

Kapitel 4

Bayes-Inferenz

4.1 Überblick

„*Definition*“ *bayesianischer Inferenz*: Anpassen eines Wahrscheinlichkeitsmodells an eine Menge von Daten.

Ergebnis: Wahrscheinlichkeitsverteilung für die Parameter des Modells (und andere unbeobachtete Größen, zum Beispiel Vorhersagen für neue Beobachtungen).

Idealisierter Prozess bayesianischer Datenanalyse:

1. Stelle ein *volles* Wahrscheinlichkeitsmodell oder eine gemeinsame Wahrscheinlichkeitsverteilung für alle beobachtbaren und unbeobachtbaren Größen auf. Dabei ist Wissen über das zugrundeliegende wissenschaftliche Problem und den datengenerierenden Prozess hilfreich.
2. Berechnung der Posterioriverteilung der unbeobachtbaren Größen (Parameter, missing data, ...): bedingte Wahrscheinlichkeitsverteilung der unbeobachtbaren Größen gegeben die beobachteten Daten.
3. Modelldiagnose: Fit, Sensitivität (bezüglich der Annahmen in 1.).

Ergebnis: „kohärentes System“

4.2 Exchangeability

Exchangeability („Austauschbarkeit“) ist ein wichtiges Konzept für die statistische Modellbildung. Es geht auf de Finetti zurück.

Definition 4.1 (Finite Exchangeability). *Die Zufallsgrößen X_1, \dots, X_n sind exchangeable bezüglich des Wahrscheinlichkeitsmaßes \mathbb{P} , wenn*

$$\mathbb{P}(x_1, \dots, x_n) = \mathbb{P}(x_{\pi(1)}, \dots, x_{\pi(n)})$$

für alle Permutationen

$$\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$$

gilt. Existiert eine Dichte f zu \mathbb{P} , so gilt entsprechend:

$$f(x_1, \dots, x_n) = f(x_{\pi(1)}, \dots, x_{\pi(n)}).$$

Definition 4.2 (Infinite Exchangeability). Die unendliche Folge X_1, X_2, \dots ist exchangeable, wenn jede endliche Teilfolge exchangeable ist.

Bemerkung. Analog zu obigen Definitionen kann auch bedingte Exchangeability definiert werden, etwa im Regressionsfall für $Y_1 | \mathbf{x}_1, \dots, Y_n | \mathbf{x}_n$.

Satz 4.3 (Darstellungssatz für 0-1 Zufallsvariablen). Sei X_1, X_2, \dots eine unendliche Folge binärer Zufallsvariablen, die exchangeable sind, mit zugrundeliegendem Wahrscheinlichkeitsmaß \mathbb{P} . Dann existiert eine Verteilungsfunktion Q , so dass die gemeinsame Dichte $f(x_1, \dots, x_n)$ folgende Gestalt hat:

$$f(x_1, \dots, x_n) = \int_0^1 \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} dQ(\theta)$$

mit

$$Q(\theta) = \lim_{n \rightarrow \infty} \mathbb{P}(y_n/n \leq \theta)$$

und

$$y_n = \sum_{i=1}^n x_i, \quad \theta = \lim_{n \rightarrow \infty} y_n/n.$$

Interpretation:

1. Bedingt auf θ sind X_1, X_2, \dots unabhängige, identisch verteilte Bernoulli-Zufallsgrößen.
2. θ wird eine Verteilung zugeordnet.
3. Q ist der „Glaube“ („Belief“) über den Grenzwert der relativen Häufigkeit der Einsen.

Konventionelle Schreibweise:

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i}.$$

Satz 4.4. Wenn die benötigten Dichten existieren und X_1, X_2, \dots eine (unendliche) Folge reellwertiger Zufallsgrößen ist, dann gilt

$$f(x_1, \dots, x_n) = \int_{\Theta} \prod_{i=1}^n f(x_i | \theta) dQ(\theta).$$

Wir betrachten nun die *a posteriori prädiktive Verteilung* oder *bedingte prädiktive Verteilung* von zukünftigen (unbeobachteten) Daten x_{m+1}, \dots, x_n gegeben die beobachteten Daten x_1, \dots, x_m :

$$\begin{aligned}
 f(x_{m+1}, \dots, x_n | x_1, \dots, x_m) &= \frac{f(x_1, \dots, x_n)}{f(x_1, \dots, x_m)} && \text{(Satz von Bayes)} \\
 &= \frac{\int \prod_{i=1}^n f(x_i | \theta) dQ(\theta)}{\int \prod_{i=1}^m f(x_i | \theta) dQ(\theta)} && \text{(Darstellungssatz)} \\
 &= \frac{\int \prod_{i=1}^m f(x_i | \theta) \prod_{i=m+1}^n f(x_i | \theta) dQ(\theta)}{\int \prod_{i=1}^m f(x_i | \theta) dQ(\theta)} \\
 &= \int \prod_{i=m+1}^n f(x_i | \theta) \cdot \frac{\prod_{i=1}^m f(x_i | \theta) dQ(\theta)}{\int \prod_{i=1}^m f(x_i | \theta) dQ(\theta)}.
 \end{aligned}$$

Dabei ist

$$\frac{\prod_{i=1}^m f(x_i | \theta) dQ(\theta)}{\int \prod_{i=1}^m f(x_i | \theta) dQ(\theta)} = dQ(\theta | x_1, \dots, x_m)$$

die Posterioriverteilung für θ gegeben Daten x_1, \dots, x_m . Hier haben wir aus „vergangenen“, beobachteten Daten für zukünftige Beobachtungen gelernt. Eine Erweiterung auf andere Zufallsgrößen ist möglich:

Satz 4.5 (Allgemeiner Darstellungssatz). *Sei X_1, X_2, \dots eine unendliche Folge reellwertiger Zufallsvariablen, die exchangeable sind, mit zugrundeliegendem Wahrscheinlichkeitsmaß \mathbb{P} . Dann existiert ein Wahrscheinlichkeitsmaß Q über \mathcal{F} , dem Raum aller Verteilungsfunktionen F auf \mathbb{R} , so dass*

$$\mathbb{P}(x_1, \dots, x_n) = \int_{\mathcal{F}} \prod_{i=1}^n F(x_i) dQ(F),$$

wobei

$$Q(F) = \lim_{n \rightarrow \infty} \mathbb{P}(F_n),$$

wobei F_n die zu x_1, \dots, x_n gehörende empirische Verteilungsfunktion bezeichnet.

Man beachte, dass obige Aussage sich (auch) auf nichtparametrische Inferenz bezieht. So steht $Q(F)$ für eine Prioriverteilung auf dem Raum aller Verteilungsfunktionen.

4.3 Bayes-Inferenz im Schnelldurchlauf

Notation:

- X : beobachtete Daten
- \tilde{X} : unbeobachtete Daten
- θ : Parameter

Ziel:

- Wahrscheinlichkeitsaussagen bedingt auf beobachtete Daten
- Vorhersage / prädiktive Inferenz

Basiskomponenten in der Bayes-Inferenz:

- $p(\theta)$ Prioriverteilung
- $f(x|\theta)$ Datenverteilung
- $f(\theta|x)$ Posterioriverteilung
- $f(\tilde{x}|x)$ prädiktive Verteilung

Nach dem Satz von Bayes ist die gemeinsame Verteilung von (θ, x) gleich

$$f(\theta, x) = f(x|\theta) \cdot p(\theta),$$

deshalb

$$f(\theta|x) = \frac{f(\theta, x)}{f(x)} = \frac{f(x|\theta)p(\theta)}{f(x)},$$

wobei

$$f(x) = \sum_{\theta \in \Theta} f(x|\theta)p(\theta), \quad \text{falls } \theta \text{ diskret,}$$
$$f(x) = \int_{\Theta} f(x|\theta)p(\theta) d\theta, \quad \text{falls } \theta \text{ stetig.}$$

Unnormalisierte Posteriori:

$$f(\theta|x) \propto f(x|\theta)p(\theta)$$

A priori prädiktive Verteilung (vor Beobachtung der Daten):

$$f(x) = \int_{\Theta} f(\theta, x) d\theta = \int_{\Theta} f(x|\theta) p(\theta) d\theta$$

A posteriori prädiktive Verteilung (nach Beobachtung der Daten x):

$$f(\tilde{x}|x) = \int_{\Theta} f(\tilde{x}, \theta|x) d\theta = \int_{\Theta} f(\tilde{x}|\theta, x) f(\theta|x) d\theta = \int_{\Theta} f(\tilde{x}|\theta) f(\theta|x) d\theta,$$

da \tilde{x} bedingt unabhängig von x gegeben θ ist.

Likelihood und Odds Ratios

Die Likelihoodfunktion ist $f(x|\theta)$ als Funktion von θ nach Beobachtung von x . Die Daten beeinflussen die Posteriori-Inferenz also nur über die Likelihood. Die Posteriori-Odds von θ_1 verglichen mit θ_2 sind

$$\frac{f_{\theta_1|x}(\theta_1|x)}{f_{\theta_2|x}(\theta_2|x)} = \frac{\frac{f(x|\theta_1)f(\theta_1)}{f(x)}}{\frac{f(x|\theta_2)f(\theta_2)}{f(x)}} = \frac{f(x|\theta_1)}{f(x|\theta_2)} \cdot \frac{f(\theta_1)}{f(\theta_2)},$$

es gilt also

$$\text{Posteriori-Odds} = \text{Priori-Odds} \times \text{Likelihoodquotient.}$$

4.4 Wiederholung: Modelle mit einem Parameter

- Gemeint ist: θ ist skalar.
- Prioriverteilung kann mehr als einen Parameter haben.
- Hier funktionieren folgende Konzepte gut:
 - **Konjugierte Prioriverteilungen**, zum Beispiel bei (einparametrischen) Exponentialfamilien.
Vorteil: Analytische Berechenbarkeit, keine Simulation nötig.
Nachteil: Für komplexe Modelle meist nicht verfügbar, deshalb eher als Baustein in komplizierteren Modellen verwendet.
 - **Referenzprioris/Referenzanalyse:**
 - * *Idee:* Priori so wählen, dass die Daten auch im Fall geringen Stichprobenumfangs die Posterioriverteilung dominieren („let the data speak for themselves“). Dies benötigt entscheidungs- und informationstheoretische Grundlagen. Suche nach nicht-informativen Prioriverteilungen: Im skalaren Fall zum Beispiel

$$0 < \theta < 1 \rightarrow \psi = \log\left(\frac{\theta}{1-\theta}\right).$$

- * *Jeffreys' Priori*

$$p(\theta) \propto \sqrt{I(\theta)}$$

ist invariant gegenüber bijektiven Transformationen von θ .

Beispiel 4.1 (Binomial- und Negative Binomialverteilung).

1. Binomialverteilung: *Die Likelihood lautet*

$$f(x|\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}.$$

Als Referenzpriori kann Jeffreys' Priori, Beta($\frac{1}{2}, \frac{1}{2}$), verwendet werden:

$$p(\theta) \propto \theta^{-1/2} (1-\theta)^{-1/2}.$$

Sei $y = \sum_{i=1}^n x_i$. Dann ist die Referenzposteriori:

$$\begin{aligned} f(\theta|x) &\propto f(x|\theta)p(\theta) \\ &\propto \theta^y(1-\theta)^{n-y}\theta^{-1/2}(1-\theta)^{-1/2} \\ &= \theta^{y-1/2}(1-\theta)^{n-y-1/2}. \end{aligned}$$

Dies entspricht dem Kern der Dichte einer Beta $(\frac{1}{2} + y, \frac{1}{2} + n - y)$ -Verteilung. $f(\theta|x)$ ist auch für die Extremfälle $y = 0$ oder $y = n$ noch „proper“. Verwendet man dagegen Haldane's Priori

$$p(\theta) \propto \theta^{-1}(1-\theta)^{-1},$$

die eine uneigentliche Priori 'Beta(0,0)' darstellt, ist die Posteriori Beta($y, n - y$) für die Extreme $y = 0$ oder $y = n$ nicht proper.

2. Negative Binomialverteilung: Sei X die Anzahl an Versuchen bis $y \geq 1$ Erfolge eintreten. Dann lautet die Likelihood

$$f(x|\theta) \propto \binom{x-1}{y-1} \theta^y (1-\theta)^{x-y} \quad \text{für } x \geq y.$$

Die Referenzpriori ist durch Jeffreys' Priori für die geometrische Verteilung gegeben (das entspricht $y = 1$):

$$p(\theta) \propto \theta^{-1}(1-\theta)^{-1/2},$$

woraus die Referenzposteriori

$$f(\theta|x) \propto \theta^{y-1}(1-\theta)^{x-y-1/2},$$

also eine Beta($y, x - y + 1/2$)-Verteilung, resultiert. Da $y \geq 1$ und $x \geq y$, ist auch diese a posteriori stets proper.

Bemerkung. Konzepte für eindimensionale Modelle sind im mehrdimensionalen Fall im Allgemeinen schwierig umzusetzen bzw. umstritten (zum Beispiel Verwendung von Referenzprioris). Man geht daher oft zu sogenannten hierarchischen Modellen über: Füge zusätzliche Stufen in das Modell ein mit dem Ziel, die Posteriori-Analyse stärker von Priori-Annahmen zu entkoppeln.

4.5 Mehr-Parameter-Modelle

4.5.1 Normalverteilung

Wir betrachten in diesem Abschnitt Daten $x_1, \dots, x_n | \mu, \sigma^2 \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$ mit μ, σ^2 unbekannt.

- (i) Gemeinsame Posterioriverteilung von $\mu, \sigma^2 | x$:

$$f(\mu, \sigma^2 | x) \propto f(x | \mu, \sigma^2) \cdot p(\mu, \sigma^2)$$

(ii) *Bedingte Posterioriverteilungen von $\mu|\sigma^2, x$ bzw. $\sigma^2|\mu, x$:*

$$f(\mu|\sigma^2, x) \quad \text{bzw.} \quad f(\sigma^2|\mu, x)$$

(iii) *Marginale Posterioriverteilung von $\mu|x$:*

$$f(\mu|x) = \int f(\mu, \sigma^2|x) d\sigma^2 = \int f(\mu|\sigma^2, x)f(\sigma^2|x) d\sigma^2$$

I. Nichtinformative Prioriverteilung

Ist nur einer der beiden Parameter unbekannt, so wählt man oft folgende **Prioriverteilungen** (Jeffreys' Prioris):

$$\begin{aligned} \sigma^2 \text{ bekannt:} & \quad p(\mu) \propto \text{const}, \\ \mu \text{ bekannt:} & \quad p(\sigma^2) \propto (\sigma^2)^{-1}. \end{aligned}$$

Eine Möglichkeit, daraus eine mehrdimensionale Priori zu konstruieren, ist:

$$p(\mu, \sigma^2) = p(\mu) \cdot p(\sigma^2) \propto (\sigma^2)^{-1},$$

d.h. wir nehmen unabhängige Prioris für μ und σ^2 an. Die **gemeinsame Posterioriverteilung** $f(\mu, \sigma^2|x)$ lautet dann:

$$\begin{aligned} f(\mu, \sigma^2|x) & \propto \text{Likelihood} \times \text{Priori} \\ & = \left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \right] \cdot (\sigma^2)^{-1} \\ & \propto \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \\ & = \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2} ((n-1)s^2 + n(\bar{x} - \mu)^2)\right) \end{aligned}$$

mit $s^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)$. Die **bedingte Posteriori** von μ , $f(\mu|\sigma^2, x)$, kann auf den Fall mit bekannter Varianz σ^2 zurückgeführt werden. Aus Statistik IV ist bekannt, dass $f(\mu|\sigma^2, x) \sim N(\bar{x}, \sigma^2/n)$. Für die **marginale Posteriori** $f(\sigma^2|x)$ hat man

$$\begin{aligned} f(\sigma^2|x) & = \int f(\mu, \sigma^2|x) d\mu \\ & \propto \int \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2} ((n-1)s^2 + n(\bar{x} - \mu)^2)\right) d\mu \\ & \propto \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2} (n-1)s^2\right) \int \exp\left(-\frac{1}{2\sigma^2} n(\bar{x} - \mu)^2\right) d\mu. \end{aligned}$$

Es gilt

$$\int \exp\left(-\frac{1}{2\sigma^2} n(\bar{x} - \mu)^2\right) d\mu = \sqrt{2\pi \frac{\sigma^2}{n}}$$

und damit

$$\begin{aligned} f(\sigma^2|x) &\propto \sigma^{-n-2} \exp\left(-\frac{1}{2\sigma^2}(n-1)s^2\right) \sqrt{2\pi \frac{\sigma^2}{n}} \\ &\propto (\sigma^2)^{-\frac{n-2}{2}} (\sigma^2)^{\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2}(n-1)s^2\right) \\ &= (\sigma^2)^{-\frac{n+1}{2}} \exp\left(-\frac{1}{2\sigma^2}(n-1)s^2\right). \end{aligned}$$

Der Kern dieser Dichte gehört zur inversen Gamma-Verteilung mit den Parametern $(n-1)/2$ und $(n-1)s^2/2$.

Wegen

$$f(\mu, \sigma^2|x_1, \dots, x_n) = f(\mu|\sigma^2, x_1, \dots, x_n) \cdot f(\sigma^2|x_1, \dots, x_n)$$

kann die gemeinsame Posterioriverteilung von $\mu, \sigma^2|x_1, \dots, x_n$ nun wie folgt simuliert werden:

Algorithmus 2 : Direkte Simulation der gemeinsamen Posterioriverteilung bei nichtinformativer Priori

Wiederhole für $s = 1, \dots, S$:

Schritt 1: Ziehe $(\sigma^2)^{(s)}$ aus $\text{IG}\left(\frac{n-1}{2}, \frac{n-1}{2}s^2\right)$.

Schritt 2: Ziehe $(\mu)^{(s)}$ aus $N(\bar{x}, \frac{1}{n}(\sigma^2)^{(s)})$.

Man erhält Paare $[(\mu^{(1)}, (\sigma^2)^{(1)}), \dots, (\mu^{(S)}, (\sigma^2)^{(S)})]$.

σ^2 als Nuisance-Parameter

Interessiert nur μ , so gibt es (mindestens) zwei Möglichkeiten zur Simulation:

1. Simuliere die gemeinsame Posteriori $f(\mu, \sigma^2|x)$ gemäß obigem Algorithmus und betrachte nur die Ziehungen $\mu^{(1)}, \dots, \mu^{(S)}$.
2. Berechne direkt die marginale Posteriori $f(\mu|x)$:

$$f(\mu|x) = \int_0^\infty f(\mu, \sigma^2|x) d\sigma^2.$$

Führt man die Substitution $z = A/(2\sigma^2)$ mit $A = (n-1)s^2 + n(\mu - \bar{x})^2$ durch, so erhält man wegen

$$\sigma^2 = \frac{1}{2}Az^{-1} \quad \text{und} \quad d\sigma^2 = -2A^{-1}\sigma^4 dz = -\frac{1}{2}Az^{-2}dz$$

für $f(\mu|x)$

$$\begin{aligned} \int_0^\infty f(\mu, \sigma^2|x) d\sigma^2 &\propto \int_0^\infty A^{-\frac{n+2}{2}} z^{\frac{n+2}{2}} \exp(-z) A z^{-2} dz \\ &= \int_0^\infty A^{-\frac{n}{2}} z^{\frac{n-2}{2}} \exp(-z) dz \\ &= A^{-\frac{n}{2}} \int_0^\infty z^{\frac{n-2}{2}} \exp(-z) dz. \end{aligned}$$

Allgemein gilt für $a > 0$ und $m > -1$:

$$\int_0^\infty x^m \exp(-ax) dx = \frac{\Gamma(m+1)}{a^{m+1}}.$$

Daraus folgt, dass das Integral konstant bezüglich μ ist und somit

$$\begin{aligned} f(\mu|x) &\propto A^{-\frac{n}{2}} \\ &= [(n-1)s^2 + n(\mu - \bar{x})^2]^{-\frac{n}{2}} \\ &= \left[1 + \frac{(\mu - \bar{x})^2}{(n-1)s^2/n} \right]^{-\frac{n}{2}} \\ &= \left[1 + \frac{1}{n-1} \left(\frac{\mu - \bar{x}}{\frac{s}{\sqrt{n}}} \right)^2 \right]^{-\frac{n}{2}} \end{aligned}$$

was der Kern einer skalierten nichtzentralen t-Verteilung mit Skalenparameter $m = s/\sqrt{n}$, Lokationsparameter $l = \bar{x}$ und $\nu = n - 1$ Freiheitsgraden ist. Allgemein hat der Kern der Dichte einer solchen allgemeinen t-Verteilung die Gestalt

$$\text{kern}(f(\theta)) = \left[1 + \frac{1}{\nu} \left(\frac{\theta - l}{m} \right)^2 \right]^{-\frac{\nu+1}{2}}.$$

Bemerkung.

1. Statt σ^2 lässt sich auch die sogenannte Präzision $\kappa = (\sigma^2)^{-1}$ verwenden. Bei Verwendung von $p(\mu, \kappa) \propto (\kappa)^{-1}$ folgt, dass $\kappa|x \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2}s^2\right)$.
2. Statt inverser Gammaverteilung wird häufig der Spezialfall einer sogenannten skalierten inversen χ^2 -Verteilung $\text{inv-}\chi^2$ verwendet (siehe unten).

II. Konjugierte Prioriverteilung

Verwende gemäß Bemerkung 2 die skalierte inverse $\chi^2(\nu_0, \sigma_0^2)$ -Verteilung als Priori.

Vorteil: Bessere Interpretation (das werden wir allerdings erst dann verstehen, wenn wir informative Prioriverteilungen in Form von (inversen) Gammaverteilungen betrachten).

Nachteil: Diese Vorgehensweise ist „non-standard“.

Zufallszahlen aus einer skalierten inversen χ^2 -Verteilung kann man wie folgt simulieren:

Algorithmus 3 : Simulation von $\theta \sim \text{inv-}\chi^2(\nu_0, \sigma_0^2)$

1. Ziehe $X^* \sim \chi^2(\nu_0)$.
 2. Setze $\theta = \frac{\nu_0 \sigma_0^2}{X^*}$.
-

Es gilt:

$$\text{inv-}\chi^2(\nu_0, \sigma_0^2) = \text{IG}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right).$$

Dies lässt sich mit dem Transformationssatz für Dichten verifizieren: Definiere $\alpha = \nu_0/2$ und $\beta = 1/2$, so dass $X^* \sim \text{Gamma}(\alpha, \beta)$. Die Umkehrfunktion der Transformation in Schritt 2 lautet

$$X^* = g^{-1}(\theta) = \frac{\nu_0 \sigma_0^2}{\theta}$$

und die zugehörige Ableitung nach θ

$$(g^{-1})'(\theta) = -\frac{\nu_0 \sigma_0^2}{\theta^2}.$$

Man erhält somit:

$$f(\theta) = f_{X^*}(g^{-1}(\theta)) \cdot |(g^{-1})'(\theta)| = \frac{(\beta \nu_0 \sigma_0^2)^\alpha}{\Gamma(\alpha)} \theta^{-\alpha-1} \exp(-\beta \nu_0 \sigma_0^2 / \theta).$$

Dies ist die Dichte einer inversen Gammaverteilung mit Parametern $(\alpha, \beta \nu_0 \sigma_0^2)$, welche gerade der gewünschten inversen χ^2 -Verteilung entspricht. Eine mögliche Parametrisierung ist nun

$$p(\mu, \sigma^2) = p(\mu | \sigma^2) \cdot p(\sigma^2)$$

mit

$$\mu | \sigma^2 \sim N\left(\mu_0, \frac{\sigma^2}{\kappa_0}\right) \quad \text{und} \quad \sigma^2 \sim \text{inv-}\chi^2(\nu_0, \sigma_0^2)$$

Man schreibt hierfür kurz: $N\text{-inv-}\chi^2\left(\mu_0, \frac{\sigma_0^2}{\kappa_0}; \nu_0, \sigma_0^2\right)$. Die Prioris sind nunmehr voneinander abhängig.

Sei nun also a priori $(\mu, \sigma^2) \sim N\text{-inv-}\chi^2\left(\mu_0, \frac{\sigma_0^2}{\kappa_0}; \nu_0, \sigma_0^2\right)$. Die **Prioridichte** lautet dann

$$\begin{aligned} p(\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2\kappa_0^{-1}}} \exp\left(-\frac{1}{2\frac{\sigma^2}{\kappa_0}}(\mu - \mu_0)^2\right) \times \frac{\left(\frac{1}{2}\nu_0\sigma_0^2\right)^{\nu_0/2}}{\Gamma(\nu_0/2)} (\sigma^2)^{-(\frac{\nu_0}{2}+1)} \exp\left(-\frac{1}{2}\frac{\nu_0\sigma_0^2}{\sigma^2}\right) \\ &\propto (\sigma^2)^{-\frac{1}{2}} (\sigma^2)^{-(\frac{\nu_0}{2}+1)} \exp\left(-\frac{1}{2\sigma^2} [\nu_0\sigma_0^2 + \kappa_0(\mu - \mu_0)^2]\right). \end{aligned}$$

Die **gemeinsame Posteriori** bei gegebenen Daten $x = (x_1, \dots, x_n)$ aus $N(\mu, \sigma^2)$ ergibt sich zu:

$$\begin{aligned} f(\mu, \sigma^2 | x) &\propto (\sigma^2)^{-\frac{1}{2}} (\sigma^2)^{-(\frac{\nu_0}{2}+1)} \exp\left(-\frac{1}{2\sigma^2} [\nu_0\sigma_0^2 + \kappa_0(\mu - \mu_0)^2]\right) \\ &\quad \times (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} ((n-1)s^2 + n(\bar{x} - \mu)^2)\right). \end{aligned}$$

Man kann zeigen (vgl. Übung), dass die Posteriori wieder $N\text{-inv-}\chi^2$ -verteilt ist mit Parametern

$$\begin{aligned} \mu_n &= \left(\frac{\kappa_0}{\kappa_0 + n}\right) \mu_0 + \left(\frac{n}{\kappa_0 + n}\right) \bar{x}, \\ \kappa_n &= \kappa_0 + n, \\ \nu_n &= \nu_0 + n, \\ \nu_n \sigma_n^2 &= \nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{x} - \mu_0)^2. \end{aligned}$$

Die Interpretation der Parameter ist wie folgt:

- μ_n ist gewichteter Mittelwert aus Stichprobenmittel und Priori-Erwartungswert. In den Grenzfällen $\kappa_0 \rightarrow \infty$ ist $\mu_n = \mu_0$ bzw. für $n \rightarrow \infty$ ist $\mu_n = \bar{x}$.
- ν_n sind die Posteriori-Freiheitsgrade als Summe von Priori-Freiheitsgraden und Stichprobenumfang.
- Die Posteriori-Quadratsumme $\nu_n \sigma_n^2$ lässt sich partitionieren in die Priori-Quadratsumme $\nu_0 \sigma_0^2$, die Quadratsumme $(n-1)s^2$ der Stichprobe und einen Term, der die Unsicherheit, die durch die Differenz zwischen Stichprobenmittel und Priori-Erwartungswert entsteht, quantifiziert.

Die **bedingte Posteriori** von $\mu|\sigma^2, x$ ist

$$\begin{aligned} \mu|\sigma^2, x &\sim N\left(\mu_n, \frac{\sigma^2}{\kappa_n}\right) \\ &\hat{=} N\left(\frac{\kappa_0}{\kappa_0+n}\mu_0 + \frac{n}{\kappa_0+n}\bar{x}, \frac{\sigma^2}{\kappa_0+n}\right) \\ &\hat{=} N\left(\frac{\frac{\kappa_0}{\sigma^2}\mu_0 + \frac{n}{\sigma^2}\bar{x}}{\frac{\kappa_0}{\sigma^2} + \frac{n}{\sigma^2}}, \frac{1}{\frac{\kappa_0}{\sigma^2} + \frac{n}{\sigma^2}}\right). \end{aligned}$$

Die Gewichte κ_0/σ^2 und n/σ^2 entsprechen der Priori- bzw. Datenpräzision. Die **marginale Posteriori** von $\sigma^2|x$ ist

$$\sigma^2|x \sim \text{inv-}\chi^2(\nu_n, \sigma_n^2).$$

Dies ermöglicht die Simulation der gemeinsamen Posteriori Verteilung:

Algorithmus 4 : Direkte Simulation der gemeinsamen Posterioriverteilung bei konjugierter Priori

Schritt 1: Ziehe $(\sigma^2)^*$ aus $\text{inv-}\chi^2(\nu_n, \sigma_n^2)$.

Schritt 2: Ziehe μ^* aus $N\left(\mu_n, \frac{(\sigma^2)^*}{\kappa_n}\right)$.

Die **marginale Posteriori** von $\mu|x$ lautet

$$f(\mu|x) \propto \left[1 + \frac{\kappa_n(\mu - \mu_n)^2}{\nu_n \sigma_n^2}\right]^{-(\nu_n+1)/2}.$$

Dies entspricht einer t-Verteilung mit ν_n Freiheitsgraden, Lokationsparameter μ_n und Skalenparameter σ_n^2/κ_n .

III. Semi-konjugierte Prioriverteilung

Die Parameter μ und σ^2 sollen nun a priori unabhängig sein. Wir wählen deshalb **a priori**

$$\mu|\sigma^2 \sim N(\mu_0, \tau_0^2) \quad \text{und} \quad \sigma^2 \sim \text{inv-}\chi^2(\nu_0, \sigma_0^2).$$

Der einzige Unterschied zum Fall der konjugierten Priori ist also, dass wir τ_0^2 statt σ_0^2/κ_0 verwenden und so die Prioris entkoppeln. Es folgt:

$$p(\mu, \sigma^2) = p(\mu|\sigma^2) \cdot p(\sigma^2) \stackrel{\text{Unabhängigkeit}}{=} p(\mu) \cdot p(\sigma^2).$$

Die resultierende gemeinsame Posteriori hat allerdings keine Form, die einer bekannten Verteilung zugeordnet werden kann. Allerdings ist $f(\mu|\sigma^2, x)$ explizit berechenbar und $f(\sigma^2|x)$ einfach zu simulieren. Die **bedingte Posteriori** ist $\mu|\sigma^2, x \sim N(\mu_n, \tau_n^2)$ mit

$$\mu_n = \frac{\frac{1}{\tau_0^2}\mu_0 + \frac{n}{\sigma^2}\bar{x}}{\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}} \quad \text{und} \quad \tau_n^2 = \frac{1}{\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}}.$$

Zur Herleitung der Posteriori $f(\sigma^2|x)$ benutzt man, dass

$$f(\sigma^2|x) = \frac{f(\mu, \sigma^2|x)}{f(\mu|\sigma^2, x)}$$

und (salopp), dass

$$f(\mu, \sigma^2|x) \propto N(\mu|\mu_0, \tau_0^2) \times \text{inv-}\chi^2(\sigma^2|\nu_0, \sigma_0^2) \times \prod_{i=1}^n N(x_i|\mu, \sigma^2).$$

Die **marginale Posteriori** hat dann die Struktur

$$f(\sigma^2|x) \propto \frac{N(\mu|\mu_0, \tau_0^2) \cdot \text{inv-}\chi^2(\sigma^2|\nu_0, \sigma_0^2) \cdot \prod_{i=1}^n N(x_i|\mu, \sigma^2)}{N(\mu|\mu_n, \tau_n^2)}.$$

Da $f(\sigma^2|x)$ nicht von μ abhängen kann, können die entsprechenden Terme ignoriert werden, die nur von μ abhängen. Man beachte jedoch, dass der Nenner über die Parameter μ_n, τ_n^2 noch von σ^2 abhängt. Setzen wir $\mu = \mu_n$, so erhalten wir

$$f(\sigma^2|x) \propto \tau_n \exp\left(-\frac{1}{2\tau_0^2}(\mu_n - \mu_0)^2\right) \cdot \text{inv-}\chi^2(\sigma^2|\nu_0, \sigma_0^2) \cdot \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_n)^2\right).$$

Diese Verteilung lässt sich beispielsweise über einen empirischen CDF-Sampler simulieren. Voraussetzung hierfür ist, dass wir die Dichte einer (univariaten) Zufallsvariable bis auf eine Proportionalitätskonstante c kennen.

Algorithmus 5 : Empirischer CDF-Sampler

- Diskretisiere den Träger der zu simulierenden Verteilung in eine Menge von N Punkten $x_1 \leq \dots \leq x_N$.
 - Evaluiere die bis auf Proportionalität bekannte Dichte an $x_1 \leq \dots \leq x_N$, um Werte f_1, \dots, f_N zu erhalten.
 - Schätze die Proportionalitätskonstante c über $c = f_1 + \dots + f_N$.
 - Ziehe Zufallszahlen aus $x_1 \leq \dots \leq x_N$ gemäß den Wahrscheinlichkeiten $f_1/c, \dots, f_N/c$.
-

Dies führt zu folgendem Algorithmus zur Simulation aus $f(\mu, \sigma^2|x)$:

1. Ziehe $(\sigma^2)^*$ aus der marginalen (approximativen) Posteriori gemäß CDF-Sampler.
2. Ziehe μ aus $f(\mu|(\sigma^2)^*, x)$.

Abschließend betrachten wir noch die **prädiktive Posterioriverteilung** für zukünftige Beobachtungen \tilde{x} gegeben Daten x . Diese lautet

$$f(\tilde{x}|x) = \int \int f(\tilde{x}|\mu, \sigma^2, x) f(\mu, \sigma^2|x) d\mu d\sigma^2 = \int \int f(\tilde{x}|\mu, \sigma^2) f(\mu, \sigma^2|x) d\mu d\sigma^2.$$

Simulation:

1. Ziehe $(\mu^*, (\sigma^2)^*)$ aus der Posteriori wie oben beschrieben.
2. Ziehe $\tilde{x} \sim N(\mu^*, (\sigma^2)^*)$.

4.5.2 Dirichlet-Multinomial Modell

Im *Dirichlet-Multinomial-Modell* wird für die Daten y_1, \dots, y_n eine *Multinomialverteilung* angenommen, also die Verallgemeinerung der Binomialverteilung auf mehr als zwei mögliche Ereignisse bei festem Stichprobenumfang n . Beispielsweise könnte eine fest vorgegebene Anzahl an Personen nach ihrer Parteipräferenz befragt werden.

Eine multinomialverteilte Zufallsvariable Y kann k mögliche Ausprägungen annehmen (zum Beispiel CDU/CSU, SPD, FDP, Grüne, Linke, andere). Die Zufallsvariable X_j , $1 \leq j \leq k$, bezeichnet die Anzahl der j -ten Ausprägung in der Stichprobe; es gilt $\sum_{j=1}^k X_j = n$. Der Parameter $\theta_j = \mathbb{P}(Y = j) \in [0, 1]$ für $Y \in \{1, \dots, k\}$ mit $\sum_{j=1}^k \theta_j = 1$ bezeichnet die Wahrscheinlichkeit für die Ausprägung j .

Die Likelihood von $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^\top$ bei Beobachtungen $\boldsymbol{x} = (x_1, \dots, x_k)^\top$ lautet

$$L(\boldsymbol{\theta}) = f(\boldsymbol{x}|\boldsymbol{\theta}) \propto \prod_{j=1}^k \theta_j^{x_j}.$$

Wegen der Restriktion $\sum_{j=1}^k \theta_j = 1$ liegen faktisch nur $k - 1$ Parameter vor, denn der k -te Parameter lässt sich deterministisch durch $\theta_k = 1 - \theta_1 - \dots - \theta_{k-1}$ berechnen. Die Likelihood lässt sich daher auch in der Form

$$L(\boldsymbol{\theta}) \propto \left(\prod_{j=1}^{k-1} \theta_j^{x_j} \right) (1 - \theta_1 - \dots - \theta_{k-1})^{x_k}$$

schreiben.

Die zur Multinomialverteilung konjugierte Verteilung ist die sogenannte *Dirichletverteilung*, eine Verallgemeinerung der Beta-Verteilung, geschrieben

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^\top \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) = D(\boldsymbol{\alpha}),$$

mit Dichtefunktion

$$p(\boldsymbol{\theta}) = \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \cdot \dots \cdot \Gamma(\alpha_k)} \theta_1^{\alpha_1-1} \cdot \dots \cdot \theta_k^{\alpha_k-1},$$

wobei $\alpha_j > 0$ für alle $j = 1, \dots, k$ und wieder $\theta_j \in [0, 1]$ mit $\sum_{j=1}^k \theta_j = 1$. Die Dirichlet-Verteilung spezifiziert also eine Verteilung auf einem $(k-1)$ -dimensionalen offenen Simplex.

Eigenschaften der Dirichletverteilung:

Definiere $\alpha_0 = \sum_{j=1}^k \alpha_j$.

- Momente:

$$\mathbb{E}(\theta_j) = \frac{\alpha_j}{\alpha_0}, \quad \text{Var}(\theta_j) = \frac{\alpha_j(\alpha_0 - \alpha_j)}{\alpha_0^2(\alpha_0 + 1)}, \quad \text{Cov}(\theta_i, \theta_j) = -\frac{\alpha_i \alpha_j}{\alpha_0^2(\alpha_0 + 1)},$$

wobei die Restriktion $\sum_{j=1}^k \theta_j = 1$ die negative Korrelation impliziert.

- Modus:

$$\text{Modus}(\boldsymbol{\theta})_j = \frac{\alpha_j - 1}{\alpha_0 - k}$$

ist die j -te Komponente des k -dimensionalen Modus.

- Jede Randverteilung ist wieder eine Dirichletverteilung, zum Beispiel

$$(\theta_i, \theta_j, 1 - \theta_i - \theta_j) \sim \text{Dirichlet}(\alpha_i, \alpha_j, \alpha_0 - \alpha_i - \alpha_j).$$

Insbesondere ist

$$\theta_j \sim \text{Beta}(\alpha_j, \alpha_0 - \alpha_j).$$

- Die bedingten Verteilungen sind ebenfalls Dirichletverteilt. Setzt man

$$\theta'_i = \frac{\theta_m}{1 - \sum_{r=1}^{m-1} \theta_r}, \quad m \leq i \leq k,$$

gegeben die Realisationen $\theta_1, \dots, \theta_{m-1}$, so ist

$$(\theta'_m, \dots, \theta'_k)^\top \sim \text{Dirichlet}(\alpha_m, \dots, \alpha_k).$$

Algorithmus 6 : Simulation aus der Dirichletverteilung

• Simulation 1:

1. Ziehe Z_1, Z_2, \dots, Z_k aus (unabhängigen) Gamma-Verteilungen mit Parametern $(\alpha_1, 1), \dots, (\alpha_k, 1)$.
2. Setze

$$\theta_j = \frac{Z_j}{\sum_{i=1}^k Z_i} .$$

• Simulation 2 („Stick Breaking Prior“):

1. Ziehe $\theta_1 \sim \text{Beta}(\alpha_1, \alpha_0 - \alpha_1)$.
2. Für $j = 2, \dots, k - 1$:
 - (i) Ziehe $Z_j \sim \text{Beta}(\alpha_j, \sum_{i=j+1}^k \alpha_i)$.
 - (ii) Setze $\theta_j = \left(1 - \sum_{i=1}^{j-1} \theta_i\right) Z_j$.
3. Setze $\theta_k = 1 - \sum_{i=1}^{k-1} \theta_i$.

Für $\mathbf{x}|\boldsymbol{\theta} \sim \text{Multinomial}(n; \theta_1, \dots, \theta_k)$ und $\boldsymbol{\theta}|\boldsymbol{\alpha} \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ lautet die **Posteriori-Verteilung** von $\boldsymbol{\theta}$:

$$\begin{aligned} f(\boldsymbol{\theta}|\mathbf{x}) &\propto L(\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}|\boldsymbol{\alpha}) \\ &\propto \prod_{j=1}^k \theta_j^{x_j} \cdot \prod_{j=1}^k \theta_j^{\alpha_j - 1} \\ &= \prod_{j=1}^k \theta_j^{x_j + \alpha_j - 1}, \end{aligned}$$

d.h. die Posteriori ist Dirichlet($x_1 + \alpha_1, \dots, x_k + \alpha_k$)-verteilt.

Interpretation der Posteriori

Der Posteriori-Erwartungswert

$$\mathbb{E}[\theta_j|\mathbf{x}] = \frac{x_j + \alpha_j}{\sum_{i=1}^k x_i + \sum_{i=1}^k \alpha_i} = \frac{x_j + \alpha_j}{n + \alpha_0}$$

lässt sich umschreiben zu

$$\mathbb{E}[\theta_j|\mathbf{x}] = \frac{\alpha_0}{\alpha_0 + n} \cdot \underbrace{\frac{\alpha_j}{\alpha_0}}_{\text{Priori-Erwartungswert}} + \frac{n}{\alpha_0 + n} \cdot \underbrace{\frac{x_j}{n}}_{\text{MLE}} .$$

Der Parameter α_0 lässt sich als „a priori Anzahl Beobachtungen“ und α_j als „a priori Erfolge“ für Kategorie j interpretieren.

Bemerkung.

1. Die Wahl $\alpha_1 = \dots = \alpha_k = 0$ entspricht einer Gleichverteilung für $\{\log \theta_j\}_{j=1}^k$. In diesem Fall ist die Posteriori nur dann proper, wenn $x_j \geq 1$, $j = 1, \dots, p$.
2. Die Wahl $\alpha_1 = \dots = \alpha_k = 1/2$ entspricht Jeffreys' Priori.
3. Die Wahl $\alpha_1 = \dots = \alpha_k = 1$ entspricht einer Priori-Gleichverteilung auf dem Simplex.

Bemerkung. Die Dirichlet-Verteilung eignet sich auch als Priori bei der Analyse von Kontingenztafeln mit multinomialem Erhebungsschema:

| | |
|-------|-------|
| x_1 | x_2 |
| x_3 | x_4 |

($n = x_1 + x_2 + x_3 + x_4$). Erweiterungen auf restringierte Multinomialverteilungen sind möglich (loglineare Modelle).

Ad-hoc Prozedur: Addiere $1/2$ zu jedem Eintrag der Kontingenztabelle und berechne dann den Maximum-Likelihood-Schätzer; das entspricht $\alpha_1 = \dots = \alpha_k = 3/2$ und der Posteriori-Modus Schätzung

$$\begin{aligned} \text{Modus}(\boldsymbol{\theta}|\mathbf{x})_j &= \frac{x_j + \alpha_j - 1}{\sum_{i=1}^k x_i + \sum_{i=1}^k \alpha_i - k} \\ &= \frac{x_j + \frac{1}{2}}{\sum_{i=1}^k x_i + \frac{1}{2}k}. \end{aligned}$$

4.5.3 Multivariate Normalverteilung

Notation:

- $\mathbf{X} = (X_1, \dots, X_p)^\top$ ist p -dimensionaler Zufallsvektor.
- $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$ ist p -dimensionaler Erwartungswertvektor.
- Die symmetrische und positiv definite (Notation: $\boldsymbol{\Sigma} > 0$) Matrix

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$$

ist $p \times p$ -dimensionale Kovarianzmatrix.

- Eine Beobachtung $\mathbf{x} = (x_1, \dots, x_p)^\top$ ist MVN($\boldsymbol{\mu}, \boldsymbol{\Sigma}$) (multivariat normalverteilt), wenn

$$f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

- Die Likelihood für n i.i.d. Realisationen $\mathbf{x}_1, \dots, \mathbf{x}_n$ lautet

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = f(\mathbf{x}_1, \dots, \mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &\propto |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right) \\ &= |\boldsymbol{\Sigma}|^{-n/2} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right) \end{aligned}$$

mit $\mathbf{S}_0 = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top$. Die zweite Identität ergibt sich über die Umformungen

$$\begin{aligned} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) &= \text{tr}\left(\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right) \\ &= \text{tr}\left(\sum_{i=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top\right) \\ &= \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0). \end{aligned}$$

I. Konjugierte Prioriverteilung bei unbekanntem $\boldsymbol{\mu}$ und bekanntem $\boldsymbol{\Sigma}$

Konjugierte **Prioriverteilung** für $\boldsymbol{\mu}$ bei bekanntem $\boldsymbol{\Sigma}$ ist

$$\boldsymbol{\mu} \sim \text{MVN}(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0)$$

mit $\boldsymbol{\Lambda}_0 > 0$. Die **Posteriori** für $\boldsymbol{\mu}$ ist

$$f(\boldsymbol{\mu} | \mathbf{x}, \boldsymbol{\Sigma}) \propto \exp\left(-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Lambda}_0^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right).$$

Der Term im Exponenten ist eine quadratische Form in $\boldsymbol{\mu}$, die sich über eine quadratische Ergänzung und Vernachlässigung von Konstanten ergibt. Man erhält

$$f(\boldsymbol{\mu} | \mathbf{x}, \boldsymbol{\Sigma}) \sim \text{MVN}(\boldsymbol{\mu}_n, \boldsymbol{\Lambda}_n)$$

mit

$$\begin{aligned} \boldsymbol{\mu}_n &= (\boldsymbol{\Lambda}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} (\boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0 + n\boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}), \\ \boldsymbol{\Lambda}_n^{-1} &= \boldsymbol{\Lambda}_0^{-1} + n\boldsymbol{\Sigma}^{-1} \end{aligned}$$

und $\bar{\mathbf{x}} = (\sum_{i=1}^n \mathbf{x}_i)/n$ in Analogie zum univariaten Fall.

Die **bedingte** und **marginale Posteriori** für *Subvektoren* von $\boldsymbol{\mu}$ folgen aus den Eigenschaften der multivariaten Normalverteilung: Betrachte die Partitionierungen

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}^{(1)} \\ \boldsymbol{\mu}^{(2)} \end{pmatrix}, \quad \boldsymbol{\mu}_n = \begin{pmatrix} \boldsymbol{\mu}_n^{(1)} \\ \boldsymbol{\mu}_n^{(2)} \end{pmatrix}, \quad \boldsymbol{\Lambda}_n = \begin{pmatrix} \boldsymbol{\Lambda}_n^{(11)} & \boldsymbol{\Lambda}_n^{(12)} \\ \boldsymbol{\Lambda}_n^{(21)} & \boldsymbol{\Lambda}_n^{(22)} \end{pmatrix}.$$

Dann gilt für die bedingten Verteilungen

$$\boldsymbol{\mu}^{(1)} | \boldsymbol{\mu}^{(2)}, \mathbf{x} \sim \text{MVN}\left(\boldsymbol{\mu}_n^{(1)} + \boldsymbol{\beta}^{1|2} (\boldsymbol{\mu}^{(2)} - \boldsymbol{\mu}_n^{(2)}), \boldsymbol{\Lambda}_n^{1|2}\right)$$

mit

$$\begin{aligned}\beta^{1|2} &= \Lambda_n^{(12)} \left(\Lambda_n^{(22)} \right)^{-1}, \\ \Lambda^{1|2} &= \Lambda_n^{(11)} - \Lambda_n^{(12)} \left(\Lambda_n^{(22)} \right)^{-1} \Lambda_n^{(21)}\end{aligned}$$

und für die marginalen Verteilungen

$$\boldsymbol{\mu}^{(1)} \sim \text{MVN} \left(\boldsymbol{\mu}_n^{(1)}, \Lambda_n^{(11)} \right).$$

Die **prädiktive Posterioriverteilung** ist (informell)

$$f(\tilde{\boldsymbol{x}}, \boldsymbol{\mu} | \boldsymbol{x}) = \text{MVN}(\tilde{\boldsymbol{x}} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \cdot \text{MVN}(\boldsymbol{\mu} | \boldsymbol{\mu}_n, \Lambda_n).$$

Im Exponenten erhält man eine quadratische Form in $(\tilde{\boldsymbol{x}}, \boldsymbol{\mu})$. $(\tilde{\boldsymbol{x}}, \boldsymbol{\mu})$ sind gemeinsam multivariat normalverteilt, und daher folgt $\tilde{\boldsymbol{x}}$ als Randverteilung der beiden Komponenten ebenfalls einer multivariaten Normalverteilung. Die erforderlichen Kenngrößen lassen sich über die iterierte Erwartung und Varianz berechnen:

$$\mathbb{E}[\tilde{\boldsymbol{x}} | \boldsymbol{x}] = \mathbb{E}[\mathbb{E}[\tilde{\boldsymbol{x}} | \boldsymbol{\mu}, \boldsymbol{x}] | \boldsymbol{x}] = \mathbb{E}[\boldsymbol{\mu} | \boldsymbol{x}] = \boldsymbol{\mu}_n$$

und

$$\begin{aligned}\text{Var}(\tilde{\boldsymbol{x}} | \boldsymbol{x}) &= \mathbb{E}[\text{Var}(\tilde{\boldsymbol{x}} | \boldsymbol{\mu}, \boldsymbol{x}) | \boldsymbol{x}] + \text{Var}[\mathbb{E}(\tilde{\boldsymbol{x}} | \boldsymbol{\mu}, \boldsymbol{x}) | \boldsymbol{x}] \\ &= \mathbb{E}[\boldsymbol{\Sigma} | \boldsymbol{x}] + \text{Var}[\boldsymbol{\mu} | \boldsymbol{x}] \\ &= \boldsymbol{\Sigma} + \Lambda_n.\end{aligned}$$

II. Konjugierte Prioriverteilung bei unbekanntem $\boldsymbol{\mu}$ und unbekanntem $\boldsymbol{\Sigma}$

In Abschnitt 4.5.1-II (Seite 94) hatten wir als konjugierte Prioriverteilungen für die Parameter μ und σ^2 der univariaten Normalverteilung

$$\mu | \sigma^2 \sim N \left(\mu_0, \frac{\sigma^2}{\kappa_0} \right) \quad \text{und} \quad \sigma^2 \sim \text{inv-}\chi^2(\nu_0, \sigma_0^2),$$

kurz

$$\mu, \sigma^2 \sim \text{N-inv-}\chi^2 \left(\mu_0, \frac{\sigma_0^2}{\kappa_0}; \nu_0, \sigma_0^2 \right),$$

verwendet. Hier nun verwenden wir die multivariaten Analoga

$$\boldsymbol{\mu} | \boldsymbol{\Sigma} \sim \text{MVN} \left(\boldsymbol{\mu}_0, \frac{1}{\kappa_0} \boldsymbol{\Sigma} \right) \quad \text{und} \quad \boldsymbol{\Sigma} \sim \text{inv-Wishart}_{\nu_0}(\Lambda_0^{-1}),$$

kurz

$$\boldsymbol{\mu}, \boldsymbol{\Sigma} \sim \text{MVN-inv-Wishart} \left(\boldsymbol{\mu}_0, \frac{1}{\kappa_0} \Lambda_0; \nu_0, \Lambda_0 \right).$$

Die **gemeinsame Prioridichte** ist dann

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-\left(\frac{\nu_0+p}{2}+1\right)} \cdot \exp\left(-\frac{1}{2}\left(\text{tr}(\boldsymbol{\Lambda}_0\boldsymbol{\Sigma}^{-1}) - \kappa_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right)\right).$$

Dabei bezeichnet ν_0 die a priori Anzahl der Freiheitsgrade, κ_0 die a priori Anzahl an Messungen auf der $\boldsymbol{\Sigma}$ -Skala. Die **gemeinsame Posterioriverteilung** von $\boldsymbol{\mu}$ und $\boldsymbol{\Sigma}$ lautet

$$\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{x} \sim \text{MVN-inv-Wishart}\left(\boldsymbol{\mu}_n, \frac{1}{\kappa_n} \boldsymbol{\Lambda}_n; \nu_n, \boldsymbol{\Lambda}_n\right).$$

mit

$$\begin{aligned}\boldsymbol{\mu}_n &= \frac{\kappa_0}{\kappa_0 + n} \boldsymbol{\mu}_0 + \frac{n}{\kappa_0 + n} \bar{\boldsymbol{x}}, \\ \kappa_n &= \kappa_0 + n, \\ \nu_n &= \nu_0 + n, \\ \boldsymbol{\Lambda}_n &= \boldsymbol{\Lambda}_0 + \boldsymbol{S} + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{\boldsymbol{x}} - \boldsymbol{\mu}_0)(\bar{\boldsymbol{x}} - \boldsymbol{\mu}_0)^\top,\end{aligned}$$

wobei $\boldsymbol{S} = \sum_{i=1}^n (\boldsymbol{x}_i - \bar{\boldsymbol{x}})(\boldsymbol{x}_i - \bar{\boldsymbol{x}})^\top$ in Analogie zum univariaten Fall. Die Interpretation der Parameter von Seite 95 lässt sich direkt übertragen: Der Posteriori-Erwartungswert ist ein gewichtetes Mittel aus Stichprobenmittelwertvektor und Priori-Erwartungswert. Die Gesamtstreuungsmatrix $\boldsymbol{\Lambda}_n$ lässt sich in Priori-Streuungsmatrix, empirische Streuungsmatrix und Streuung zwischen Priori-Erwartungswert und Stichprobenmittel partitionieren.

Die **marginale Posteriori** für $\boldsymbol{\mu}$ folgt einer multivariaten t-Verteilung mit Parametern $\boldsymbol{\mu}_n$ und $\boldsymbol{\Lambda}_n / (\kappa_n \cdot (\nu_n - p + 1))$, die marginale Posteriori für $\boldsymbol{\Sigma}$ einer inversen Wishart-Verteilung mit Parametern ν_n und $\boldsymbol{\Lambda}_n^{-1}$. Zur Simulation aus der gemeinsamen Posteriori oder aus der prädiktiven Verteilung ist folgender Algorithmus anwendbar:

Algorithmus 7 : Simulation aus der gemeinsamen Posteriori und der prädiktiven Verteilung bei konjugierter Priori

Für $s = 1, \dots, S$:

1. Ziehe $\boldsymbol{\Sigma}^{(s)} | \boldsymbol{x} \sim \text{inv-Wishart}_{\nu_n}(\boldsymbol{\Lambda}_n^{-1})$.
2. Ziehe $\boldsymbol{\mu}^{(s)} | \boldsymbol{\Sigma}^{(s)}, \boldsymbol{x} \sim \text{MVN}\left(\boldsymbol{\mu}_n, \frac{1}{\kappa_n} \boldsymbol{\Sigma}^{(s)}\right)$.
3. Ziehe $\tilde{\boldsymbol{x}}^{(s)} | \boldsymbol{\mu}^{(s)}, \boldsymbol{\Sigma}^{(s)}, \boldsymbol{x} \sim \text{MVN}\left(\boldsymbol{\mu}^{(s)}, \boldsymbol{\Sigma}^{(s)}\right)$.

Dann ist $(\boldsymbol{\mu}^{(s)}, \boldsymbol{\Sigma}^{(s)})$ eine Ziehung aus der gemeinsamen Posterioridichte, $\tilde{\boldsymbol{x}}$ eine Ziehung aus der prädiktiven Verteilung.

III. Nichtinformative Prioriverteilung bei unbekanntem $\boldsymbol{\mu}$ und unbekanntem $\boldsymbol{\Sigma}$

Als nichtinformative Prioriverteilung bei unbekanntem $\boldsymbol{\mu}$ und $\boldsymbol{\Sigma}$ eignet sich die multivariate Jeffreys' Priori

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(p+1)/2}.$$

Diese entspricht dem Grenzfall $\kappa_0 \rightarrow 0, \nu_0 \rightarrow 1, |\mathbf{\Lambda}_0| \rightarrow 0$ bei der konjugierten Priori. Für die Posteriori-Kenngrößen ergibt sich in diesem Fall für die **bedingte Verteilung** von $\boldsymbol{\mu}$

$$\boldsymbol{\mu} | \boldsymbol{\Sigma}, \mathbf{x} \sim \text{MVN} \left(\bar{\mathbf{x}}, \frac{1}{n} \boldsymbol{\Sigma} \right),$$

und für die **marginalen Verteilungen**

$$\begin{aligned} \boldsymbol{\Sigma} | \mathbf{x} &\sim \text{inv-Wishart}_{n-1}(\mathbf{S}), \\ \boldsymbol{\mu} | \mathbf{x} &\sim \text{mv-t}_{n-p} \left(\bar{\mathbf{x}}, \frac{1}{n(n-p)} \mathbf{S} \right). \end{aligned}$$

Als Beispiel wird im Folgenden die bivariate Normalverteilung betrachtet. Es dient auch dazu, die folgende für die Bayes-Inferenz wichtige Simulationsstrategie, das sogenannte Gibbs-Sampling, zu illustrieren.

Algorithmus 8 : Gibbs-Sampler

Gegeben: Ein mehrdimensionaler stetiger Zufallsvektor \mathbf{X} mit Verteilung π . Der Einfachheit halber seien alle Komponenten von \mathbf{X} stetig.

Wir erzeugen im Folgenden eine Markovkette $\mathbf{X}^{(0)}, \mathbf{X}^{(1)}, \dots$ mit Startwert $\mathbf{X}^{(0)}$ und stationärer Verteilung π . Sei $\mathbf{X}^{(t)}$ der aktuelle Zustand der Markovkette. $\mathbf{X}^{(t)}$ lasse sich in k Subvektoren $\mathbf{X}^{(t)} = (\mathbf{X}_{\bullet 1}^{(t)}, \mathbf{X}_{\bullet 2}^{(t)}, \dots, \mathbf{X}_{\bullet k}^{(t)})$ partitionieren. Definiere

$$\mathbf{X}_{-s}^{(t)} = \left(\mathbf{X}_{\bullet 1}^{(t)}, \mathbf{X}_{\bullet 2}^{(t)}, \dots, \mathbf{X}_{\bullet (s-1)}^{(t)}, \mathbf{X}_{\bullet (s+1)}^{(t-1)}, \dots, \mathbf{X}_{\bullet k}^{(t-1)} \right)$$

für $s = 1, \dots, k$. Ferner seien die *vollständig bedingten Verteilungen* („full conditionals“)

$$\pi_{\bullet s | -s}^{(t)} = \pi \left(\mathbf{X}_{\bullet s}^{(t)} | \mathbf{X}_{-s}^{(t)} \right) = \frac{\pi \left(\mathbf{X}_{\bullet s}^{(t)}, \mathbf{X}_{-s}^{(t)} \right)}{\int \pi \left(\mathbf{X}_{\bullet s}^{(t)}, \mathbf{X}_{-s}^{(t)} \right) d\mathbf{X}_{\bullet s}^{(t)}}$$

gegeben und simulierbar.

Dann wird der nächste Zustand $\mathbf{X}^{(t+1)}$ komponentenweise wie folgt erzeugt:

Schritt 1: Ziehe $\mathbf{X}_{\bullet 1}^{(t+1)} \sim \pi_{\bullet 1 | -1}^{(t+1)}$.

Schritt 2: Ziehe $\mathbf{X}_{\bullet 2}^{(t+1)} \sim \pi_{\bullet 2 | -2}^{(t+1)}$.

⋮

Schritt k: Ziehe $\mathbf{X}_{\bullet k}^{(t+1)} \sim \pi_{\bullet k | -k}^{(t+1)}$.

Wiederhole diese Schritte ausreichend oft.

Nach einer gewissen Zahl von Wiederholungen kann $\mathbf{X}^{(t)}$ als Ziehung aus π angesehen werden. Im Gegensatz zu obigen „direkten“ Simulationsalgorithmen liegen nun allerdings abhängige Realisationen vor.

Beispiel 4.2. (Bivariate Normalverteilung) Sei \mathbf{x} bivariat normalverteilt mit Erwartungswertvektor $(\mu_1, \mu_2)^\top$ und Kovarianzmatrix

$$\begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho \\ \rho & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \text{ und } \rho \text{ bekannt.}$$

Bei einer nichtinformativen Priori $p(\mu_1, \mu_2) \propto \text{const}$ für μ_1, μ_2 reduziert sich die Posteriori auf die Likelihood bei gegebenen Daten $\mathbf{x} = ((x_{11}, x_{12})^\top, \dots, (x_{n1}, x_{n2})^\top)$:

$$L(\mu_1, \mu_2) = \left(\frac{1}{2\pi}\right)^n (1 - \rho^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2(1 - \rho^2)} A\right),$$

wobei

$$A = \sum_{i=1}^n [(x_{i1} - \mu_1)^2 - 2\rho(x_{i1} - \mu_1)(x_{i2} - \mu_2) + (x_{i2} - \mu_2)^2]$$

Wir möchten nun die vollständig bedingten Verteilungen $\mu_1|\mu_2, \mathbf{x}$ und $\mu_2|\mu_1, \mathbf{x}$ berechnen. Natürlich ist es aus Symmetriegründen ausreichend, nur $\mu_1|\mu_2, \mathbf{x}$ zu ermitteln. Wegen

$$f(\mu_1|\mu_2, \mathbf{x}) = \frac{f(\mu_1, \mu_2|\mathbf{x})}{f(\mu_2|\mathbf{x})} = \frac{f(\mu_1, \mu_2, \mathbf{x})}{f(\mu_2, \mathbf{x})} \propto f(\mu_1, \mu_2|\mathbf{x})$$

genügt es, aus der gemeinsamen Posteriori lediglich die Terme zu betrachten, die von der jeweiligen Variablen in der bedingten Verteilung abhängen. Man erhält dann

$$f(\mu_1|\mu_2, \mathbf{x}) \propto \exp\left(-\frac{1}{2(1 - \rho^2)} n [\mu_1^2 - 2\mu_1(\bar{x}_1 + \rho(\mu_2 - \bar{x}_2))]\right)$$

mit $\bar{x}_j = (\sum_{i=1}^n x_{ij})/n$ für $j = 1, 2$. Eine quadratische Ergänzung des Terms in eckigen Klammern um $\bar{x}_1 + \rho(\mu_2 - \bar{x}_2)$ liefert schließlich das Endresultat

$$p(\mu_1|\mu_2, \mathbf{x}) \propto \exp\left(-\frac{1}{2\frac{1-\rho^2}{n}} \left(\mu_1 - [\bar{x}_1 + \rho(\mu_2 - \bar{x}_2)]\right)^2\right).$$

Dies entspricht dem Kern einer $N(\bar{x}_1 + \rho(\mu_2 - \bar{x}_2), (1 - \rho^2)/n)$ -Verteilung. Der zugehörige Gibbs-Sampler hat die Gestalt:

1. Wähle einen Startwert $\mu_2^{(0)}$.
2. Für $s = 1, \dots, S$:
 - (a) Ziehe $\mu_1^{(s)}|\mu_2^{(s-1)} \sim N\left(\bar{x}_1 + \rho\left(\mu_2^{(s-1)} - \bar{x}_2\right), \frac{1-\rho^2}{n}\right)$.
 - (b) Ziehe $\mu_2^{(s)}|\mu_1^{(s)} \sim N\left(\bar{x}_2 + \rho\left(\mu_1^{(s)} - \bar{x}_1\right), \frac{1-\rho^2}{n}\right)$.