

Methoden der Statistik

Markov Chain Monte Carlo–Methoden

Thorsten Dickhaus

Humboldt-Universität zu Berlin

07.02.2012



Problemstellung

Wie kann eine Zufallsstichprobe am Computer simuliert werden, deren Verteilung aus einem komplizierten Bildungsgesetz (z. B. hierarchisch Bayes) herrührt?

Definition

Eine Markov Chain Monte Carlo (MCMC)–Methode zur Simulation einer Wahrscheinlichkeitsverteilung mit Dichte f produziert eine **ergodische Markoffkette**, welche die Verteilung mit Dichte f als stationäre Verteilung aufweist.

Übersicht

- 1 Grundlegendes über Markoffketten
- 2 Der Metropolis-Hastings-Algorithmus
- 3 Anwendung auf das Logit-Modell
- 4 Der Gibbs-Sampler

Im Folgenden bezeichnet (Ω, \mathcal{F}) einen messbaren Raum.

Definition

Ein Übergangskern ist eine Abbildung $\mathcal{K} : \Omega \times \mathcal{F} \rightarrow [0, 1]$ mit

- $\mathcal{K}(x, \cdot)$ ist ein Wahrscheinlichkeitsmaß auf $\mathcal{F} \quad \forall x \in \Omega$
- $\mathcal{K}(\cdot, B)$ ist messbar $\forall B \in \mathcal{F}$

Definition

Ein stochastischer Prozess in diskreter Zeit X_0, \dots, X_n, \dots heißt Markoff-Kette (MC) (X_n) auf Ω , falls $\forall B \in \mathcal{F}$

$$\begin{aligned} \mathbb{P}(X_{k+1} \in B | x_0, \dots, x_k) &= \mathbb{P}(X_{k+1} \in B | x_k) \\ &= \mathcal{K}(x_k, B) = \int_B \mathcal{K}(x_k, dx) \end{aligned}$$

für $k \in \mathbb{N}$ mit einem Übergangskern \mathcal{K} gilt.

Definition (Sei Ω abzählbar.)

- a) (X_n) heißt **zeithomogen**, falls $\forall k \in \mathbb{N}$:

$$\mathbb{P}(X_{k+1} \in B | X_k = x) = \mathbb{P}(X_1 \in B | X_0 = x), x \in \Omega.$$

- b) **Ersteintrittszeit** in $x \in \Omega$: $T^x := \min \{n > 0 | X_n = x\}$,

$$f(x, y) := \mathbb{P}(T^y < \infty | X_0 = x) \text{ (Kette startet in } x).$$

- c) Ein Zustand $x \in \Omega$ heißt **rekurrent**, falls $f(x, x) = 1$.

Eine Markoffkette (X_n) heißt rekurrent, falls jeder ihrer möglichen Zustände rekurrent ist. Ansonsten heißt (X_n) transient.

- d) Eine Markoffkette heißt **irreduzibel**, falls jeder Zustand y von jedem Zustand x aus erreichbar ist.

Satz

- 1 Falls MC (X_n) irreduzibel:
Ein Zustand rekurrent $\Leftrightarrow (X_n)$ rekurrent
- 2 Ist eine Markoffkette **irreduzibel und rekurrent**, dann existiert ein **eindeutig bestimmtes invariantes Maß**
 $\mu \equiv \mu_{\mathcal{K}}$.
Ist μ endlich, heißt die Markoffkette positiv-rekurrent.

Korollar

MC irreduzibel und positiv-rekurrent, dann gilt für $x \in \Omega$:

$$\mathbb{E}[T^x] < \infty \text{ und } \mu(x) = (\mathbb{E}[T^x])^{-1}.$$

Ergodizität

Definition

Periode: $d_x := \text{ggT}\{n \mid \mathcal{K}^n(x, x) > 0\}$

MC heißt **aperiodisch**, wenn $d_x = 1 \forall x \in \Omega$.

Satz (Konvergenz ins Gleichgewicht)

Ist (X_n) eine irreduzible, positiv-rekurrenente, aperiodische (**ergodische**) Markoffkette mit stationärer Verteilung μ , so gilt für jede beliebige Startverteilung μ_0 von X_0 :

$$\mu_0 \mathcal{K}^n \rightarrow \mu \text{ für } n \rightarrow \infty .$$

Markov Chain Monte Carlo-Algorithmen

Prinzip von MCMC-Algorithmen:

Für beliebige Startverteilungen wird eine ergodische Markoffkette generiert.

Dabei wird der **Übergangskern** so gewählt, dass Konvergenz gegen eine festgelegte invariante Verteilung (**die Zielverteilung**) stattfindet.

Anwendungsbeispiel:

Näherungsweise Berechnung von $\int h(x)f(x)dx$ durch $T^{-1} \sum_{i=n_0}^{n_0+T} h(X_i)$, wobei f die stationäre Dichte ist und die sogenannte „burn-in-Phase“ $\{1, \dots, n_0 - 1\}$ ausgenommen wird.

Übersicht

- 1 Grundlegendes über Markoffketten
- 2 Der Metropolis-Hastings-Algorithmus**
- 3 Anwendung auf das Logit-Modell
- 4 Der Gibbs-Sampler

Aufgabe: Simulation einer Stichprobe, die der Zielverteilung mit Dichte f folgt.

Gegeben sei $X^{(t)} = x^{(t)}$.

- 1 Generiere $Y_t \sim q(y|x^{(t)})$ gemäß einer **proposal-Verteilung q** .
- 2 Acceptance-Rejection Schritt:

$$X^{(t+1)} = \begin{cases} Y_t & \text{mit Wahrscheinlichkeit } \rho(x^{(t)}, Y_t) \\ x^{(t)} & \text{mit Wahrscheinlichkeit } 1 - \rho(x^{(t)}, Y_t), \end{cases}$$

wobei

$$\rho(x, y) = \min \left(1, \frac{f(y) q(x|y)}{f(x) q(y|x)} \right).$$

Wahl der proposal-Verteilung q

Metropolis-Hastings erzeugt eine irreduzible MC, falls

$$q(y|x) > 0 \quad \forall (x, y) \in \text{supp}(f) \times \text{supp}(f).$$

Aperiodizität gilt gerade, falls

$$\mathbb{P} \left(f(X^{(t)})q(Y_t|X^{(t)}) \leq f(Y_t)q(X^{(t)}|Y_t) \right) < 1.$$

Zur Anwendung von MH sollte q leicht zu simulieren sein und von f muss $f(y)/q(y|x)$ bis auf eine Konstante bekannt sein.

Übergangskern des Metropolis-Hastings-Algorithmus'

Explizite Angabe von \mathcal{K} :

$$\mathcal{K}(x, y) = \rho(x, y)q(y|x) + \left(1 - \int \rho(x, y)q(y|x)dy\right) \delta_x(y).$$

Damit erfüllt (X_n) die **detailed balance**-Bedingung

$$\mathcal{K}(x, y)f(x) = \mathcal{K}(y, x)f(x) \quad \forall x, y.$$

Satz

Besitzt eine Markoffkette die „detailed balance“ Eigenschaft, so ist sie irreduzibel und f **Dichte der invarianten Verteilung**.

Satz (Konvergenz von Metropolis-Hastings)

Ist die Metropolis-Hastings-Markoffkette $(X^{(t)})$ f -irreduzibel, dann gilt

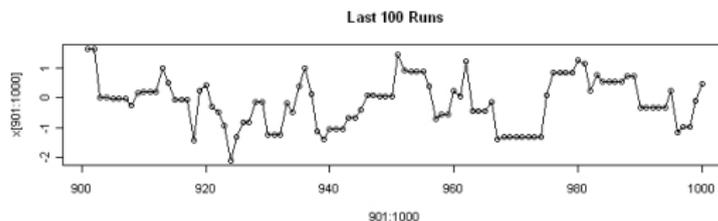
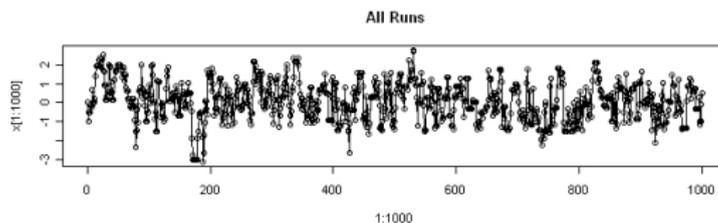
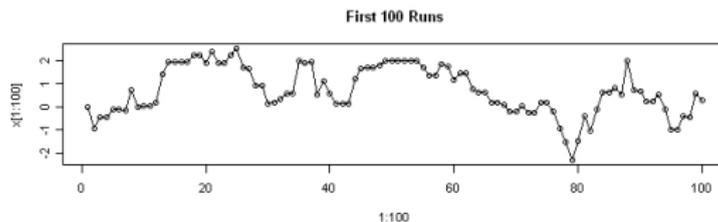
$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T h(X^{(t)}) = \int h(x) f(x) dx \quad \forall h \in L^1(f), \quad f - \text{f. ü.}$$

Ist $(X^{(t)})$ aperiodisch, so gilt außerdem

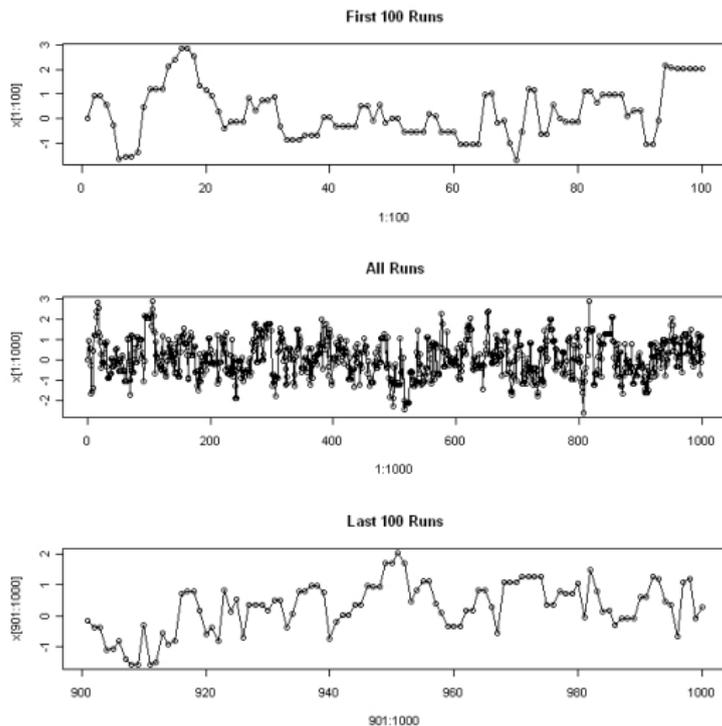
$$\lim_{n \rightarrow \infty} \left\| \int \mathcal{K}^n(x, \cdot) \mu(dx) - f \right\|_{TV} = 0$$

für jede Startverteilung μ .

Beispiel: Normalverteilung

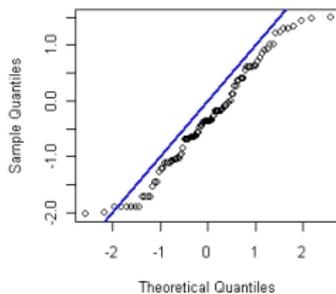


Beispiel: Normalverteilung

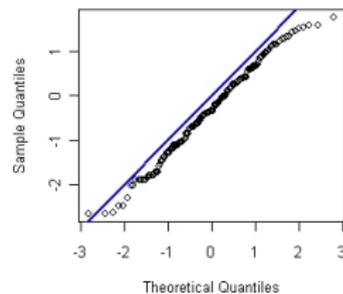


Beispiel: Normalverteilung

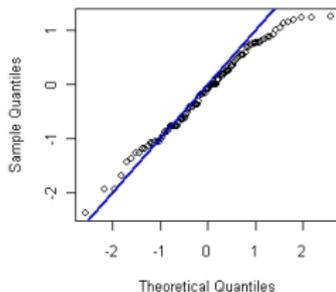
first 100 qq-plot



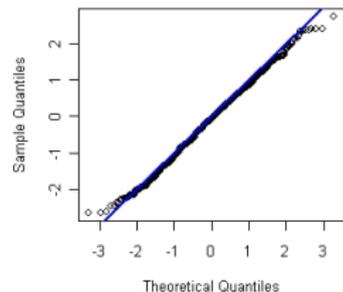
next 100 qq-plot



last 100 qq-plot



qq-plot



Übersicht

- 1 Grundlegendes über Markoffketten
- 2 Der Metropolis-Hastings-Algorithmus
- 3 Anwendung auf das Logit-Modell**
- 4 Der Gibbs-Sampler

Challenger-Beispiel

Am 28. Januar 1986 explodierte die Raumfähre Challenger mit sieben Astronauten an Bord nach dem Start aufgrund von Materialermüdungserscheinungen an den O-Ringen.

Zum Startzeitpunkt herrschte eine ungewöhnlich niedrige Außentemperatur von 31°F (ca. 0°C).

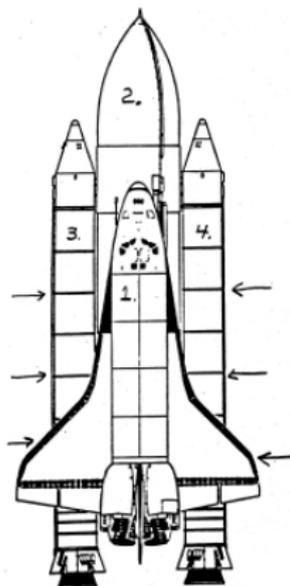
Nachfolgend soll der Zusammenhang zwischen der Temperatur und der Überbeanspruchung der Dichtungsringe analysiert werden.

Dazu im Folgenden:

X : Außentemperatur in $^{\circ}\text{F}$,

$Y = 1$ bzw. $Y = 0$: Dichtungsringausfall ja bzw. nein

Challenger-Beispiel: Datensatz, n=23 Messpunkte



X	Y	X	Y
66	0	67	0
70	1	53	1
69	0	67	0
68	0	75	0
67	0	70	0
72	0	81	0
73	0	76	0
70	0	79	0
57	1	75	0
63	1	76	0
70	1	58	1
78	0	—	—

Logistische Regression

Modell: Logistische Regression (Logit-Modell)

Sei Y binäre, Bernoulli-verteilte Zufallsvariable mit Parameter $p(X)$, wobei X eine reellwertige Zufallsvariable. Es gelte:

$$\begin{aligned}\mathbb{P}_\beta(Y = 1|X = x) &=: p \equiv p(x) = \frac{\exp(\alpha + \beta x)}{1 + \exp(\alpha + \beta x)} \\ &= 1 - \mathbb{P}_\beta(Y = 0|X = x).\end{aligned}$$

Die **Logit-Funktion** $\text{logit}(p) = \log(p/(1-p)) = \alpha + \beta x$ hänge also **linear** von $X = x$ ab.

Logistische Regression

Bayes-Ansatz zur Schätzung der Parameter α und β :

$$\begin{aligned} \ell((\alpha, \beta), \mathbf{Daten}) &= \prod_{i=1}^{23} \mathbb{P}(Y_i = y_i | X_i = x_i) \\ &= \prod_i \left(\frac{\exp(\alpha + \beta x_i)}{1 + \exp(\alpha + \beta x_i)} \right)^{y_i} \left(\frac{1}{1 + \exp(\alpha + \beta x_i)} \right)^{1-y_i}. \end{aligned}$$

A priori-Verteilung:

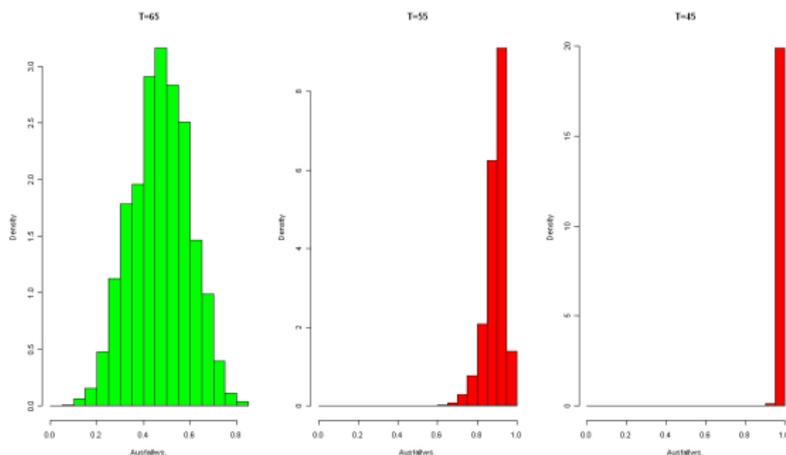
$$f_{(\tilde{\alpha}, \tilde{\beta})}(\alpha, \beta) = f_{\tilde{\alpha} | \tilde{\beta} = \beta}(\alpha) f_{\tilde{\beta}}(\beta)$$

(Hierarchisches Bayes-Verfahren)

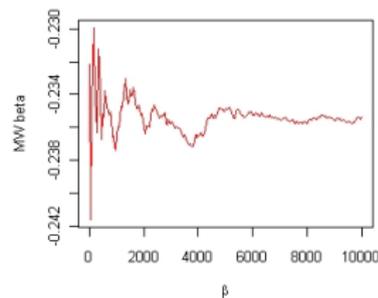
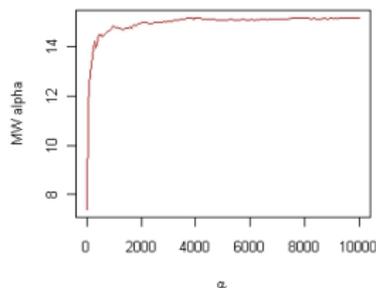
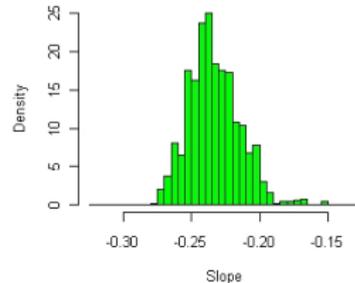
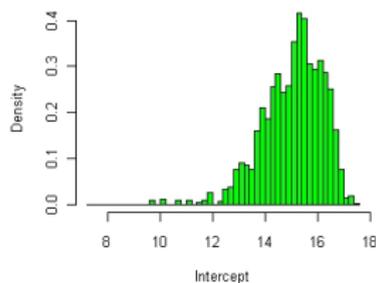
Metropolis-Hastings im Logit-Modell

Um von der Zielverteilung (A posteriori-Verteilung) zu sampeln, wird nun ein Metropolis-Hastings-Algorithmus implementiert.

Resultierende Histogramme für die Ausfallwahrscheinlichkeit bei verschiedenen Temperaturen:



Konvergenz Metropolis-Hastings



Übersicht

- 1 Grundlegendes über Markoffketten
- 2 Der Metropolis-Hastings-Algorithmus
- 3 Anwendung auf das Logit-Modell
- 4 Der Gibbs-Sampler**

Aufgabenstellung Gibbs-Sampler

Es sei $X = (X_1, \dots, X_p)^T$ ein \mathbb{R}^p -wertiger Zufallsvektor.

Annahmen:

- 1 Weder von der gemeinsamen Verteilung mit Dichte f noch von den univariaten Randdichten f_1, \dots, f_p kann direkt simuliert werden.
- 2 Jedoch kann von den bedingten Verteilungen der

$$X_j | X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p \sim f_j(\cdot | X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p)$$

gesamlet werden.

Algorithmus: Gibbs-Sampler

Gegeben sei $\mathbf{x}^{(t)} = (x_1^{(t)}, \dots, x_p^{(t)})$. Generiere sukzessive

$$X_1^{(t+1)} \sim f_1(\cdot | x_2^{(t)}, \dots, x_p^{(t)}),$$

$$X_2^{(t+1)} \sim f_2(\cdot | x_1^{(t+1)}, x_3^{(t)}, \dots, x_p^{(t)}),$$

$$\vdots$$

$$X_p^{(t+1)} \sim f_p(\cdot | x_1^{(t+1)}, \dots, x_{p-1}^{(t+1)}).$$

⇒ Zusammensetzung von p Metropolis-Hastings-Algorithmen mit Akzeptanzwahrscheinlichkeiten 1 und

$$q_i(\tilde{\mathbf{x}} | \mathbf{x}^{(t)}) = \delta_{(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_p)} \\ \times (\tilde{x}_1, \dots, \tilde{x}_{i-1}, \tilde{x}_{i+1}, \dots, \tilde{x}_p) f_i(\tilde{x}_i | x_1^{(t)}, \dots, x_{i-1}^{(t)}, x_{i+1}^{(t)}, \dots, x_p^{(t)}).$$

Satz (Konvergenz des Gibbs-Samplers)

Ist die vom Gibbs-Sampler erzeugte Markoffkette $(\mathbf{X}^{(t)})$ ergodisch, so ist f die Dichte der stationären Verteilung und es gilt für jede Startverteilung μ :

$$\lim_{n \rightarrow \infty} \left\| \int \mathcal{K}^n(x, \cdot) \mu(dx) - f \right\|_{TV} = 0.$$

Bemerkung:

Ergodizität von $(\mathbf{X}^{(t)})$ kann charakterisiert werden.

Beispiel (bedingte Verteilungen herleitbar)

Bivariate Normalverteilung:

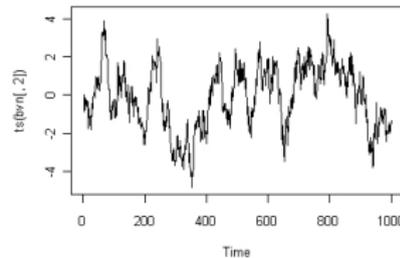
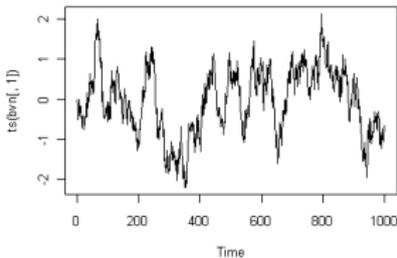
$$\text{Es sei } (X, Y)^T \sim \mathbf{N}_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \right).$$

Die bedingten univariaten Verteilungen sind gegeben durch

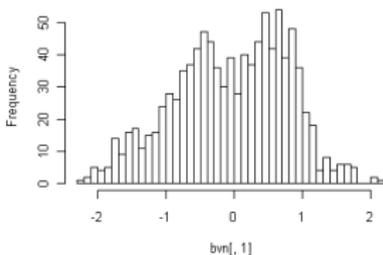
$$\mathcal{L}(X|Y=y) = \mathbf{N}\left(\rho y \sigma_1 / \sigma_2, (1 - \rho^2) \sigma_1^2\right),$$

$$\mathcal{L}(Y|X=x) = \mathbf{N}\left(\rho x \sigma_2 / \sigma_1, (1 - \rho^2) \sigma_2^2\right).$$

Illustration: Gibbs-Sampler bivariate Normalverteilung



Histogram of $bvn[, 1]$



Histogram of $bvn[, 2]$

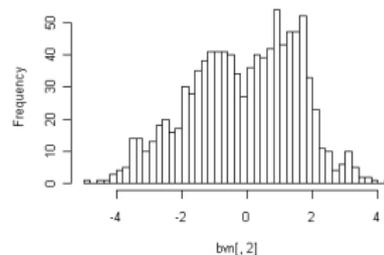
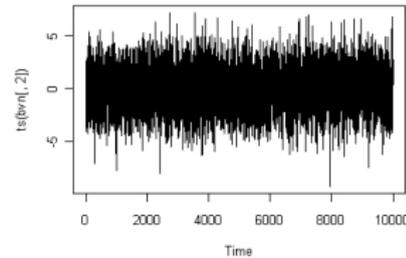
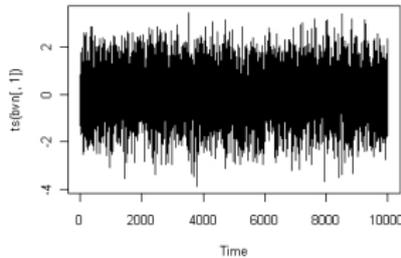
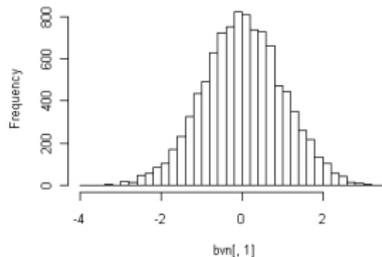


Illustration: Gibbs-Sampler bivariate Normalverteilung



Histogram of $bvn[, 1]$



Histogram of $bvn[, 2]$

