linearen Regressionsmodell mit unabhängigen Fehlern führt die Berücksichtigung dieser Korrelation

- zu einer **verbesserten Schätzgenauigkeit** (kleineren wahren Standardfehlern) → **Übung**
- zu validen modellbasierten Standardfehlern und damit Konfidenzintervallen und Tests → **Übung**

Vorteile der Analyse mit gemischten Modellen

Im Gegensatz zum linearen Regressionsmodell mit festen Effekten für die Beobachtungseinheiten

- können Effekte für Kovariablen (z.B. Geschlecht), die nur zwischen Beobachtungseinheiten variieren, geschätzt werden.
- werden die festen Effekte wegen der kleineren Anzahl von Modellparametern effizienter geschätzt

Im Gegensatz zum linearen Regressionsmodell mit allgemeiner Kovarianz

- erlauben die Vorhersagen für die zufälligen Effekte **individuelle Prognosen**.

Weitere Annahmen in gemischten Modellen

Auch wenn die zufälligen Effekte die Effekte von unbeobachteten Kovariablen auffangen, so können Sie **nicht Confounding** durch solche Kovariablen verhindern.

Dies liegt daran, dass die zufälligen Effekte (marginal, siehe später), die Kovarianzstruktur beeinflussen, nicht jedoch den Erwartungswert.

Streng genommen, lautet in (7) die Annahme an ε und \mathbf{b} :

$$E(\varepsilon|\mathbf{X}, \mathbf{b}) = \mathbf{0}, \quad E(\mathbf{b}|\mathbf{X}) = \mathbf{0}$$

(Regressionsannahme und Random effects-Annahme), so dass $E(\mathbf{y}|\mathbf{X}) = \mathbf{X}\beta$. Dies ist insbesondere verletzt, wenn Confounding vorliegt.

Beispiel Confounding

Betrachte das Modell

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + u_{ij}\gamma + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2).$$

y_{ij} misst den Schulerfolg von Kind j in Schule i ; u_{ij} ist a) Förderbedarf des Kindes oder b) Familieneinkommen. Beispiele für mögliches Confounding wären

- a) Bei Förderbedarf wählen Eltern Schulen mit unterstützender Schulkultur, die sich auch im Schulerfolg widerspiegelt (Random effects-Annahme verletzt).
- b) Familieneinkommen korreliert mit Sprachkenntnissen, Ziele der Eltern für ihre Kinder etc., die mit Schulerfolg korrelieren (Regressionsannahme verletzt).

Ein Modell mit festen Effekten b_i würde den Bias durch a), jedoch nicht den Bias durch b) verhindern. Ziel muss auf jeden Fall sein, in \mathbf{x}_{ij} mögliche Confounder möglichst gut abzubilden.

Frei nach Clarke, Crawford, Steele & Vignoles (2010): *The Choice Between Fixed and Random Effects Models: Some Considerations for Educational Research*. IZA Discussion Paper No. 5287

Beispiel 100-Meter-Lauf

Das Modell für die 100-Meter-Lauf-Daten

$$y_{ij} = \beta_{g_i} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2), \quad b_i \perp \varepsilon_{ij},$$

lässt sich in die Form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}), \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}), \quad \mathbf{b} \perp \boldsymbol{\varepsilon}$$

bringen mit

$$\mathbf{y} = (y_{11}, \dots, y_{20,3})'$$

$$\mathbf{b} = (b_1, \dots, b_{20})$$

$$\mathbf{X} = \begin{pmatrix} 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 \end{pmatrix}'$$

$$\mathbf{G} = \tau^2 \mathbf{I}_{20}$$

$$\boldsymbol{\beta} = (\beta_1, \beta_2)'$$

$$\boldsymbol{\varepsilon} = (\varepsilon_{11}, \dots, \varepsilon_{20,3})$$

$$\mathbf{Z} = \begin{pmatrix} 1 & 1 & 1 & \dots & 0 & 0 & 0 \\ & & & \ddots & & & \\ 0 & 0 & 0 & \dots & 1 & 1 & 1 \end{pmatrix}'$$

$$\mathbf{R} = \sigma^2 \mathbf{I}_{60}.$$

Spezialfall Longitudinal- und Clusterdaten

Wiederholte Beobachtungen y_{ij} der Zielvariablen → **Beispiel in Übung**

- von Subjekt i zum Zeitpunkt t_{ij} bei **Longitudinaldaten**
- für das j -te Objekt aus dem Cluster i bei **Clusterdaten**

mit jeweils Kovariablenvektor $(\mathbf{x}'_{ij}, \mathbf{z}'_{ij})'$, $i = 1, \dots, N$, $j = 1, \dots, n_i$.

Verschiedene **Variabilitätsquellen** in den Daten:

- Zwischen den Subjekten/Clustern, Abweichungen vom Populationsmittel.
- Innerhalb des Subjekts/Clusters, Abweichungen einer Messung vom Mittelwert des entsprechenden Subjekts/Clusters.

Spezialfall Longitudinal- und Clusterdaten

Das lineare gemischte Modell auf Beobachtungs- bzw. Cluster/Subjekt-Ebene ist

$$\begin{aligned}y_{ij} &= \mathbf{x}'_{ij}\boldsymbol{\beta} + \mathbf{z}'_{ij}\mathbf{b}_i + \varepsilon_{ij}, & j = 1, \dots, n_i, i = 1, \dots, N & \quad \text{bzw.} \\ \mathbf{y}_i &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\varepsilon}_i, & i = 1, \dots, N.\end{aligned}$$

wobei \mathbf{X}_i und \mathbf{Z}_i die n_i Zeilen \mathbf{x}'_{ij} bzw. \mathbf{z}'_{ij} enthalten und $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})$.

Annahmen:

- Unabhängig und identisch normalverteilte zufällige Effekte $\mathbf{b}_i \stackrel{iid}{\sim} N(\mathbf{0}, \mathbf{D})$,
- unabhängig und normalverteilte Fehler $\varepsilon_i \sim N(\mathbf{0}, \boldsymbol{\Sigma}_i)$,
- $\mathbf{b}_1, \dots, \mathbf{b}_N, \varepsilon_1, \dots, \varepsilon_N$ unabhängig,
- \mathbf{D} positiv semi-definit und $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_N$ positiv definit.

Spezialfall Longitudinal- und Clusterdaten

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \varepsilon_i, \quad i = 1, \dots, N.$$

- $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ entspricht den (festen) **Populationseffekten**.

- $\mathbf{b}_i = (b_{i1}, \dots, b_{iq})'$ entspricht den (zufälligen) **subjekt- / clusterspezifischen Effekten**.

- $\mathbf{X}_i = \begin{pmatrix} \mathbf{x}'_{i1} \\ \vdots \\ \mathbf{x}'_{in_i} \end{pmatrix}$, $\mathbf{Z}_i = \begin{pmatrix} \mathbf{z}'_{i1} \\ \vdots \\ \mathbf{z}'_{in_i} \end{pmatrix}$ sind die **Designmatrizen** für die p populationsspezifischen bzw. q subjektspezifischen Kovariablen.

Dabei können die Kovariablen mit j (bzw. t_{ij} , *zeitvariierend*) variieren oder nicht.

Falls die Variablen \mathbf{z}_{ij} in \mathbf{x}_{ij} enthalten sind, lassen sich die \mathbf{b}_i mit $E(\mathbf{b}_i) = \mathbf{0}$ als **individuelle Abweichungen** vom Populationsmittel interpretieren.

Spezialfall Longitudinal- und Clusterdaten

Das lineare gemischte Modell für Longitudinal- und Clusterdaten ist ein Spezialfall des allgemeinen LMMs

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}), \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}), \quad \mathbf{b} \perp \boldsymbol{\varepsilon}$$

mit

- $\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_N)'$ und $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}'_1, \dots, \boldsymbol{\varepsilon}'_N)'$ der Länge $n = \sum_{i=1}^N n_i$,
- $\boldsymbol{\beta}$ der Länge p ,
- $\mathbf{b} = (\mathbf{b}'_1, \dots, \mathbf{b}'_N)'$ der Länge $r = Nq$,
- $\mathbf{X} = (\mathbf{X}'_1 | \dots | \mathbf{X}'_N)'$ der Dimension $n \times p$,
- $\mathbf{Z} = \text{blockdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_N)$ der Dimension $n \times Nq$,
- $\mathbf{G} = \text{blockdiag}(\mathbf{D}, \dots, \mathbf{D}, \dots, \mathbf{D})$ der Dimension $Nq \times Nq$,
- $\mathbf{R} = \text{blockdiag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_i, \dots, \boldsymbol{\Sigma}_N)$ der Dimension $n \times n$.

Die blockdiagonalen Kovarianzmatrizen resultieren aus der Unabhängigkeitsannahme für Beobachtungen an verschiedenen Individuen / Clustern.

Spezialfall hierarchische Struktur: Beispiel

Die Lesefähigkeit von 875 achtjährigen Schülern in 29 Klassen in 11 Schulen wird anhand eines standardisierten Scores y gemessen. Mögliches Modell:

$$y_{ijk} = \beta_0 + b_i + b_{ij} + \varepsilon_{ijk},$$

$$b_i \stackrel{iid}{\sim} N(0, \tau_1^2), \quad b_{ij} \stackrel{iid}{\sim} N(0, \tau_2^2), \quad \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2), \quad b_i, b_{ij}, \varepsilon_{ijk} \text{ unabh.}$$

wobei

- y_{ijk} : der Lesescore des k -ten Kindes in der j -ten Klasse der i -ten Schule
- β_0 : die allgemeine mittlere Lesefähigkeit von Achjährigen
- b_i : die Abweichung der mittleren Lesefähigkeit in Schule i vom allgemeinen Mittel
- b_{ij} : die Abweichung der mittleren Lesefähigkeit der Klasse j vom Mittel ihrer Schule i
- ε_{ijk} : die Abweichung der Lesefähigkeit von Kind k von der mittleren Lesefähigkeit seiner Klasse.

Spezialfall hierarchische Struktur: Beispiel

Das Modell lässt sich wieder in die allgemeine LMM-Form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}), \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}), \quad \mathbf{b} \perp \boldsymbol{\varepsilon}$$

bringen mit

- $\mathbf{y} = (y_{1,1,1}, \dots, y_{11,3,30})'$ und $\boldsymbol{\varepsilon} = (\varepsilon_{1,1,1}, \dots, \varepsilon_{11,3,30})'$ der Länge $n = 875$,
- $\mathbf{X} = (1, \dots, 1)'$ der Dimension 875×1 ,
- $\boldsymbol{\beta} = \beta_0$,
- $\mathbf{b} = (b_1, \dots, b_{11} | b_{1,1}, \dots, b_{11,3})'$ der Länge $r = 11 + 29 = 40$,
- $\mathbf{Z} = (\text{blockdiag}(\mathbf{Z}_1, \dots, \mathbf{Z}_{11}) | \text{blockdiag}(\mathbf{Z}_{1,1}, \dots, \mathbf{Z}_{11,3}))$ der Dimension 875×40 , wobei \mathbf{Z}_i bzw. \mathbf{Z}_{ij} Einservektoren sind mit Länge gleich der Anzahl der Schüler (bzw. der Schüler in Klasse j) in Schule i ,
- $\mathbf{G} = \text{blockdiag}(\tau_1^2 \mathbf{I}_{11}, \tau_2^2 \mathbf{I}_{29})$ der Dimension 40×40 ,
- $\mathbf{R} = \sigma^2 \mathbf{I}_{875}$ der Dimension 875×875 .

Spezialfall gekreuzte Struktur: Beispiel

In einem Phonetik-Experiment sprechen 9 Subjekte 140 Worte, in denen *s*- und *sch*-Laute vorkommen, je 5 mal. Ein akustischer Index y misst, ob der gesprochene Laut einem *s* oder *sch* näher ist.

Mögliches Modell für y_{ijk} (i te Person, j tes Wort, k te Wiederholung):

$$y_{ijk} = \mathbf{x}'_j \boldsymbol{\beta} + b_i + c_j + d_{ij} + \varepsilon_{ijk},$$

$$b_i \stackrel{iid}{\sim} N(0, \tau_b^2), \quad c_j \stackrel{iid}{\sim} N(0, \tau_c^2), \quad d_{ij} \stackrel{iid}{\sim} N(0, \tau_d^2), \quad \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2),$$

$$b_i, c_j, d_{ij}, \varepsilon_{ijk} \text{ unabhängig, } i = 1, \dots, 9; j = 1, \dots, 140; k = 1, \dots, 5,$$

mit

- $\boldsymbol{\beta}$ Effekte von Wortmerkmalen wie Betonung etc.
- zufällige Effekte b_i für Person i , c_j für Wort j und d_{ij} für deren Interaktion.

Ein einfacheres Modell wäre das Modell ohne Interaktion d_{ij} .

Spezialfall gekreuzte Struktur: Beispiel

Das Modell lässt sich wieder in die allgemeine LMM-Form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}), \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}), \quad \mathbf{b} \perp \boldsymbol{\varepsilon}$$

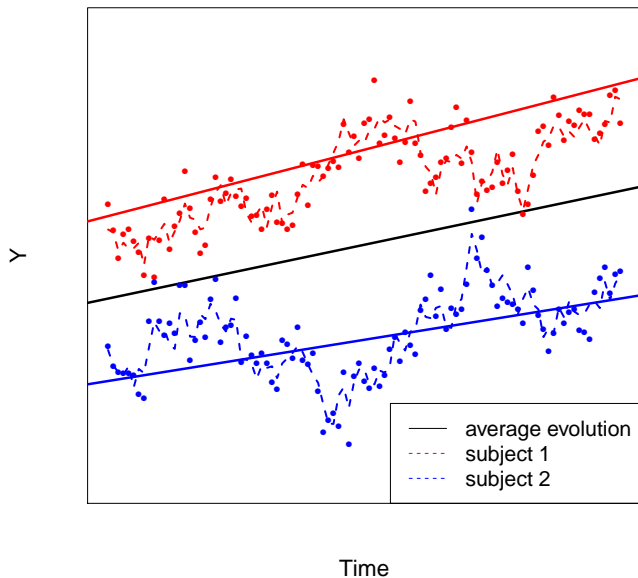
bringen mit

- $\mathbf{y} = (y_{1,1,1}, \dots, y_{9,140,5})'$ und $\boldsymbol{\varepsilon} = (\varepsilon_{1,1,1}, \dots, \varepsilon_{9,140,5})'$ der Länge $n = 9 \cdot 140 \cdot 5 = 6300$,
- \mathbf{X} der Dimension $6300 \times p$ enthält \mathbf{x}_j zeilenweise,
- $\mathbf{Z} = (\mathbf{I}_9 \otimes \mathbf{1}_{140 \cdot 5} | \mathbf{1}_9 \otimes \mathbf{I}_{140} \otimes \mathbf{1}_5 | \mathbf{I}_{9 \cdot 140} \otimes \mathbf{1}_5)$ der Dimension 6300×1409 , wobei $\mathbf{1}_l$ der Einservektor der Länge l ist und \otimes das Kroneckerprodukt.
- $\mathbf{b} = (b_1, \dots, b_9 | c_1, \dots, c_{140} | d_{1,1}, \dots, d_{9,140})'$ der Länge $r = 9 + 140 + 9 \cdot 140 = 1409$,
- $\mathbf{G} = \text{blockdiag}(\tau_b^2 \mathbf{I}_9, \tau_c^2 \mathbf{I}_{140}, \tau_d^2 \mathbf{I}_{9 \cdot 140})$ der Dimension 1409×1409 ,
- $\mathbf{R} = \sigma^2 \mathbf{I}_{6300}$ der Dimension 6300×6300 .

Die Kovarianzstruktur

- \mathbf{G} und \mathbf{R} modellieren die Abhängigkeitsstruktur von \mathbf{b} bzw. ε . Zusätzliche Annahmen (z.B. Diagonalmatrix) - unabhängig für \mathbf{G} und \mathbf{R} - ergeben Modelle verschiedener Komplexität und Flexibilität.
- \mathbf{G} , \mathbf{R} und \mathbf{Z} implizieren zusammen die Kovarianzstruktur für \mathbf{y} ,
$$\mathbf{V} = \text{Cov}(\mathbf{y}) = \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R}.$$
- $\mathbf{Z}\mathbf{b}$ zusammen mit \mathbf{G} modelliert Unterschiede zwischen Beobachtungseinheiten - z.B. zwischen Schülern und Klassen.
- ε ist der Fehlerterm. \mathbf{R} fängt möglicherweise verbleibende Autokorrelation auf, die nicht durch $\mathbf{Z}\mathbf{b}$ erklärt wird.

Die Kovarianzstruktur - longitudinales Beispiel



Conditional Independence Model

Die stärkste Annahme für die Fehler ist, dass sie unabhängig und identisch normalverteilt sind, $\mathbf{R} = \sigma^2 \mathbf{I}_n$ oder

$$\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2), i = 1, \dots, n \Leftrightarrow \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

Aus der Unabhängigkeit der Fehler folgt die bedingte Unabhängigkeit der y_i gegeben \mathbf{b} , also von $y_i | \mathbf{b}, \dots, y_n | \mathbf{b}$.

Die **Korrelation** zwischen den Beobachtungen y_i wird im **Conditional Independence Model** nur durch den Vektor \mathbf{b} der **zufälligen Effekten** erzeugt.

Werden die zufälligen Effekte zusätzlich unabhängig angenommen, (bei Longitudinal-/Clusterdaten $\mathbf{D} = \text{diag}(\tau_1^2, \dots, \tau_q^2)$ diagonal) so spricht man von einem **Varianzkomponentenmodell**.

Spezialfall Random Intercept Modell

Bei Designvektor $\mathbf{z}'_{ij} = 1$ ergibt sich das **Random Intercept Modell**

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2)$$

mit individuellen Interzepten (Beispiel 100-Meter-Lauf mit $\mathbf{x}'_{ij}\boldsymbol{\beta} = \beta_{g_i}$).

In Kombination mit $\varepsilon_i \sim N(0, \sigma^2 \mathbf{I}_{n_i})$ führt dies zur marginalen Kovarianzstruktur

$$\begin{aligned} \text{Cov}(y_{ij}, y_{ik}) &= \tau^2 + \sigma^2 \delta_{jk} \\ \Rightarrow \text{Corr}(y_{ij}, y_{ik}) &= \frac{\tau^2}{\sigma^2 + \tau^2} =: \rho \geq 0, \quad j \neq k \end{aligned}$$

(mit Kronecker-Delta $\delta_{jk} = 1$ für $j = k$, $\delta_{jk} = 0$ sonst.)

Block-konstante Korrelationsstruktur der Zielvariablen (**Compound Symmetry**).

ρ groß wenn interindividuelle Varianz groß relativ zur intraindividuellen Varianz.

Random Intercept-Random Slope Modell

Beim *Random Intercept-Random Slope Modell*

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + b_{0i} + b_{1i}t_{ij} + \varepsilon_{ij}, \quad (b_{0i}, b_{1i})' \stackrel{iid}{\sim} N\left(\mathbf{0}, \begin{pmatrix} \tau_1^2 & \tau_{12} \\ \tau_{12} & \tau_2^2 \end{pmatrix}\right)$$

unterscheiden sich individuelle Interzepte und Steigungen, z.B. über die Zeit. Bei $\varepsilon_i \sim N(0, \sigma^2 \mathbf{I}_{n_i})$ ergibt sich eine quadratische Varianzfunktion:

$$\begin{aligned} \text{Var}(y_{ij}) &= \tau_1^2 + 2\tau_{12}t_{ij} + \tau_2^2 t_{ij}^2 + \sigma^2 \quad \text{und} \\ \text{Cov}(y_{ij}, y_{ik}) &= \tau_1^2 + \tau_{12}t_{ij} + \tau_{12}t_{ik} + \tau_2^2 t_{ij}t_{ik}, \quad j \neq k. \end{aligned}$$

Bei zusätzlichem quadratischen Term $b_{2i}t_{ij}^2$ ergibt sich ein Polynom 4. Ordnung.

Nicht-longitudinales Beispiel: t_{ij} Dosis eines Medikaments - b_{1i} berücksichtigt individuelle Unterschiede in der Reaktion auf das Medikament.

Allgemeines R

Manchmal ist die Annahme $R = \sigma^2 I_n$ zu vereinfachend.

Autokorrelation

Bei Longitudinaldaten (\rightarrow ALD) z.B. wird R block-diagonal gewählt, mit (bei balanzierten Daten) unstrukturierten Kovarianzen Σ_j oder

$$\text{Cov}(\varepsilon_{ij}, \varepsilon_{ik}) = \underbrace{\sigma_1^2 \delta_{jk}}_{\text{„Messfehler“}} + \underbrace{\sigma_2^2 g(|t_{ij} - t_{ik}|)}_{\text{Autokorrelation}}$$

für eine monoton fallende Funktion $g(\cdot)$ mit $g(0) = 1$ und $\lim_{u \rightarrow \infty} g(u) = 0$.

Häufig nur entweder unabhängiger oder autokorrelierter Teil gut schätzbar.

Heteroskedastizität

Bei den 100-Meter-Lauf-Daten könnte man z.B. zulassen, dass die Varianz σ_{gi}^2 der individuellen Zeiten vom Geschlecht abhängt.

Konditionale und marginale Perspektive

Konditionale oder bedingte Perspektive auf das gemischte Modell:

$$\mathbf{y}|\mathbf{b} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \mathbf{R}), \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}). \quad (8)$$

Interpretation: Die zufälligen Effekte sind individuelle Effekte (von Beobachtungseinheiten), die in der Population variieren und unter Normalverteilungsannahme geschätzt werden.

In dieser hierarchischen Formulierung des LMM wird der Erwartungswert von y_i als Funktion von Populationseffekten und individuellen Effekten modelliert.

Marginale Perspektive auf das gemischte Modell:

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}) \quad \text{mit} \quad \mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R} \quad (9)$$

Interpretation: Die zufälligen Effekte induzieren eine Korrelationsstruktur und ermöglichen so eine valide statistische Analyse korrelierter Daten.

In der marginalen Formulierung des LMM wird der marginale, über die Population gemittelte Erwartungswert von y_i als Funktion von Populationseffekten modelliert.

Konditionale und marginale Perspektive

- Aus der hierarchischen Darstellung (8) folgt die marginale Darstellung (9).
- Bezeichne mit p die Dichten der entsprechenden Verteilungen. Dann ist einfach zu zeigen, dass die Dichte der marginalen Verteilung

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{b})p(\mathbf{b})d\mathbf{b}$$

die Dichte einer Normalverteilung (NV) mit EW $\mathbf{X}\beta$ und Kovarianz $\mathbf{V} = \mathbf{ZGZ}' + \mathbf{R}$ ist.

- Im LMM, wenn $\mathbf{y}|\mathbf{b}$ normalverteilt ist, lässt sich diese Integration analytisch durchführen. Dass dies für andere Verteilungen der Exponentialfamilie nicht geht, ist ein wesentlicher Grund dafür, dass die Inferenz für **generalisierte lineare gemischte Modelle** (GLMMs) schwieriger ist als für LMMs.

Konditionale und marginale Perspektive

- Aus der marginalen Verteilung von \mathbf{y} alleine folgt nicht die bedingte Verteilung für \mathbf{y} gegeben \mathbf{b} und die Verteilung von \mathbf{b} .
- Das marginale Modell für sich betrachtet nimmt keine zufälligen Effekte an, um Heterogenitäten darzustellen.
- Nicht jede Kovarianz \mathbf{V} erlaubt hierarchische Interpretation, $\mathbf{V} = \mathbf{ZGZ}' + \mathbf{R}$.

Hierarchisches und **marginale Modell** sind also **nicht äquivalent**. Beispiel:

- Random Intercept Modell hierarchisch: τ^2 als Varianz muss nicht-negativ sein.
- Marginal Compound Symmetry: τ^2 als Kovarianz könnte negativ werden.

Dennoch gleiche Interpretation der festen Effekte β in hierarchischer und marginaler Formulierung des LMM. Dies ist aber i.A. für GLMMs nicht der Fall! (Mehr dazu in Kapiteln 6+7.)

Konsequenzen für die Schätzung

- Hauptinteresse nur an festen Effekten β
⇒ Verwendung des marginalen Modells.
- Interesse an festen und zufälligen Effekten β , \mathbf{b} sowie den Varianz-/Kovarianzkomponenten in der Kovarianzmatrix \mathbf{D}
⇒ Verwendung der hierarchischen Darstellung.

Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
 - Schätzung der festen und Vorhersage der zufälligen Effekte
 - Schätzung der Kovarianzstruktur
 - Numerische Berechnung der Schätzer
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Schätzung der festen Effekte

Für die Schätzung der festen Effekte β verwenden wir die marginale Verteilung (9), $\mathbf{y} \sim N(\mathbf{X}\beta, \mathbf{V})$.

Wir nehmen zunächst an, dass die Kovarianzen $\text{Cov}(\mathbf{b}) = \mathbf{G}$ und $\text{Cov}(\varepsilon) = \mathbf{R}$, und damit $\text{Cov}(\mathbf{y}) = \mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R}$, bekannt sind. Dann ist (9) ein allgemeines lineares Modell. Das verallgemeinerte Kleinste-Quadrate (KQ)-Kriterium

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta)$$

für β ergibt den verallgemeinerten KQ-Schätzer (Aitken-Schätzer)

$$\hat{\beta} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y}. \quad (10)$$

(Wir nehmen hier an, dass die Inversen \mathbf{V}^{-1} und $(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1}$ existieren. Verallgemeinerungen mit generalisierten Inversen existieren.)

Schätzung der festen Effekte

Dies ist gleichzeitig der Maximum-Likelihood (ML)-Schätzer, wenn wir die Normalverteilungsannahme treffen: Die log-Likelihood für β aus dem marginalen Modell (9) ist

$$l(\beta) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta). \quad (11)$$

$$\frac{d}{d\beta} l(\beta) = \mathbf{X}' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) \stackrel{!}{=} \mathbf{0} \quad \Rightarrow \quad \hat{\beta} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y}.$$

Optimalität des Schätzers für β

$\hat{\beta}$ ist der BLUE, der beste lineare erwartungstreue Schätzer (best linear unbiased estimator) für β . **Beweis:**

Erwartungstreue $E(\hat{\beta}) = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\beta = \beta$ ist klar.

Sei $\tilde{\beta} = \mathbf{a} + \mathbf{B}\mathbf{y}$ ein anderer linearer erwartungstreuer Schätzer für β .

$\Rightarrow E(\tilde{\beta}) = \mathbf{a} + \mathbf{B}\mathbf{X}\beta \stackrel{!}{=} \beta \quad \forall \beta \Rightarrow \mathbf{a} = \mathbf{0}$ und $\mathbf{B}\mathbf{X} = \mathbf{I}_p$.

\Rightarrow Für alle $\mathbf{c} \in \mathbb{R}^p$ gilt

$$\begin{aligned} & \text{Var}(\mathbf{c}'\tilde{\beta}) - \text{Var}(\mathbf{c}'\hat{\beta}) \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}\mathbf{B}'\mathbf{c} - \mathbf{c}'(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{V}\mathbf{V}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{c} \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}\mathbf{B}'\mathbf{c} - \mathbf{c}'\mathbf{B}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{B}'\mathbf{c} \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}^{1/2}[\mathbf{I}_n - \mathbf{V}^{-1/2}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1/2}]\mathbf{V}^{1/2}\mathbf{B}'\mathbf{c} \geq 0, \end{aligned}$$

da die mittlere Matrix eine Projektionsmatrix (idempotent) und damit positiv semi-definit ist. Damit ist $\hat{\beta}$ der lineare erwartungstreue Schätzer mit der kleinsten Varianz. □

Optimalität des Schätzers für β

- Unter der Normalverteilungsannahme kann man sogar zeigen, dass $\hat{\beta}$ der beste erwartungstreue Schätzer ist.
- Mit der Transformation $\mathbf{X}^* = \mathbf{V}^{-1/2}\mathbf{X}$, $\mathbf{y}^* = \mathbf{V}^{-1/2}\mathbf{y}$, $\varepsilon^* = \mathbf{V}^{-1/2}(\mathbf{Z}\mathbf{b} + \varepsilon)$ gilt $\mathbf{y}^* = \mathbf{X}^*\beta + \varepsilon^*$ mit $\varepsilon^* \sim N(\mathbf{0}, \mathbf{I}_n)$.

Somit kann das marginale Modell $\mathbf{y} \sim N(\mathbf{X}\beta, \mathbf{V})$ auf das gewöhnliche lineare Modell $\mathbf{y}^* \sim N(\mathbf{X}^*\beta, \mathbf{I}_n)$ zurückgeführt werden.

Die Optimalitätseigenschaften von $\hat{\beta}$ ergeben sich dann unmittelbar aus dem Gauß-Markov-Theorem für das gewöhnliche lineare Modell.

Prädiktion der zufälligen Effekte

Oft ist es von Interesse, auch Vorhersagen für die zufälligen Effekte \mathbf{b} zu erhalten. Dazu benötigen wir die hierarchische Modellformulierung (8).

Unter der Normalverteilungsannahme ist der bedingte Erwartungswert $\hat{\mathbf{b}} := E(\mathbf{b}|\mathbf{y})$ von \mathbf{b} , gegeben die Daten \mathbf{y} , die **beste lineare erwartungstreue Vorhersage** (BLUP, best linear unbiased prediction) für \mathbf{b} .

- $\hat{\mathbf{b}}$ ist wegen des Satzes vom iterierten Erwartungswert **unverzerrt**:

$$E(\hat{\mathbf{b}}) = E(E(\mathbf{b}|\mathbf{y})) = E(\mathbf{b}) = \mathbf{0}.$$

Für die zufälligen Effekte wird hierbei Unverzerrtheit als $E(\hat{\mathbf{b}}) = E(\mathbf{b})$ definiert und nicht etwa als $E(\hat{\mathbf{b}}|\mathbf{b}) = \mathbf{b}$ für alle \mathbf{b} .

Prädiktion der zufälligen Effekte

- $\hat{\mathbf{b}} = E(\mathbf{b}|\mathbf{y})$ minimiert $E[(\tilde{\mathbf{b}} - \mathbf{b})'(\tilde{\mathbf{b}} - \mathbf{b})]$ in der Klasse der unverzerrten, linearen Schätzer $\tilde{\mathbf{b}}$:

$$\begin{aligned}
 E[(\tilde{\mathbf{b}} - \mathbf{b})'(\tilde{\mathbf{b}} - \mathbf{b})] &= E[(\tilde{\mathbf{b}} - \hat{\mathbf{b}} + \hat{\mathbf{b}} - \mathbf{b})'(\tilde{\mathbf{b}} - \hat{\mathbf{b}} + \hat{\mathbf{b}} - \mathbf{b})] \\
 &= E[(\tilde{\mathbf{b}} - \hat{\mathbf{b}})'(\tilde{\mathbf{b}} - \hat{\mathbf{b}})] + 2 E[(\tilde{\mathbf{b}} - \hat{\mathbf{b}})'(\hat{\mathbf{b}} - \mathbf{b})] \\
 &\quad + E[(\hat{\mathbf{b}} - \mathbf{b})'(\hat{\mathbf{b}} - \mathbf{b})] \\
 &= E[(\tilde{\mathbf{b}} - \hat{\mathbf{b}})'(\tilde{\mathbf{b}} - \hat{\mathbf{b}})] + E[(\hat{\mathbf{b}} - \mathbf{b})'(\hat{\mathbf{b}} - \mathbf{b})] \\
 &\geq E[(\hat{\mathbf{b}} - \mathbf{b})'(\hat{\mathbf{b}} - \mathbf{b})],
 \end{aligned}$$

da $E[(\tilde{\mathbf{b}} - \hat{\mathbf{b}})'(\hat{\mathbf{b}} - \mathbf{b})] = E_{\mathbf{y}} E_{\mathbf{b}|\mathbf{y}}[(\tilde{\mathbf{b}} - \hat{\mathbf{b}})'(E(\mathbf{b}|\mathbf{y}) - \mathbf{b})] = \mathbf{0}$.

Gleichheit gilt für $\tilde{\mathbf{b}} = \hat{\mathbf{b}}$.

Prädiktion der zufälligen Effekte

Um $\hat{\mathbf{b}}$ explizit zu bestimmen, betrachte die gemeinsame Verteilung von \mathbf{y} und \mathbf{b} . Diese ergibt sich aus dem zweistufigen hierarchischen Modell für $\mathbf{y}|\mathbf{b}$ und \mathbf{b} zu

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{b} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V} & \mathbf{Z}\mathbf{G} \\ \mathbf{G}\mathbf{Z}' & \mathbf{G} \end{pmatrix} \right). \quad (12)$$

Herleitung für Interessierte: \rightarrow **Extrablatt**.

Klar: Aus (12) folgt wiederum (Satz B.4, Fahrmeir et al., 2009)

- $\mathbf{b} \sim N(\mathbf{0}, \mathbf{G})$
- $\mathbf{y}|\mathbf{b} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \mathbf{V} - \mathbf{Z}\mathbf{G}\mathbf{G}^{-1}\mathbf{G}\mathbf{Z}' = \mathbf{R})$.

Aus (12) folgt (Satz B.4, Fahrmeir et al., 2009) der bedingte Erwartungswert

$$E(\mathbf{b}|\mathbf{y}) = \mathbf{G}\mathbf{Z}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (13)$$

In der Praxis wird der unbekanntes Vektor $\boldsymbol{\beta}$ in (13) durch den verallgemeinerten KQ-Schätzer $\hat{\boldsymbol{\beta}}$ aus (10) ersetzt,

$$\hat{\mathbf{b}} = \mathbf{G}\mathbf{Z}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = (\mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1})^{-1}\mathbf{Z}'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

Henderson's mixed model equations

BLUE und BLUP (gemeinsam auch mit BLUP bezeichnet) $\hat{\beta}$ und $\hat{\mathbf{b}}$ sind die Lösung der sogenannten Henderson's mixed model equations

$$\begin{aligned} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X}\hat{\beta} + \mathbf{X}'\mathbf{R}^{-1}\mathbf{Z}\hat{\mathbf{b}} &= \mathbf{X}'\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{Z}'\mathbf{R}^{-1}\mathbf{X}\hat{\beta} + (\mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1})\hat{\mathbf{b}} &= \mathbf{Z}'\mathbf{R}^{-1}\mathbf{y} \end{aligned} \quad (14)$$

Für diese Schätzgleichungen gibt es mehrere Herleitungen, u.a. folgende drei:

- 1 $\hat{\beta}$ und $\hat{\mathbf{b}}$ haben minimale erwartete quadratische Abweichung unter den linearen erwartungstreuen Schätzern/Vorhersagen.
(Analog: Linearkombinationen, z.B. $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} + \mathbf{Z}\hat{\mathbf{b}}$.)
- 2 Eine bayesianische Herleitung (mehr in Kapitel 4).
- 3 Die Schätzgleichungen ergeben sich bei simultaner Maximierung (bzgl. β und \mathbf{b}) der sogenannten **gemeinsamen Likelihood** $L(\beta, \mathbf{b})$ basierend auf der gemeinsamen Dichte von \mathbf{y} und \mathbf{b} .

Gemeinsame / penalisierte Likelihood

Wegen $p(\mathbf{y}, \mathbf{b}) = p(\mathbf{y}|\mathbf{b})p(\mathbf{b})$ ergibt sich (bis auf additive Konstanten)

$$l(\boldsymbol{\beta}, \mathbf{b}) = \log L(\boldsymbol{\beta}, \mathbf{b}) = -\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})' \mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) - \frac{1}{2}\mathbf{b}' \mathbf{G}^{-1}\mathbf{b}. \quad (15)$$

Diese kann als **penalisierte Likelihood** aufgefasst werden mit **Strafterm** $\mathbf{b}' \mathbf{G}^{-1}\mathbf{b}$.
Die Maximierung ist äquivalent zur Minimierung des **penalisierten KQ-Kriteriums**

$$KQ_{pen}(\boldsymbol{\beta}, \mathbf{b}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})' \mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) + \mathbf{b}' \mathbf{G}^{-1}\mathbf{b}. \quad (16)$$

- Der erste Term entspricht einem verallgemeinerten KQ-Kriterium.
- Der zweite Term berücksichtigt die **Verteilungsannahme** für \mathbf{b} und bestraft ($\mathbf{b}' \mathbf{G}^{-1}\mathbf{b} \geq 0$ wegen $\mathbf{G} \geq 0$ als Kovarianz) Abweichungen von $E(\mathbf{b}) = \mathbf{0}$.
 - Für $\mathbf{G} = \mathbf{0}$ gilt $\hat{\mathbf{b}} = \mathbf{0}$.
 - Für $\mathbf{G}^{-1} \rightarrow \mathbf{0}$ geht $\mathbf{b}' \mathbf{G}^{-1}\mathbf{b} \rightarrow 0$ und \mathbf{b} wird wie ein fester Effekt geschätzt.

Differenzieren von $KQ_{pen}(\boldsymbol{\beta}, \mathbf{b})$ und Nullsetzen der Ableitungen ergibt (14).

Definiert man $\mathbf{C} = (\mathbf{X}|\mathbf{Z})$ und die partitionierte Matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^{-1} \end{pmatrix}, \quad (17)$$

so lässt sich die Lösung der Schätzgleichung (14) auch in der folgenden kompakten Form schreiben

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{b}} \end{pmatrix} = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1}\mathbf{C}'\mathbf{R}^{-1}\mathbf{y}. \quad (18)$$

In dieser Form erkennt man die enge Beziehung zur [Ridge-Schätzung](#). Ebenso besteht ein enger Zusammenhang zur (empirischen) Bayes-Schätzung (mehr dazu später).

Shrinkage

$\hat{\mathbf{b}}$ ist ein **Shrinkage-Schätzer** für \mathbf{b} , d.h. seine Komponenten haben eine geringere Streuung als sie hätten, wenn \mathbf{b} als feste Effekte behandelt würden.

Dies scheint auf den ersten Blick der Eigenschaft der Unverzerrtheit des BLUP zu widersprechen. Erwartungstreue bedeutet hier jedoch

$$E(\hat{\mathbf{b}}) = E(\mathbf{b}) = \mathbf{0},$$

und nicht

$$E(\hat{\mathbf{b}}|\mathbf{b}) = \mathbf{b} \text{ für alle } \mathbf{b}.$$

Shrinkage

$$\begin{aligned}
 \hat{\mathbf{y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\mathbf{b}} \\
 &= \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\mathbf{G}\mathbf{Z}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\
 &= (\mathbf{V} - \mathbf{Z}\mathbf{G}\mathbf{Z}')\mathbf{V}^{-1}\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\mathbf{G}\mathbf{Z}'\mathbf{V}^{-1}\mathbf{y} \\
 &= \mathbf{R}\mathbf{V}^{-1}\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\mathbf{G}\mathbf{Z}'\mathbf{V}^{-1}\mathbf{y}.
 \end{aligned}$$

Der BLUP für \mathbf{y} ist also ein **gewichtetes Mittel** des Populationsmittels $\mathbf{X}\hat{\boldsymbol{\beta}}$ und der individuellen Daten \mathbf{y} . Die beobachteten Daten werden gegen das Populationsmittel „geschrumpft“ (“borrowing of strength”).

Beachte, dass $\mathbf{V} = \mathbf{R} + \mathbf{Z}\mathbf{G}\mathbf{Z}'$.

- Wenn die Varianz der Residuen \mathbf{R} im Verhältnis zu $\mathbf{Z}\mathbf{G}\mathbf{Z}'$ groß ist, wird viel Gewicht auf das Populationsmittel $\mathbf{X}\hat{\boldsymbol{\beta}}$ gelegt.
- Wenn \mathbf{R} klein ist, gilt das Gegenteil.

Shrinkage: Varianz-Bias-Tradeoff

Beispiel Random-Intercept-Modell (ohne Kovariablen)

$$y_{ij} = b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \tau^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2), \\ j = 1, \dots, n_i, i = 1, \dots, N.$$

Für den BLUP $\hat{\mathbf{b}}^{BLUP}$ und den KQ-Schätzer $\hat{\mathbf{b}}^{KQ}$, der \mathbf{b} als feste Effekte schätzt, gilt

	$\hat{\mathbf{b}}^{BLUP}$	$\hat{\mathbf{b}}^{KQ}$	
\hat{b}_i			
$E(\hat{b}_i b_i)$			
$\text{Var}(\hat{b}_i b_i)$			
$\text{MSE}(\hat{b}_i)$			

$$\text{MSE}(\hat{b}_i) = E_b E_{y|b}[(\hat{b}_i - b_i)^2] = E_b[\text{Var}(\hat{b}_i | b_i) + (E(\hat{b}_i | b_i) - b_i)^2].$$

Schätzung der Kovarianzstruktur

Sei ϑ der Parametervektor, der alle unbekannt Parameter in den Kovarianzmatrizen $\mathbf{R} = \mathbf{R}(\vartheta)$, $\mathbf{G} = \mathbf{G}(\vartheta)$ und $\mathbf{V} = \mathbf{V}(\vartheta)$ enthält.

Indem man einen (konsistenten) Schätzer $\hat{\vartheta}$ einsetzt, erhält man die geschätzten Kovarianzmatrizen

$$\hat{\mathbf{R}} = \mathbf{R}(\hat{\vartheta}), \quad \hat{\mathbf{G}} = \mathbf{G}(\hat{\vartheta}), \quad \hat{\mathbf{V}} = \mathbf{V}(\hat{\vartheta}),$$

die für die Berechnung der BLUPs von β und \mathbf{b} verwendet werden können.

Maximum-Likelihood (ML)-Schätzung von ϑ

Die **ML-Schätzung** von ϑ basiert auf der **Likelihood des marginalen Modells**

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}(\vartheta)).$$

Die log-Likelihood für $\boldsymbol{\beta}$ und ϑ ist bis auf additive Konstanten gegeben durch

$$l(\boldsymbol{\beta}, \vartheta) = -\frac{1}{2} \left\{ \log |\mathbf{V}(\vartheta)| + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}(\vartheta)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

Maximiert man $l(\boldsymbol{\beta}, \vartheta)$ für festes ϑ bezüglich $\boldsymbol{\beta}$, so erhält man (wieder)

$$\widehat{\boldsymbol{\beta}}(\vartheta) = \arg \max_{\boldsymbol{\beta}} l(\boldsymbol{\beta}, \vartheta) = (\mathbf{X}' \mathbf{V}(\vartheta)^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}(\vartheta)^{-1} \mathbf{y}.$$

Einsetzen von $\widehat{\boldsymbol{\beta}}(\vartheta)$ in $l(\boldsymbol{\beta}, \vartheta)$ liefert die nur von ϑ abhängige **Profil-log-Likelihood**

$$l_P(\vartheta) = l(\widehat{\boldsymbol{\beta}}(\vartheta), \vartheta) = -\frac{1}{2} \left\{ \log |\mathbf{V}(\vartheta)| + (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\vartheta))' \mathbf{V}(\vartheta)^{-1} (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\vartheta)) \right\}.$$

Maximierung von $l_P(\vartheta)$ bezüglich ϑ liefert den **ML-Schätzer** $\widehat{\vartheta}_{ML}$.

Motivation REML-Schätzung

ML-Schätzer für Varianzen sind (nach unten) verzerrt, da der Verlust von Freiheitsgraden durch die Schätzung von $\hat{\beta}$ unberücksichtigt bleibt.

Beispiel lineares Modell $\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$

- ML-Schätzer ist $\hat{\sigma}^2 = \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon}$, aber $E(\hat{\sigma}^2) = \frac{n-p}{n} \sigma^2$.
- Benutzt wird üblicherweise der erwartungstreue Schätzer $\hat{\sigma}^2 = \frac{1}{n-p} \hat{\varepsilon}' \hat{\varepsilon}$.

Dies ist der **restringierte Maximum Likelihood**-Schätzer. Er berücksichtigt den Verlust an Freiheitsgraden durch Schätzung von β .

REstringierte Maximum-Likelihood (REML) Schätzung von ϑ

Grundidee: Basiere Likelihood nicht auf \mathbf{y} , sondern auf $(n - p)$ linear unabhängigen Fehlerkontrasten \mathbf{Ay} , deren Verteilung von β unabhängig ist.

Teile die log-Likelihood in zwei Teile auf für \mathbf{Ay} (für ϑ) und \mathbf{By} (für β) mit

- 1 \mathbf{A} von Rang $n - p$ und \mathbf{B} von Rang p .
- 2 Die zwei Teile sind statistisch unabhängig, also hier $\text{Cov}(\mathbf{Ay}, \mathbf{By}) = \mathbf{0}$.
- 3 \mathbf{Ay} sind Fehlerkontraste, d.h. $E(\mathbf{Ay}) = \mathbf{0}$.
- 4 \mathbf{BX} hat Rang p ($\Rightarrow \mathbf{By}$ schätzt eine eindeutige Funktion von β).

Geeignete Matrizen: $\mathbf{A} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}$, $\mathbf{B} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}$.

\rightarrow Englischer Name **residual** oder **restricted maximum likelihood estimation**.

Aus den ersten zwei Punkten folgt, dass die log-Likelihood $l = l' + l^*$ ist (bis auf additive Konstanten), l' für $\mathbf{A}\mathbf{y}$ und l^* für $\mathbf{B}\mathbf{y}$.

Maximierung von l^* (aus $\mathbf{B}\mathbf{y}$) bzgl. β ergibt den MLE/BBLUE $\hat{\beta}$ (abhängig von ϑ).

Maximierung von l' (aus $\mathbf{A}\mathbf{y}$) bzgl. ϑ : Die restringierte log-Likelihood ist (bis auf additive Konstanten) unabhängig von dem gewählten Fehlerkontrast

$$\begin{aligned} l_R(\vartheta) &= -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}| - \frac{1}{2} (\mathbf{y} - \mathbf{X} \hat{\beta})' \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \hat{\beta}) \\ &= l(\hat{\beta}(\vartheta), \vartheta) - \frac{1}{2} \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}| \end{aligned} \quad (19)$$

Herleitung für Interessierte: → **Extrablatt**

Maximierung ergibt den restringierten ML- (REML)-Schätzer $\hat{\vartheta}_{REML}$ für ϑ .

Alternative Herleitung des REML-Schätzers

Alternativ maximiert der REML-Schätzer $\hat{\vartheta}_{REML}$ die **marginale log-Likelihood** für ϑ , bei der β aus der Likelihood $L(\beta, \vartheta)$ des marginalen Modells herausintegriert wird

$$l_R(\vartheta) = \log \left(\int L(\beta, \vartheta) d\beta \right). \quad (20)$$

Diese ergibt sich bis auf additive Konstanten wieder als (19). → **Übung**

Auch beim LMM ist die **Reduktion der Verzerrung** von $\hat{\vartheta}_{ML}$ der Hauptgrund für die Verwendung von $\hat{\vartheta}_{REML}$ als Schätzer für ϑ . Es ist jedoch allgemein nicht gesichert, ob auch der mean squared error (MSE) geringer wird.

Empirische Versionen von BLUP und BLUE

Mit den durch **(RE)ML** geschätzten Kovarianzmatrizen $\hat{\mathbf{R}}$, $\hat{\mathbf{G}}$ und $\hat{\mathbf{V}}$ ergeben sich empirische Versionen von BLUP und BLUE, **EBLUP** und **EBLUE**, für \mathbf{b} und β ,

$$\hat{\beta} = (\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{y}, \quad \hat{\mathbf{b}} = \hat{\mathbf{G}}\mathbf{Z}'\hat{\mathbf{V}}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta}).$$

- Im Gegensatz zum linearen Modell sind die Schätzer der festen Effekte β im LMM abhängig von der Kovarianzmatrix \mathbf{V} . Daher sind die Schätzer $\hat{\beta}(\hat{\vartheta}_{REML})$ und $\hat{\beta}(\hat{\vartheta}_{ML})$ **nicht identisch**.
- Nach Einsetzen von $\hat{\vartheta}$ gelten die Optimalitätseigenschaften nicht mehr exakt.
- Ebenso sind die Varianzen bzw. Kovarianzmatrizen der Schätzer analytisch nicht exakt zugänglich (mehr später).

Numerische Berechnung der Schätzer

- Die Maximierung von $l_R(\vartheta)$ [$l_P(\vartheta)$] zur Berechnung von $\hat{\vartheta}_{REML}$ [$\hat{\vartheta}_{ML}$] erfolgt wegen Nichtlinearität in ϑ meist per [Newton-Raphson](#) oder [Fisher-Scoring](#).
- Dabei können diese Algorithmen keine Restriktion von Parametern auf $[0, \infty)$ berücksichtigen, sodass negative Schätzer für Varianzen entstehen könnten.
- Daher maximiert `lme()` in R-Paket `nlme` die (restringierte) log-Likelihood bzgl. der skalierten log-Varianzen. Allerdings kann diese Methode ein Maximum in der Null (kein zufälliger Effekt) nicht finden.
- `lmer()` in R-Paket `lme4` verwendet eine Parametrisierung, die über eine Nebenbedingung positiv (semi-)definites \mathbf{G} sicherstellt und singuläres \mathbf{G} erlaubt. Die Funktion verwendet Optimierung mit Nebenbedingungen statt Newton-Raphson.

Numerische Berechnung der Schätzer

- Die Bedingungen an ϑ unterscheiden sich zwischen marginalem (\mathbf{V} positiv (semi-)definit) und konditionalem Modell (\mathbf{G} und \mathbf{R} positiv (semi-)definit). Softwarepakete maximieren bzgl. ϑ zum Teil über einen größeren Parameterraum als das konditionale Modell implizieren würde.
- Die historisch erste Alternative ist der EM-Algorithmus, der die zufälligen Effekte als fehlende Daten behandelt. Die Konvergenz kann langsam sein.

Likelihood-Schätzung in R

Das klassische Paket zur LMM-Schätzung in R ist `nlme`. Eine ausführliche Beschreibung findet man im Buch der Autoren, Pinheiro & Bates (2000).

Der iterative Optimierungsalgorithmus ist eine Hybridversion aus zunächst **EM**- und im Anschluss **Newton-Raphson(NR)-Algorithmus**.

- **EM-Algorithmus**: Iterationen schnell und einfach zu berechnen. Schnell nahe ans Optimum, Konvergenz nahe des Optimums potentiell sehr langsam.
- **NR-Algorithmus**: Iterationen sehr rechenintensiv. (Berechnung 1. und 2. Ableitung, Scorefunktion, Hessematrix, in jedem Schritt.) Möglicherweise instabil entfernt vom Optimum. Sehr schnelle Konvergenz nahe des Optimums.

Alternativen

- Für additive und generalisierte additive gemischte Modelle (mehr dazu in Kapitel 5) eignet sich das R-Paket `mgcv`, siehe auch Wood (2006).
- `lme4` ist eine neuere Weiterentwicklung von Douglas Bates und Koautoren, das lineare und generalisierte lineare gemischte Modelle fiten kann (unter Benutzung von S4-Klassen). Siehe auch Bates, Maechler, Bolker & Walker (2015). *Fitting Linear Mixed-Effects Models using lme4*. Journal of Statistical Software, angenommen. <http://arxiv.org/abs/1406.5823>.
- Viele Beispiele für lineare gemischte Modelle mit SAS `proc mixed` finden sich in Verbeke & Molenberghs, 2000.

Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell**
 - Kovarianzmatrizen für die Schätzer
 - Tests für die festen Effekte
 - Tests für die zufälligen Effekte oder Varianzkomponenten
 - Modellselektion
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Schätzung der Kovarianzmatrix der Schätzer

Bei bekannter Kovarianzstruktur ist wegen (18) → **Übung**

$$\text{Cov} \begin{pmatrix} \hat{\beta} \\ \hat{\mathbf{b}} - \mathbf{b} \end{pmatrix} = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1}. \quad (21)$$

Beachte, dass $\text{Cov}(\hat{\mathbf{b}} - \mathbf{b})$ statt $\text{Cov}(\hat{\mathbf{b}})$ verwendet wird, was den Bias durch Shrinkage und die Variabilität in \mathbf{b} berücksichtigt.

Geschätzte Versionen mit eingesetztem $\hat{\mathbf{R}}$ und $\hat{\mathbf{A}}$ können zur Konstruktion approximativer Konfidenz- oder Vorhersageintervalle für Ausdrücke in β und \mathbf{b} verwendet werden (z.B. für neue Beobachtungen \mathbf{y}_i); die Diagonalelemente ergeben geschätzte Varianzen für die Komponenten von $\hat{\beta}$ bzw. $\hat{\mathbf{b}}$.

Beachte, dass die Variabilität durch die Schätzung der Kovarianzparameter nicht berücksichtigt wird und die tatsächlichen Varianzen daher unterschätzt werden.

Schätzung der Kovarianzmatrix der Schätzer

Manchmal wird auch eine Alternative verwendet, bei der \mathbf{b} wie ein fester Effekt behandelt wird, indem man auf \mathbf{b} bedingt. Dann erhält man \rightarrow Übung

$$\text{Cov} \left(\left(\begin{array}{c} \hat{\beta} \\ \hat{\mathbf{b}} \end{array} \right) \middle| \mathbf{b} \right) = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1} \mathbf{C}'\mathbf{R}^{-1}\mathbf{C} (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1}. \quad (22)$$

Diese Kovarianzmatrix ist von der **Sandwich-Matrix**-Form, da sie wie ein Sandwich aus $(\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1}$ und $\mathbf{C}'\mathbf{R}^{-1}\mathbf{C}$ zusammengesetzt ist.

Nach Einsetzen der Schätzer $\hat{\mathbf{R}}$, $\hat{\mathbf{A}}$ erhält man die geschätzte Kovarianzmatrix.

Wir präferieren (21). (22) führt zu kleineren Standardfehlern,

$$\begin{aligned} \text{Cov} \left(\begin{array}{c} \hat{\beta} \\ \hat{\mathbf{b}} - \mathbf{b} \end{array} \right) &= \text{Cov} \left(\left(\begin{array}{c} \hat{\beta} \\ \hat{\mathbf{b}} \end{array} \right) \middle| \mathbf{b} \right) \\ &= (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{A})^{-1} \geq \mathbf{0}. \end{aligned}$$

Tests für die festen Effekte

Mit den geschätzten Kovarianzmatrizen lassen sich **Wald-Tests** für

$$H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0} \text{ gegen } H_A : \mathbf{L}\boldsymbol{\beta} \neq \mathbf{0} \quad (23)$$

oder Konfidenzintervalle konstruieren. Dazu benutzen wir, dass

$$T_W = (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{L}' [\mathbf{L} \text{Cov}(\hat{\boldsymbol{\beta}}) \mathbf{L}']^{-1} \mathbf{L} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$$

asymptotisch $\chi^2_{\text{Rang}(\mathbf{L})}$ -verteilt ist. Es können auch robuste Kovarianzmatrix-Schätzer verwendet werden, siehe z.B. Liang & Zeger (1986). *Longitudinal Data Analysis Using Generalized Linear Models*. Biometrika, 73 (1), 13-22.

Wichtig: Alle asymptotischen Ergebnisse in diesem Kapitel gelten in LMMs für **Longitudinal- oder Cluster-Daten**, in denen \mathbf{y} aus N unabhängigen Teilvektoren \mathbf{y}_i besteht, für $N \rightarrow \infty$. Resultate für allgemeine LMMs sind schwierig zu zeigen (vgl. die Diskussion in Ruppert, Wand & Carroll (2003), Kapitel 4.8).

Die Wald-Teststatistik basiert auf Standardfehlern (21), die die wahre Variabilität in $\hat{\beta}$ durch Vernachlässigung der Schätzung von ϑ unterschätzen.

Dieses Problem wird häufig durch approximative *t-* oder *F-Statistiken* verringert:

- Für einen einzelnen Parameter β_j in β wird die Verteilung von $(\hat{\beta}_j - \beta_j)/\widehat{SE}(\hat{\beta}_j)$ durch eine *t*-Verteilung approximiert.
- Für generelle Hypothesen (23) wird die Verteilung von

$$F = \frac{(\hat{\beta} - \beta)' \mathbf{L}' [\mathbf{L} \text{Cov}(\hat{\beta}) \mathbf{L}']^{-1} \mathbf{L} (\hat{\beta} - \beta)}{\text{Rang}(\mathbf{L})}$$

durch eine *F*-Verteilung mit $\text{Rang}(\mathbf{L})$ Zähler-Freiheitsgraden approximiert.

Die Freiheitsgrade für die t -Verteilung bzw. die Nenner-Freiheitsgrade für die F -Verteilung werden aus den Daten geschätzt. Häufig wird die Satterthwaite-Approximation verwendet, die die Momente der Verteilungen matcht.

Für kleine Stichproben gibt die Methode nach Kenward und Rogers bessere Ergebnisse. Diese berücksichtigt die zusätzliche Unsicherheit durch Schätzung von ϑ .

Die Berechnung von Satterthwaite-artigen Freiheitsgraden ist recht rechenintensiv.

In R: R-Paket `pbkrtest`, siehe Halekoh & Højsgaard (2014), A Kenward-Roger Approximation and Parametric Bootstrap Methods for Tests in Linear Mixed Models – The R Package `pbkrtest`, Journal of Statistical Software, 59 (9).

Alternativ können wir den **Likelihood-Quotienten-Test** (LQT) verwenden. Bei nicht so großem n gibt der LQT meist genauere Ergebnisse als der Wald-Test.

Die Teststatistik des LQT

$$T_{LQT} = 2 \sup_{H_A} l(\boldsymbol{\beta}, \boldsymbol{\vartheta}) - 2 \sup_{H_0} l(\boldsymbol{\beta}, \boldsymbol{\vartheta})$$

ist asymptotisch unter der Nullhypothese $\chi^2_{\text{Rang}(L)}$ -verteilt. Voraussetzungen:

- genestete Modelle
- gleiche Kovarianzstruktur in jedem Modell
- **keine REML-Schätzung!** (Bei verschiedenen festen Effekten unter H_0 und H_A unterscheiden sich die Fehlerkontraste $\mathbf{A}_0 \mathbf{y}$ und $\mathbf{A}_A \mathbf{y} \Rightarrow$ die Likelihoods für $\mathbf{A}_0 \mathbf{y}$ und $\mathbf{A}_A \mathbf{y}$ sind nicht vergleichbar.)

Beachte außerdem, dass die Approximation bei kleinem $n - p$ schlecht sein kann.

Tests für die (Ko-)Varianzkomponenten in ϑ

Bei Longitudinal- oder Clusterdaten mit

$$\mathbf{D} = \begin{pmatrix} d_{11} & \dots & d_{1q} \\ \vdots & \ddots & \vdots \\ d_{1q} & \dots & d_{qq} \end{pmatrix} = \left(\begin{array}{c|c} \mathbf{D}_1 & \begin{matrix} d_{1q} \\ \vdots \\ d_{qq} \end{matrix} \\ \hline d_{1q} & \dots & d_{qq} \end{array} \right)$$

könnte z.B. ein Test für

$$H_0: \mathbf{D} = \begin{pmatrix} \mathbf{D}_1 & 0 \\ \hline 0 & \dots & 0 \end{pmatrix} \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \text{ mit } \mathbf{D}_1 \text{ positiv-definit } (q-1 \times q-1) \text{ gegen}$$

$$H_A: \mathbf{D} \text{ positiv-semidefinit } (q \times q) \quad (24)$$

von Interesse sein. Dies entspricht einem Test für den q ten zufälligen Effekt, z.B.

- H_0 : Random Intercept gegen H_A : Random Intercept und Slope ($q = 2$)
- H_0 : keine zufälligen Effekte gegen H_A : Random Intercept ($q = 1$).

Asymptotische Verteilung des LQT

(24) verletzt die Standardannahme der getesteten Parameter im **Innern des Parameterraums**. Z.B. $d_{qq} \geq 0$ als Varianz und $d_{qq} = 0$ (unter H_0) daher am **Rand des Parameterraums**. (Generell: komplizierte Restriktionen an $\mathbf{D} \geq 0$ und $\mathbf{D}_1 \geq 0$.)

Für Longitudinal- oder Cluster-Daten mit $N \rightarrow \infty$ hat T_{LQT} für (24) asymptotisch unter H_0 eine 0.5:0.5-Mischungsverteilung aus zwei χ_{q-1}^2 - und χ_q^2 -Verteilungen (χ_0^2 -Verteilung = Punktmasse an der Null).

Für komplexere Hypothesen für \mathbf{D} können komplexere χ^2 -Mischungen entstehen.

Diese Mischungsverteilung gilt nicht, wenn ein Nuisance-Parameter (Störparameter) am Rand des Parameterraums liegt, also z.B. die Varianz der Random Intercepts (nahe) Null ist, während die Varianz der Random Slopes getestet wird.

Exakte Verteilung des LQT

Kann \mathbf{y} nicht in viele unabhängige Teilvektoren $\mathbf{y}_i, i = 1, \dots, N$ unterteilt werden, so approximiert die Mischungsverteilung die Verteilung von T_{LQT} meist schlecht.

Die exakte Verteilung von T_{LRT} für $H_0: \tau_1^2 = 0$ im allgemeinen LMM mit $\mathbf{G} = \tau_1^2 \mathbf{I}_{q_1}$ und $\mathbf{R} = \sigma^2 \mathbf{I}_n$ ist im R-Paket `RLRsim` (F. Scheipl) implementiert. Approximationen für $\mathbf{G} = \text{diag}(\tau_1^2 \mathbf{I}_{q_1}, \dots, \tau_s^2 \mathbf{I}_{q_s})$, $s > 1$, beruhen auf

- Greven, S., Crainiceanu, C., Küchenhoff, H. & Peters, A. (2008). *Restricted Likelihood Ratio Testing for Zero Variance Components in Linear Mixed Models*. JCGS, 17 (4): 870-891.
- Scheipl, F., Greven, S. & Küchenhoff, H. (2008). *Size and Power of Tests for a Zero Random Effect Variance or Polynomial Regression in Additive and Linear Mixed Models*. CSDA, 52 (7): 3283-3299.

Wenn nur Parameter in ϑ getestet werden, kann auch ein restringierter LQT basierend auf der REML-Schätzung verwendet werden. Die Asymptotik ist entsprechend, die exakte Verteilung ist ebenfalls in `RLRsim` implementiert.

Modellselektion

AIC- und BIC-Werte in Modelloutputs, z.B. von `lme()` in R, basieren meist auf der marginalen Likelihood (9), z.B. das Akaike Informationskriterium (AIC)

$$AIC = -2 \log l(\hat{\beta}, \hat{\vartheta}) + 2(p + d),$$

mit d die Anzahl der Parameter in ϑ . Ein AIC basierend auf der marginalen restringierten log-Likelihood kann analog definiert werden.

Die Herleitung des AIC nimmt unabhängig und identisch verteilte Beobachtungen sowie einen Parameterraum \mathbb{R}^{p+d} an. Bei der Selektion von Parametern in ϑ (mit Restriktionen, z.B. Varianzen aus $[0, \infty)$) führt dies zu einem Bias und tendenziell zu kleineren Modellen ohne zufällige Effekte.

Modellselektion

Neuere Ansätze verwenden für die Selektion von zufälligen Effekten ein AIC basierend auf der konditionalen log-Likelihood für $\mathbf{y}|\mathbf{b} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \mathbf{R})$. Die korrekten Freiheitsgrade werden hergeleitet in

- Greven, S. und Kneib, T. (2010). *On the Behaviour of Marginal and Conditional AIC in Linear Mixed Models*. *Biometrika*, 97(4): 773-789

und das konditionale AIC ist implementiert im R-Paket `cAIC4`.

Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 **Bayes-Schätzung für lineare gemischte Modelle**
 - Wiederholung: Bayes-Inferenz
 - Bayesianisches LMM
 - Empirische Bayes-Schätzung
 - Volle Bayes-Schätzung
 - Erweiterungen: Flexiblere Verteilung der zufälligen Effekte
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Wiederholung: Bayes-Inferenz

- Parameter $\theta \in \Theta$ nicht deterministisch, sondern als zufällig angenommen.
- Volles Wahrscheinlichkeitsmodell für alle beobachteten und unbeobachteten Größen bestehend aus:
 - 1 **Beobachtungsmodell:** bedingte Verteilung der Daten gegeben unbekannte Parameter θ , $p(y|\theta)$.
 - 2 **Priori-Verteilung** $p(\theta)$: drückt Vorwissen/Annahmen über θ aus.
- Mathematisch günstig: zum Beobachtungsmodell **konjugierte Prioris**. (Prioris und Posterioris in der gleichen Verteilungsfamilie.)
- Informationsgehalt der Priori: **nicht** oder **schwach informative Prioris**.

Wiederholung: Bayes-Inferenz

- **Statistische Schlüsse** basieren auf der **Posteriori-Verteilung** $p(\theta|\mathbf{y})$, $\theta \in \Theta$, der bedingten Verteilung der unbeobachteten Größen gegeben die beobachteten Daten.
- **Berechnung** mit dem Satz von Bayes:

$$p(\theta|\mathbf{y}) = \frac{p(\mathbf{y}|\theta) p(\theta)}{\int_{\Theta} p(\mathbf{y}|\theta) p(\theta) d\theta} \propto p(\mathbf{y}|\theta) p(\theta)$$

mit der **Normierungskonstanten** der Posteriori-Dichte im Nenner.

Wiederholung: Bayes-Inferenz

- Übliche **Punktschätzer** der Bayesianischen Inferenz:
 - Posteriori-Erwartungswert $\hat{\theta} = E(\theta|\mathbf{y})$,
 - Posteriori-Median $\hat{\theta} = \inf \{\theta : F(\theta|\mathbf{y}) \geq 0.5\}$,
 - Posteriori-Modus $\hat{\theta} = \arg \max \{p(\theta|\mathbf{y})\}$.
- **Problem:** Posteriori-Verteilung meist analytisch unzugänglich
- **Lösung:** Verwendung von **MCMC-Verfahren**, mit denen (abhängige) Zufallszahlen aus der Posteriori Verteilung gezogen werden können.

Grundidee MCMC

- Konstruiere Markov-Kette (MK), deren stationäre Verteilung mit der Posteriori Verteilung übereinstimmt.
- Zustände der MK entsprechen gezogenen Zufallszahlen, die (nach entsprechender Konvergenzzeit (Burn-In-Phase) der MK) abhängige Stichprobe aus Posteriori darstellen.
- Abhängigkeit kann durch geeignetes Ausdünnen der Stichprobe reduziert werden.
- Interessierende Größen (z.B. P.-Erwartungswert), werden dann aus dieser (ausgedünnten) Stichprobe durch die **empirischen Analog**a geschätzt.
- Bekanntester Algorithmus: **Metropolis-Hastings-Algorithmus**, Spezialfall: **Gibbs-Sampler**

Beobachtungsmodell

Beobachtungsmodell:

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{b}, \boldsymbol{\vartheta} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \mathbf{R}(\boldsymbol{\vartheta}))$$

entspricht $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$ mit $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}(\boldsymbol{\vartheta}))$.

Priori für β

- Jetzt auch „feste“ Effekte β als Zufallsgrößen
- kein Vorwissen über $\beta \Rightarrow$ **nichtinformative Priori**, d.h.

$$p(\beta) \propto \text{const},$$

sonst

$\beta \sim N(\mathbf{m}, \mathbf{M})$ mit bekanntem EWert \mathbf{m} , Kovarianz \mathbf{M} .

- nichtinformative Priori $p(\beta) \propto \text{const}$ ergibt sich als Grenzfall der Priori-Normalverteilung (NV) für **Präzisionsmatrix** $\mathbf{M}^{-1} \rightarrow \mathbf{0}$.
- zum Beobachtungsmodell **konjugierte Normalverteilung** für β
 \Rightarrow Posteriori-Inferenz vergleichsweise einfach.

Prioris für \mathbf{b} , ε

Üblicherweise:

$$\mathbf{b} \sim N(\mathbf{0}, \mathbf{G}(\vartheta)); \quad \varepsilon \sim N(\mathbf{0}, \mathbf{R}(\vartheta))$$

- Kovarianzmatrizen $\mathbf{G} = \text{Cov}(\mathbf{b})$ und $\mathbf{R} = \text{Cov}(\varepsilon)$ hängen i.A. von unbekanntem **Hyperparametern** im Vektor ϑ ab.
- **Voller Bayes-Ansatz:**
 ϑ ebenfalls **Zufallsvariable**, mit (Hyper-)Priori $p(\vartheta)$, die in Ermittlung der Posteriori mit einfließt.
- **empirischer Bayes-Ansatz:**
 ϑ als unbekannter, aber **fester Parameter**.

Weitere Annahme: Zufallsgrößen β , \mathbf{b} und ε a priori **unabhängig**.

Gemeinsame Posteriori bei NV-Priori für β und b

$$\begin{aligned}
 p(\beta, \mathbf{b} | \mathbf{y}) &\propto p(\mathbf{y} | \beta, \mathbf{b}) p(\beta) p(\mathbf{b}) \\
 &\propto \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}) \right. \\
 &\quad \left. - \frac{1}{2} (\beta - \mathbf{m})' \mathbf{M}^{-1} (\beta - \mathbf{m}) - \frac{1}{2} \mathbf{b}' \mathbf{G}^{-1} \mathbf{b} \right).
 \end{aligned}$$

$$\Rightarrow \begin{pmatrix} \beta \\ \mathbf{b} \end{pmatrix} \Big| \mathbf{y} \sim N(\boldsymbol{\mu}_{\beta, \mathbf{b}}, \boldsymbol{\Sigma}_{\beta, \mathbf{b}}) \text{ (Übung) mit}$$

$$\boldsymbol{\Sigma}_{\beta, \mathbf{b}} = (\mathbf{C}' \mathbf{R}^{-1} \mathbf{C} + \mathbf{A})^{-1}; \quad \mathbf{C} = [\mathbf{X} | \mathbf{Z}]; \quad \mathbf{A} = \begin{bmatrix} \mathbf{M}^{-1} & 0 \\ 0 & \mathbf{G}^{-1} \end{bmatrix}$$

$$\boldsymbol{\mu}_{\beta, \mathbf{b}} = \boldsymbol{\Sigma}_{\beta, \mathbf{b}} (\tilde{\mathbf{m}} + \mathbf{C}' \mathbf{R}^{-1} \mathbf{y}); \quad \tilde{\mathbf{m}} = \begin{pmatrix} \mathbf{M}^{-1} \mathbf{m} \\ \mathbf{0} \end{pmatrix}$$

Posteriori bei nichtinformativer Priori für β

Nichtinformativer Priori $p(\beta) \propto \text{const}$ entspricht Präzisionsmatrix $\mathbf{M}^{-1} = \mathbf{0}$:

$$p(\beta, \mathbf{b} | \mathbf{y}) \propto p(\mathbf{y} | \beta, \mathbf{b}) p(\mathbf{b}) \\ \propto \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}) - \frac{1}{2} \mathbf{b}' \mathbf{G}^{-1} \mathbf{b} \right)$$

$$\Rightarrow \text{Posteriori-EW für } \begin{pmatrix} \beta \\ \mathbf{b} \end{pmatrix}: (\mathbf{C}' \mathbf{R}^{-1} \mathbf{C} + \mathbf{A})^{-1} \mathbf{C}' \mathbf{R}^{-1} \mathbf{y}$$

$$\Rightarrow \text{Posteriori-Kov. für } \begin{pmatrix} \beta \\ \mathbf{b} \end{pmatrix}: (\mathbf{C}' \mathbf{R}^{-1} \mathbf{C} + \mathbf{A})^{-1}$$

- **Posteriori** äquivalent zu **penalisiertem KQ-Kriterium** $KQ_{pen}(\beta, \mathbf{b})$ in (16)
- **Posteriori-Modus** als Maximierer identisch mit **BLUP-Schätzern**, die $KQ_{pen}(\beta, \mathbf{b})$ minimieren. **Posteriori-EW** = Posteriori-Modus wegen NV.
- Posteriori-Kov. identisch mit Kovarianz in (21).

Empirische Bayes-Schätzung

- **Empirische Bayes-Schätzer** für β , \mathbf{b} durch Einsetzen der geschätzten Kovarianzmatrizen $\hat{\mathbf{G}} = \mathbf{G}(\hat{\vartheta})$ und $\hat{\mathbf{R}} = \mathbf{R}(\hat{\vartheta})$ in vorige Ausdrücke.
- Schätzung von $\hat{\vartheta}$ z.B. durch Maximieren der **marginalen Likelihood** für ϑ :

$$\hat{\vartheta} = \arg \max p(\mathbf{y}|\vartheta) = \arg \max \int p(\mathbf{y}|\beta, \mathbf{b}, \vartheta)p(\mathbf{b}|\vartheta)p(\beta)d\beta d\mathbf{b}.$$

- Für nicht-informative $p(\beta) \propto \text{const}$ ist die **marginale Likelihood**

$$p(\mathbf{y}|\vartheta) = \int p(\mathbf{y}|\beta, \mathbf{b}, \vartheta)p(\mathbf{b}|\vartheta)d\beta d\mathbf{b} = \int p(\mathbf{y}|\beta, \vartheta)d\beta$$

proportional zur **restringierten Likelihood** $\exp\{l_R(\vartheta)\}$, siehe (20).

\Rightarrow emp. Bayes-Schätzer in diesem Fall äquivalent zu **REML-Schätzer** $\hat{\vartheta}_{REML}$ und den dazugehörigen EBLUPs für β und \mathbf{b} .

Volle Bayes-Inferenz

- **Voller Bayes-Ansatz:** Priori-Verteilung $p(\vartheta)$ auch für unbekannte Parameter ϑ ; β , $\mathbf{b}|\vartheta$ und ϑ als unabhängig angenommen.
- Inferenz basiert auf **Posteriori-Verteilung**

$$p(\beta, \mathbf{b}, \vartheta | \mathbf{y}) \propto p(\mathbf{y} | \beta, \mathbf{b}, \vartheta) p(\beta) p(\mathbf{b} | \vartheta) p(\vartheta)$$

- $p(\beta, \mathbf{b}, \vartheta | \mathbf{y})$ echte Posteriori-Dichte, wenn zur **Normierung** gilt:

$$p(\mathbf{y}) = \int p(\mathbf{y} | \beta, \mathbf{b}, \vartheta) p(\beta) p(\mathbf{b} | \vartheta) p(\vartheta) d\beta d\mathbf{b} d\vartheta < \infty.$$

- Bei echter, **informativer** Priori mit $\int p(\vartheta) d\vartheta = 1$ existiert auch $p(\beta, \mathbf{b}, \vartheta | \mathbf{y})$.
- Für **nichtinformativ** Priori mit $\int p(\vartheta) d\vartheta = \infty$ Existenz der Posteriori nicht allgemein gesichert.

Beziehungen zur Likelihood-Inferenz

Bei nichtinformativer Priori $p(\vartheta) \propto \text{const}$ und Existenz der Posteriori:

- **REML-Schätzer** $\hat{\vartheta}_{REML}$ als Maximierer der marginalen Likelihood = Posteriori-Modus der marginalen Posteriori von ϑ wegen

$$p(\vartheta|\mathbf{y}) = p(\mathbf{y}|\vartheta) \frac{p(\vartheta)}{p(\mathbf{y})} \propto p(\mathbf{y}|\vartheta).$$

- Bei zusätzlich nichtinformativer Priori $p(\beta) \propto \text{const}$: **ML-Schätzer** $\hat{\vartheta}_{ML}$ als Maximierer der Likelihood = ϑ -Komponente des Posteriori-Modus der gemeinsamen Posteriori von β und ϑ wegen

$$p(\beta, \vartheta|\mathbf{y}) = p(\mathbf{y}|\vartheta, \beta) \frac{p(\vartheta)p(\beta)}{p(\mathbf{y})} \propto p(\mathbf{y}|\vartheta, \beta).$$

Inferenz

- Normierungskonstante der Posteriori i.A. nicht analytisch zugänglich
⇒ Posteriori-Dichte $p(\beta, \mathbf{b}, \vartheta | \mathbf{y})$ **nicht in geschlossener Form** darstellbar
- volle Bayes-Inferenz daher üblicherweise mittels ***MCMC-Simulation***
(Details siehe z.B. Fahrmeir et al. (2007, Abschnitt B.5.3))
- moderne (approximative) Alternativen:
INLA (Integrated Nested Laplace Approximation),
Variational Bayes-Verfahren

MCMC mit blockweisem Gibbs-Sampling

Vorgehen: Teile Parametervektor $\theta = (\beta, \mathbf{b}, \vartheta)$ in **Teilvektoren** **zusammengehöriger Parameter**, d.h. üblicherweise β , \mathbf{b} und ϑ auf.

- Wähle Startwerte $\beta^{(0)}$, $\mathbf{b}^{(0)}$, $\vartheta^{(0)}$ und Anzahl der Iterationen T
- Bilde **vollständig bedingte Dichten** (full conditionals) gegeben der restlichen Parameter und \mathbf{y}
- Ziehe **sequentiell Zufallszahlen** $\beta^{(t)}$, $\mathbf{b}^{(t)}$, $\vartheta^{(t)}$ aus diesen (geg. jeweils die momentan aktuellen Zustände) bis T erreicht.

Nach einer gewissen Konvergenzphase können die Zufallszahlen als Ziehungen aus den Marginalverteilungen von $\beta|\mathbf{y}$, $\mathbf{b}|\mathbf{y}$ und $\vartheta|\mathbf{y}$ angesehen werden.

Vollständig bedingte Dichten

$$\begin{aligned}
 p(\boldsymbol{\beta} | \mathbf{b}, \boldsymbol{\vartheta}, \mathbf{y}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{b}, \boldsymbol{\vartheta}) p(\boldsymbol{\beta}) \\
 &\propto \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})' \mathbf{R}(\boldsymbol{\vartheta})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) \right. \\
 &\qquad \qquad \qquad \left. -\frac{1}{2} (\boldsymbol{\beta} - \mathbf{m})' \mathbf{M}^{-1} (\boldsymbol{\beta} - \mathbf{m}) \right)
 \end{aligned}$$

$$\begin{aligned}
 p(\mathbf{b} | \boldsymbol{\beta}, \boldsymbol{\vartheta}, \mathbf{y}) &\propto \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})' \mathbf{R}(\boldsymbol{\vartheta})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) \right. \\
 &\qquad \qquad \qquad \left. -\frac{1}{2} \mathbf{b}' \mathbf{G}(\boldsymbol{\vartheta})^{-1} \mathbf{b} \right)
 \end{aligned}$$

$$\begin{aligned}
 p(\boldsymbol{\vartheta} | \boldsymbol{\beta}, \mathbf{b}, \mathbf{y}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{b}, \boldsymbol{\vartheta}) p(\mathbf{b} | \boldsymbol{\vartheta}) p(\boldsymbol{\vartheta}) \\
 &\propto |\mathbf{R}(\boldsymbol{\vartheta})|^{-1/2} |\mathbf{G}(\boldsymbol{\vartheta})|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{b}' \mathbf{G}(\boldsymbol{\vartheta})^{-1} \mathbf{b} \right. \\
 &\qquad \qquad \qquad \left. -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})' \mathbf{R}(\boldsymbol{\vartheta})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) \right) p(\boldsymbol{\vartheta})
 \end{aligned}$$

Vollständig bedingte Dichten von β , b

$$\beta|\cdot \sim N(\mu_\beta, \Sigma_\beta) \text{ mit (Übung)}$$

$$\Sigma_\beta = (\mathbf{X}'\mathbf{R}(\vartheta)^{-1}\mathbf{X} + \mathbf{M}^{-1})^{-1}$$

$$\mu_\beta = \Sigma_\beta (\mathbf{M}^{-1}\mathbf{m} + \mathbf{X}'\mathbf{R}(\vartheta)^{-1}(\mathbf{y} - \mathbf{Z}\mathbf{b}))$$

$$b|\cdot \sim N(\mu_b, \Sigma_b) \text{ mit (analog)}$$

$$\Sigma_b = (\mathbf{Z}'\mathbf{R}(\vartheta)^{-1}\mathbf{Z} + \mathbf{G}(\vartheta)^{-1})^{-1}$$

$$\mu_b = \Sigma_b (\mathbf{Z}'\mathbf{R}(\vartheta)^{-1}(\mathbf{y} - \mathbf{X}\beta))$$

- **nichtinformative Priori** mit $\mathbf{M}^{-1} = \mathbf{0}$: Erwartungswert $\mu_\beta = (\mathbf{X}'\mathbf{R}(\vartheta)^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{R}(\vartheta)^{-1}(\mathbf{y} - \mathbf{Z}\mathbf{b})$ ist gewichteter KQ-Schätzer angewandt auf die um $\mathbf{Z}\mathbf{b}$ bereinigten Daten
- **informative Priori** \Rightarrow Erwartungswert μ_β ist gewichtetes Mittel aus KQ-Schätzer und Priori-Erwartungswert.
- Analog für μ_b .

Vollständig bedingte Dichten von ϑ im Spezialfall

LMM für Clusterdaten mit $\text{Cov}(\varepsilon) = \sigma^2 I_n \Rightarrow \vartheta$ enthält σ^2 und Parameter in \mathbf{D} .

- nichtinformative Jeffreys Prioris $p(\sigma^2) \propto \sigma^{-2}$, $p(\mathbf{D}) \propto |\mathbf{D}|^{-\frac{q+1}{2}}$ führen i.A. zu uneigentlichen (d.h. nicht normierbaren) Posteriori-Verteilungen

\Rightarrow Schwach informative inverse Gammaverteilung $\sigma^2 \sim IG(a_\sigma, b_\sigma)$; a_σ, b_σ klein

- Für \mathbf{D} oft inverse Wishart-Verteilung. Bei $\mathbf{D} = \text{diag}(\tau_1^2, \dots, \tau_q^2)$ mit unabh. τ_j^2 ergibt diese ein Produkt von IGs mit $\tau_j^2 \sim IG(a_{\tau_j}, b_{\tau_j}), j = 1, \dots, q$.

Dann full conditionals bei zusätzlich $p(\beta) \propto \text{const}$ (nichtinformative Priori):

$$\sigma^2 | \cdot \sim IG(\tilde{a}_\sigma, \tilde{b}_\sigma) \text{ mit } \tilde{a}_\sigma = a_\sigma + \frac{1}{2}, \tilde{b}_\sigma = b_\sigma + \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}\|_2^2$$

$$\tau_j^2 | \cdot \sim IG(\tilde{a}_{\tau_j}, \tilde{b}_{\tau_j}) \text{ mit } \tilde{a}_{\tau_j} = a_{\tau_j} + \frac{N}{2}, \tilde{b}_{\tau_j} = b_{\tau_j} + \frac{1}{2} \sum_{i=1}^N b_{ij}^2$$

Probleme der NV-Annahme für \mathbf{b}

NV-Annahme für zufällige Effekte $\mathbf{b} \sim N(\mathbf{0}, \mathbf{G})$ mathematisch günstig. Aber wie sensitiv ist die Schätzung bei Fehlspezifikation?

- Die Schätzer der festen Effekte β und der Kovarianzparameter ϑ sind meist sehr robust gegenüber Fehlspezifikation der Verteilung der \mathbf{b} . Die Standardfehler können jedoch über/unterschätzt werden.
- Durch den Shrinkage-Effekt können die EBLUPs der zufälligen Effekte \mathbf{b} normalverteilt aussehen, selbst wenn die Verteilung der \mathbf{b} z.B. bimodal / schief ist / hohe Wahrscheinlichkeitsmasse an den Rändern hat (heavy tails) \Rightarrow Schlechte Vorhersagen $\hat{\mathbf{b}}$. q-q-Plots der $\hat{\mathbf{b}}$ eignen sich nicht zur Diagnose. Diagnose durch Fitten eines flexibleren Modells.

Alternative Prioris für b : Skalenmischungen

- Verwendung von Prioris mit mehr Masse auf den Rändern (heavy-tailed): z.B. t-Verteilung mit niedrigen Freiheitsgraden, Laplace-Verteilung. Diese sind oft darstellbar als Skalenmischung von Normalverteilungen

$$p(b_i) = \int \phi(b_i | \mu, \sigma^2) p(\sigma^2 | \theta) d\sigma^2$$

- Darstellbar als Skalenmischung \Rightarrow sehr leicht in Modellhierarchie für LMM einzubauen
- Beispiel: t-Verteilung mit $df = \nu$ ist Skalenmischung aus $N(0, \sigma^2)$ mit $\sigma^{-2} \sim \Gamma(\nu/2, \nu/2)$

Alternative Prioris für b : Finite Mischungen

- Zur Aufdeckung von Clustern (z.B. durch unbeobachtete Kovariablen oder latente Subpopulationen) können multimodale Verteilungen für die zufälligen Effekte verwendet werden, z.B. **finite Mischverteilungsmodelle**:

$$p(\mathbf{b}_i | \boldsymbol{\pi}, \boldsymbol{\phi}) = \sum_{k=1}^K \pi_k p(\mathbf{b}_i | \boldsymbol{\phi}_k),$$

mit Gewichten $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$, $\sum_k \pi_k = 1$, und parametrischer Verteilungsfamilie $p(\mathbf{b}_i | \boldsymbol{\phi}_k)$ mit Parametern $\boldsymbol{\phi} = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_K)$ (z.B. multivariate NV $p(\mathbf{b}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$ mit Erwartungswerten $\boldsymbol{\mu}_k$ und homogenen Kovarianzen $\boldsymbol{\Sigma}$)

- nichtparametrische Erweiterung: K nicht fest/konstant, wird mitgeschätzt.
- Noch flexibler: nichtparametrische Bayes-Ansätze, z.B. **Dirichlet-Prozess**, Dirichlet-Prozess-Mischungs-Prioris.
- Inferenz für diese Modelle basiert i.d.R. auf MCMC-Techniken.

Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 **Additive gemischte Modelle**
 - Penalisierte Splines
 - Penalisierte Regression und LMMs
 - Bayesianische Sicht auf Penalisierungsansätze
 - Weitere Beispiele
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Motivation

Die gemeinsame log-Likelihood (basierend auf der gemeinsamen Dichte von \mathbf{y} und \mathbf{b}) im LMM, aus der sich die BLUPs für β und \mathbf{b} ergeben, hat die Form (s. (15))

$$l(\beta, \mathbf{b}) = \log L(\beta, \mathbf{b}) = -\frac{1}{2}(\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b})' \mathbf{R}^{-1}(\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}) - \frac{1}{2}\mathbf{b}' \mathbf{G}^{-1} \mathbf{b}.$$

Diese kann als **penalisierte Likelihood** aufgefasst werden mit **Strafterm** $\mathbf{b}' \mathbf{G}^{-1} \mathbf{b}$.

Dies öffnet eine Verbindung zu **penalisierten Regressionsmethoden** (mit quadratischen Straftermen) und damit zu einer sehr allgemeinen Modellklasse.

Wir beginnen mit dem ausführlichen Beispiel der penalisierten Splines; weitere knappe Beispiele zeigen dann die Größe der flexiblen Modellklasse auf.

Das nichtparametrische Regressionsproblem

Modell: Die metrische Zielvariable y lässt sich durch eine deterministische Funktion $f(z)$ der metrischen Kovariablen z und einen additiven Fehler ε erklären:

$$y_i = f(z_i) + \varepsilon_i, i = 1, \dots, n.$$

- **Annahmen über f :**

Glattheitsanforderungen, z.B. via Stetigkeit oder Differenzierbarkeit von f .
Deutlich flexibler als Linearitätsannahme $f(z) = \beta_0 + \beta_1 z$ (lineares Modell).

- **Annahmen über ε :**

Wie im linearen Modell $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \Rightarrow E(y_i) = f(z_i)$.

D.h. die Funktion f modelliert den Erwartungswert der Zielvariablen.

Splines

Idee: Modelliere f stückweise durch Polynome, ohne dass Sprungstellen / Kanten entstehen \rightarrow **Splines** auf dem Definitionsbereich $[a, b]$ von z .

Definition von Splines:

Wähle **Knoten** $a = \kappa_1 < \dots < \kappa_m = b$. f ist ein Spline vom Grad $\ell \geq 0$, falls gilt

- $f(z)$ ist auf den Intervallen $[\kappa_j, \kappa_{j+1})$ ein Polynom vom Grad ℓ .
- $f(z)$ ist stetig (falls $\ell > 0$) und $(\ell - 1)$ -mal stetig differenzierbar. (**Glattheit**)

Konstruktive Darstellung durch Basisfunktionen:

Jeder Spline f vom Grad ℓ mit Knoten $\kappa_1 < \dots < \kappa_m$ kann eindeutig als Linearkombination $f(z) = \sum_{j=1}^d \gamma_j B_j(z)$ von $d = \ell + m - 1$ Basisfunktionen B_j dargestellt werden, die den d dimensionalen Spliner Raum $S_\ell(\kappa_1, \dots, \kappa_m)$ aufspannen.

Bekannte Basen: **Truncated-Power-Series** (TP) und **Basic-Splines** (B-Splines).

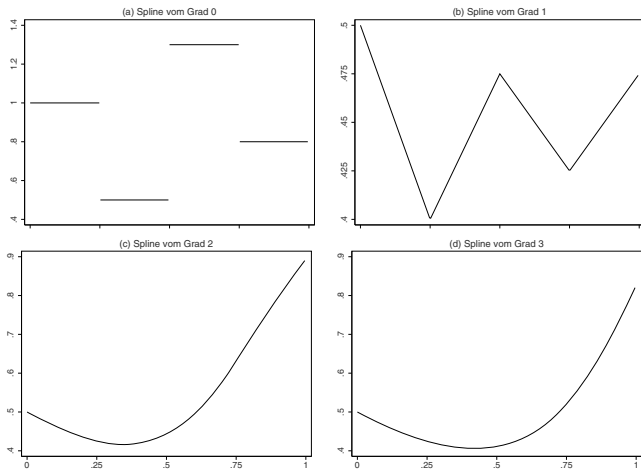


Abbildung 1: Beispiele für Splines vom Grad 0, 1, 2 und 3 zu den Knoten $\kappa_1 = 0$, $\kappa_2 = 0.25$, $\kappa_3 = 0.5$, $\kappa_4 = 0.75$ und $\kappa_5 = 1$.

Quelle: Fahrmeir et al (2009, Abb. 7.5).

Definition der d Basisfunktionen

Für eine Knotenmenge $\{\kappa_1, \dots, \kappa_m\}$ und einen Grad ℓ :

TP-Basisfunktionen

$$B_j(z) = z^{j-1}, \quad j = 1, \dots, \ell + 1 \quad (\text{Polynome})$$

$$B_{\ell+j}(z) = (z - \kappa_j)_+^\ell, \quad j = 2, \dots, m - 1 \quad (\text{abgeschnittene Potenzen}).$$

mit $(u)_+^\ell = [\max(u, 0)]^\ell$.

B-Spline Basisfunktionen (rekursiv)

$$B_j^0(z) = \mathbb{1}_{[\kappa_j, \kappa_{j+1})}(z), \quad j = 1, \dots, d,$$

$$B_j^\ell(z) = \frac{z - \kappa_{j-\ell}}{\kappa_j - \kappa_{j-\ell}} B_{j-1}^{\ell-1}(z) + \frac{\kappa_{j+1} - z}{\kappa_{j+1} - \kappa_{j+1-\ell}} B_j^{\ell-1}(z), \quad j = 1, \dots, d$$

mit zusätzlichen Knoten $\kappa_{1-\ell}, \dots, \kappa_0 < a$ und $b < \kappa_{m+1}, \dots, \kappa_{m+\ell}$.

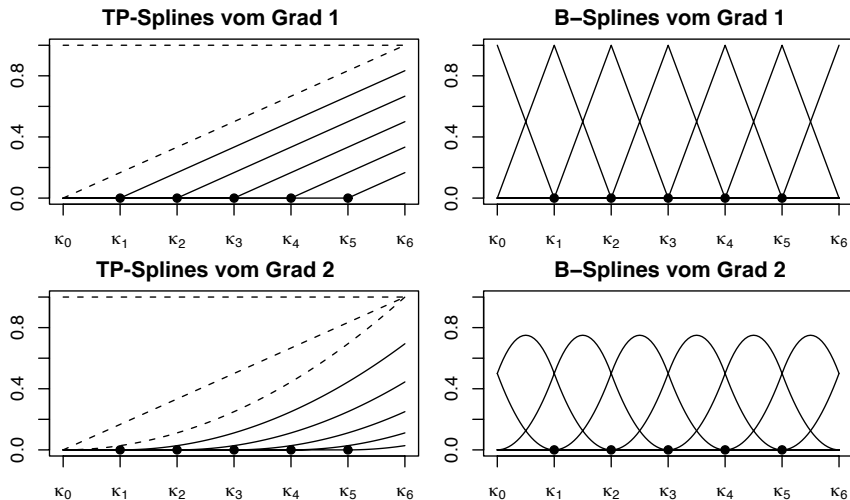


Abbildung 2: Lineare (oben) und quadratische (unten) TP- (links) und B-Spline-Basifunktionen (rechts) für $K = 5$ äquidistante Knoten mit $\kappa_0 = 0$. Für die TP-Basis sind die Polynomterme gestrichelt dargestellt.

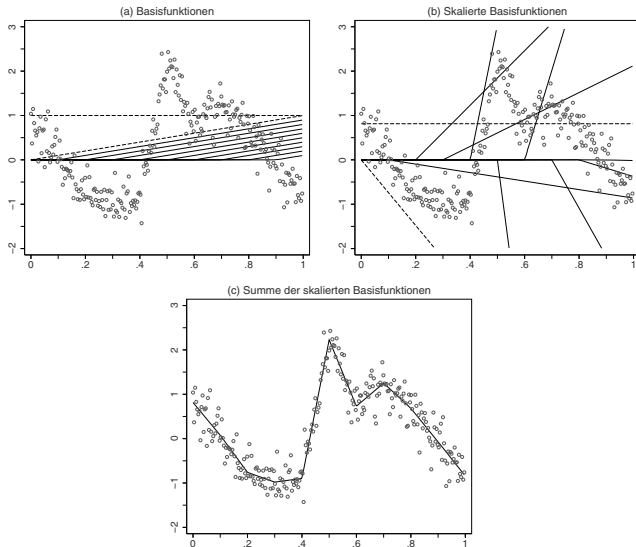


Abbildung 3: Konstruktion einer Funktion aus skalierten TP-Basisfunktionen für lineare Splines. Quelle: Fahrmeir et al (2009, Abb. 7.6).

Basisdarstellung

Damit ergibt sich die Darstellung des nichtparametrischen Regressionsproblems als

$$y_i = f(z_i) + \varepsilon_i \approx \sum_{j=1}^d \gamma_j B_j(z_i) + \varepsilon_i.$$

Man erhält ein **lineares Modell** in dem Parametervektor $\gamma = (\gamma_1, \dots, \gamma_d)'$

$$\mathbf{y} = \mathbf{B}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}$$

mit Designmatrix aus TP- bzw. B-Spline-Basis

$$\mathbf{B} = \begin{pmatrix} B_1(z_1) & \dots & B_d(z_1) \\ \vdots & & \vdots \\ B_1(z_n) & \dots & B_d(z_n) \end{pmatrix}.$$

Mit KQ-Schätzer ergibt sich $\hat{f}(z) = (B_1(z), \dots, B_d(z))(\mathbf{B}'\mathbf{B})^{-1}\mathbf{B}'\mathbf{y}$.

Einfluss der Knotenzahl

Es gilt der **Varianz-Bias-Trade Off**:

- große Knotenzahl: flexible Funktionsschätzung, aber große Variabilität (Overfitting)
- kleine Knotenzahl: glattere Schätzung, aber möglicher Bias

⇒ Wesentliche Frage nach der **Anzahl der zu verwendenden Knoten!**

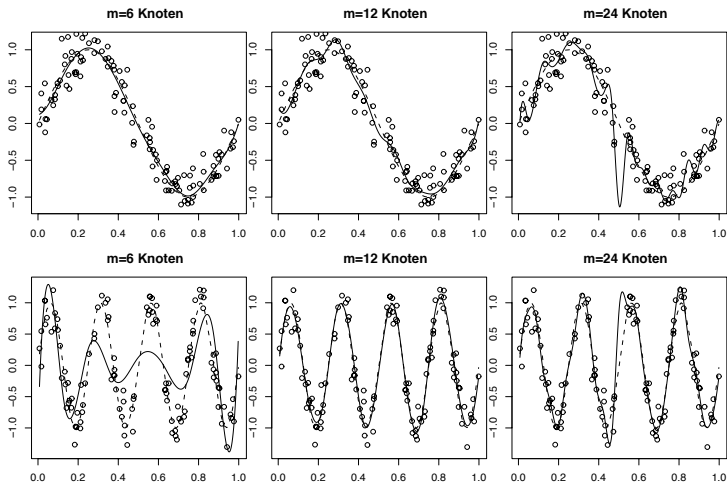


Abbildung 4: Abhängigkeit der Glättung mit Splines von der Knotenwahl. Die gestrichelte Linie gibt jeweils die wahre Funktion wieder, die durchgezogene Linie die Schätzung mit B-Splines mit 6, 12 bzw. 24 Knoten.

Grundidee penalisierte Splines

- Approximiere $f(z)$ durch einen Spline mit einer **großen Zahl von Knoten** (üblicherweise ca. 20 bis 40) \Rightarrow prinzipiell flexibel genug, auch stark variierende Funktionen $f(z)$ darzustellen.
- Führe zusätzlich einen **Strafterm** ein, der eine zu große **Variabilität der Schätzung bestraft**. (**Regularisierung**)
- Minimiere anstelle des üblichen KQ-Kriteriums ein um diesen Strafterm erweitertes **penalisiertes KQ-Kriterium** (äquivalent: Maximierung einer penalisierten log-Likelihood).

Penalisierte Splines auf TP-Basis

$$f(z) = \underbrace{\gamma_1 + \gamma_2 z + \dots + \gamma_{\ell+1} z^\ell}_{\text{1. Teil: globale polynomiale Form}} + \underbrace{\gamma_{\ell+2} (z - \kappa_2)_+^\ell + \dots + \gamma_d (z - \kappa_{m-1})_+^\ell}_{\text{2. Teil: Abweichung von globaler polynomialer Form}} .$$

Rauhe Funktionsschätzungen bei großer Variabilität des 2. Modellteils.

Da die ℓ te Ableitung von $(z - \kappa_j)_+^\ell$ gleich $\ell!$ für $z > \kappa_j$ und 0 sonst ist, $j \geq \ell + 2$, sind die γ_j proportional zu **Sprüngen in der ℓ ten Ableitung**.

⇒ Strafterm für betragsmäßig große Koeffizienten der abgeschnittenen Potenzen:

$$\text{pen}(\gamma, \mathbf{K}) = \lambda \sum_{j=\ell+2}^d \gamma_j^2 = \lambda \gamma' \mathbf{K} \gamma$$

mit $\gamma = (\gamma_1, \dots, \gamma_d)'$ und $\mathbf{K} = \text{diag}(\underbrace{0, \dots, 0}_{(\ell+1)}, \underbrace{1, \dots, 1}_{(m-2)})$.

Penalisierte Splines auf B-Spline-Basis: P-Splines

- Ein geeigneter Strafterm lässt sich über das Integral der (quadrierten) k -ten Ableitung definieren. Häufig wird dabei $k = 2$ verwendet:

$$\int (f''(z))^2 dz = \sum_{i=1}^d \sum_{j=1}^d \gamma_i \gamma_j \int B_i''(z) B_j''(z) dz = \gamma' \mathbf{K} \gamma$$

mit Strafmatrix $\mathbf{K} = (\int B_i''(z) B_j''(z) dz)_{i,j}$.

- Für die k -ten Ableitungen lässt sich der Strafterm über **Differenzen k -ter Ordnung** Δ^k der Parameter γ approximieren. Diese sind rekursiv definiert

$$\begin{aligned} \Delta^1 \gamma_j &= \gamma_j - \gamma_{j-1}, \\ \Delta^2 \gamma_j &= \Delta^1 \Delta^1 \gamma_j = \Delta^1 \gamma_j - \Delta^1 \gamma_{j-1} = \gamma_j - 2\gamma_{j-1} + \gamma_{j-2}, \\ &\vdots \\ \Delta^k \gamma_j &= \Delta^{k-1} \gamma_j - \Delta^{k-1} \gamma_{j-1}. \end{aligned}$$

P-Splines

Der Strafterm für B-Splines mit k ten Differenzen $\Delta^k \gamma_j$ hat die Gestalt

$$\text{pen}(\gamma, \mathbf{K}_k) = \lambda \gamma' \mathbf{K}_k \gamma = \lambda \gamma' \mathbf{D}'_k \mathbf{D}_k \gamma$$

mit der $(d - k) \times d$ **Differenzenmatrix** \mathbf{D}_k und der $d \times d$ Strafmatrix \mathbf{K}_k

$$\mathbf{D}_1 = \begin{pmatrix} -1 & 1 & & & & & \\ & -1 & 1 & & & & \\ & & & \ddots & \ddots & & \\ & & & & -1 & 1 & \\ & & & & & & \end{pmatrix}, \quad \mathbf{K}_1 = \begin{pmatrix} 1 & -1 & & & & & \\ -1 & 2 & -1 & & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & -1 & 2 & -1 & \\ & & & & -1 & 1 & \end{pmatrix},$$

$$\mathbf{D}_k = \mathbf{D}_1 \mathbf{D}_{k-1}, \quad \mathbf{K}_k = \mathbf{D}'_k \mathbf{D}_k$$

(bei entsprechend angepassten Dimensionen).

Penalisiertes KQ-Kriterium

Minimiere das **penalisierte KQ-Kriterium** über γ in Abhängigkeit vom $\lambda \geq 0$,

$$PKQ(\lambda, \gamma) = (\mathbf{y} - \mathbf{B}\gamma)'(\mathbf{y} - \mathbf{B}\gamma) + \lambda\gamma' \mathbf{K}\gamma.$$

- Der **Glättungsparameter** $\lambda \geq 0$ steuert dabei den **Kompromiss zwischen Datentreue** und **Glattheit** der Funktionsschätzung.
- $\lambda \rightarrow 0$: Geringes Gewicht für den Strafterm. $PKQ(\lambda, \gamma) \approx KQ(\gamma)$ und $\hat{\gamma}$ nahe beim KQ-Schätzer.
- $\lambda \rightarrow \infty$: Schätzung vollständig durch den Strafterm dominiert. Resultat:
 - Bei der TP-Basis ein Polynom vom Grad ℓ als Schätzung für $f(z)$.
 - Bei P-Splines des Grades $\ell \geq k - 1$, k te Differenzen, ein Polynom vom Grad $k - 1$.
- Durch Wahl von $\lambda \geq 0$ Kontinuum zwischen diesen beiden Extremen.

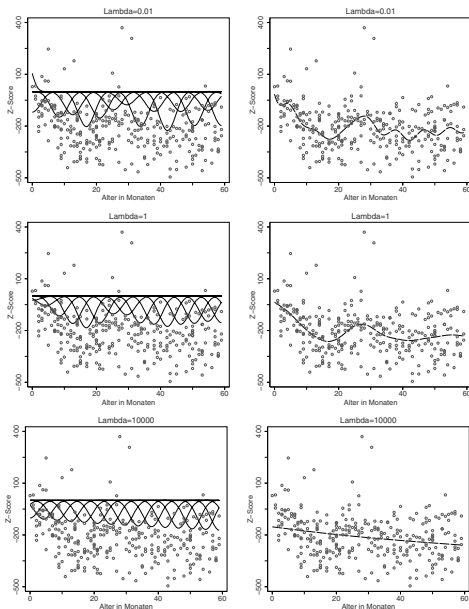


Abbildung 5: Einfluss des Glättungsparameters auf die Schätzung eines P-Splines basierend auf zweiten Differenzen. Quelle: Fahrmeir (2009, Abb. 7.15).

Die Rolle des Glättungsparameters

- Vorteil: Die **Glätteit** der Schätzung wird nicht mehr über Anzahl und Position der Knoten, sondern nur durch den Glättungsparameter λ gesteuert.
⇒ Frage nach der Wahl **geeigneter Werte für den Glättungsparameter λ** !

Optimalitätskriterien zur Bestimmung von **Glättungsparametern**, z.B.:

- Überlegungen zum **Bias-Varianz-Trade Off**. Einen solchen Kompromiss erhält man durch die Betrachtung des mittleren quadratischen Fehlers.
- Approximation des quadratischen Prognosefehlers durch **Kreuzvalidierung**.
- Verwendung von anderen aus der **Modellwahl** bekannten Kriterien, z.B. AIC.
- Hier: **Darstellung der Penalisierungsansätze als gemischtes Modell**.

TP-Penalierungsansatz als gemischtes Modell

Für penalisierte Splines mit TP-Basis ist das penalisierte KQ-Kriterium

$$PKQ(\lambda, \gamma) = (\mathbf{y} - \mathbf{B}\gamma)'(\mathbf{y} - \mathbf{B}\gamma) + \lambda \sum_{j=l+2}^d \gamma_j^2.$$

Zerlege die Designmatrix $\mathbf{B} = (\mathbf{X}|\mathbf{Z})$ und die Koeffizienten γ in

- $\beta = (\gamma_1, \dots, \gamma_{l+1})'$, die **nicht penalisierten** Polynomkoeffizienten,
 - $\mathbf{b} = (\gamma_{l+2}, \dots, \gamma_d)'$, die **penalisierten** Koeff. der abgeschnittenen Potenzen.
- $$\Rightarrow PKQ(\lambda, \gamma)/\sigma^2 = (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b})' \left(\frac{1}{\sigma^2} \mathbf{I}_n \right) (\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}) + \mathbf{b}' \left(\frac{\lambda}{\sigma^2} \mathbf{I}_{m-2} \right) \mathbf{b}.$$

Dies ist das penalisierte KQ-Kriterium (16) für ein LMM

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \varepsilon, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad \mathbf{b} \sim N(\mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}_{m-2}). \quad (25)$$

- β und \mathbf{b} können als Vektoren **fester / zufälliger Effekte** aufgefasst werden.
- Die **Varianz** τ^2 der zufälligen Effekte ist dabei $\tau^2 = \sigma^2/\lambda$.

Schätzung über gemischte Modelle

- Im gemischten Modell (25) können die Varianzen σ^2 und τ^2 über **ML-** oder **REML-Schätzung** bestimmt werden.
- Den **optimalen Glättungsparameter** erhält man dann als $\hat{\lambda} = \hat{\sigma}^2 / \hat{\tau}^2$.
- Die Hypothese

$H_0 : f$ Polynom vom Grad ℓ bzw. $k - 1$ gegen $H_A : f$ flexiblere Funktion

lässt sich durch einen Test überprüfen für

$$H_0 : \tau^2 = 0 \quad \text{gegen} \quad H_A : \tau^2 > 0.$$

(Beachte die Hinweise zum Testen von Varianzkomponenten, Folien 66 ff.! Hier lässt sich \mathbf{y} nicht in N unabhängige Teilvektoren zerlegen, sodass die asymptotischen Mischungsverteilungen für T_{LQT} nicht gültig sind. Die exakte Verteilung oder Bootstrap-Verfahren bieten Alternativen.)

Punktweise Konfidenzbänder

Der BLUP für $f(z) = \mathbf{C}_z(\boldsymbol{\beta}', \mathbf{b}')'$ mit $\mathbf{C}_z = (B_1(z), \dots, B_d(z))$ ergibt sich aus (18) mit $\mathbf{C} = \mathbf{B} = [\mathbf{X}|\mathbf{Z}]$ als

$$\hat{f}(z) = \mathbf{C}_z(\hat{\boldsymbol{\beta}}', \hat{\mathbf{b}}')' = \mathbf{C}_z(\mathbf{C}'\mathbf{C} + \lambda \text{diag}(\mathbf{0}, \mathbf{I}_{m-2}))^{-1} \mathbf{C}'\mathbf{y},$$

mit der (bayesianischen) Kovarianz aus (21)

$$\text{Cov}(\hat{f}(z) - f(z)) = \sigma^2 \mathbf{C}_z(\mathbf{C}'\mathbf{C} + \lambda \text{diag}(\mathbf{0}, \mathbf{I}_{m-2}))^{-1} \mathbf{C}'_z.$$

Diese Kovarianz berücksichtigt möglichen Bias in $\hat{f}(z)$. → Übung

Unter der Annahme asymptotischer Normalverteilung können hiermit approximative punktweise Konfidenzbänder konstruiert werden. → Übung

Für simultane Konfidenzbänder siehe Ruppert, Wand und Carroll (2003, 6.5).

Penalisierte Regression und LMMs

Für allgemeine Penalisierungsansätze mit penalisierter Residuenquadratsumme

$$PKQ(\lambda) = (\mathbf{y} - \mathbf{B}\boldsymbol{\gamma})'(\mathbf{y} - \mathbf{B}\boldsymbol{\gamma}) + \lambda\boldsymbol{\gamma}'\mathbf{K}\boldsymbol{\gamma} \quad (26)$$

könnte man nun versuchen, mit $\tau^2 = \sigma^2/\lambda$ ein LMM zu definieren der Form

$$\mathbf{y} = \mathbf{B}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I}_n), \quad \boldsymbol{\gamma} \sim N(\mathbf{0}, \tau^2\mathbf{K}^{-1}).$$

Problem: Die Inverse \mathbf{K}^{-1} existiert nicht immer, $0 < \text{Rang}(\mathbf{K}) < \dim(\boldsymbol{\gamma})$.

- **Beispiel:** Für penalisierte B-Splines hat $\mathbf{K}_k = \mathbf{D}_k' \mathbf{D}_k$ keinen vollen Rang.

Konsequenz: Die aus (26) resultierende Dichte der zufälligen Effekte

$$p(\boldsymbol{\gamma}) \propto \exp\left(-\frac{1}{2\tau^2}\boldsymbol{\gamma}'\mathbf{K}\boldsymbol{\gamma}\right)$$

ist damit **teilweise uneigentlich**, d.h. sie lässt sich nicht normieren.

Konstruktion einer geeigneten Zerlegung

Sei $h = \text{Rang}(\mathbf{K})$ und $d = \dim(\gamma)$. Zerlege γ in Teilvektoren β und \mathbf{b} , so dass

$$\gamma = \begin{matrix} \tilde{\mathbf{X}} & \beta & + & \tilde{\mathbf{Z}} & \mathbf{b} \\ d \times (d-h) & (d-h) \times 1 & & d \times h & h \times 1 \end{matrix}.$$

Wähle $\tilde{\mathbf{X}}$ und $\tilde{\mathbf{Z}}$, so dass

- $\mathbf{K}\tilde{\mathbf{X}} = \mathbf{0}$, so dass β nicht durch \mathbf{K} penalisiert wird,
- $\tilde{\mathbf{Z}}' \mathbf{K} \tilde{\mathbf{Z}} = \mathbf{I}_h$, so dass sich der Strafterm als $\lambda \mathbf{b}' \mathbf{b}$ schreiben lässt:

$$\gamma' \mathbf{K} \gamma = (\tilde{\mathbf{X}}\beta + \tilde{\mathbf{Z}}\mathbf{b})' \underbrace{\mathbf{K}(\tilde{\mathbf{X}}\beta + \tilde{\mathbf{Z}}\mathbf{b})}_{\mathbf{0}} = \mathbf{b}' \underbrace{\tilde{\mathbf{Z}}' \mathbf{K} \tilde{\mathbf{Z}}}_{\mathbf{I}_h} \mathbf{b} = \mathbf{b}' \mathbf{b}.$$

β , \mathbf{b} können dann als feste und (i.i.d.) zufällige Effekte aufgefasst werden im LMM

$$\mathbf{y} = \mathbf{B}\gamma + \varepsilon = \mathbf{B}\tilde{\mathbf{X}}\beta + \mathbf{B}\tilde{\mathbf{Z}}\mathbf{b} + \varepsilon =: \mathbf{X}\beta + \mathbf{Z}\mathbf{b} + \varepsilon, \quad \mathbf{b} \sim N(\mathbf{0}, \tau^2 \mathbf{I}_h).$$

$\lambda = \sigma^2 / \tau^2$ lässt sich wieder über REML oder ML schätzen.

Details zur Konstruktion

- Verwende eine Basis des Nullraums von \mathbf{K} als Spalten von $\tilde{\mathbf{X}} \Rightarrow \mathbf{K}\tilde{\mathbf{X}} = \mathbf{0}$.
Z.B. Polynome vom Grad 0 bis $k - 1$ für B-Splines mit k ten Differenzen.
- Für $\tilde{\mathbf{Z}}$, verwende Spektralzerlegung $\mathbf{K} = \underset{d \times h}{\mathbf{\Gamma}} \underset{h \times h}{\mathbf{\Omega}_+} \underset{h \times d}{\mathbf{\Gamma}'}$ mit $\mathbf{\Omega}_+$ Diagonalmatrix der positiven Eigenwerte, $\mathbf{\Gamma}$ orthonormale Matrix der Eigenvektoren.
Definiere $\tilde{\mathbf{Z}} = \mathbf{L}(\mathbf{L}'\mathbf{L})^{-1}$ mit $\mathbf{L} = \mathbf{\Gamma}\mathbf{\Omega}_+^{1/2}$ (also $\mathbf{K} = \mathbf{L}\mathbf{L}'$).
 $\Rightarrow \tilde{\mathbf{Z}}'\mathbf{K}\tilde{\mathbf{Z}} = (\mathbf{L}'\mathbf{L})^{-1} \mathbf{L}'\mathbf{L}\mathbf{L}'\mathbf{L}(\mathbf{L}'\mathbf{L})^{-1} = \mathbf{I}_h$.
- Die Spektralzerlegung ist nicht immer notwendig. So kann beispielsweise für P-Splines auch $\mathbf{L} = \mathbf{D}'$ mit der Differenzenmatrix \mathbf{D} gewählt werden.

\Rightarrow Die Zerlegung von γ ist nicht eindeutig.

Frequentistische und Bayesianische Sicht

- **Frequentistisch** betrachtet sind die Parameter γ feste, unbekannte Parameter.
- Im umformulierten gemischten Modell enthält der **b** -Teil von γ zufällige Effekte - keine (festen) Parameter, sondern Zufallsgrößen!
- Streng genommen wäre die Darstellung als LMM nur als algorithmischer Trick zu betrachten, nicht als tatsächliche Umformulierung des Modells.
- **Bayesianisch** betrachtet ist dies kein Problem, da eh alle Parameter Zufallsgrößen sind. Die zwei Darstellungen sind äquivalente Formulierungen der gleichen Priori-Annahmen.

Bayesianische Sicht auf Penalisierungsansätze

Beobachtungsmodell

$$\mathbf{y} = \mathbf{B}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

mit Designmatrix \mathbf{B} , die entsprechende Basisfunktionen enthält.

Penalisierung der Form $pen(\lambda, \mathbf{K}) = \lambda \boldsymbol{\gamma}' \mathbf{K} \boldsymbol{\gamma}$ entspricht (improperer) multivariater **Priori-Normalverteilung** $p(\boldsymbol{\gamma} | \mathbf{K}, \tau^2) \propto \exp\left(-\frac{1}{2\tau^2} \boldsymbol{\gamma}' \mathbf{K} \boldsymbol{\gamma}\right)$ für $\tau^2 = \sigma^2 / \lambda$ oder

$$\boldsymbol{\gamma} \sim N(\mathbf{0}, \tau^2 \mathbf{K}^{-1}).$$

Beispiel: TP-Basis

- **Priori-Normalverteilung** für die Koeffizienten der abgeschnittenen Polynome: $\gamma_{l+2}, \dots, \gamma_d \stackrel{iid}{\sim} N(0, \tau^2)$.
- **Nichtinformative Priori** für die Koeffizienten der globalen Polynome: $p(\gamma_j) \propto \text{const}, j = 1, \dots, l + 1$.

Die beiden verschiedenen Typen von Priori-Verteilungen spiegeln die Unterscheidung in unpenalisiert und penalisiert zu schätzende Parameter wieder.

Beispiel: B-Spline-Basis

- Priori-Annahme für γ : **Random Walks** (Irrfahrten) **der Ordnung k** als stochastisches Analogon zur Differenzenbestrafung. Z.B. $k = 1$ (RW1):

$$\gamma_j = \gamma_{j-1} + u_j, \quad u_j \sim N(0, \tau^2), \quad j = 2, \dots, d.$$

- **nichtinformative** Priori-Verteilung $p(\gamma_1) \propto \text{const}$
- \Rightarrow **bedingte Verteilungen** $\gamma_j | \gamma_{j-1}, \dots, \gamma_1 \sim N(\gamma_{j-1}, \tau^2)$
- **gemeinsame Verteilung** für γ : NV mit EW $\mathbf{0}$ und Präzisionsmatrix \mathbf{K}/τ^2 , $\mathbf{K} = \mathbf{D}'_1 \mathbf{D}_1$
- **Präzisionsmatrix** hat **keinen vollen Rang** \Rightarrow Kovarianzmatrix $\tau^2 \mathbf{K}^{-1}$ existiert nicht bzw. gemeinsame Priori ist eine **teilweise uneigentlichen Verteilung**
- trotzdem **eigentliche Posteriori-Verteilung** (NV) mit

$$E(\gamma | \mathbf{y}) = \left(\mathbf{B}' \mathbf{B} + \frac{\sigma^2}{\tau^2} \mathbf{K} \right)^{-1} \mathbf{B}' \mathbf{y}, \quad \text{Cov}(\gamma | \mathbf{y}) = \left(\mathbf{B}' \mathbf{B} + \frac{\sigma^2}{\tau^2} \mathbf{K} \right)^{-1}.$$

Variierende Koeffizienten

- Glatte Interaktion $g(z)x$ zwischen metrischer Variable z , binärer Variable x .
- Das Modell sollte dann auch die Haupteffekte $x\beta_1$ und $f(z)$ enthalten,

$$y_i = \beta_0 + x_i\beta_1 + f(z_i) + g(z_i)x_i + \varepsilon_i.$$

- $f(z)$ ist der nichtlineare Effekt von z , falls $x = 0$,
- $f(z) + g(z) + \beta_1$ ist der nichtlineare Effekt von z , falls $x = 1$.

Zentriere $g(z)$ um Null, damit das Modell identifizierbar ist:

$$g(z)x + \beta_1x = (g(z) + c)x + (\beta_1 - c)x.$$

- $g(z)$ lässt sich wieder über Basisfunktionen approximieren. Die Designmatrix enthält das Produkt der Werte der Basisfunktionen und der x -Werte.
- Eine geeignete Strafmatrix ist wie zuvor z.B. $\mathbf{K}_k = \mathbf{D}'_k \mathbf{D}_k$ für B-Splines.
- Schätzung über das lineare gemischte Modell analog zu penalisierten Splines. Analog für x stetig oder kategorial. Details in Fahrmeir et al (2009, 8.3-8.4).

Bivariate Glättung und räumliche Effekte

Für die Darstellung von räumlichen Effekten oder nichtparametrischen Interaktionsoberflächen lässt sich die Methodik der P-Splines verallgemeinern:

$$y_i = f(z_{i1}, z_{i2}) + \varepsilon_i \approx \sum_{j=1}^{d_1} \sum_{k=1}^{d_2} \gamma_{jk} B_{jk}(z_{i1}, z_{i2}) + \varepsilon_i$$

mit z.B.

$$B_{jk}(z_1, z_2) = B_j^{(1)}(z_1) B_k^{(2)}(z_2), \quad j = 1, \dots, d_1, \quad k = 1, \dots, d_2$$

eine **Tensorprodukt-Basis** (Produkte aller univariater B-Spline-Basisfunktionen).
Ein geeigneter Strafterm ist hier z.B. (mit Kronecker-Produkt \otimes , analog für \mathbf{D}_k)

$$\begin{aligned} \lambda \gamma' \mathbf{K} \gamma &= \lambda \left[\sum_{j=1}^{d_1} \sum_{k=2}^{d_2} (\gamma_{jk} - \gamma_{j,k-1})^2 + \sum_{k=1}^{d_2} \sum_{j=2}^{d_1} (\gamma_{jk} - \gamma_{j-1,k})^2 \right] \\ &= \lambda \gamma' [(\mathbf{I}_{d_2} \otimes \mathbf{D}_1)' (\mathbf{I}_{d_2} \otimes \mathbf{D}_1) + (\mathbf{D}_1 \otimes \mathbf{I}_{d_1})' (\mathbf{D}_1 \otimes \mathbf{I}_{d_1})] \gamma. \end{aligned}$$

Die Schätzung über gemischte Modelle läuft wie für allgemeine Penaliserungsansätze beschrieben. Für Details, weitere (z.B. radiale) Basisfunktionen und Beziehungen zum Kriging in der räumlichen Statistik, siehe Fahrmeir et al (2009, Abschnitt 7.2).

Für höherdimensionale Glättung, siehe Fahrmeir et al (2009, Abschnitt 7.3).

Markov-Zufallsfelder

Bei räumlichen Daten in Regionen oder auf einem Gitter wird räumliche Glättung meist über Nachbarschaften umgesetzt:

$$r \in N(s) \Leftrightarrow r \text{ ist zu } s \text{ benachbart.}$$

Verwenden wir ein Intercept γ_s pro Region s , so ist ein möglicher Strafterm

$$\lambda \sum_{s=1}^d \sum_{r \in N(s), r < s} (\gamma_r - \gamma_s)^2 = \lambda \gamma' \mathbf{K} \gamma$$

mit $\mathbf{K} = (K_{sr})_{sr}$, $K_{sr} = -1$ für $r \in N(s)$, $K_{sr} = |N(s)|$ für $r = s$, $K_{sr} = 0$, sonst. Eine Schätzung über gemischte Modelle ist wieder möglich, für Details und die Beziehung zu Markov-Zufallsfeldern siehe Fahrmeir et al (2009, Abschnitt 7.2.4).

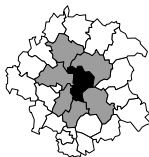
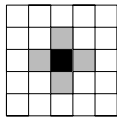


Abbildung 6: Nachbarschaften 1. Ordnung für reguläres Gitter und irreguläre Regionen. Die Nachbarn der schwarz gekennzeichneten Region sind in grau wiedergegeben. Quelle: Fahrmeir et al (2009, Abb. 7.45).

Ordinale Kovariablen

Bei einer Kovariablen mit geordneten Kategorien $k = 0, \dots, K$, bei der ein glatter Anstieg/Abfall der Koeffizienten $\gamma_0 = 0, \gamma_1, \dots, \gamma_K$ angenommen wird, schlagen

- Gertheiss, J. and G. Tutz (2009). *Penalized regression with ordinal predictors*. International Statistical Review 77, 345–365.

den Strafterm

$$\lambda \gamma' \mathbf{K} \gamma = \lambda \sum_{j=1}^K (\gamma_j - \gamma_{j-1})^2 = \lambda \gamma' \mathbf{D}'_1 \mathbf{D}_1 \gamma$$

vor. Die Schätzung über ein LMM läuft analog zu penalierten Splines. Der Test

$$H_0 : \gamma_1 = \dots = \gamma_K = \gamma_0 = 0 \quad \text{gegen} \quad H_A : \gamma_k \neq 0 \text{ für mindestens ein } k \quad (27)$$

auf Einfluss der Kovariablen über einen (R)LRT

$$H_0 : \tau^2 = 0 \quad \text{gegen} \quad H_A : \tau^2 > 0$$

im zugehörigen LMM hat größere Power als ein F-Test für (27) im linearen Modell:

- Gertheiss, J. and F. Oehrlin (2011). *Testing linearity and relevance of ordinal predictors*. Electronic Journal of Statistics 5, 1935-1959.

Lineare Gemischte Modelle als Modulare Modelle

Die Beispiele zeigen, dass LMMs ein mächtiges Werkzeug sind, mit dem flexibel Modellbausteine in **strukturiert-additiven Regressionsmodellen** modular kombiniert werden können:

- zufällige Effekte für z.B. Longitudinal- oder Clusterdaten
- räumliche Terme
- glatte Terme oder Interaktionen, variierende Koeffizienten
- ...

Damit sind z.B. Modelle der Form

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + b_{i0} + u_{ij1}b_{i1} + f_1(z_{ij1}) + f_2(z_{ij2}, z_{ij3}) + f_3(\mathbf{s}_{ij}) + f_4(z_{ij4})x_{ij,p+1} + \varepsilon_{ij}$$

mit Kovariablen \mathbf{x} und \mathbf{u} , metrischen Kovariablen \mathbf{z} , räumlichen Koordinaten \mathbf{s} , festen Effekten $\boldsymbol{\beta}$ und zufälligen Effekten \mathbf{b} möglich.

Software

Penalisierungsansätze können formal als LMM aufgeschrieben werden, weisen jedoch nicht die sonst typische Gruppierungsstruktur auf.

⇒ Verwende zur Schätzung spezialisierte Software, die an diese Struktur angepasst ist (und die nötigen Designmatrizen \mathbf{X} und \mathbf{Z} automatisch konstruiert).

Hier ist insbesondere das Paket `mgcv` von Simon Wood zu nennen.

Die `gamm()`-Funktion oder `gam()` mit Option `method = "(RE)ML"` fitten (generalisierte) additive gemischte Modelle mit (RE)ML-Schätzung.

- Feste und zufällige Effekte werden wie in `lme()` spezifiziert.
- Glatte Effekte $f(z)$ über `s(z)`.
- Glatte räumliche oder Interaktions-Oberflächen über `s(z1, z2)` oder `te(z1, z2)`.
- Für variierende Koeffizienten ist das `by`-Argument vorgesehen, `s(z, by=x)`.
- Weitere Möglichkeiten und Optionen, siehe die Manual sowie Wood (2006).

Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell**
 - Wiederholung: GLM und LMM
 - $GLMM = GLM + LMM$
 - GLMM für Longitudinal-/Clusterdaten
 - GLMM in allgemeiner Form
 - GLMM: Marginale und konditionale Verteilung
 - Generalisierte Additive Gemischte Modelle (GAMMs)
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle

Generalisiertes Lineares Modell (GLM): Strukturannahme

- **Strukturannahme:** Verknüpfung von Erwartungswert und linearem Prädiktor $\eta_i = \mathbf{x}_i' \boldsymbol{\beta}$ durch die Linkfunktion $g(\cdot)$:

$$E(y_i | \mathbf{x}_i) = g^{-1}(\eta_i) \Leftrightarrow g(E(y_i | \mathbf{x}_i)) = \eta_i$$

- Binäre Daten: z.B. Logitmodell: $E(y_i | \mathbf{x}_i) = P(y_i = 1 | \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i' \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i' \boldsymbol{\beta})}$
- Zähldaten: z.B. log-lineares Modell: $E(y_i | \mathbf{x}_i) = \exp(\mathbf{x}_i' \boldsymbol{\beta})$

Generalisiertes Lineares Modell (GLM): Verteilungsannahme

- **Verteilungsannahme:** Gegeben η_i sind die y_i unabhängig verteilt und ihre Dichte $f(y_i)$ gehört zu einer Exponentialfamilie:

$$f(y_i|\theta_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i} - c(y_i, \phi_i)\right)$$

$$\text{mit } E(y_i|\mathbf{x}_i) = b'(\theta_i) = g^{-1}(\mathbf{x}_i'\boldsymbol{\beta}); \text{ Var}(y_i|\mathbf{x}_i) = \phi b''(\theta_i);$$

- θ_i : natürlicher/kanonischer Parameter
- ϕ_i : Skalenparameter ($\phi_i = 1$ für Binär-/Poisson-Daten)

Die Exponentialfamilie - Beispiele

Dichten / Wahrscheinlichkeitsfunktionen:

$$\text{Normal : } f(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

$$\text{Binomial : } f(y|\pi) = \binom{n}{y} \pi^y (1 - \pi)^{n-y}$$

$$\text{Poisson : } f(y|\lambda) = \frac{\lambda^y}{y!} \exp(-\lambda)$$

	μ	θ	ϕ	$b(\theta)$
Normal	μ	μ	σ^2	$\frac{1}{2}\theta^2$
Binomial	π	$\log\left(\frac{\pi}{1-\pi}\right)$	1	$n \log(1 + \exp(\theta))$
Poisson	λ	$\log(\lambda)$	1	$\exp(\theta)$

GLM

- Strukturannahme:

$$E(y_i | \mathbf{x}_i) = g^{-1}(\mathbf{x}_i' \boldsymbol{\beta})$$

- Verteilungsannahme:

$$f(y_i | \theta_i) = \exp \left(\frac{y_i \theta_i - b(\theta_i)}{\phi_i} - c(y_i, \phi_i) \right)$$

⇒ Verknüpfung von Erwartungswertstruktur und Varianzstruktur

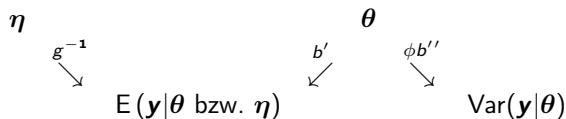
- $E(y_i | \mathbf{x}_i) = g^{-1}(\mathbf{x}_i' \boldsymbol{\beta}) = b'(\theta_i)$
- $\text{Var}(y_i | \mathbf{x}_i) = \phi_i b''(\theta_i)$

GLM

$$f(y_i|\theta_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i} - c(y_i, \phi_i)\right)$$

- $E(y_i|\mathbf{x}_i) = g^{-1}(\mathbf{x}_i'\boldsymbol{\beta}) = b'(\theta_i)$
- $\text{Var}(y_i|\mathbf{x}_i) = \phi_i b''(\theta_i)$

Also:

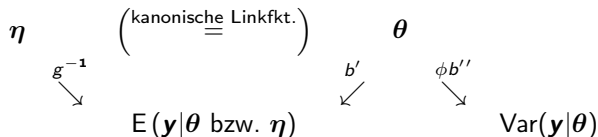


GLM mit kanonischem Link

$$f(y_i|\theta_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi} - c(y_i, \phi)\right)$$

- $E(y_i|\mathbf{x}_i) = \mathbf{g}^{-1}(\mathbf{x}'_i\boldsymbol{\beta}) = b'(\theta_i)$
- $\text{Var}(y_i|\mathbf{x}_i) = \phi b''(\theta_i)$
- kanonische Linkfunktion: $\mathbf{g}^{-1}(\cdot) = b'(\cdot)$

Also:



- Beispiele: $g(\cdot) = \text{id}(\cdot)$ für Normalverteilung, $g(\cdot) = \text{logit}(\cdot)$ für Bernoulliverteilung, $g(\cdot) = \text{log}(\cdot)$ für Poissonverteilung

LM \Rightarrow LMM

- Was ist anders/neu im gemischten Modell?
- Wozu gemischte Modelle?
- Warum sind gemischte Modelle so beliebt?

Formal: Was ist anders/neu im LMM?

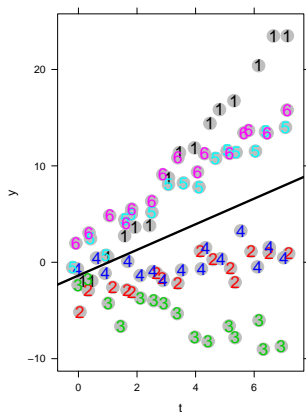
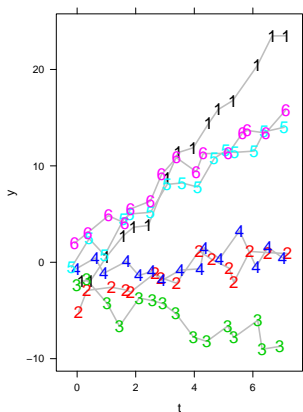
- Erweiterung des linearen Prädiktors um zufällige/penalisierte Effekte \mathbf{b} :
 $\eta = \mathbf{X}\beta + \mathbf{Zb}$
- Verteilungsannahme über \mathbf{b} : $\mathbf{b} \sim \mathcal{P}(\vartheta)$
üblich: $\mathbf{b} \sim N(\mathbf{0}, \mathbf{G}(\vartheta))$
- Inferenz über gemeinsame / penalisierte (Log-)Likelihood:
 $l_{pen}(\mathbf{y}|\beta, \mathbf{b}, \vartheta) = l(\mathbf{y}|\beta, \mathbf{b}, \vartheta) + l(\mathbf{b}|\vartheta)$
- konditionales Modell / marginales Modell:
 \mathbf{Z} und $\mathbf{G}(\vartheta)$ bestimmen marginale Kovarianzstruktur von \mathbf{y} .

$$\mathbf{y}|\mathbf{b} \sim N(\mathbf{X}\beta + \mathbf{Zb}, \mathbf{R}(\vartheta)); \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}(\vartheta))$$
$$\Rightarrow \mathbf{y} \sim N(\mathbf{X}\beta, \mathbf{ZG}(\vartheta)\mathbf{Z}' + \mathbf{R}(\vartheta))$$

Inhaltlich: Wozu LMMs?

Konditional betrachtet:

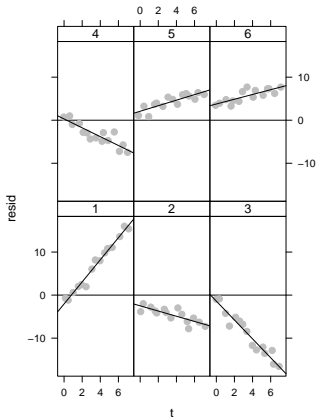
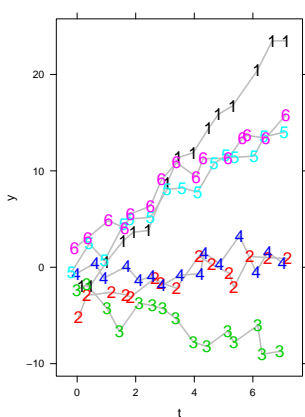
(zeitliche/räumliche/hierarchische) Struktur der Daten legt bestimmte Struktur der Abweichungen der Beobachtungen y von $X\beta$ nahe:



Inhaltlich: Wozu LMMs?

Konditional betrachtet:

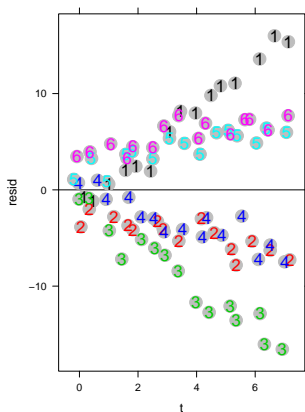
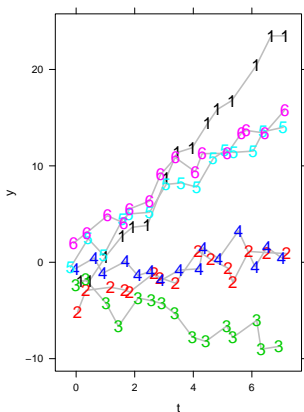
(zeitliche/räumliche/hierarchische) Struktur der Daten legt bestimmte Struktur der Abweichungen der Beobachtungen \mathbf{y} von $\mathbf{X}\beta$ nahe:



Inhaltlich: Wozu LMMs?

Marginal betrachtet:

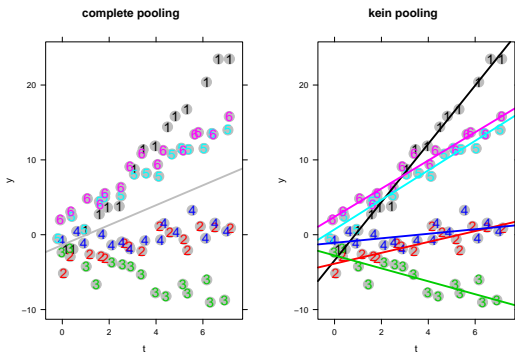
(zeitliche/räumliche/hierarchische) Struktur der Daten legt bestimmte Kovarianzstruktur der Beobachtungen nahe:



Warum sind LMMs so beliebt?

LMM bilden *Kompromiss* zwischen *komplettem pooling* und *keinem pooling*

- *komplettes pooling*: gewöhnliches LM ohne Einbezug der Gruppierungsstruktur
- *kein pooling*: separates LM (oder ein fester Effekt) für jede Ausprägung der Gruppierung



Warum sind LMMs so beliebt?

- LMM für hierarchische Daten als Kompromiss zwischen *komplettem pooling* und *keinem pooling*: Schätzung der Varianzparameter ϑ erfasst Heterogenität zwischen Gruppen auf Basis des gesamten Datensatzes → Shrinkage der Koeffizienten zum Populationsmittel entsprechend stärker oder schwächer.
- penalisiertes Likelihoodkriterium stabilisiert Parameterschätzung (aber: Bias-Varianz-Tradeoff!)
- bayesianisch: priori-Annahmen über \mathbf{b} bilden Vorwissen über Struktur der Daten ab. Einführung dieser zusätzlichen Information erlaubt stabile Schätzung von Modellen mit sehr vielen Parametern.
- sehr vielseitig: Modelle für zeitliche, räumliche, hierarchische und gruppierte Datenstrukturen sowie nichtlineare Effekte lassen sich als (G)LMM auffassen.

LMM \Rightarrow GLMM: Flexibilisierung der Verteilungsannahme

- LMM: nur normalverteilter Response: $\mathbf{y}|\boldsymbol{\eta} \sim N(\boldsymbol{\eta}, \mathbf{R})$
- GLMM: Response aus beliebiger Exponentialfamilie:

$$E(\mathbf{y}|\boldsymbol{\eta}) = g^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b})$$
$$\mathbf{y}|\boldsymbol{\eta} \sim \text{Expo.fam.}$$

GLMM = GLM + LMM

- 1. Stufe:
 y_i bedingt auf den erweiterten linearen Prädiktor $\mathbf{X}\beta + \mathbf{Zb}$
unabhängig verteilt nach Exponentialfamilie (analog GLM)
- 2. Stufe:
Verteilungsannahme über $\mathbf{b} \Rightarrow$ regularisierte Schätzung (analog LMM)
Interpretation der Verteilungsannahme als Strafterm oder Eigenschaft der Grundgesamtheit (freq.) oder Priori-Annahme (bayes.)

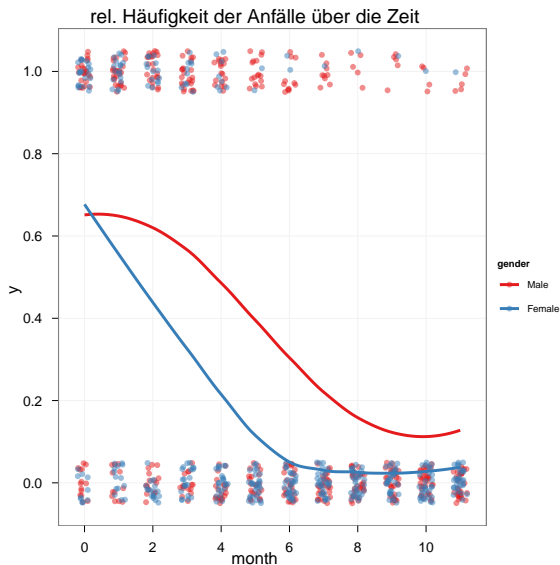
GLMM für Longitudinal-/Clusterdaten

- y_{ij} ; $i = 1, \dots, N$; $j = 1, \dots, n_i$ mit je n_i Messwiederholungen an N Beobachtungseinheiten (z.B. Binär-/Zählraten)
- Bedingte Verteilung von $y_{ij} | \mathbf{b}_i$ ist aus Exponentialfamilie (z.B. Bernoulli/Poisson)
- Bedingter Erwartungswert $\mu_{ij} = E(y_{ij} | \mathbf{b}) = g^{-1}(\mathbf{x}'_{ij}\boldsymbol{\beta} + \mathbf{z}'_{ij}\mathbf{b}_i)$
- Die zufälligen Effekte \mathbf{b}_i sind normalverteilt,

$$\mathbf{b}_i \stackrel{iid}{\sim} N(\mathbf{0}, \mathbf{D}).$$

Beispiel: Madras-Daten

- 69 stationäre Schizophrenie-Patienten
- 1 Jahr lang nach Ersteinlieferung beobachtet
- 12 monatliche Beobachtungen pro Patient ob Schub (ja/nein)



Beispiel: Madras-Daten

- Response: $y_{it} \in 0, 1$: Patient/in i hat Schub im Monat t
- Kovariablen:
 - month: Zeit in Monaten nach Einlieferung (0-11)
 - gender: Geschlecht
 - id: Patient
- Logit-Modell mit zufälligem Intercept: → Übung

$$\text{logit}(P(y_{it} = 1|b_i)) = \beta_0 + b_i + \beta_M \text{month}_t + \beta_G \text{gender}_i$$

Beispiel: Madras-Daten - Ergebnisse

Random effects:

Groups Name	Variance	Std.Dev.
id (Intercept)	4.91	2.22

Number of obs: 828, groups: id, 69

Fixed effects:

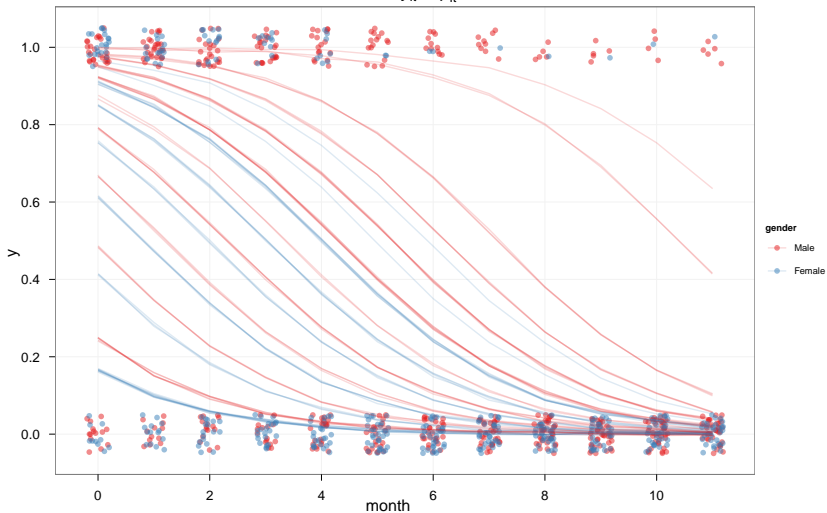
	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	1.933	0.437	4.42	9.8e-06	***
month	-0.573	0.045	-12.72	< 2e-16	***
genderFemale	-1.574	0.595	-2.65	0.0081	**

```
R> quantile(exp(ranef(m1)$id$(Intercept)'), p=c(.2,.4,.6,.8))
```

20%	40%	60%	80%
0.138	0.634	2.167	5.828

⇒ relativ starke individuelle Heterogenität

Beispiel: Madras-Daten - Ergebnisse

Modell 1: y_{it} & \hat{p}_{it} 

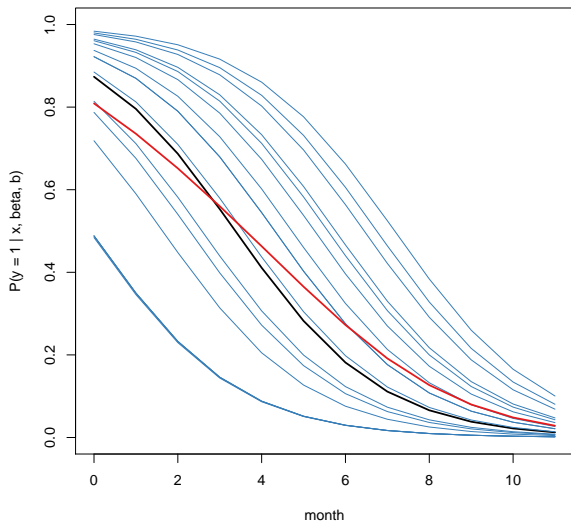
Beispiel: Madras-Daten - Interpretation

Wie sind die Parameter in diesem Logit-Modell mit Random Intercept zu interpretieren?

	beta	exp(beta)
month	-0.573	0.564

Graphisch:

- $\frac{\exp(\beta_0 + \beta_1 \text{month}_t)}{1 + \exp(\beta_0 + \beta_1 \text{month}_t)}$
- $P(y = 1 | x, \beta, b)$ für Quantile von b
- $E_b[P(y = 1 | x, \beta, b)] = P(y = 1 | x, \beta)$



Interpretation Random-Intercept-Logitmodell

Konditional: Odds-Ratio

$$\frac{\frac{P(y_{it}=1|x_{it}+1, b_i)}{P(y_{it}=0|x_{it}+1, b_i)}}{\frac{P(y_{it}=1|x_{it}, b_i)}{P(y_{it}=0|x_{it}, b_i)}} = \frac{\exp((x_{it} + 1)\beta + b_i)}{\exp(x_{it}\beta + b_i)} = \exp(\beta)$$

Marginal: Odds-Ratio

$$\frac{\frac{P(y_{it}=1|x_{it}+1)}{P(y_{it}=0|x_{it}+1)}}{\frac{P(y_{it}=1|x_{it})}{P(y_{it}=0|x_{it})}} = \frac{E_{b_i}[P(y_{it}=1|x_{it}+1, b_i)]}{E_{b_i}[P(y_{it}=0|x_{it}+1, b_i)]} \neq \exp(\beta)$$

$$\text{mit } E_{b_i}[P(y_{it} = 1|x_{it}, b_i)] = \int \frac{\exp(x_{it}\beta + b_i)}{1 + \exp(x_{it}\beta + b_i)} p(b_i) db_i \quad \text{etc.}$$

⇒ Interpretation von β nur **bedingt auf zufällige Effekte b_i zulässig**

⇒ β ist **individuenspezifischer Parameter, nicht Populationsparameter**

Beispiel: Madras-Daten - Interpretation

Wie sind die Parameter in diesem Logit-Modell mit Random Intercept zu interpretieren?

	beta	exp(beta)
month	-0.573	0.564

GLMM in allgemeiner Form

- **Verteilungsannahme I:** Gegeben die zufälligen Effekte \mathbf{b} sind die y_i bedingt unabhängig und die bedingte Dichte $f(y_i|\mathbf{b})$ gehört zur Exponentialfamilie
- **Strukturannahme:** Der **bedingte** Erwartungswert $\mu_i = E(y_i|\mathbf{b}_i)$ ist mit dem linearen Prädiktor

$$\eta_i = \mathbf{x}'_i\boldsymbol{\beta} + \mathbf{z}'_i\mathbf{b}$$

durch $\mu_i = g^{-1}(\eta_i)$ verknüpft mit bekannter Linkfunktion $g(\cdot)$.

- **Verteilungsannahme II:** Die zufälligen Effekte \mathbf{b} sind normalverteilt,

$$\mathbf{b} \sim N(\mathbf{0}, \mathbf{G})$$

(auch andere Verteilungen möglich).

Hierarchisches GLMM Spezialfall mit $\mathbf{X} = (\mathbf{X}'_1 | \dots | \mathbf{X}'_N)'$, $\mathbf{Z} = \text{diag}(\mathbf{Z}_1, \dots, \mathbf{Z}_N)$ und $\mathbf{G} = \text{diag}(\mathbf{D}, \dots, \mathbf{D})$ analog zum hierarchischen LMM. ▶ Hierarch. GLMM

GLMM: Marginale und konditionale Verteilung

Erinnerung:

$$E(y_i|\eta_i) = g^{-1}(\mathbf{x}'_i\boldsymbol{\beta} + \mathbf{z}'_i\mathbf{b})$$

Für $g^{-1}(x) \neq x$ gilt:

$$\begin{aligned} E_b(E(y_i|\eta_i)) &= \int E(y_i|\eta_i)p(\mathbf{b}|\boldsymbol{\vartheta})d\mathbf{b} \\ &= \int g^{-1}(\mathbf{x}'_i\boldsymbol{\beta} + \mathbf{z}'_i\mathbf{b})p(\mathbf{b}|\boldsymbol{\vartheta})d\mathbf{b} \\ &\neq g^{-1}(\mathbf{x}'_i\boldsymbol{\beta}) \end{aligned}$$

GLMM mit $g^{-1}(x) \neq x$

⇒ marginaler Erwartungswert $E_b(E(y_i|\eta_i)) \neq g^{-1}(\mathbf{x}'_i\boldsymbol{\beta})$

Parameter $\boldsymbol{\beta}$ haben nur **konditionale** Interpretation.

GLMM: Marginale und konditionale Verteilung

- **Konditionale Kovarianzstruktur:** Gegeben die zufälligen Effekte \mathbf{b} sind die y_i bedingt unabhängig

$$\text{Cov}(y_i, y_{i'} | \mathbf{b}) = 0, \quad i \neq i'.$$

(Es gibt keine Residuenkovarianzmatrix \mathbf{R} wie im LMM.)

- **Marginale Kovarianzstruktur:** Es folgt für die marginale Kovarianz

$$\begin{aligned} \text{Cov}(y_i, y_{i'}) &= \text{Cov}(E(y_i | \mathbf{b}), E(y_{i'} | \mathbf{b})) + E(\text{Cov}(y_i, y_{i'} | \mathbf{b})) \\ &= \text{Cov}(\mu_i, \mu_{i'}) + 0 \\ &= \text{Cov}(g^{-1}(\mathbf{x}_i' \boldsymbol{\beta} + \mathbf{z}_i' \mathbf{b}), g^{-1}(\mathbf{x}_{i'}' \boldsymbol{\beta} + \mathbf{z}_{i'}' \mathbf{b})), i \neq i'. \end{aligned}$$

Im Allgemeinen hängt die marginale Korrelation von den Kovariablen \mathbf{x}_i ab.

GLMM + AMM = GAMM

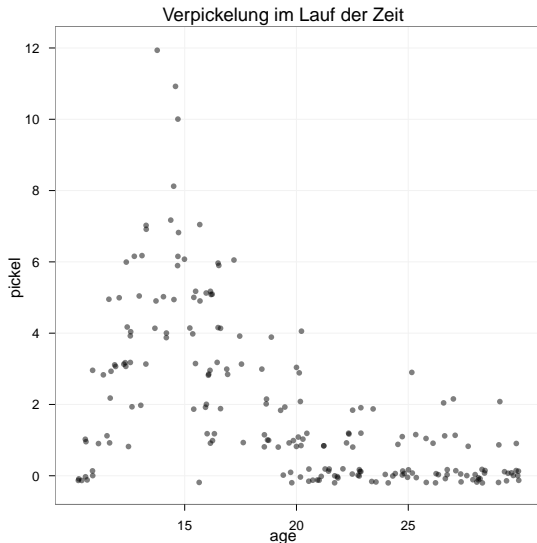
- Additive gemischte Modelle (AMMs) und
- Generalisierte lineare gemischte Modelle (GLMMs)

lassen sich zu generalisierten additiven gemischten Modellen (GAMMs) verbinden.

Hierzu werden die Designmatrizen \mathbf{X} und \mathbf{Z} und die Kovarianz \mathbf{G} der zufälligen Effekte \mathbf{b} im GLMM analog zu den additiven gemischten Modellen gewählt, siehe Kapitel 5.

Beispiel Verpickelung

- keine echten Daten
- Pickel auf den Gesichtern von 200 Teens und Twens wurden gezählt (1 Beob./Person)
- Verlauf der Verpickelung über die Zeit



Beispiel Verpickelung: Modell

- Response pickel_i : Person i hat pickel_i Pickel im Gesicht
- Kovariable alter_i : Alter in Jahren (10 – 30)
- Modell für $\mathbf{pickel} = (\text{pickel}_1, \dots, \text{pickel}_n)$ mit Splines aus der TP(2)-Basis für $\mathbf{alter} = (\text{alter}_1, \dots, \text{alter}_n)$:

$$\mathbf{pickel} | \mathbf{b} \sim \text{Poisson}(\exp(f(\mathbf{alter})))$$

$$f(\mathbf{alter}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}$$

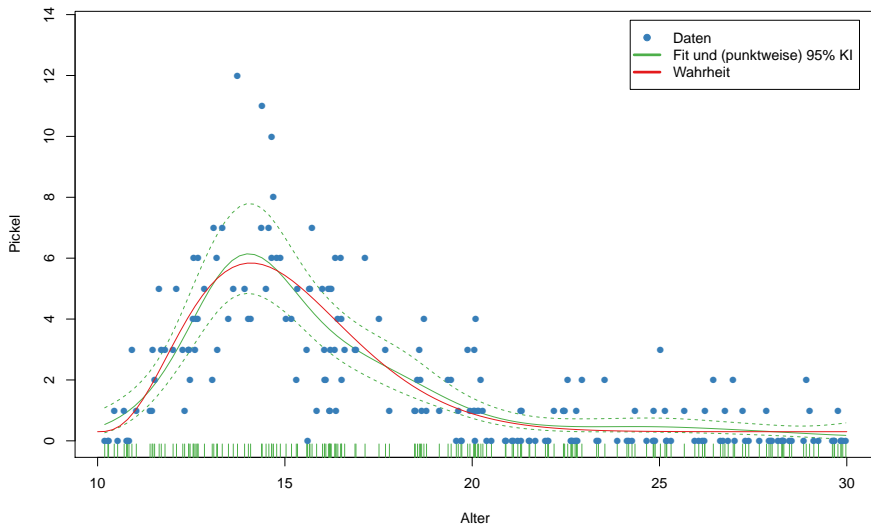
$$\mathbf{X} = (\mathbf{1}, \mathbf{alter}, \mathbf{alter}^2)$$

$$\mathbf{Z} = ((\text{alter}_i - \kappa_k)_+^2)_{i=1, \dots, n; k=1, \dots, m}$$

$$\mathbf{b} \sim N(\mathbf{0}, \tau^2 \mathbf{I}_m)$$

- Wir verwenden für den Fit auf der nächsten Folie quadratische B-Splines mit Differenzen 2. Ordnung. → **Übung**

Beispiel Verpickelung: Ergebnis



Inhalt

- 1 Das lineare gemischte Modell
- 2 Likelihood-Schätzung für lineare gemischte Modelle
- 3 Likelihood-Inferenz im linearen gemischten Modell
- 4 Bayes-Schätzung für lineare gemischte Modelle
- 5 Additive gemischte Modelle
- 6 Das generalisierte lineare gemischte Modell
- 7 Likelihood-Schätzung für generalisierte lineare gemischte Modelle
 - Laplace-Approximation und P-IRLS
 - Adaptive Gauss-Hermite Quadratur (AGQ)
 - Penalized Quasi-Likelihood (PQL)
 - Inferenz in GLMMs

Likelihood eines GLMM

Verteilungsannahme für $\mathbf{y}|\mathbf{b}$: $y_i|\mathbf{b} \sim \text{Expo.fam.}(\boldsymbol{\theta}, \phi)$ unabhängig

Verteilungsannahme für $\mathbf{b}|\boldsymbol{\vartheta}$: $\mathbf{b}|\boldsymbol{\vartheta} \sim N(\mathbf{0}, \mathbf{G}(\boldsymbol{\vartheta}))$

⇒ gemeinsame / penalisierte Likelihood:

$$L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}) = f(\mathbf{y}, \mathbf{b}|\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}) = \left(\prod_{i=1}^n f(y_i|\boldsymbol{\beta}, \phi, \mathbf{b}, \boldsymbol{\vartheta}) \right) f(\mathbf{b}|\boldsymbol{\vartheta})$$

bzw. marginale Likelihood:

$$L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}) = f(\mathbf{y}|\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}) = \int \left(\prod_{i=1}^n f(y_i|\boldsymbol{\beta}, \phi, \mathbf{b}, \boldsymbol{\vartheta}) \right) f(\mathbf{b}|\boldsymbol{\vartheta}) d\mathbf{b}.$$

Erinnerung: im LMM ist das Integral analytisch lösbar.

$\int f(\mathbf{y}, \mathbf{b}|\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}) d\mathbf{b}$ ist die Dichte einer $N(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\mathbf{G}(\boldsymbol{\vartheta})\mathbf{Z}' + \mathbf{R}(\phi, \boldsymbol{\vartheta}))$ -Verteilung.

Im GLMM: ??

ML-Inferenz

Interessante Parameter: primär β, ϑ , evtl. ϕ

⇒ Problem: finde $\operatorname{argmax} L(\beta, \vartheta, \phi)$, wobei

$$\begin{aligned}
 L(\beta, \vartheta, \phi) &= \int \left(\prod_{i=1}^n f(y_i | \beta, \phi, \mathbf{b}, \vartheta) \right) f(\mathbf{b} | \vartheta) d\mathbf{b} \\
 &= \int \left(\prod_{i=1}^n \exp \left(\frac{y_i \theta_i - b(\theta_i)}{\phi} - c(y_i, \phi) \right) \right) |\mathbf{G}(\vartheta)|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{b}' \mathbf{G}(\vartheta)^{-1} \mathbf{b} \right) d\mathbf{b} \\
 &\quad \text{mit } \theta_i = (b')^{-1}(g^{-1}(\eta_i)), \quad \eta_i = \mathbf{x}'_i \beta + \mathbf{z}'_i \mathbf{b}.
 \end{aligned}$$

Hoch-dimensionales Integral, i.A. nicht analytisch lösbar

⇒ iterative Optimierung einer Approximation der Likelihood

Laplace-Approximation

Problem: Löse r -dimensionales Integral $H = \int \exp(h(\theta)) d\theta$

Ansatz:

- bestimme $\theta_0 = \arg \max h(\theta)$

- quadratische Taylor-Entwicklung von $h(\theta)$ um θ_0 :

$$h(\theta) \approx h(\theta_0) + \frac{1}{2}(\theta - \theta_0)' \underbrace{\left(\frac{\partial^2}{\partial \theta \partial \theta'} h(\theta_0) \right)}_{=-\mathbf{P}} (\theta - \theta_0)$$

- $\Rightarrow \int \exp(h(\theta)) d\theta \approx \int \exp(h(\theta_0) - \frac{1}{2}(\theta - \theta_0)' \mathbf{P} (\theta - \theta_0)) d\theta$
wie bei $N(\theta_0, \mathbf{P}^{-1})$

- $\Rightarrow H \approx \exp(h(\theta_0)) \underbrace{((2\pi)^{r/2} |\mathbf{P}|^{-1/2})}_{1/\text{Normierung der } N(\theta_0, \mathbf{P}^{-1})}$

GLMM-Likelihood: Laplace-Approximation

Verwende eine Laplace-Approximation, mit Entwicklung um den Maximierer $\hat{\mathbf{b}}$ von $L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}) = f(\mathbf{y}|\boldsymbol{\beta}, \phi, \mathbf{b})f(\mathbf{b}|\boldsymbol{\vartheta})$,

$$\begin{aligned} \log(L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta})) &= \log \left(\int f(\mathbf{y}|\boldsymbol{\beta}, \phi, \mathbf{b})f(\mathbf{b}|\boldsymbol{\vartheta})d\mathbf{b} \right) \\ &= \log \left(\int L(\boldsymbol{\beta}, \phi, \mathbf{b})(2\pi)^{-r/2}|\mathbf{G}(\boldsymbol{\vartheta})|^{-\frac{1}{2}} \exp \left(-\frac{1}{2}\mathbf{b}'\mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b} \right) d\mathbf{b} \right) \\ &\approx \log(L(\boldsymbol{\beta}, \hat{\mathbf{b}}, \phi)) - \frac{r}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{G}(\boldsymbol{\vartheta})| - \frac{1}{2} \hat{\mathbf{b}}' \mathbf{G}(\boldsymbol{\vartheta})^{-1} \hat{\mathbf{b}} \\ &\quad + \log \left(\int \exp \left(-\frac{1}{2}(\mathbf{b} - \hat{\mathbf{b}})' \mathcal{I}(\hat{\mathbf{b}})(\mathbf{b} - \hat{\mathbf{b}}) \right) d\mathbf{b} \right) \\ &= l(\boldsymbol{\beta}, \hat{\mathbf{b}}, \phi) - \frac{1}{2} \log |\mathcal{I}(\hat{\mathbf{b}})| - \frac{1}{2} \log |\mathbf{G}(\boldsymbol{\vartheta})| - \frac{1}{2} \hat{\mathbf{b}}' \mathbf{G}(\boldsymbol{\vartheta})^{-1} \hat{\mathbf{b}}. \end{aligned}$$

Dabei ist $\mathcal{I}(\mathbf{b}) = -E \left(\frac{\partial^2}{\partial \mathbf{b} \partial \mathbf{b}'} l(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}) \right)$ (mit zusätzlichem Erwartungswert), Herleitung siehe S. 166.

Schaukel-Algorithmus

In der Laplace-Approximation wird der Maximierer $\hat{\mathbf{b}}$ von $L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b})$ benötigt.
⇒ iterativer, zweistufiger Schaukel-Algorithmus:

- 1 Für gegebene $\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}$ bestimme $\hat{\mathbf{b}} = \arg \max_{\mathbf{b}} L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b})$ über penalisierten IRLS-Algorithmus (P-IRLS).
- 2 Maximiere Laplace-Approximation $\tilde{L}(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta})$ von $L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta})$ in $\hat{\mathbf{b}}$ mit numerischer Optimierung (Pseudo-Newton-Algorithmen wie BFGS).

Iteriere bis zur Konvergenz der Devianz $-2 \log L(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b})$.

Dieser Schaukel-Algorithmus ist im R-Paket [lme4](#) implementiert.

Grundidee: IRLS-Algorithmus

- IRLS = Fisher-Scoring für GLM
 - IRLS: Iteratively **R**e-Weighted **L**east **S**quares
 - Führt GLM-Schätzproblem auf iterierte gewichtete KQ-Schätzung zurück.
- Fisher-Scoring zur Lösung des Score-Gleichungssystems $s(\boldsymbol{\theta}) \stackrel{!}{=} \mathbf{0}$
 - lineare Taylor-Entwicklung $s(\boldsymbol{\theta}) \approx s(\boldsymbol{\theta}_0) - \mathcal{J}(\boldsymbol{\theta}_0)(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \stackrel{!}{=} \mathbf{0}$
mit der **beobachteten Informationsmatrix** $\mathcal{J}(\boldsymbol{\theta})$.
 - Ersetzen von $\mathcal{J}(\boldsymbol{\theta})$ (Newton-Raphson) durch die **erwartete Informationsmatrix** $\mathcal{I}(\boldsymbol{\theta})$ (Fisher-Scoring) liefert

$$\mathcal{I}(\boldsymbol{\theta}_0)\boldsymbol{\theta} = \mathcal{I}(\boldsymbol{\theta}_0)\boldsymbol{\theta}_0 + s(\boldsymbol{\theta}_0).$$

$s(\mathbf{b})$ und $\mathcal{I}(\mathbf{b})$ für den kanonischen Link

Für den kanonischen Link ist $\boldsymbol{\theta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}$.

$$\begin{aligned}\Rightarrow s(\mathbf{b}) &= \frac{\partial}{\partial \mathbf{b}} l(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}) = \frac{\partial}{\partial \mathbf{b}} \left(\text{const} + \frac{\boldsymbol{\theta}'\mathbf{y} - b(\boldsymbol{\theta})'\mathbf{1}}{\phi} - \frac{1}{2}\mathbf{b}'\mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b} \right) \\ &= \frac{1}{\phi}(\mathbf{Z}'\mathbf{y} - \mathbf{Z}'\text{diag}(b'(\boldsymbol{\theta}))\mathbf{1}) - \mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b} \\ &= \frac{1}{\phi}\mathbf{Z}'(\mathbf{y} - \boldsymbol{\mu}) - \mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b} \quad \text{und}\end{aligned}$$

$$\begin{aligned}\mathcal{I}(\mathbf{b}) &= -\text{E} \left(\frac{\partial^2}{\partial \mathbf{b} \partial \mathbf{b}'} l(\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}) \right) = - \left(\frac{\partial}{\partial \mathbf{b}'} s(\mathbf{b}) \right) \\ &= \frac{1}{\phi} \mathbf{Z} \text{diag}(b''(\boldsymbol{\theta})) \mathbf{Z}' + \mathbf{G}(\boldsymbol{\vartheta})^{-1} =: \mathbf{Z}\mathbf{W}\mathbf{Z}' + \mathbf{G}(\boldsymbol{\vartheta})^{-1}\end{aligned}$$

P-IRLS

Fisher-Scoring ausgehend von Startwert \mathbf{b}_0 :

$$\mathcal{I}(\mathbf{b}_0)\mathbf{b} = \mathcal{I}(\mathbf{b}_0)\mathbf{b}_0 + s(\mathbf{b}_0)$$

$$\text{mit im GLMM: } s(\mathbf{b}) = \frac{1}{\phi} \mathbf{Z}'(\mathbf{y} - \boldsymbol{\mu}) - \mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b};$$

$$\mathcal{I}(\mathbf{b}) = \mathbf{Z}'\mathbf{W}\mathbf{Z} + \mathbf{G}(\boldsymbol{\vartheta})^{-1}$$

liefert mit $\mathbf{W}_0 = \mathbf{W}(\mathbf{b}_0)$, $\boldsymbol{\mu}_0 = \boldsymbol{\mu}(\mathbf{b}_0)$:

$$(\mathbf{Z}'\mathbf{W}_0\mathbf{Z} + \mathbf{G}(\boldsymbol{\vartheta})^{-1})\mathbf{b} = (\mathbf{Z}'\mathbf{W}_0\mathbf{Z} + \mathbf{G}(\boldsymbol{\vartheta})^{-1})\mathbf{b}_0 + \frac{1}{\phi} \mathbf{Z}'(\mathbf{y} - \boldsymbol{\mu}_0) - \mathbf{G}(\boldsymbol{\vartheta})^{-1}\mathbf{b}_0$$

$$\Leftrightarrow (\mathbf{Z}'\mathbf{W}_0\mathbf{Z} + \mathbf{G}(\boldsymbol{\vartheta})^{-1})\mathbf{b} = \underbrace{\mathbf{Z}'\mathbf{W}_0(\mathbf{Z}\mathbf{b}_0 + \frac{1}{\phi} \mathbf{W}_0^{-1}(\mathbf{y} - \boldsymbol{\mu}_0))}_{\text{working response } \tilde{\mathbf{y}}}$$

\Rightarrow Schätzugleichung eines LMM mit bekanntem \mathbf{W}_0 , $\mathbf{G}(\boldsymbol{\vartheta})$, vgl. (14):

$$\tilde{\mathbf{y}}|\mathbf{b} \sim N(\mathbf{Z}\mathbf{b}, \mathbf{W}_0^{-1}); \mathbf{b} \sim N(\mathbf{0}, \mathbf{G}(\boldsymbol{\vartheta}))$$

P-IRLS

P-IRLS-Algorithmus iteriert folgende Schritte bis zur Konvergenz von $\hat{\mathbf{b}}$:

- i. setze Iterationswert $\mathbf{b}_0 = \hat{\mathbf{b}}^{(k)}$, berechne mit β , \mathbf{b}_0 die working responses $\tilde{\mathbf{y}}$ und Gewichte \mathbf{W}_0
- ii. berechne $\hat{\mathbf{b}}^{(k+1)}$ als Lösung des daraus abgeleiteten gewichteten, penalisierten KQ-Problems

(Adaptive) Gauss'sche Quadratur

- Laplace-Approximation recht schnell, aber ungenau besonders für kleine Clustergrößen und starke „Diskretheit“ (logistische Regression schlechter als Poisson-Regression).
- Gauss-Quadratur genauere Methode, um $\int f(\mathbf{y}|\boldsymbol{\beta}, \phi, \mathbf{b})f(\mathbf{b}|\boldsymbol{\vartheta})d\mathbf{b}$ zu approximieren, aber wesentlich rechenaufwändiger.
- Benutzt orthonormalisierte zufällige Effekte $\mathbf{b}^* = \mathbf{G}(\boldsymbol{\vartheta})^{-1/2}\mathbf{b}$.
- $\int f(\mathbf{y}|\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}^*)f(\mathbf{b}^*)d\mathbf{b}^* \approx \sum_{k=1}^Q w_k f(\mathbf{y}|\boldsymbol{\beta}, \phi, \boldsymbol{\vartheta}, \mathbf{b}_k^*)$
- Die Stützstellen \mathbf{b}_k^* ergeben sich aus den Nullstellen des Q-ten Hermite-Polynoms und bestimmen auch die Gewichte w_k .

(Adaptive) Gauss'sche Quadratur

- Genauigkeit der Approximation steigt mit wachsendem Q
⇒ erhöhe Q solange bis keine Änderung in Schätzung mehr zu beobachten
- in `lme4` nur implementiert für Modelle mit *einer* einzigen Gruppierungsvariable (option: `nAGQ`)
- Laplace-Approximation ergibt sich als Spezialfall $Q = 1$.

Penalized Quasi-Likelihood (PQL):

Ähnliche Idee wie P-IRLS: Approximiere \mathbf{y} durch $\boldsymbol{\mu} = E(\mathbf{y}|\mathbf{b})$ plus Fehler mit Varianz $\text{Var}(\mathbf{y}|\mathbf{b})$. Eine Taylor-Approximation von $\boldsymbol{\mu}$ um die aktuellen Werte $\boldsymbol{\beta}_0$ und \mathbf{b}_0 und Umsortieren ergibt

$$\tilde{\mathbf{y}} := \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\mathbf{b}_0 + \frac{1}{\phi} \mathbf{W}_0^{-1}(\mathbf{y} - \boldsymbol{\mu}_0) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}.$$

Algorithmus:

- 1 Für gegebene Werte $\boldsymbol{\beta}_0 = \hat{\boldsymbol{\beta}}^{(k)}$, $\mathbf{b}_0 = \hat{\mathbf{b}}^{(k)}$ bestimme *working responses* $\tilde{\mathbf{y}}$.
- 2 Bestimme $\hat{\boldsymbol{\beta}}^{(k+1)}$, $\hat{\mathbf{b}}^{(k+1)}$ und $\hat{\boldsymbol{\vartheta}}^{(k+1)}$, $\hat{\phi}^{(k+1)}$ aus dem LMM $\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$ mit $\mathbf{b} \sim N(\mathbf{0}, \mathbf{G}(\boldsymbol{\vartheta}))$ und $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{W}_0^{-1})$

Iteriere bis zur Konvergenz.

PQL

- Alternative Herleitung für PQL: ergibt sich ebenfalls bei Fisher-Scoring für die gemeinsame/penalisierte log-Likelihood (die im LMM die BLUPs ergibt). Daher auch mögliche Schätzart bei GAMMs.
- in R implementiert in Funktion `g1mmPQL` im Paket MASS (`g1mmPQL` wird auch von Funktion `gamm` im Paket `mgcv` benutzt)
- benutzt iterierte ML- oder REML-Schätzung
- Schätzung der Varianzparameter nach unten verzerrt, v.a. falls n_i klein. Bessere Approximation, je näher Daten an Normalverteilung (z.B. Poisson-Verteilung mit größeren λ 's)
- Konvergenz nicht garantiert
- Kein AIC, BIC oder Devianz berechenbar

Vergleich glmmPQL/lme vs. (g)lmer

	MASS::glmmPQL / nlme::lme	lme4::(g)lmer
Daten	nur genestete Daten; grosse Datensätze/Modelle oft nicht zu fiten	genestete & gekreuzte Datenstrukturen; auch riesige Datensätze werden gefittet
(G)AMs Cov(\mathbf{b}_i)	via mgcv::gamm sehr breite, erweiterbare Klasse von Kovarianzstrukturen (s. nlme::pdMat)	via gamm4 nur unstrukturierte oder diagonale Kovarianzen (in gamm4 beliebige Präzisionsmatrizen \mathbf{P} mit $\mathbf{G}(\vartheta) = \vartheta \mathbf{P}^{-1}$)
Cov(ε)	sehr breite, erweiterbare Klasse von Kovarianzstrukturen (s. nlme::varFunc)	Cov(ε) = $\sigma^2 \mathbf{I}_n$ oder $\sigma^2 \text{diag}(\mathbf{w}^{-1})$ mit bekannten Gewichten \mathbf{w}
Stabilität	sehr instabil für komplexe Strukturen von \mathbf{b}	sehr stabil für LMMs; für GLMMs ab > 3 zuf. Effekten oft kritisch
Speed	relativ langsam	LA sehr schnell (benutzt Sparse-Matrix-Algorithmen); AGQ deutlich langsamer

Vorhersage der zufälligen Effekte

Die beste Vorhersage für \mathbf{b} (minimaler mean squared error of prediction $E(\mathbf{b} - \hat{\mathbf{b}})^2$) wäre wieder

$$\hat{\mathbf{b}} = E(\mathbf{b}|\mathbf{y}).$$

In der Praxis würde man $\hat{\beta}$ und $\hat{\vartheta}$ einsetzen. Allerdings erfordert

$$\hat{\mathbf{b}} = \int \mathbf{b}f(\mathbf{b}|\mathbf{y}, \beta, \vartheta)d\mathbf{b} = \frac{\int \mathbf{b}f(\mathbf{y}|\mathbf{b}, \beta, \phi)f(\mathbf{b}|\vartheta)d\mathbf{b}}{\int f(\mathbf{y}|\mathbf{b}, \beta, \phi)f(\mathbf{b}|\vartheta)d\mathbf{b}}$$

wieder numerische Integration.

Alternativ zum Posteriori-EW wird häufig (z.B. in `lme4` als Teil von P-IRLS) der Posteriori-Modus berechnet, der $f(\mathbf{b}|\mathbf{y}, \beta, \vartheta) \propto f(\mathbf{y}|\mathbf{b}, \beta, \phi)f(\mathbf{b}|\vartheta)$ maximiert. Die beiden unterscheiden sich i.A. außer im LMM unter NV.

PQL schätzt \mathbf{b} direkt im Schätzalgorithmus.

Hypothesen-Tests für β

- Wegen der Maximum-Likelihood-Schätzung von β können prinzipiell Wald-, Likelihood-Quotienten- (LQT) oder Score-Tests verwendet werden mit einer entsprechenden χ^2 -Verteilung als Referenzverteilung.
- Die Güte der Approximation hängt ab von
 - der Güte der Approximation an die Likelihood in der Schätzung
 - der Güte der asymptotischen Approximation, die nur für Longitudinal/Clusterdaten bei $N \rightarrow \infty$ greift.
- PQL-Schätzung beruht auf der Likelihood von Pseudodaten und daher ist kein Likelihood-Quotiententest möglich. Inferenz für PQL wird meist basiert auf dem LMM in der PQL-Schätzung. Allerdings ist $\hat{\beta}$ i.A. nicht konsistent.

Inferenz für D

Hypothesen-Tests

- Für das longitudinale/Cluster-GLMM gelten die gleichen asymptotischen Ergebnisse für Tests von Parametern in D wie im LMM (Rand des Parameterraums).
- Es ist keine exakte Verteilung vorhanden.

Modellselektion

- Ein konditionales AIC für Poisson-, binomial- oder normalverteilte Zielgrößen ist im R-Paket `cAIC4` implementiert, siehe auch Saefken, Kneib, van Waveren & Greven (2014). A unifying approach to the estimation of the conditional Akaike information in generalized linear mixed models. *Electronic Journal Statistics*, 8, 201-225.

Konfidenzintervalle

- Generell gilt für einen Parameter θ :
Ein Wert θ_0 ist im $(1 - \alpha)\%$ -Konfidenzintervall für $\theta \Leftrightarrow$ Die Hypothese $H_0 : \theta = \theta_0$ gegen $H_A : \theta \neq \theta_0$ wird zum Level α nicht verworfen.
(Sofern Test und Konfidenzintervall auf der gleichen Statistik beruhen.)
- Im GLMM: Für Parameter θ in β oder \mathbf{D} , betrachte ein Gitter von θ_0 -Werten um den Schätzer $\hat{\theta}$.
- Konstruiere das Konfidenzintervall für θ so, dass alle Werte auf dem Gitter enthalten sind, für die ein LQT mit Referenzverteilung χ_1^2 H_0 nicht ablehnt.
- Dieser Ansatz funktioniert auch bei nicht-symmetrischen Verteilungen von Schätzern, jedoch nicht bei Parametern θ_0 in \mathbf{D} nahe des Randes des Parameterraums.
- Der Ansatz ist in `lme4` implementiert in der Funktion `confint`. Alternativen in `confint` sind Wald- (für β) und Bootstrap-basierte Konfidenzintervalle.