

Dies funktioniert (meist) unter Annahme von Fisher-Regularität. Nur in einfachen Fällen ist die Lösung analytisch zugänglich. Die numerische Lösung geschieht über Verfahren wie Newton-Raphson, Fisher-Scoring, Quasi-Newton oder über den EM-Algorithmus. Erstere drei Verfahren arbeiten mit der Hesse-Matrix der Log-Likelihood bzw. Approximationen an diese:

$$J(\theta; x) = -\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top} = \left(-\frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} \right)$$

heißt *beobachtete Informationsmatrix*. Bildet man den Erwartungswert bezüglich allen möglichen Stichproben X aus \mathcal{X} , so erhält man die *erwartete Informationsmatrix*

$$I(\theta) = \mathbb{E}_\theta[J(\theta; X)].$$

Unter Fisher-Regularität gilt (vgl. Abschnitt 2):

$$\mathbb{E}_\theta[s(\theta)] = 0 \quad \text{und} \quad \text{Cov}_\theta(s(\theta)) = \mathbb{E}_\theta[s(\theta)s(\theta)^\top] = I(\theta).$$

Beispiel 3.3 (Lineares Modell). *Betrachte*

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

- *Likelihood:*

$$L(\boldsymbol{\beta}, \sigma^2) \propto (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}\|^2\right)$$

- *Log-Likelihood:*

$$\ell(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \underbrace{\|\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}\|^2}_{\text{KQ-Kriterium}}$$

- *Score-Funktion:*

$$\begin{aligned} s_{\boldsymbol{\beta}}(\boldsymbol{\beta}, \sigma^2) &= \frac{\partial \ell(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} = \frac{1}{\sigma^2} \mathbf{Z}^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) \\ s_{\sigma^2}(\boldsymbol{\beta}, \sigma^2) &= -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \|\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}\|^2 \end{aligned}$$

Man verifiziert leicht, dass $\mathbb{E}[s_{\boldsymbol{\beta}}] = \mathbb{E}[s_{\sigma^2}] = 0$ ist. Aus den ML-Gleichungen, die sich durch Nullsetzen der Score-Funktionen ergeben, folgt:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{ML} &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y}, \\ \sigma_{ML}^2 &= \frac{1}{n} \|\mathbf{y} - \mathbf{Z}\hat{\boldsymbol{\beta}}_{ML}\|^2. \end{aligned}$$

Der ML-Schätzer für $\boldsymbol{\beta}$ entspricht also dem KQ-Schätzer. Der ML-Schätzer für σ^2 ist verzerrt, aber asymptotisch erwartungstreu. Der Restricted Maximum Likelihood (REML) Schätzer

$$\sigma_{REML}^2 = \frac{1}{n-p} \|\mathbf{y} - \mathbf{Z}\hat{\boldsymbol{\beta}}_{ML}\|^2$$

ist erwartungstreu für σ^2 . Dabei ist p die Dimension von $\boldsymbol{\beta}$.

- Informationsmatrizen:

$$\begin{aligned}
-\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} &= -\frac{\partial s_{\boldsymbol{\beta}}}{\partial \boldsymbol{\beta}^\top} = \frac{1}{\sigma^2} \mathbf{Z}^\top \mathbf{Z} = \left(\text{Cov}(\widehat{\boldsymbol{\beta}}) \right)^{-1} && \text{(von } \mathbf{y} \text{ unabhängig)} \\
-\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \sigma^2} &= \frac{1}{\sigma^4} \mathbf{Z}^\top (\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}) && \Rightarrow \mathbb{E} \left[-\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \sigma^2} \right] = 0 \\
-\frac{\partial^2 \ell}{\partial \sigma^2 \partial \sigma^2} &= -\frac{n}{2(\sigma^2)^2} + \frac{\|\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}\|^2}{(\sigma^2)^3} && \Rightarrow \mathbb{E} \left[-\frac{\partial^2 \ell}{\partial \sigma^2 \partial \sigma^2} \right] = \frac{n}{2\sigma^4}
\end{aligned}$$

Der letzte Erwartungswert folgt aus $\|\mathbf{y} - \mathbf{Z}\boldsymbol{\beta}\|^2 = \sum_{i=1}^n \varepsilon_i^2 \sim \sigma^2 \chi^2(n)$.

Beispiel 3.4 (GLM). Seien $y_i \stackrel{\text{unabh.}}{\sim} f(y_i|\mu_i)$ für $i = 1, \dots, n$ mit $\mu_i = h(\mathbf{x}_i^\top \boldsymbol{\beta})$, etwa $y_i \sim \text{Po}(\lambda_i)$ und $\lambda_i = \exp(\mathbf{x}_i^\top \boldsymbol{\beta})$ (loglineares Poisson-Modell, vgl. Übung/generalisierte Regression).

Beispiel 3.5 (GLMM für Longitudinaldaten). Sei $\mathbf{y}_i = (y_{i1}, \dots, y_{it}, \dots, y_{iT})$ mit bedingt unabhängigen Komponenten $y_{it} \sim f(y_{it}|\mu_{it})$ und $\mu_{it} = h(\mathbf{z}_i^\top \boldsymbol{\beta} + \mathbf{w}_i^\top \boldsymbol{\gamma}_i)$. Die $\boldsymbol{\gamma}_i$ sind individualspezifische Intercepts ($\mathbf{w}_i \equiv \mathbf{1}$) mit Priorverteilung $\boldsymbol{\gamma}_i \stackrel{\text{i.i.d.}}{\sim} N(0, \tau^2)$. Die Likelihood des Parameters $\theta = (\boldsymbol{\beta}, \tau^2)$ lautet

$$L(\boldsymbol{\beta}, \tau^2) = \int \prod_{i=1}^n f(y_{it}|\boldsymbol{\beta}, \tau^2, \boldsymbol{\gamma}_i) p(\boldsymbol{\gamma}_i) d\boldsymbol{\gamma}_i.$$

Lösungsansätze für die Maximierung der Likelihood: EM-Algorithmus mit REML bzw. Bayes-Inferenz.

3.2.2 Iterative numerische Verfahren zur Berechnung des ML-Schätzers

EM (Expectation-Maximization)-Algorithmus

Der EM-Algorithmus ist eine Alternative zu Newton-Raphson, Fisher-Scoring usw., vor allem in Modellen mit unvollständigen Daten oder latenten (nicht direkt beobachtbaren) Variablen oder Faktoren (vgl. Computerintensive Methoden).

Notation:

- \mathbf{x} beobachtbare („unvollständige“) Daten
- \mathbf{z} unbeobachtbare Daten/latente Variablen
- (\mathbf{x}, \mathbf{z}) vollständige Daten
- $L(\theta; \mathbf{x}) = f(\mathbf{x}|\theta)$ Likelihood der beobachtbaren Daten
- $L(\theta; \mathbf{x}, \mathbf{z}) = f(\mathbf{x}, \mathbf{z}|\theta)$ Likelihood der vollständigen Daten

Der EM-Algorithmus ist insbesondere nützlich, wenn $L(\theta; \mathbf{x})$ schwierig zu berechnen und $L(\theta; \mathbf{x}, \mathbf{z})$ leichter zu handhaben ist.

Algorithmus 1 : EM-Algorithmus

Startwert: $\theta^{(0)}$

- **E-Schritt:** Berechne

$$Q(\theta) = Q(\theta; \theta^{(0)}) = \mathbb{E}_{\mathbf{z}|\mathbf{x}}[\ell(\theta; \mathbf{x}, \mathbf{Z})|\mathbf{x}, \theta^{(0)}].$$

- **M-Schritt:** Berechne $\theta^{(1)}$, so dass $Q(\theta)$ maximiert wird:

$$\theta^{(1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta).$$

Iteriere **E/M-Schritte**: $\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(k)}$ bis zur Konvergenz.

Satz 3.3. *Unter relativ allgemeinen Annahmen gilt $\theta^{(k)} \rightarrow \hat{\theta}_{ML}$ für $k \rightarrow \infty$.*

Eigenschaften des EM-Algorithmus:

- Monotonie: $\ell(\theta^{(k+1)}; \mathbf{x}) \geq \ell(\theta^{(k)}; \mathbf{x})$.
- Langsame Konvergenz.
- Der Standardfehler des resultierenden Schätzers ist schwierig zu bestimmen, die Informationsmatrix ist nicht direkt zugänglich wie beim Fisher-Scoring.

Eine Alternative bietet die Bayes-Inferenz.

Beispiel 3.6 (Mischverteilungen). *Seien X_1, \dots, X_n i.i.d. wie $X \sim f(x|\theta)$. Betrachte die Mischverteilung*

$$f(x|\theta) = \sum_{j=1}^J \pi_j f_j(x|\theta_j) \quad \text{mit} \quad \theta = (\{\theta\}_{j=1}^J, \{\pi_j\}_{j=1}^J). \quad (3.1)$$

Dabei sind

- π_j unbekannte Mischungsanteile, $\sum_{j=1}^J \pi_j = 1$,
- $f_j(x|\theta_j)$ die j -te Mischungskomponente,
- θ_j der unbekannte Parameter(-vektor) .

Speziell: Bei einer Mischung von Normalverteilungen erhalten wir

$$f_j(x|\theta_j) \propto |\Sigma_j|^{-1/2} \exp\left(-\frac{1}{2}(x - \mu_j)^\top \Sigma_j^{-1}(x - \mu_j)\right)$$
$$X \sim \pi_1 N(\mu_1, \Sigma_1) + \pi_2 N(\mu_2, \Sigma_2) + \dots + \pi_J N(\mu_J, \Sigma_J).$$

Im univariaten Fall mit zwei Mischungskomponenten also:

$$X \sim \pi_1 N(\mu_1, \sigma_1^2) + \pi_2 N(\mu_2, \sigma_2^2).$$

Interpretation des Mischungsmodells (3.1): x_i entstammt einer von J Subpopulationen, wobei in Subpopulation j gilt:

$$X_i|j \sim f_j(x_i|\theta_j).$$

Definiere die unbeobachtete (latente) Indikatorvariable Z_i für $j = 1, \dots, J$ durch

$$Z_i = j \Leftrightarrow x_i \text{ ist aus Population } j.$$

Die Randverteilung sei $\mathbb{P}(Z_i = j) = \pi_j$, $j = 1, \dots, J$. Dann lautet die bedingte Verteilung von $x_i|Z_i$:

$$x_i|Z_i = j \sim f_j(x_i|\theta_j).$$

Die Log-Likelihood der beobachteten Daten x ist

$$\ell(\theta; x) = \sum_{i=1}^n \log \left(\sum_{j=1}^J \pi_j f_j(x_i|\theta_j) \right),$$

die der vollständigen Daten (x, z)

$$\ell(\theta; x, z) = \sum_{i=1}^n \log f(x_i, z_i|\theta) = \sum_{i=1}^n \log (f(x_i|z_i; \theta) \cdot f(z_i)) = \sum_{i=1}^n (\log f_{z_i}(x_i|\theta_{z_i}) + \log \pi_{z_i}).$$

E-Schritt:

$$\begin{aligned} Q(\theta) &= \mathbb{E}_{z|\mathbf{x}}[\ell(\theta; \mathbf{x}, \mathbf{Z})|\mathbf{x}, \theta^{(k)}] \\ &= \sum_{i=1}^n \sum_j^J p_{ij}^{(k)} \left\{ \log \pi_j - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j) \right\} \end{aligned}$$

wobei wir nur

$$p_{ij}^{(k)} = \mathbb{P}(Z_i = j|x_i, \theta^{(k)}) \stackrel{\text{Bayes}}{=} \frac{\pi_j^{(k)} f_j(x_i|\theta_j^{(k)})}{\sum_{s=1}^J \pi_s^{(k)} f_j(x_i|\theta_s^{(k)})}.$$

für $i = 1, \dots, n$, $j = 1, \dots, J$ tatsächlich in der Praxis berechnen müssen.

M-Schritt: Berechne

$$\begin{aligned} \pi_j^{(k+1)} &= \operatorname{argmax}_{\pi_j} Q(\theta) \stackrel{1.}{=} \frac{1}{n} \sum_{i=1}^n p_{ij}^{(k)} \\ \mu_j^{(k+1)} &= \operatorname{argmax}_{\mu_j} Q(\theta) \stackrel{2.}{=} \sum_{i=1}^n w_{ij}^{(k)} x_i \\ \Sigma_j^{(k+1)} &= \operatorname{argmax}_{\Sigma_j} Q(\theta) \stackrel{3.}{=} \sum_{i=1}^n w_{ij}^{(k)} (x_i - \mu_j^{(k+1)})(x_i - \mu_j^{(k+1)})^T \end{aligned}$$

mit $w_{ij}^{(k)} = \frac{p_{ij}^{(k)}}{\sum_{s=1}^J p_{is}^{(k)}}$. 1. folgt für $J = 2$ als Maximierer der binomialen Likelihood (für $J > 2$ braucht man Lagrange). 2.+3. folgt als Maximierer der gewichteten Normalverteilungslikelihood.