

4.6 Bayesianisches lineares Modell

Modell:

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

wobei $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\boldsymbol{\beta} \in \mathbb{R}^p$, $\boldsymbol{\varepsilon} \in \mathbb{R}^n$

Annahmen und Notation:

$$\begin{aligned} p &= \text{rang}(\mathbf{X}) \\ \boldsymbol{\varepsilon} &= (\varepsilon_1, \dots, \varepsilon_n)^\top, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2) \end{aligned}$$

Bayesianisch:

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

Likelihood:

$$\begin{aligned} f(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2) &\propto |\sigma^2 \mathbf{I}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \\ &= (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \end{aligned}$$

4.6.1 Nichtinformative Prioriverteilung

Die **nichtinformative Priori**

$$p(\boldsymbol{\beta}, \sigma^2) \propto (\sigma^2)^{-1}$$

ist insbesondere im Fall $p \ll n$ nützlich. Für die **gemeinsame Posteriori** folgt:

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-(\frac{n}{2}+1)} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right).$$

Sei

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \\ \hat{\mathbf{y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{y} = \mathbf{H}\mathbf{y}, \\ \hat{\boldsymbol{\varepsilon}} &= (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{y} - \hat{\mathbf{y}}. \end{aligned}$$

Aus der Theorie linearer Modelle ist bekannt, dass

$$\mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = 0, \quad \hat{\mathbf{y}}^\top \hat{\boldsymbol{\varepsilon}} = 0.$$

Daraus ergeben sich folgende Umformungen:

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) &= [(\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})]^\top [(\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})] \\ &= \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} + (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + 2(\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top \hat{\boldsymbol{\varepsilon}} \\ &= \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} + (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top (\hat{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) \\ \hat{\mathbf{y}} &\stackrel{=}{=} \mathbf{X}\hat{\boldsymbol{\beta}} \quad \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}), \end{aligned}$$

so dass sich die Posteriori schreiben lässt als

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-(\frac{n}{2}+1)} \exp\left(-\frac{1}{2\sigma^2} \left(\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right)\right).$$

Die **bedingte Posteriori** von $\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}$ ist

$$p(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{1}{2\sigma^2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right),$$

da $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ nicht von $\boldsymbol{\beta}$ abhängt. Man identifiziert obigen Ausdruck als Kern einer multivariaten Normalverteilung, genauer ist

$$p(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \sim \text{MVN}(\hat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}),$$

ein bekanntes Resultat aus der Theorie linearer Modelle.

Die **marginale Posteriori** von σ^2 erhält man über Herausintegrieren von $\boldsymbol{\beta}$ bzw. einfacher über den Satz von Bayes

$$f(\sigma^2 | \mathbf{y}, \mathbf{X}) = \frac{f(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X})}{f(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X})}.$$

Die Normalisierungskonstante für die bedingte Posteriori von $\boldsymbol{\beta}$ ist $\sigma^{-p/2}$, also

$$\begin{aligned} f(\sigma^2 | \mathbf{y}, \mathbf{X}) &\propto (\sigma^2)^{-(\frac{n}{2}+1)} \sigma^{\frac{p}{2}} \exp\left(-\frac{1}{2\sigma^2} \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}\right) \\ &= (\sigma^2)^{-(\frac{n-p}{2}+1)} \exp\left(-\frac{1}{2\sigma^2} \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}\right). \end{aligned}$$

Dies ist der Kern einer $\text{inv-}\chi^2(n-p, \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-p})$ bzw. $\text{IG}(\frac{n-p}{2}, \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{2})$ -Verteilung. Es gilt:

$$\mathbb{E}[\sigma^2 | \mathbf{y}, \mathbf{X}] = \frac{n-p}{n-p-2} \cdot \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-p} = \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-p-2}.$$

Algorithmus 9 : Direkte Simulation von $\boldsymbol{\beta}$ und σ^2 im bayesianischen linearen Modell

Für $t = 1, \dots, T$:

1. Ziehe $(\sigma^2)^{(t)}$ aus $f(\sigma^2 | \mathbf{y}, \mathbf{X})$, d.h. aus $\text{inv-}\chi^2(n-p, \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-p})$.
2. Ziehe $\boldsymbol{\beta}^{(t)}$ aus $f(\boldsymbol{\beta} | (\sigma^2)^{(t)}, \mathbf{y}, \mathbf{X})$, d.h. aus $\text{MVN}(\hat{\boldsymbol{\beta}}, (\sigma^2)^{(t)} (\mathbf{X}^\top \mathbf{X})^{-1})$, wobei $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

Eine Alternative zur direkten Simulation besteht in der Verwendung von Gibbs-Sampling, indem zusätzlich zur vollständig bedingten Dichte von $\boldsymbol{\beta}$ die vollständig bedingte Dichte von σ^2 ,

$$f(\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X}) \propto p(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-(\frac{n}{2}+1)} \exp\left(-\frac{1}{2\sigma^2} \left(\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right)\right),$$

zur Simulation verwendet wird. Dies ist für festes β der Kern einer skalierten $\text{inv-}\chi^2\left(n, \frac{\hat{\epsilon}^\top \hat{\epsilon} + (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)}{n}\right)$ -Verteilung. Damit lässt sich auch die **marginale Posteriori** von β herleiten:

$$\begin{aligned} f(\beta | \mathbf{y}, \mathbf{X}) &= \frac{f(\beta, \sigma^2 | \mathbf{y}, \mathbf{X})}{f(\sigma^2 | \mathbf{y}, \mathbf{X}, \beta)} = \frac{f(\beta | \sigma^2, \mathbf{y}, \mathbf{X}) \cdot f(\sigma^2 | \mathbf{y}, \mathbf{X})}{f(\sigma^2 | \beta, \mathbf{y}, \mathbf{X})} \\ &\propto \frac{\exp\left(-\frac{1}{2\sigma^2} (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)\right)}{\left[\frac{\hat{\epsilon}^\top \hat{\epsilon} + (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)}{n}\right]^{n/2} \exp\left(-\frac{1}{2\sigma^2} [\hat{\epsilon}^\top \hat{\epsilon} + (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)]\right)}. \end{aligned}$$

Damit:

$$f(\beta | \mathbf{y}, \mathbf{X}) \propto \left[\hat{\epsilon}^\top \hat{\epsilon} + (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)\right]^{-n/2}.$$

Setzt man

$$\hat{\sigma}_\epsilon^2 = \frac{\hat{\epsilon}^\top \hat{\epsilon}}{n-p} \Leftrightarrow \hat{\epsilon}^\top \hat{\epsilon} = (n-p)\hat{\sigma}_\epsilon^2,$$

so ist

$$\begin{aligned} f(\beta | \mathbf{y}, \mathbf{X}) &\propto \left[(n-p)\hat{\sigma}_\epsilon^2 + (\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)\right]^{-n/2} \\ &= \left((n-p)\hat{\sigma}_\epsilon^2 \cdot \left[1 + \frac{(\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)}{(n-p)\hat{\sigma}_\epsilon^2}\right]\right)^{-\frac{n}{2}} \\ &\propto \left[1 + \frac{(\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta} - \beta)}{(n-p)\hat{\sigma}_\epsilon^2}\right]^{-\frac{(n-p)+p}{2}}. \end{aligned}$$

Dies entspricht dem Kern einer multivariaten t-Verteilung mit $n-p$ Freiheitsgraden, Lokationsparameter $\hat{\beta}$ und Skalenparameter $\hat{\sigma}_\epsilon^2 (\mathbf{X}^\top \mathbf{X})^{-1}$, also

$$\beta | \mathbf{y}, \mathbf{X} \sim \text{mv-t}_{n-p}\left(\hat{\beta}, \hat{\sigma}_\epsilon^2 (\mathbf{X}^\top \mathbf{X})^{-1}\right).$$

Abschließend betrachten wir noch die **prädiktive Verteilung** für $\tilde{\mathbf{y}} | \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}$. Seien

- m die Anzahl neuer Beobachtungen,
- $\tilde{\mathbf{X}}$ neue Beobachtungen von Regressoren der Dimension $m \times p$,
- $\tilde{\mathbf{y}}$ der Vektor der Prognosen der Dimension $m \times 1$.

Zur Simulation können wir Algorithmus 9 wie folgt erweitern:

Algorithmus 10 : Direkte Simulation der prädiktiven Verteilung im bayesianischen linearen Modell

Für $t = 1, \dots, T$:

1. Ziehe $(\sigma^2)^{(t)}$ aus $f(\sigma^2 | \mathbf{y}, \mathbf{X})$, d.h. aus $\text{inv-}\chi^2\left(n - p, \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n - p}\right)$.
 2. Ziehe $\boldsymbol{\beta}^{(t)}$ aus $f(\boldsymbol{\beta} | (\sigma^2)^{(t)}, \mathbf{y}, \mathbf{X})$, d.h. aus $\text{MVN}\left(\hat{\boldsymbol{\beta}}, (\sigma^2)^{(t)} (\mathbf{X}^\top \mathbf{X})^{-1}\right)$, wobei $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.
 3. Für $i = 1, \dots, m$:
 Ziehe $\tilde{\mathbf{y}}_i \sim \text{MVN}\left(\tilde{\mathbf{x}}_i^\top \boldsymbol{\beta}^{(t)}, (\sigma^2)^{(t)}\right)$, wobei $\tilde{\mathbf{x}}_i^\top$ die i -te Zeile von $\tilde{\mathbf{X}}$ bezeichnet.
-

Es ist sogar eine analytische Berechnung möglich:

$$f(\tilde{\mathbf{y}} | \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \sim \text{mv-t}_{n-p} \left[\tilde{\mathbf{X}} \hat{\boldsymbol{\beta}}, \hat{\sigma}_\varepsilon^2 \left(\tilde{\mathbf{X}} (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{X}}^\top + \mathbf{I} \right) \right]$$

in Analogie zur Berechnung von Prognose und Prognoseintervallen für lineare Modelle aus frequentistischer Sicht.

4.6.2 Konjugierte Prioriverteilung

Im Falle der **konjugierten Priori**

$$\sigma^2 \sim \text{inv-}\chi^2(\kappa_0, \sigma_0^2), \quad \boldsymbol{\beta} | \sigma^2 \sim \text{MVN}(\boldsymbol{\beta}_0, \sigma^2 \boldsymbol{\Sigma}_0)$$

bzw.

$$\boldsymbol{\beta}, \sigma^2 \sim \text{MVN-inv-}\chi^2(\boldsymbol{\beta}_0, \sigma_0^2 \boldsymbol{\Sigma}_0; \kappa_0, \sigma_0^2)$$

ergibt sich die **gemeinsame Posteriori**

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \text{inv-}\chi^2(\kappa_n, \sigma_n^2), \quad \boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim \text{MVN}(\boldsymbol{\beta}_n, \sigma^2 \boldsymbol{\Sigma}_n)$$

bzw.

$$\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X} \sim \text{MVN-inv-}\chi^2(\boldsymbol{\beta}_n, \sigma_n^2 \boldsymbol{\Sigma}_n; \kappa_n, \sigma_n^2),$$

wobei

$$\begin{aligned} \boldsymbol{\beta}_n &= (\boldsymbol{\Sigma}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} (\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^\top \mathbf{y}), \\ \boldsymbol{\Sigma}_n &= (\boldsymbol{\Sigma}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}, \\ \kappa_n &= \kappa_0 + n, \\ \sigma_n^2 &= (\boldsymbol{\beta}_0^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^\top \boldsymbol{\Sigma}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}^\top \mathbf{y} + \kappa_0 \sigma_0^2) / (\kappa_0 + n). \end{aligned}$$

Als **bedingte Posteriori** von $\boldsymbol{\beta}$ ergibt sich

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim \text{MVN}(\boldsymbol{\beta}_n, \sigma^2 \boldsymbol{\Sigma}_n),$$

als **marginale Posteriori** von σ^2

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \text{inv-}\chi^2(\kappa_n, \sigma_n^2).$$

4.6.3 Spezialfälle und Erweiterungen

1. Ridge-Regression

Hinweis: Es ist im Allgemeinen sinnvoll, die „echten“ Kovariablen (ohne Intercept) zu standardisieren, um Unterschiede in der Skala zu beseitigen. Ferner geht man zum zentrierten Response über, so dass der Intercept entfällt. Man erhält \mathbf{X}^* , \mathbf{y}^* . Betrachte nun

$$\mathbf{y}^* = \mathbf{X}^* \boldsymbol{\beta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}).$$

Der Ridge-Schätzer ist durch

$$[(\mathbf{X}^*)^\top \mathbf{X}^* + \lambda \mathbf{I}]^{-1} (\mathbf{X}^*)^\top \mathbf{y}^*$$

mit $\lambda > 0$ gegeben. Dieser lässt sich wie folgt bayesianisch interpretieren: Sei

$$p(\boldsymbol{\beta}^*) \sim N(0, \tau^2 \mathbf{I}),$$

d.h. die Komponenten von $\boldsymbol{\beta}^*$ sind a priori unkorreliert (also wegen der Normalverteilung auch unabhängig). Dann ist die bedingte Posteriori $f(\boldsymbol{\beta}^* | \mathbf{y}^*, \mathbf{X}^*, \sigma^2, \tau^2)$ gleich

$$\text{MVN} \left(\left[(\mathbf{X}^*)^\top \mathbf{X}^* + \frac{\sigma^2}{\tau^2} \mathbf{I} \right]^{-1} (\mathbf{X}^*)^\top \mathbf{y}^*, \sigma^2 \left((\mathbf{X}^*)^\top \mathbf{X}^* + \frac{\sigma^2}{\tau^2} \mathbf{I} \right)^{-1} \right).$$

Der Parameter λ entspricht dabei dem Quotienten σ^2/τ^2 .

2. Ungleiche Varianzen der Störvariablen / abhängige Störvariablen

Allgemein:

$$\begin{aligned} \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Sigma}_\varepsilon) \\ \mathbf{y} &\sim \text{MVN}(\mathbf{X} \boldsymbol{\beta}, \boldsymbol{\Sigma}_\mathbf{y}), \quad \boldsymbol{\Sigma}_\mathbf{y} = \boldsymbol{\Sigma}_\varepsilon \end{aligned}$$

Problem: Spezifikation der Prioriverteilung für $\boldsymbol{\Sigma}_\varepsilon$.

Mögliche Auswege sind:

(a) Parametrisiere

$$\boldsymbol{\Sigma}_\mathbf{y} = \sigma^2 \mathbf{Q}_\mathbf{y}$$

mit $\mathbf{Q}_\mathbf{y}$ bekannt und $p(\boldsymbol{\beta}, \sigma^2) \propto (\sigma^2)^{-1}$. Dieser Fall ist auf das Modell aus Abschnitt 4.6.1 reduzierbar, indem man das Modell

$$\underbrace{\mathbf{Q}^{-1/2} \mathbf{y}}_{\mathbf{y}^*} = \underbrace{\mathbf{Q}^{-1/2} \mathbf{X}}_{\mathbf{X}^*} \boldsymbol{\beta} + \underbrace{\mathbf{Q}^{-1/2} \boldsymbol{\varepsilon}}_{\boldsymbol{\varepsilon}^*}$$

betrachtet. Man erhält dann wieder ein homoskedastisches Regressionsmodell in den Größen \mathbf{y}^* , \mathbf{X}^* , $\boldsymbol{\varepsilon}^*$.

(b) Gewichtete Regression:

$$\boldsymbol{\Sigma}_\mathbf{y} = \text{diag}(\sigma^2 w_i^{-1})_{1 \leq i \leq n}$$

lässt sich als Spezialfall von (a) auffassen.

(c) Korrelationen: Schreibe

$$\Sigma_{\mathbf{y}} = \mathbf{SRS} \text{ mit } \mathbf{S} = \text{diag}(\sigma_1, \dots, \sigma_p)$$

mit beispielweise

$$p(\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2) = \prod_{j=1}^p p(\sigma_j^2) \quad \text{und} \quad p(\sigma_j^2) \sim \text{inv-}\chi^2(\nu_j, \sigma_{0j}^2).$$

Es bleibt die Spezifikation der Korrelationsmatrix. Priori-Spezifikationen müssen insbesondere positive Definitheit gewährleisten. Eine einfache Variante besteht in der Verwendung von (positiver) „Äqui-Korrelation“, was zum Beispiel bei Clusterdaten eine vernünftige Annahme darstellt:

$$\mathbf{R} = \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix},$$

wobei $\rho \sim U(0, 1)$ eine positive Korrelation erzwingt. Bei Messwiederholungen greift man oft auf eine autoregressive Kovarianzstruktur zurück. Für in 1. Ordnung autokorrelierte Residuen

$$\varepsilon_t = \rho\varepsilon_{t-1} + Z_t, \quad Z_t \sim N(0, \sigma^2),$$

erhält man

$$\mathbf{R} = \begin{bmatrix} 1 & \rho & \rho^2 & \cdots & \rho^{p-1} \\ \rho & 1 & \rho & \cdots & \rho^{p-2} \\ \rho^2 & \rho & 1 & \cdots & \rho^{p-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^{p-1} & \rho^{p-2} & \rho^{p-3} & \cdots & 1 \end{bmatrix}.$$

Andere Zerlegungen basieren auf der Cholesky- oder Spektralzerlegung und sind relativ komplex.

(d) Übergang zu Modellen mit Zufallseffekten: Modelle mit Zufallseffekten (linear mixed models, generalized linear mixed models) führen zu strukturierten, meist parametersparsamen Kovarianzmatrizen.

Aber: Die Modellgleichung ändert sich und $\Sigma_{\mathbf{y}} \neq \Sigma_{\varepsilon}$, d.h. man kommt in eine andere Modellklasse.

4.7 Bayesianisches generalisiertes lineares Modell

Struktur von GLMs: Der Response folgt einer Verteilung aus einer einfachen Exponentialfamilie ($i = 1, \dots, n$)

$$f(y_i|\theta_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i}\right) \cdot c(y_i, \phi_i) \quad (4.1)$$

oder

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

wobei in vielen Fällen $\phi_i \equiv \phi$ (Bernoulli-, Poissonverteilung). Es ist

$$\mu_i = \mathbb{E}[y_i|\theta_i] = b'(\theta_i), \quad \text{Var}(y_i|\theta_i) = b''(\theta_i)\phi_i$$

und θ_i der kanonische Parameter. Mit einer Linkfunktion g bzw. Responsefunktion $h = g^{-1}$ gelte

$$g(\mu_i) = \eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}. \quad (4.2)$$

Beispiel 4.3 (Logit-Modell). Mit $\mu_i = \mathbb{P}(y_i = 1)$ ist

$$\begin{aligned} f(y_i|\mu_i) &= \mu_i^{y_i} (1 - \mu_i)^{1-y_i} \\ &= \exp(y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)) \\ &= \exp\left(\underbrace{y_i \log\left(\frac{\mu_i}{1 - \mu_i}\right)}_{\theta_i} + \log(1 - \mu_i)\right) \end{aligned}$$

mit

$$\theta_i = \log\left(\frac{\mu_i}{1 - \mu_i}\right) \Leftrightarrow \mu_i = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)}.$$

Dies entspricht (4.1) mit $\phi_i = 1$, $c(y_i, \phi_i) = 1$,

$$b(\theta_i) = -\log\left(1 - \frac{\exp(\theta_i)}{1 + \exp(\theta_i)}\right) = -\log\left(\frac{1}{1 + \exp(\theta_i)}\right) = \log(1 + \exp(\theta_i))$$

und

$$b'(\theta_i) = \frac{1}{1 + \exp(\theta_i)} \cdot \exp(\theta_i) = \mu_i.$$

Als Prioriverteilung für $\boldsymbol{\beta}$ in (4.2) eignet sich

$$\boldsymbol{\beta} \sim \text{MVN}(\boldsymbol{\beta}_0, \mathbf{B}_0)$$

mit $\mathbf{B}_0 > 0$ (vgl. Abschnitt 4.5.3 zur multivariaten Normalverteilung bei bekannter Kovarianzmatrix). $\boldsymbol{\beta}$ beeinflusst $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top$ über den Prädiktor $\boldsymbol{\mu}(\boldsymbol{\beta}) = h(\mathbf{X}\boldsymbol{\beta})$, wobei die Auswertung komponentenweise zu verstehen ist. Für $\mathbf{B}_0^{-1} \rightarrow \mathbf{0}$ erhält man eine nichtinformative Priori.

Über die Darstellung (4.1) als Exponentialfamilie mit kanonischen Parametern erhält man als Posterioriverteilung

$$\begin{aligned} f(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) &\propto \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \cdot \prod_{i=1}^n \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i}\right) \\ &= \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \exp\left(\sum_{i=1}^n \frac{y_i\theta_i - b(\theta_i)}{\phi_i}\right) \\ &= \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0) + \sum_{i=1}^n \frac{y_i\theta_i - b(\theta_i)}{\phi_i}\right). \end{aligned}$$

Beispiel 4.4 (Logit-Modell). *Im Falle des Logit-Modells erhält man als Posteriori*

$$\begin{aligned} f(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) &\propto \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \times \prod_{i=1}^n \mu_i(\boldsymbol{\beta})^{y_i} (1 - \mu_i(\boldsymbol{\beta}))^{1-y_i} \\ &= \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \times \prod_{i=1}^n h(\mathbf{x}_i^\top \boldsymbol{\beta})^{y_i} (1 - h(\mathbf{x}_i^\top \boldsymbol{\beta}))^{1-y_i} \\ &= \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \times \prod_{i=1}^n \left(\frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}\right)^{y_i} \left(\frac{1}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}\right)^{1-y_i}. \end{aligned}$$

Problem: Der Posteriori-Kern entspricht keinem Kern einer bekannten Verteilung. Die Posteriori-Verteilung ist demnach nicht analytisch zugänglich. Mögliche Auswege sind

1. Approximation oder
2. Exploration der Posteriori durch Generierung von Samples aus der Posteriori.

Wir betrachten im Folgenden Lösung 2. Hier gibt es mehrere Möglichkeiten; sehr etabliert ist ein Vorschlag von Gamerman (1997)¹, eine Variante des *Metropolis-Hastings-Algorithmus*.

4.7.1 Ein MCMC-Algorithmus: Metropolis-Hastings

Zunächst folgt eine Darstellung des Grundproblems, ohne näher auf die zugrundeliegende mathematische Theorie einzugehen. Bekannt ist, dass für $X_i \stackrel{i.i.d.}{\sim} \pi$, $i = 1, \dots, n$, wobei π eine Verteilung bezeichnet, interessierende Kennzeichen dieser Verteilung (Momente, Dichte etc.) — Existenz vorausgesetzt — durch Simulation von Zufallszahlen gemäß π als Monte-Carlo-Schätzung gewonnen werden können, z.B.

$$\widehat{\mathbb{E}[X]} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Dies ist wenig „spannend“, da, wenn π bekannt ist, in der Regel auch der Erwartungswert zugänglich ist. Angenommen jedoch, man betrachtet eine (nichtlineare) Funktion von X , zum Beispiel $g(X) = X^2$. Dann ist möglicherweise die Dichte der Transformation $g(X)$ noch analytisch bestimmbar, aber der Erwartungswert komplex zu berechnen. Im Fall, dass X mehrdimensional ist, kann die analytische Bestimmung derartiger Kenngrößen analytisch unmöglich und bei höherer Dimension mittels numerischer Integration zu instabil sein.

Unter geeigneten Voraussetzungen lässt sich obige Monte-Carlo Schätzung erweitern zu

$$\widehat{\mathbb{E}[g(X)]} = \frac{1}{n} \sum_{i=1}^n g(x_i).$$

(Dies ist ein allgemeines Prinzip, also nicht notwendigerweise bayesianisch, solange es sich bei π nicht zum Beispiel um eine Posterioriverteilung handelt.) Es sei allerdings bemerkt,

¹Gamerman (1997): *Sampling from the posteriori distribution in generalized linear models*. Statistics and Computing **7**, pp. 57-68.

dass dieses Vorgehen im Vergleich zur exakten Lösung mit einem Monte-Carlo Fehler behaftet ist. Wesentliche Voraussetzung ist zudem, dass Zufallszahlen aus π gezogen werden können. Verfahren zur Generierung von i.i.d. Zufallszahlen sind beispielsweise das Inversionsverfahren, Rejection Sampling oder Importance Sampling (vgl. Vorlesung Computerintensive Methoden). Gerade bei höherer Dimension sind diese jedoch zum Teil nicht oder nur sehr kompliziert anwendbar.

Eine Alternative stellen *Markov Chain Monte Carlo (MCMC)*-Verfahren dar. Ziel ist die Generierung einer Markov-Kette (X_0, \dots, X_n) von (abhängigen!) Zufallszahlen, deren Verteilung gegen die interessierende Verteilung konvergiert, d.h. π ist die stationäre oder invariante Verteilung der Markov-Kette. Der *Ergodensatz* erlaubt dann Schätzungen der Form

$$\widehat{\mathbb{E}[X]} = \frac{1}{n - \text{burnin}} \sum_{i=\text{burnin}+1}^n x_i \quad \text{bzw.} \quad \widehat{\mathbb{E}[g(X)]} = \frac{1}{n - \text{burnin}} \sum_{i=\text{burnin}+1}^n g(x_i),$$

wobei $x_0, \dots, x_{\text{burnin}}$ Werte am Anfang der Sequenz bezeichnen, bevor sich die Kette in der stationären Verteilung befindet, und die deshalb „weggeworfen“ werden.

Praktische Umsetzung: Starte mit einem Startwert x_0 und ziehe dann für $i = 1, \dots, n$ Werte $X_i \sim P(\cdot | X_{i-1})$, wobei P den Markov-Übergangskern bezeichnet, der nur vom aktuellen Zustand der Kette abhängt. An ihn bzw. die Markov-Kette werden die folgenden Anforderungen gestellt:

1. Die Markov-Kette ist homogen.
2. Die Markov-Kette ist irreduzibel.
3. Die Markov-Kette ist aperiodisch.
4. Die Markov-Kette ist positiv rekurrent.

Wir betrachten hier Markov-Ketten in diskreter Zeit bei diskretem oder stetigem Zustandsraum, gewöhnlich eine Teilmenge des \mathbb{R}^p . Für allgemeine Zustandsräume ist „mehr Technik“ erforderlich, aber keine „neuen Ideen“. Für den hier betrachteten Fall ist die Zielverteilung π immer gegeben.

Univariater Metropolis-Hastings

Wir beschreiben nun den *Metropolis-Hastings-Algorithmus* (kurz: *MH*) zur Generierung einer wie oben beschriebenen Markov-Kette für den univariaten Fall; dieser Algorithmus enthält den Gibbs-Sampler als Spezialfall.

Sei π die Dichte der Zielverteilung, aus der wir simulieren möchten, und q eine geeignete Vorschlagsdichte, aus der neue Zustände der Kette generiert werden, d.h.

$$X_i \sim q(\cdot | x_{i-1}),$$

zum Beispiel $q_{x_i|x_{i-1}} = N(x_{i-1}, 1)$ oder $q_{x_i|x_{i-1}} = U(x_{i-1} - c, x_{i-1} + c)$. Die Vorschläge werden nicht immer, sondern nur mit einer gewissen *Akzeptanzwahrscheinlichkeit* $\alpha(x_{i-1}, x_i)$

akzeptiert. Für den MH-Algorithmus hat diese die Gestalt

$$\alpha(x_{i-1}, x_i) = \min \left(1, \frac{\pi(x_i) \cdot q(x_{i-1}|x_i)}{\pi(x_{i-1}) \cdot q(x_i|x_{i-1})} \right).$$

Wird x_i nicht akzeptiert, so setzt man $x_i \leftarrow x_{i-1}$, d.h. der alte Zustand wird beibehalten.

Ein wesentlicher Vorteil dieses Verfahrens besteht darin, dass sich die (meist unbekannte) Normalisierungskonstante von π herauskürzt, d.h. der MH-Algorithmus kann auch (bzw. gerade) für diese Fälle angewendet werden. Die Konstruktion von α gewährleistet, dass die Bedingungen 1. bis 4. eingehalten werden.

Für $q(x_{i-1}|x_i) = q(x_i|x_{i-1})$ reduziert sich der MH-Algorithmus auf den *Metropolis-Algorithmus* mit

$$\alpha(x_{i-1}, x_i) = \min \left(1, \frac{\pi(x_i)}{\pi(x_{i-1})} \right),$$

d.h. wenn die Zieldichte an der Stelle x_i größer als an x_{i-1} ist, wird der neue Vorschlag stets akzeptiert, andernfalls nur im Verhältnis $\pi(x_i)/\pi(x_{i-1})$. Setzt man die Akzeptanzwahrscheinlichkeit konstant gleich eins, erhält man den Gibbs-Sampler.

Der MH-Algorithmus akzeptiert tendenziell neue Werte in Bereichen mit hoher Dichte (relevante Bereiche). Die Akzeptanzwahrscheinlichkeit sollte nicht zu gering sein, um regelmäßige Zustandsänderungen in der Kette zu erhalten. Sie sollte allerdings auch nicht zu hoch sein, d.h. die Varianz der Vorschlagsverteilung sollte nicht zu niedrig sein, damit der Träger von π ausreichend gut exploriert wird.

Algorithmus 11 : Univariater Metropolis-Hastings-Algorithmus

Setze Startwert X_0 .

Für $i = 1, \dots, n$:

1. Ziehe X_i aus $q(\cdot|x_{i-1})$.
2. Ziehe $U \sim U(0, 1)$; akzeptiere, wenn

$$U \leq \alpha(x_{i-1}, x_i),$$

ansonsten setze $x_i \leftarrow x_{i-1}$.

Multivariater Metropolis-Hastings

Die Verallgemeinerung auf den multivariaten Fall ist im Prinzip einfach, zum Beispiel mit

$$q(\mathbf{x}_i|\mathbf{x}_{i-1}) = \text{MVN}(\mathbf{x}_{i-1}, \mathbf{\Sigma}).$$

Ein Problem stellt hier die Wahl der „Tuning-Matrix“ $\mathbf{\Sigma}$ dar, die die Akzeptanzwahrscheinlichkeit steuert. Meist ist $\mathbf{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$; man startet mehrere Läufe und berechnet die Akzeptanzrate. Die Varianzen der Vorschlagsdichte werden dann solange variiert, bis „angemessene“ Akzeptanzraten erreicht werden.

Im Fall bayesianischer GLMs existiert eine Variante, die automatisch brauchbare Vorschlagsdichten berechnet, wie im folgenden Abschnitt beschrieben wird.

4.7.2 Metropolis-Hastings mit IWLS-Vorschlagsdichte

Aus der Vorlesung Generalisierte Regression ist das *Fisher-Scoring* bekannt:

Beispiel 4.5 (Fisher-Scoring beim Logit-Modell). Die Scorefunktion im Logit-Modell (bei kanonischer Linkfunktion) hat die Form

$$\mathbf{s}(\boldsymbol{\beta}) = \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})).$$

Die Fisher-Information ist

$$\mathbf{F}(\boldsymbol{\beta}) = \mathbf{X}^\top \text{diag}(\mu_i(\boldsymbol{\beta})(1 - \mu_i(\boldsymbol{\beta}))) \mathbf{X}.$$

Bezeichnet $\hat{\boldsymbol{\beta}}$ den ML-Schätzer, so ist

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = \left[\mathbf{X}^\top \text{diag}(\mu_i(\hat{\boldsymbol{\beta}})(1 - \mu_i(\hat{\boldsymbol{\beta}}))) \mathbf{X} \right]^{-1}.$$

Der Fisher-Scoring Algorithmus hat dann die Form

$$\begin{aligned} \hat{\boldsymbol{\beta}}^{(k+1)} &= \hat{\boldsymbol{\beta}}^{(k)} + \mathbf{F}^{-1}(\hat{\boldsymbol{\beta}}^{(k)}) \mathbf{s}(\hat{\boldsymbol{\beta}}^{(k)}) \\ &= \hat{\boldsymbol{\beta}}^{(k)} + \left[\mathbf{X}^\top \text{diag}(\mu_i(\hat{\boldsymbol{\beta}}^{(k)})(1 - \mu_i(\hat{\boldsymbol{\beta}}^{(k)}))) \mathbf{X} \right]^{-1} \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}^{(k)})). \end{aligned}$$

Allgemein lässt sich das Fisher-Scoring wie folgt umschreiben: Definiere Pseudo-Beobachtungen $\tilde{\mathbf{y}} = (\tilde{y}_1(\boldsymbol{\beta}), \dots, \tilde{y}_n(\boldsymbol{\beta}))^\top$, wobei

$$\tilde{y}_i(\boldsymbol{\beta}) = \mathbf{x}_i^\top \boldsymbol{\beta} + D_i^{-1}(y_i - \mu_i)$$

mit

$$D_i(\boldsymbol{\beta}) = \frac{\partial h(\eta_i)}{\partial \eta_i} = \frac{\partial h(\mathbf{x}_i^\top \boldsymbol{\beta})}{\partial \mathbf{x}_i^\top \boldsymbol{\beta}} \quad \text{und} \quad \eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}.$$

Im Spezialfall des Logit-Modells ist

$$D_i(\boldsymbol{\beta}) = \mu_i(1 - \mu_i) = \mu_i(\boldsymbol{\beta})(1 - \mu_i(\boldsymbol{\beta})).$$

Fasse diese Einträge zu $\mathbf{D} = \text{diag}(D_1, \dots, D_n)$ zusammen. Definiere weiter

$$w_i(\boldsymbol{\beta}) = D_i^2(\boldsymbol{\beta}) [\text{Var}(y_i)]^{-1} \quad \text{und} \quad \mathbf{W} = \text{diag}(w_1(\boldsymbol{\beta}), \dots, w_n(\boldsymbol{\beta})).$$

Im Logit-Modell:

$$w_i(\boldsymbol{\beta}) = \frac{[\mu_i(1 - \mu_i)]^2}{\mu_i(1 - \mu_i)} = \mu_i(1 - \mu_i).$$

Dann lässt sich das Fisher-Scoring als *iterierte kleinste Quadrate-Schätzung (IWLS, iteratively (re)-weighted least squares)* schreiben:

$$\hat{\boldsymbol{\beta}}^{(k+1)} = (\mathbf{X}^\top \mathbf{W}^{(k)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(k)} \tilde{\mathbf{y}}^{(k)}.$$

Aus der Analogie von kleinster Quadrate- und Maximum-Likelihood-Schätzung im Normalverteilungsfall lässt sich dies interpretieren als

$$\tilde{\mathbf{y}}^{(k)} \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, (\mathbf{W}^{-1})^{(k)}).$$

Bayesianische Version: Kombiniere das Ganze mit der Priorverteilung $\boldsymbol{\beta} \sim \text{MVN}(\boldsymbol{\beta}_0, \mathbf{B}_0)$. Iteriere dazu:

1. Aktueller Zustand sei $\beta^{(t)}$; berechne

$$\tilde{\mathbf{y}}^{(t)} = \mathbf{X}^\top \beta^{(t)} + \mathbf{D}^{-1}(\beta^{(t)})(\mathbf{y} - \boldsymbol{\mu}(\beta^{(t)})).$$

2. Ziehe $\beta^* \sim \text{MVN}(\beta^{(t+1)}, \mathbf{C}^{(t+1)})$ mit

$$\begin{aligned} \beta^{(t+1)} &= (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \mathbf{X})^{-1} \cdot [\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \tilde{\mathbf{y}}(\beta^{(t)})], \\ \mathbf{C}^{(t+1)} &= (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \mathbf{X})^{-1}. \end{aligned}$$

3. Akzeptiere β^* mit Wahrscheinlichkeit

$$\alpha(\beta^{(t)}, \beta^*) = \min \left(1, \frac{f(\beta^* | \mathbf{X})}{f(\beta^{(t)} | \mathbf{X})} \times \frac{q(\beta^{(t)} | \beta^*)}{q(\beta^* | \beta^{(t)})} \right),$$

wobei $q(\beta^{(t)} | \beta^*)$ dem Wert der Dichte von

$$\text{MVN} \left(\left(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^*) \mathbf{X} \right)^{-1} \left(\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{W}(\beta^*) \tilde{\mathbf{y}}(\beta^*) \right), \left(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^*) \mathbf{X} \right)^{-1} \right)$$

an der Stelle $\beta^{(t)}$ entspricht und analog $q(\beta^* | \beta^{(t)})$ dem Wert der Dichte von

$$\text{MVN} \left(\left(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \mathbf{X} \right)^{-1} \left(\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \tilde{\mathbf{y}}(\beta^{(t)}) \right), \left(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{W}(\beta^{(t)}) \mathbf{X} \right)^{-1} \right)$$

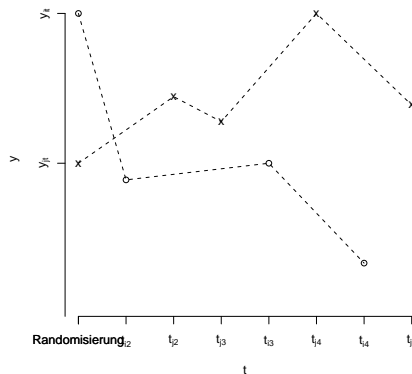
an der Stelle β^* .

4.8 Bayesianische generalisierte lineare gemischte Modelle

Der Prädiktor des GLM aus dem vorherigen Abschnitt wird hier erweitert. Im Folgenden konzentrieren wir uns auf Cluster- und Longitudinaldaten. Bei letzteren lassen sich die Daten für ein Individuum i wie folgt strukturieren:

Response	Kovariablen			Zeitpunkt der Beobachtung
y_{i1}	x_{i11}	...	x_{ip1}	t_{i1}
\vdots				\vdots
y_{iT_i}	x_{i1T_i}	...	x_{ipT_i}	t_{iT_i}

Dabei sind y_{i1}, \dots, y_{iT_i} korreliert, T_i kann variieren und die Beobachtungszeitpunkte können von Individuum zu Individuum variieren. Die Beobachtungszeitpunkte sollten jedoch nicht informativ für den Response sein. Die folgende Abbildung zeigt schematisch eine solche Datensituation, wie sie zum Beispiel bei einer kontrollierten Studie auftauchen könnte.



Diese Art der Daten stellt nicht nur wegen der Abhängigkeit der Beobachtungen eine Herausforderung dar. Häufig sind hier zum Beispiel Drop-Outs und damit fehlende Werte. Oft treten Longitudinaldaten zudem in Kombination mit Survival-Daten auf, zum Beispiel kann $y_{it_1}, \dots, y_{iT_i}$ der (mit Messfehlern behaftete) Verlauf eines Biomarkers sein. Die Frage ist dann, ob der Biomarkerverlauf prognostisch für die Überlebenszeit ist. Dies führt zu sogenannten *joint models* (siehe Modelle für Longitudinal- und Überlebenszeitdaten). GLMMs können mit diesem Datentyp gut umgehen, wenn man die sogenannte „bedingte Unabhängigkeitsannahme“ akzeptiert:

1. Erweitere den Prädiktor zu

$$\eta_{it} = \underbrace{\mathbf{x}_{it}^\top \boldsymbol{\beta}}_{\text{feste Effekte}} + \underbrace{\mathbf{z}_{it}^\top \boldsymbol{\alpha}_i}_{\text{zufällige Effekte}}$$

in der Annahme, dass

$$\boldsymbol{\alpha}_i \sim \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}).$$

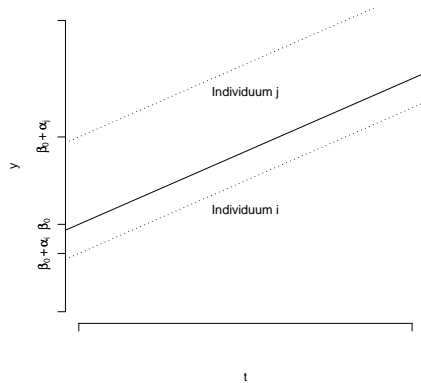
Dabei ist $\mathbf{x}_{it}^\top = (x_{i1}, \dots, x_{ipt})$ der Kovariablenvektor, und $\mathbf{z}_{it}^\top \in \mathbb{R}^{1 \times q}$ kann Kovariablen aus \mathbf{x}_{it} enthalten und zum Beispiel die Zeit t selbst.

Beispiel 4.6 (Random Intercept Modell). Sei

$$\eta_{it_i} = \beta_0 + \beta_1 t_i + \alpha_i, \quad \alpha_i \sim N(0, \sigma_\alpha^2),$$

dann haben wir für ein Individuum i :

$$\begin{pmatrix} \eta_{it_{i_1}} \\ \vdots \\ \eta_{it_{T_i}} \end{pmatrix} = \begin{pmatrix} 1 & t_{i_1} \\ \vdots & \vdots \\ 1 & t_{iT_i} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \alpha_i.$$



2. Wir treffen die bedingte Unabhängigkeitsannahme

$$y_{it} \perp\!\!\!\perp y_{it'} \mid \boldsymbol{\alpha}_i, \boldsymbol{\beta}$$

für alle $t \neq t'$. Diese erlaubt die Darstellung der gemeinsamen Verteilung von $(y_{i1}, \dots, y_{iT_i})$ als Produkt der bedingten Verteilungen

$$f(y_{i1}, \dots, y_{iT_i}) = \prod_{t=1}^{T_i} f(y_{it} \mid \boldsymbol{\alpha}_i).$$

Ohne diese bedingte Unabhängigkeitsannahme verlieren GLMMs deutlich an Attraktivität.

Das volle Setup bei n Individuen sieht wie folgt aus:

$$\mathbf{y}_i = \begin{pmatrix} y_{i1} \\ \vdots \\ y_{iT_i} \end{pmatrix}, \quad \mathbf{X}_i = \begin{pmatrix} x_{i11} & \dots & x_{ip1} \\ \vdots & & \\ x_{i1T_i} & \dots & x_{ipT_i} \end{pmatrix}, \quad \mathbf{Z}_i = \begin{pmatrix} z_{i11} & \dots & z_{iq1} \\ \vdots & & \\ z_{i1T_i} & \dots & z_{iqT_i} \end{pmatrix}, \quad \boldsymbol{\alpha} = \begin{pmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_n \end{pmatrix}$$

und

$$\boldsymbol{\eta}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\alpha}_i.$$

Zusammenfassend in Matrixnotation ergibt sich:

$$\begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_n \end{pmatrix} \boldsymbol{\beta} + \begin{pmatrix} \mathbf{Z}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{Z}_n \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_n \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1 & \mathbf{Z}_1 & \dots & \mathbf{0} \\ \mathbf{X}_2 & \mathbf{0} & \dots & \vdots \\ \vdots & \vdots & \ddots & \mathbf{0} \\ \mathbf{X}_n & \mathbf{0} & \dots & \mathbf{Z}_n \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_n \end{pmatrix}.$$

Die bayesianische Herangehensweise ist hier im Prinzip wie im GLM mit

$$\boldsymbol{\beta} \sim \text{MVN}(\boldsymbol{\beta}_0, \mathbf{B}_0) \quad \text{und} \quad \boldsymbol{\alpha} \sim \text{MVN}(\mathbf{0}, \underbrace{\text{diag}(\boldsymbol{\Sigma}, \dots, \boldsymbol{\Sigma})}_{(n \cdot q) \times (n \cdot q)})$$

bzw.

$$\begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\alpha} \end{pmatrix} \sim \text{MVN} \left(\begin{pmatrix} \boldsymbol{\beta}_0 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \underbrace{\text{diag}(\mathbf{B}_0, \boldsymbol{\Sigma}, \dots, \boldsymbol{\Sigma})}_{(p+nq) \times (p+nq)} \right).$$

Dabei bezeichnen $\text{diag}(\boldsymbol{\Sigma}, \dots, \boldsymbol{\Sigma})$ bzw. $\text{diag}(\mathbf{B}_0, \boldsymbol{\Sigma}, \dots, \boldsymbol{\Sigma})$ Blockdiagonalmatrizen.

Bemerkung.

- (i) Bei komplexeren Situationen, zum Beispiel bei nicht unabhängigen Individuen, kann $\text{diag}(\mathbf{B}_0, \boldsymbol{\Sigma}, \dots, \boldsymbol{\Sigma})$ durch eine nicht-blockweise Matrix ersetzt werden.
- (ii) GLMMs sind hochdimensional, wenn n groß ist. Spezielle Algorithmen zur Optimierung sind notwendig.

Zusätzlich wird eine (Hyper-) Prioriverteilung für $\boldsymbol{\Sigma}$ konstruiert, weil die $\boldsymbol{\alpha}_i$ unbeobachtete, latente Variablen sind, d.h. die $\boldsymbol{\alpha}_i$ sind zu behandeln wie die ε_i im linearen Modell; auch dort haben wir für die Varianz eine Priori angenommen.

Die Priori könnte zum Beispiel $\boldsymbol{\Sigma} \sim \text{inv-Wishart}_{\nu_0}(\boldsymbol{\Lambda}_0^{-1})$, also

$$p(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(\nu_0+q+1)/2} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Lambda}_0)\right),$$

mit resultierender Posteriori

$$\begin{aligned} f(\boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{\Sigma} | \mathbf{y}, \mathbf{X}) &\propto \left[\prod_{i=1}^n \prod_{j=1}^{T_i} f(y_{it_j} | \cdot) \right] \times \\ &\exp\left(-\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right) \times \\ &|\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n \boldsymbol{\alpha}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}_i\right) \times \\ &|\boldsymbol{\Sigma}|^{-(\nu_0+q+1)/2} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Lambda}_0)\right). \end{aligned}$$

Ein möglicher Algorithmus zur Simulation der Posteriori ist dann ein blockweiser Gibbs-Sampler.

- (i) *Full-Conditional für den $\boldsymbol{\beta}$ -Block:*

$$f(\boldsymbol{\beta} | \boldsymbol{\alpha}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}) \propto \left[\prod_{i=1}^n \prod_{j=1}^{T_i} f(y_{it_j} | \cdot) \right] \cdot \exp\left(-\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)\right)$$

lässt sich wie im bayesianischen GLM behandeln, wenn man zusätzlich einen Offset $\mathbf{z}_i^\top \boldsymbol{\alpha}_i$ verwendet:

$$\tilde{y}_i(\boldsymbol{\beta}^{(t-1)} | \boldsymbol{\alpha}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}^{(t-1)} + \mathbf{z}_i^\top \boldsymbol{\alpha}_i + D_i^{-1} [y_i - \mu_i(\boldsymbol{\beta}^{(t-1)}, \boldsymbol{\alpha}_i)]$$

bzw.

$$\tilde{y}_i(\boldsymbol{\beta}^{(t-1)}|\boldsymbol{\alpha}_i) - \mathbf{z}_i^\top \boldsymbol{\alpha}_i = \mathbf{x}_i^\top \boldsymbol{\beta}^{(t-1)} + D_i^{-1}[y_i - \mu_i(\boldsymbol{\beta}^{(t-1)}, \boldsymbol{\alpha}_i)]$$

Definiert man $\tilde{\tilde{y}}_i(\boldsymbol{\beta}^{(t-1)}|\boldsymbol{\alpha}_i) = \tilde{y}_i(\boldsymbol{\beta}^{(t-1)}|\boldsymbol{\alpha}_i) - \mathbf{z}_i^\top \boldsymbol{\alpha}_i$, so lässt sich der IWLS-Metropolis-Hastings-Algorithmus aus 4.7.2 zum Ziehen aus dieser Full-Conditional anwenden.

(ii) *Full-Conditional für den $\boldsymbol{\alpha}_i$ -Block:* Für $i = 1, \dots, n$ erhält man

$$f(\boldsymbol{\alpha}_i|\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X}) \propto \left[\prod_{i=1}^n \prod_{j=1}^{T_i} f(y_{it_j}|\cdot) \right] \exp\left(-\boldsymbol{\alpha}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha}_i\right).$$

Dies lässt sich wieder wie ein GLM mit Offset $\mathbf{x}_i^\top \boldsymbol{\beta}$ interpretieren und mit dem IWLS-MH-Algorithmus mit Proposal-Verteilung

$$\text{MVN}\left([\boldsymbol{\Sigma}^{-1} + \mathbf{z}_i \mathbf{W}_i(\boldsymbol{\alpha}_i^{(t-1)}) \mathbf{z}_i^\top]^{-1} \mathbf{z}_i \mathbf{W}_i(\boldsymbol{\alpha}_i^{(t-1)}) [\tilde{\tilde{y}}_i(\boldsymbol{\alpha}_i^{(t-1)}) - \mathbf{x}_i^\top \boldsymbol{\beta}], [\boldsymbol{\Sigma}^{-1} + \mathbf{z}_i \mathbf{W}_i(\boldsymbol{\alpha}_i^{(t-1)}) \mathbf{z}_i^\top]^{-1}\right)$$

behandeln.

(iii) *Full-Conditional für $\boldsymbol{\Sigma}$:*

$$f(\boldsymbol{\Sigma}|\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{y}, \mathbf{X}) \propto |\boldsymbol{\Sigma}|^{-(n+\nu_0+q+1)/2} \exp\left(-\frac{1}{2} \text{tr}\left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Lambda}_0 + \sum_{i=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_i^\top\right)\right)$$

entspricht dem Kern einer inv-Wishart $_{\nu_0+n}(\boldsymbol{\Lambda}_0 + \sum_{i=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_i^\top)$ -Verteilung (implizite Dimension $q \times q$); diese lässt sich direkt mit einem geeigneten Zufallszahlengenerator simulieren.

Insgesamt hat man also einen blockweisen Gibbs-Sampler mit

$$\begin{array}{ccccccc} 1 & + & n & + & 1 & = & n + 2 \\ \boldsymbol{\beta} & & \{\boldsymbol{\alpha}_i\}_{i=1}^n & & \boldsymbol{\Sigma} & & \end{array}$$

Blöcken. Verwendet man beim Update von $\boldsymbol{\beta}$ und $\boldsymbol{\alpha}_i$ einen Akzeptanzmechanismus, dann handelt es sich bei dem Gibbs-Sampler genauer gesagt um einen *Metropolis-Hastings-within-Gibbs*-Algorithmus, da die einzelnen Blöcke jeweils mit Metropolis-Hastings erzeugt, am Ende des Durchgangs aber die $n+2$ Blöcke noch einmal (mit Wahrscheinlichkeit 1) formal akzeptiert werden.

4.9 Hierarchische Modelle

Dieser Abschnitt behandelt ein Beispiel aus Gelman, Carlin, Stern und Rubin (2003) zum Tumorrisiko bei Ratten. Durchgeführt wurden 71 Experimente i mit jeweils n_i Ratten j :

$$y_{ij} = \begin{cases} 1, & \text{Ratte entwickelt Tumor,} \\ 0, & \text{Ratte entwickelt keinen Tumor.} \end{cases}$$

Also ist $y_i = \sum_{j=1}^{n_i} y_{ij}$ die Anzahl an Ratten in Experiment i , die einen Tumor entwickeln.

Ideen:

1. Experiment-spezifische Wahrscheinlichkeiten θ_i , $i = 1, \dots, n$, für Tumorentwicklung betrachten, potentiell zurückzuführen auf Heterogenität der Ratten, unterschiedliche experimentelle Bedingungen usw.
2. Alle θ_i stammen aus einer Population, zum Beispiel einer $\text{Beta}(\alpha, \beta)$ -Verteilung.
3. Anstatt 71 Parameter $\theta_1, \dots, \theta_{71}$ direkt aus den Daten zu schätzen, nehmen wir eine Verteilung für die θ_i an.
4. Ohne weitere Information sind $\theta_1, \dots, \theta_{71}$ als exchangeable zu betrachten, d.h. für die gemeinsame Priori von $\theta = (\theta_1, \dots, \theta_{71})$ gilt

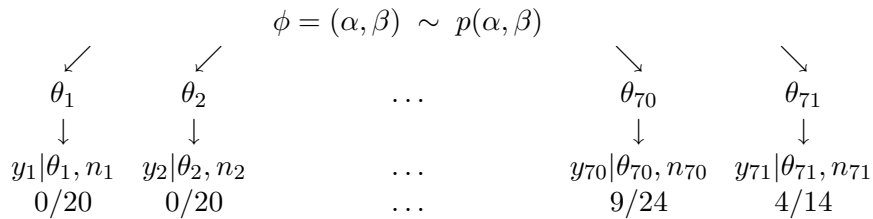
$$p(\theta|\phi) = \prod_{i=1}^{71} p(\theta_i|\phi),$$

wobei ϕ die *Hyperparameter* bezeichnet. „Averaging“ über ϕ liefert die marginale Priori

$$p(\theta) = \int \left(\prod_{i=1}^n p(\theta_i|\phi) \right) p(\phi) d\phi.$$

Da ϕ nicht bekannt ist, erhält es eine eigene (Hyper-) Prioriverteilung $p(\phi)$, im Beispiel eine Priori für α und β .

Struktur des hierarchischen Modells:



Hierarchisches Modell des Beispiels in der top-down Schreibweise:

$$\begin{array}{l}
 \text{tauchen in der Likelihood auf} \\
 \text{taucht nicht in der Likelihood auf}
 \end{array}
 \left\{ \begin{array}{l}
 y_{ij}|n_i, \theta_i \sim \text{Binomial}(n_i, \theta_i) \\
 \theta_i|\alpha, \beta \sim \text{Beta}(\alpha, \beta) \\
 (\alpha, \beta) \sim p(\alpha, \beta)
 \end{array} \right.$$

Die Posteriori für alle Parameter lautet

$$\begin{aligned}
 & f(\theta, \alpha, \beta | y_1, \dots, y_{71}) \\
 \propto & \left(\prod_{i=1}^{71} \theta_i^{y_i} (1 - \theta_i)^{n_i - y_i} \right) \left(\prod_{i=1}^{71} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta_i^{\alpha-1} (1 - \theta_i)^{\beta-1} \right) \cdot p(\alpha, \beta).
 \end{aligned}$$

Für feste α, β ist die Posteriori von $(\theta_1, \dots, \theta_{71})$ das Produkt unabhängiger Posterioris, die jeweils einer $\text{Beta}(\tilde{\alpha}_i, \tilde{\beta}_i)$ -Verteilung folgen mit $\tilde{\alpha}_i = \alpha + y_i$, $\tilde{\beta}_i = \beta + (n_i - y_i)$. Genauer:

$$f(\theta|\alpha, \beta, \{y_{ij}\}) = \prod_{i=1}^{71} \frac{\Gamma(\alpha + \beta + n_i)}{\Gamma(\alpha + y_i)\Gamma(\beta + n_i - y_i)} \cdot \theta_i^{\alpha+y_i-1} (1 - \theta_i)^{\beta+n_i-y_i-1}.$$

Die marginale Posteriori $f(\theta|\alpha, \beta, \{y_{ij}\})$ von (α, β) ergibt sich über die bereits vielfach verwandte Formel

$$f(\alpha, \beta|\{y_{ij}\}) = \frac{f(\theta, \alpha, \beta|\{y_{ij}\})}{f(\theta|\alpha, \beta, \{y_{ij}\})}.$$

Im Beispiel erhält man

$$f(\alpha, \beta|\{y_{ij}\}) \propto p(\alpha, \beta) \cdot \prod_{i=1}^{71} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + y_i)\Gamma(\beta + n_i - y_i)}{\Gamma(\alpha + \beta + n_i)}.$$

„Knackpunkt“ ist die nicht-triviale Wahl der Hyperpriori $p(\alpha, \beta)$: Oft sind keine Anhaltspunkte für die Verteilung der Hyperparameter vorhanden. Die Hyperpriori sollte in diesem Fall möglichst wenig informativ sein, allerdings gleichzeitig so, dass die Posteriori proper ist. Für das Beispiel schlagen Gelman, Carlin, Stern und Rubin (2003) vor, die Hyperparameter in den (kompletten) \mathbb{R}^2 zu transformieren, zum Beispiel durch

$$(\alpha, \beta) \mapsto \left[\text{logit} \left(\frac{\alpha}{\alpha + \beta} \right), \log(\alpha + \beta) \right],$$

wobei hier $\text{logit}(\alpha/(\alpha + \beta)) = \log(\alpha/\beta)$. Zur Interpretation sei daran erinnert, dass $\alpha/(\alpha + \beta)$ dem (Priori-) Erwartungswert einer Beta-Verteilung entspricht und $\alpha + \beta$ sich als Priori-Stichprobengröße auffassen lässt. Eine Gleichverteilung auf der transformierten Skala führt jedoch zu einer uneigentlichen Posteriori. Alternativ schlagen obige Autoren folgende Priori vor:

$$p \left(\log \left(\frac{\alpha}{\beta} \right), \log(\alpha + \beta) \right) \propto \alpha\beta(\alpha + \beta)^{-5/2}. \quad (4.3)$$

Diese ergibt sich aus einer Gleichverteilung für eine Approximation der Priori-Standardabweichung $(\alpha + \beta)^{-1/2}$, die unabhängig mit einer Gleichverteilung auf dem Priori-Erwartungswert kombiniert wird, d.h.

$$p \left(\frac{\alpha}{\alpha + \beta}, (\alpha + \beta)^{-\frac{1}{2}} \right) \propto 1.$$

Die Simulation der Posterioriverteilung $f(\theta|\alpha, \beta, \{y_{ij}\})$ ist dann wie folgt:

1. Ziehe Zufallszahlen aus $f(\log(\alpha/\beta), \log(\alpha + \beta)|\{y_{ij}\})$: Dazu wird $f(\alpha, \beta|\{y_{ij}\})$ gemäß Dichtetransformationssatz transformiert (unter anderem wird also $p(\alpha, \beta)$ durch die Hyperpriori (4.3) ersetzt) und dann auf einem feinen Gitter berechnet. $(\log(\alpha/\beta), \log(\alpha + \beta))|\{y_{ij}\}$ kann nun unter Verwendung des CDF-Samplers (Algorithmus 5 auf Seite 97) simuliert werden.
2. Transformiere die in Schritt 1 gezogenen Zufallszahlen auf die ursprüngliche (α, β) -Skala zurück.
3. Ziehe dann für $i = 1, \dots, 71$ $\theta_i|\alpha, \beta, \{y_{ij}\}$ gemäß einer $\text{Beta}(\alpha + y_i, \beta + n - y_i)$ -Verteilung.

4.10 Konvergenzdiagnostik

Schwierigkeiten bei statistischer Inferenz entstehen durch iterative Simulation:

1. Simulierte Zufallszahlen aus der Posteriori repräsentieren die Zielverteilung eventuell unzureichend (Problem der Startwertwahl, Einfluss der Startwerte auf spätere Ziehungen).
2. MCMC: Zufallszahlen sind korreliert — die Inferenz ist ungenauer, als wenn die gleiche Anzahl unabhängiger Zufallszahlen verwandt würde. Stichwort: „effektive Stichprobengröße“.

Strategien:

1. Multiple Sequenzen, die stark über den Parameterraum streuen.
2. Konvergenz-Monitoring.
3. Falls das „Mixing“ schlecht ist (der Parameterraum wird unzureichend exploriert, vorwiegend Bewegung entlang weniger lokaler Maxima der Zielverteilung usw.), sollte der Algorithmus geändert werden.

Bezüglich Punkt 2 schlagen Gelman, Carlin, Stern und Rubin (2003) für skalaren Parameter vor:

- Generiere m parallele Sequenzen der Länge $n \geq 2$ nach Entfernen der Burnin-Werte, also

$$\{\psi_{ij}\} \text{ für } i = 1, \dots, n, j = 1, \dots, m.$$

- Berechne die Varianz zwischen den Sequenzen und innerhalb jeder Sequenz: Setze

$$\bar{\psi}_{\cdot j} = \frac{1}{n} \sum_{i=1}^n \psi_{ij} \quad , \quad \bar{\psi}_{\cdot\cdot} = \frac{1}{m} \sum_{j=1}^m \bar{\psi}_{\cdot j}$$

und definiere

$$B = n \cdot \frac{1}{m-1} \sum_{j=1}^m (\bar{\psi}_{\cdot j} - \bar{\psi}_{\cdot\cdot})^2$$

sowie

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2 \quad \text{mit} \quad s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\psi_{ij} - \bar{\psi}_{\cdot j})^2.$$

Betrachte dann folgende Schätzung für die marginale Posteriori-Varianz:

$$\widehat{\text{Var}}^+(\psi|y) = \frac{n-1}{n} W + \frac{1}{n} B.$$

Die Varianz für ψ wird in der Regel überschätzt. Die Schätzung ist jedoch für $n \rightarrow \infty$ unverzerrt. W allein unterschätzt in der Regel die Varianz, aber

$$\lim_{n \rightarrow \infty} \text{E}[W] = \text{Var}(\psi|y).$$

- Führe nun ein Monitoring für die Größe

$$\hat{R} = \sqrt{\frac{\widehat{\text{Var}}^+(\psi|y)}{W}}$$

mit $R \rightarrow 1$ für $n \rightarrow \infty$ durch, d.h. wenn \hat{R} groß ist, erhält man potentiell eine Verbesserung der Inferenz, wenn weitere Simulationen durchgeführt werden.

- Also:
 - (a) Wenn \hat{R} groß: Lasse Simulation weiter laufen.
 - (b) Wenn \hat{R} „nahe“ 1: Verwende alle $m \cdot n$ Werte für Posteriori-Inferenz. Die optimale Größe von \hat{R} ist wiederum ein Tuning-Parameter, zum Beispiel $\hat{R} < 1.1$.

Erfolg ist allerdings nicht garantiert. Die Methode funktioniert gut bei approximativ normaler marginaler Posteriori, ist jedoch weniger geeignet, wenn Interesse an den extremen Quantilen besteht.

- Effektive Anzahl unabhängiger Ziehungen:

$$n_{\text{eff}} = m \cdot n \cdot \frac{\widehat{\text{Var}}^+(\psi)}{W}.$$

Dabei sollte m nicht zu klein sein, da sonst B wiederum schlecht geschätzt wird.

Obige Methode ist im R-Paket `coda` (convergence diagnostics and output analysis) implementiert.

4.11 Modellwahl und Modellkritik

Modellwahl ist ein weites Feld, auch in Likelihood-basierter Inferenz. Generelle Strategien sind schwierig zu finden, da in der Regel eine Abhängigkeit vom Kontext (Substanzwissenschaft) und Randbedingungen (zum Beispiel randomisierte Studie oder Beobachtungsstudie) oder auch dem *signal-to-noise-ratio* (kontrolliertes experimentelles Umfeld oder starke Heterogenität) besteht.

Sensitivitätsanalyse:

- Einfluss der Priori auf Ergebnisse.
- Einfluss des Modells (Likelihood) auf Güte der Vorhersagen. Welche Fälle werden schlecht durch das Modell beschrieben?

Modellwahl:

Typischerweise, zum Beispiel bei einem Regressionsmodell, erfolgt eine Auswahl der Kovariablen; „nested versus non-nested“ Situation.

Ein populärer Vorschlag zur Modellwahl ist das *DIC* (Deviance Information Criterion). Wie das AIC und BIC ist das DIC eine asymptotische Approximation und nur anwendbar, falls die Posteriori approximativ multivariat normal ist.

Allgemein ist die Devianz für Daten y und Parameter(-vektor) θ definiert als

$$D(y, \theta) = -2 \log f(y|\theta) + C(y).$$

Beurteilt man Modelle nach der Devianz, so ist ein Modell umso besser, je kleiner diese ist. Das DIC kann aus den generierten Samples der MCMC-Simulation berechnet werden. Sei $\theta_1, \dots, \theta_L$ eine generierte Sequenz. Die erwartete Devianz bezüglich der Posterioriverteilung von θ ist

$$\mathbb{E}[D(y, \theta)|y]$$

und wird durch

$$\bar{D} = \frac{1}{L} \sum_{l=1}^L D(y, \theta^l)$$

geschätzt. Dies ist ein Maß dafür, wie gut das Modell an die Daten angepasst ist (je kleiner, umso besser).

Die effektive Anzahl der Parameter in einem bayesianischen Modell ist

$$p_D = \bar{D} - D(\bar{\theta}),$$

wobei zum Beispiel

$$\bar{\theta} = \frac{1}{L} \sum_{l=1}^L \theta^l$$

die Schätzung des Posteriori-Mittelwerts von θ und $D(\bar{\theta})$ die Devianz, ausgewertet an $\bar{\theta}$, ist. Für ein lineares Modell mit Normalverteilungsannahme entspricht p_D der Anzahl unrestringierter Parameter im Modell. Das DIC ist dann definiert als

$$\text{DIC} = 2\bar{D} - D(\bar{\theta}) = \bar{D} + p_D.$$

Hinweis:

$$\begin{aligned} \bar{D}^{\text{prediction}} &= \mathbb{E}[D(y^{\text{rep}}, \bar{\theta})] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (y_i^{\text{rep}} - \mathbb{E}[y_i^{\text{rep}}|y])^2 \right], \end{aligned}$$

wobei der Erwartungswert bezüglich der a posteriori prädiktiven Verteilung zu verstehen ist. Es ist $\bar{D}^{\text{prediction}} \approx \text{DIC}$.