# Advanced Simulation - Lecture 7

George Deligiannidis

February 10th, 2019

# Markov chains - continuous space

- The state space $\mathbb{X}$ is now continuous, e.g. $\mathbb{R}^d$.

- $(X_t)_{t \geq 1}$ is a Markov chain if for any (measurable) set $A$,

$$\mathbb{P}(X_t \in A | X_1 = x_1, X_2 = x_2, ..., X_{t-1} = x_{t-1})$$
$$= \mathbb{P}(X_t \in A | X_{t-1} = x_{t-1}).$$

  *The future is conditionally independent of the past given the present.*

- We have

$$\mathbb{P}(X_t \in A | X_{t-1} = x) = \int_A K(x, y) \, dy = K(x, A),$$

  that is conditional on $X_{t-1} = x$, $X_t$ is a random variable which admits a probability density function $K(x, \cdot)$.

- $K : \mathbb{X}^2 \to \mathbb{R}$ is the **kernel** of the Markov chain.

# Markov chains - continuous space

- Denoting $\mu_1$ the pdf of $X_1$, we obtain directly

$$\mathbb{P}(X_1 \in A_1, ..., X_t \in A_t)$$

$$= \int_{A_1 \times \cdots \times A_t} \mu_1(x_1) \prod_{k=2}^{t} K(x_{k-1}, x_k) \, dx_1 \cdots dx_t.$$

- Denoting by $\mu_t$ the distribution of $X_t$, Chapman-Kolmogorov equation reads

$$\mu_t(y) = \int_{\mathbb{X}} \mu_{t-1}(x) K(x, y) \, dx$$

and similarly for $m > 1$

$$\mu_{t+m}(y) = \int_{\mathbb{X}} \mu_t(x) K^m(x, y) \, dx$$

where

$$K^m(x_t, x_{t+m}) = \int_{\mathbb{X}^{m-1}} \prod_{k=t+1}^{t+m} K(x_{k-1}, x_k) \, dx_{t+1} \cdots dx_{t+m-1}.$$

## Example

- Consider the autoregressive (AR) model

$$X_t = \rho X_{t-1} + V_t$$

where $V_t \overset{\text{i.i.d.}}{\sim} \mathcal{N}\left(0, \tau^2\right)$. This defines a Markov chain such that

$$K\left(x, y\right) = \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{1}{2\tau^2}\left(y - \rho x\right)^2\right).$$

- We also have

$$X_{t+m} = \rho^m X_t + \sum_{k=1}^{m} \rho^{m-k} V_{t+k}$$

so in the Gaussian case

$$K^m\left(x, y\right) = \frac{1}{\sqrt{2\pi\tau_m^2}} \exp\left(-\frac{1}{2}\frac{\left(y - \rho^m x\right)^2}{\tau_m^2}\right)$$

with $\tau_m^2 = \tau^2 \sum_{k=1}^{m}\left(\rho^2\right)^{m-k} = \tau^2 \frac{1 - \rho^{2m}}{1 - \rho^2}$.

# Irreducibility and aperiodicity

## Definition

Given a probability measure $\mu$ over $\mathbb{X}$, a Markov chain is
$\mu$-irreducible if

$$\forall x \in \mathbb{X} \quad \forall A : \mu(A) > 0 \quad \exists t \in \mathbb{N} \quad K^t(x, A) > 0.$$

A $\mu$-irreducible Markov chain of transition kernel $K$ is
periodic if there exists some partition of the state space
$\mathbb{X}_1, ..., \mathbb{X}_d$ for $d \geq 2$, such that

$$\forall i, j, t, s : \ \mathbb{P}\left( X_{t+s} \in \mathbb{X}_j \,\middle|\, X_t \in \mathbb{X}_i \right) = \left\{ \begin{array}{ll} 1 & j = i + s \ \mathrm{mod} \ d \\ 0 & \mathrm{otherwise.} \end{array} \right. .$$

Otherwise the chain is aperiodic.

# Recurrence and Harris Recurrence

For any measurable set $A$ of $\mathbb{X}$, let

$$\eta_A = \sum_{k=1}^{\infty} \mathbb{1}_A(X_k),$$

*the number of visits to the set $A$.*

## Definition

A $\mu$-irreducible Markov chain is recurrent if for any measurable set $A \subset \mathbb{X} : \mu(A) > 0$, then

$$\forall x \in A \quad \mathbb{E}_x(\eta_A) = \infty.$$

A $\mu$-irreducible Markov chain is Harris recurrent if for any measurable set $A \subset \mathbb{X} : \mu(A) > 0$, then

$$\forall x \in \mathbb{X} \quad \mathbb{P}_x(\eta_A = \infty) = 1.$$

Harris recurrence is stronger than recurrence.

# Invariant Distribution and Reversibility

## Definition

A distribution of density $\pi$ is invariant or *stationary* for a Markov kernel $K$, if

$$\int_{\mathbb{X}} \pi(x) K(x, y) \, dx = \pi(y).$$

A Markov kernel $K$ is $\pi$-reversible if

$$\forall f \quad \iint f(x, y) \pi(x) K(x, y) \, dx \, dy$$

$$= \iint f(y, x) \pi(x) K(x, y) \, dx \, dy$$

where $f$ is a bounded measurable function.

# Detailed balance

In practice it is easier to check the detailed balance condition:

$$\forall x, y \in \mathbb{X} \quad \pi(x)K(x, y) = \pi(y)K(y, x)$$

### Lemma

*If detailed balance holds, then $\pi$ is invariant for $K$ and $K$ is $\pi$-reversible.*

Example: the Gaussian AR process is $\pi$-reversible, $\pi$-invariant for

$$\pi(x) = \mathcal{N}\left(x; 0, \frac{\tau^2}{1 - \rho^2}\right)$$

when $|\rho| < 1$.

# Law of Large Numbers

## Theorem

*Suppose the Markov chain $\{X_i; i \geq 0\}$ is $\pi$−irreducible, with invariant distribution $\pi$, and suppose that $X_0 = x$.*
*Then for any $\pi$-integrable function $\varphi : \mathbb{X} \to \mathbb{R}$:*

$$\lim_{t \to \infty} \frac{1}{t} \sum_{i=1}^{t} \varphi(X_i) = \int_{\mathbb{X}} \varphi(w) \pi(w) \, dw$$

*almost surely, for $\pi$−almost every $x$.*

*If the chain in addition is Harris recurrent then this holds for **every** starting value $x$.*

# Convergence

## Theorem

*Suppose the kernel $K$ is $\pi$-irreducible, $\pi$-invariant, aperiodic. Then, we have*

$$\lim_{t \to \infty} \int_{\mathbb{X}} |K^t(x, y) - \pi(y)| \, dy = 0$$

*for $\pi$-almost all starting values $x$.*

Under some additional conditions, one can prove that there exists a $\rho < 1$ and a function $M : \mathbb{X} \to \mathbb{R}^+$ such that for all measurable sets $A$ and all $n$

$$|K^n(x, A) - \pi(A)| \leq M(x)\rho^n.$$

The chain is then said to be **geometrically ergodic**.

# Central Limit Theorem

## Theorem

*Under regularity conditions, for a Harris recurrent, $\pi$-invariant Markov chain, we can prove*

$$\sqrt{t}\left[\frac{1}{t}\sum_{i=1}^{t}\varphi(X_i) - \int_{\mathbb{X}}\varphi(x)\,\pi(x)\,\mathrm{d}x\right] \xrightarrow[t\to\infty]{\mathscr{D}} \mathscr{N}\left(0, \sigma^2(\varphi)\right),$$

*where the asymptotic variance can be written*

$$\sigma^2(\varphi) = \mathbb{V}_\pi[\varphi(X_1)] + 2\sum_{k=2}^{\infty}\mathrm{Cov}_\pi[\varphi(X_1), \varphi(X_k)].$$

This formula shows that (positive) correlations increase the asymptotic variance, compared to i.i.d. samples for which the variance would be $\mathbb{V}_\pi(\varphi(X))$.

# Central Limit Theorem

Example: for the AR Gaussian model,
$\pi(x) = \mathcal{N}\left(x; 0, \tau^2/(1-\rho^2)\right)$ for $|\rho| < 1$ and

$$\mathbb{C}\text{ov}(X_1, X_k) = \rho^{k-1} \mathbb{V}[X_1] = \rho^{k-1} \frac{\tau^2}{1-\rho^2}.$$

Therefore with $\varphi(x) = x$,

$$\sigma^2(\varphi) = \frac{\tau^2}{1-\rho^2}\left(1 + 2\sum_{k=1}^{\infty} \rho^k\right) = \frac{\tau^2}{1-\rho^2}\frac{1+\rho}{1-\rho} = \frac{\tau^2}{(1-\rho)^2},$$

which increases when $\rho \to 1$.

# Markov chain Monte Carlo

- We are interested in sampling from a distribution $\pi$, for instance a posterior distribution in a Bayesian framework.

- Markov chains with $\pi$ as invariant distribution can be constructed to approximate expectations with respect to $\pi$.

- For example, the Gibbs sampler generates a Markov chain targeting $\pi$ defined on $\mathbb{R}^d$ using the full conditionals

$$\pi(x_i \mid x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_d).$$

# Gibbs Sampling

- Assume you are interested in sampling from

$$\pi(x) = \pi(x_1, x_2, ..., x_d), \quad x \in \mathbb{R}^d.$$

- Notation: $x_{-i} := (x_1, ..., x_{i-1}, x_{i+1}, ..., x_d)$.

  **Systematic scan Gibbs sampler**. Let $\left(X_1^{(1)}, ..., X_d^{(1)}\right)$ be the initial state then iterate for $t = 2, 3, ...$

  **1.** Sample $X_1^{(t)} \sim \pi_{X_1|X_{-1}}\left(\cdot \,|\, X_2^{(t-1)}, ..., X_d^{(t-1)}\right)$.

  $\vdots$

  **j.** Sample $X_j^{(t)} \sim \pi_{X_j|X_{-j}}\left(\cdot \,|\, X_1^{(t)}, ..., X_{j-1}^{(t)}, X_{j+1}^{(t-1)}, ..., X_d^{(t-1)}\right)$.

  $\vdots$

  **d.** Sample $X_d^{(t)} \sim \pi_{X_d|X_{-d}}\left(\cdot \,|\, X_1^{(t)}, ..., X_{d-1}^{(t)}\right)$.

# Gibbs Sampling

A few questions one can ask about this algorithm:

- Is the joint distribution $\pi$ uniquely specified by the conditional distributions $\pi_{X_i|X_{-i}}$?
- A: Not in general![1]
- Does the Gibbs sampler provide a Markov chain with the correct stationary distribution $\pi$?
- A: Not in general!
- If yes, does the Markov chain converge towards this invariant distribution?
- It will turn out to be the case under some mild conditions.

---

[1] J.P. Hobert, C.P. Robert, C. Goutis, Connectedness conditions for the convergence of the Gibbs sampler (1997)

# Hammersley-Clifford Theorem

## Theorem

*Consider a distribution with continuous density $\pi(x_1, x_2, ..., x_d)$ such that*

$$supp(\pi) = supp\left(\bigotimes_{i=1}^{d} \pi_{X_i}\right).$$

*Then for any $(z_1, ..., z_d) \in supp(\pi)$, we have*

$$\pi(x_1, x_2, ..., x_d) \propto \prod_{j=1}^{d} \frac{\pi_{X_j|X_{-j}}\left(x_j \,|\, x_{1:j-1}, z_{j+1:d}\right)}{\pi_{X_j|X_{-j}}\left(z_j \,|\, x_{1:j-1}, z_{j+1:d}\right)}.$$

The condition above is known as the **positivity condition**.

Equivalently, if $\pi_{X_i}(x_i) > 0$ for $i = 1, ..., d$, then

$$\pi(x_1, ..., x_d) > 0.$$

# Proof of Hammersley-Clifford Theorem

### Proof.

We have

$$\pi(x_{1:d-1}, x_d) = \pi_{X_d|X_{-d}}(x_d | x_{1:d-1})\pi(x_{1:d-1}),$$
$$\pi(x_{1:d-1}, z_d) = \pi_{X_d|X_{-d}}(z_d | x_{1:d-1})\pi(x_{1:d-1}).$$

Therefore

$$\begin{aligned}
\pi(x_{1:d}) &= \pi(x_{1:d-1}, z_d)\frac{\pi(x_{1:d-1}, x_d)}{\pi(x_{1:d-1}, z_d)} \\
&= \pi(x_{1:d-1}, z_d)\frac{\pi(x_{1:d-1}, x_d)/\pi(x_{1:d-1})}{\pi(x_{1:d-1}, z_d)/\pi(x_{1:d-1})} \\
&= \pi(x_{1:d-1}, z_d)\frac{\pi_{X_d|X_{1:d-1}}(x_d | x_{1:d-1})}{\pi_{X_d|X_{1:d-1}}(z_d | x_{1:d-1})}.
\end{aligned}$$

## Proof.

Similarly, we have

$$\pi(x_{1:d-1}, z_d) = \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi(x_{1:d-1}, z_d)}{\pi(x_{1:d-2}, z_{d-1}, z_d)}$$

$$= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi(x_{1:d-1}, z_d)/\pi(x_{1:d-2}, z_d)}{\pi(x_{1:d-2}, z_{d-1}, z_d)/\pi(x_{1:d-2}, z_d)}$$

$$= \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi_{X_{d-1}|X^{-(d-1)}}(x_{d-1} \mid x_{1:d-2}, z_d)}{\pi_{X_{d-1}|X^{-(d-1)}}(z_{d-1} \mid x_{1:d-2}, z_d)}$$

hence

$$\pi(x_{1:d}) = \pi(x_{1:d-2}, z_{d-1}, z_d) \frac{\pi_{X_{d-1}|X_{-(d-1)}}(x_{d-1} \mid x_{1:d-2}, z_d)}{\pi_{X_{d-1}|X_{-(d-1)}}(z_{d-1} \mid x_{1:d-2}, z_d)}$$

$$\times \frac{\pi_{X_d|X_{-d}}(x_d \mid x_{1:d-1})}{\pi_{X_d|X_{-d}}(z_d \mid x_{1:d-1})}$$

## Proof.

By $z \in \operatorname{supp}(\pi)$ we have that $\pi_{X_i}(z_i) > 0$ for all $i$. Also, we are allowed to suppose that $\pi_{X_i}(x_i) > 0$ for all $i$. Thus all the conditional probabilities we introduce are positive since

$$\pi_{X_j | X^{-j}}(x_j \mid x_1, \ldots, x_{j-1}, z_{j+1}, \ldots, z_d)$$
$$= \frac{\pi(x_1, \ldots, x_{j-1}, x_j, z_{j+1}, \ldots, z_d)}{\pi(x_1, \ldots, x_{j-1}, z_j, z_{j+1}, \ldots, z_d)} > 0.$$

By iterating we have the theorem. $\qquad\square$

# Example: Non-Integrable Target

- Consider the following conditionals on $\mathbb{R}^+$

$$\pi_{X_1|X_2}(x_1|x_2) = x_2 \exp(-x_2 x_1)$$
$$\pi_{X_2|X_1}(x_2|x_1) = x_1 \exp(-x_1 x_2).$$

  We might expect that these full conditionals define a joint probability density $\pi(x_1, x_2)$.

- Hammersley-Clifford would give

$$\begin{aligned}
\pi(x_1, x_2, ..., x_d) &\propto \frac{\pi_{X_1|X_2}(x_1|z_2)}{\pi_{X_1|X_2}(z_1|z_2)} \frac{\pi_{X_2|X_1}(x_2|x_1)}{\pi_{X_2|X_1}(z_2|x_1)} \\
&= \frac{z_2 \exp(-z_2 x_1) x_1 \exp(-x_1 x_2)}{z_2 \exp(-z_2 z_1) x_1 \exp(-x_1 z_2)} \propto \exp(-x_1 x_2).
\end{aligned}$$

- However $\iint \exp(-x_1 x_2)\, dx_1\, dx_2 = \infty$ so
  $\pi_{X_1|X_2}(x_1|x_2) = x_2 \exp(-x_2 x_1)$ and
  $\pi_{X_2|X_1}(x_1|x_2) = x_1 \exp(-x_1 x_2)$ are not compatible.

# Example: Positivity condition violated



Figure: Gibbs sampling targeting
$\pi(x, y) \propto \mathbb{1}_{[-1,0] \times [-1,0] \cup [0,1] \times [0,1]}(x, y)$.

Positivity condition violated: any density of the form

$$f(x) = \alpha \mathbb{1}_{[-1,0]^2} + (1 - \alpha) \mathbb{1}_{[0,1]^2},$$

has same conditionals.

# Invariance of the Gibbs sampler I

The kernel of the Gibbs sampler (case $d = 2$) is

$$K(x^{(t-1)}, x^{(t)}) = \pi_{X_1|X_2}(x_1^{(t)} \mid x_2^{(t-1)}) \pi_{X_2|X_1}(x_2^{(t)} \mid x_1^{(t)})$$

Case $d > 2$:

$$K(x^{(t-1)}, x^{(t)}) = \prod_{j=1}^{d} \pi_{X_j|X_{-j}}(x_j^{(t)} \mid x_{1:j-1}^{(t)}, x_{j+1:d}^{(t-1)})$$

### Proposition

*The systematic scan Gibbs sampler kernel admits $\pi$ as invariant distribution.*

# Invariance of the Gibbs sampler II

### Proof for $d = 2$.

Let $x = (x_1, x_2)$ and $y = (y_1, y_2)$. Then we have

$$\int K(x,y)\pi(x)dx = \int \pi(y_2 \mid y_1)\pi(y_1 \mid x_2)\pi(x_1, x_2)dx_1\,dx_2$$

$$= \pi(y_2 \mid y_1)\int \pi(y_1 \mid x_2)\pi(x_2)dx_2$$

$$= \pi(y_2 \mid y_1)\pi(y_1) = \pi(y_1, y_2) = \pi(y).$$

$\square$

# Irreducibility and Recurrence

## Proposition

*Assume $\pi$ satisfies the positivity condition, then the Gibbs sampler yields a $\pi$–irreducible and recurrent Markov chain.*

## Proof.

**Recurrence.** Will follow from irreducibility and the fact that $\pi$ is invariant, [a]

**(One step)Irreducibility.** Let $\mathbb{X} \subset \mathbb{R}^d$, such that $\pi(\mathbb{X}) = 1$. Write $K$ for the kernel and let $A \subset \mathbb{X}$ such that $\pi(A) > 0$. Then for any $x \in \mathbb{X}$

$$K(x, A) = \int_A K(x, y) \mathrm{d}y$$
$$= \int_A \pi_{X_1 | X_{-1}}(y_1 \mid x_2, \ldots, x_d) \times \cdots \times \pi_{X_d | X_{-d}}(y_d \mid y_1, \ldots, y_{d-1}) \mathrm{d}y.$$

---

[a] Meyn and Tweedie, Markov chains and stochastic stability, Prop'n 10.1.1.

## Proof.

Thus if for some $x \in \mathbb{X}$ and $A$ with $\pi(A) > 0$ we have $K(x, A) = 0$, we must have that

$$\pi_{X_1 | X_{-1}}(y_1 \mid x_2, \ldots, x_d) \times \cdots \times \pi_{X_d | X_{-d}}(y_d \mid y_1, \ldots, y_{d-1}) = 0,$$

for almost all $y = (y_1, \ldots, y_d) \in A$.

Therefore, by the Hammersley-Clifford theorem, we must also have that

$$\pi(y_1, y_2, \ldots, y_d) \propto \prod_{j=1}^{d} \frac{\pi_{X_j | X_{-j}} \left( y_j \mid y_{1:j-1}, x_{j+1:d} \right)}{\pi_{X_j | X_{-j}} \left( x_j \mid y_{1:j-1}, x_{j+1:d} \right)} = 0,$$

for almost all $y = (y_1, \ldots, y_d) \in A$ and thus $\pi(A) = 0$ obtaining a contradiction.

Note: Positivity not necessary for irreducibility; e.g. $f \propto \mathbb{1}_{|x| \le 1}$.

# LLN for Gibbs Sampler

### Theorem

*If the positivity condition is satisfied then for any $\pi$-integrable function $\varphi : \mathbb{X} \to \mathbb{R}$:*

$$\lim \frac{1}{t} \sum_{i=1}^{t} \varphi\left(X^{(i)}\right) = \int_{\mathbb{X}} \varphi(x)\,\pi(x)\,\mathrm{d}x$$

*for $\pi$–almost all starting values $X^{(1)}$.*

# Example: Bivariate Normal Distribution

- Let $X := (X_1, X_2) \sim \mathcal{N}(\mu, \Sigma)$ where $\mu = (\mu_1, \mu_2)$ and

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho \\ \rho & \sigma_2^2 \end{pmatrix}.$$

- The Gibbs sampler proceeds as follows in this case

(a) Sample $X_1^{(t)} \sim \mathcal{N}\left(\mu_1 + \rho / \sigma_2^2 \left(X_2^{(t-1)} - \mu_2\right), \sigma_1^2 - \rho^2 / \sigma_2^2\right)$

(b) Sample $X_2^{(t)} \sim \mathcal{N}\left(\mu_2 + \rho / \sigma_1^2 \left(X_1^{(t)} - \mu_1\right), \sigma_2^2 - \rho^2 / \sigma_1^2\right).$

- By proceeding this way, we generate a Markov chain $X^{(t)}$ whose successive samples are correlated. If successive values of $X^{(t)}$ are strongly correlated, then we say that the Markov chain mixes slowly.

# Bivariate Normal Distribution



Figure: Case where $\rho = 0.1$, first 100 steps.

# Bivariate Normal Distribution



Figure: Case where $\rho = 0.99$, first 100 steps.

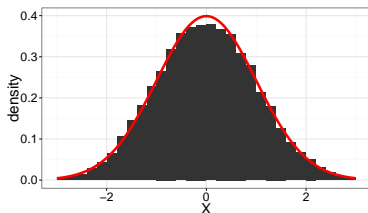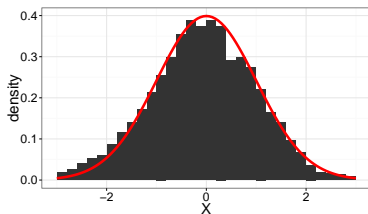# Bivariate Normal Distribution



(a) Figure A

(b) Figure B

Figure: Histogram of the first component of the chain after 1000 iterations. Small $\rho$ on the left, large $\rho$ on the right.
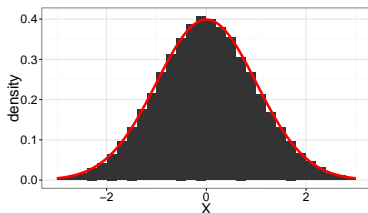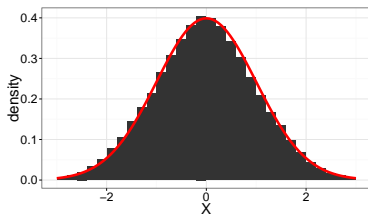
# Bivariate Normal Distribution



(a) b



(b) b

Figure: Histogram of the first component of the chain after 10000 iterations. Small $\rho$ on the left, large $\rho$ on the right.

# Bivariate Normal Distribution



(a) Figure A



(b) Figure B

Figure: Histogram of the first component of the chain after 100000 iterations. Small $\rho$ on the left, large $\rho$ on the right.

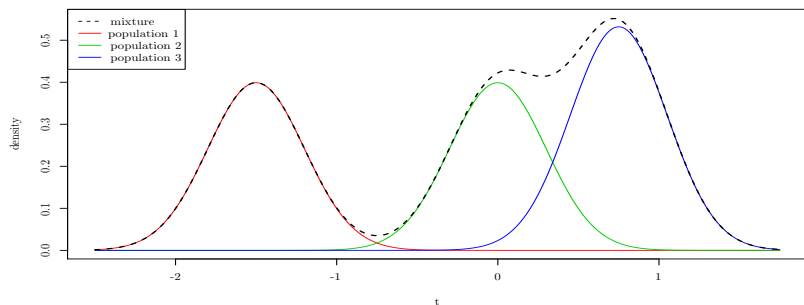# Gibbs Sampling and Auxiliary Variables

- Gibbs sampling requires sampling from $\pi_{X_j | X_{-j}}$.

- In many scenarios, we can include a set of auxiliary variables $Z_1, ..., Z_p$ and have an "extended" distribution of joint density $\overline{\pi}(x_1, ..., x_d, z_1, ..., z_p)$ such that

$$\int \overline{\pi}(x_1, ..., x_d, z_1, ..., z_p) \, dz_1 ... dz_d = \pi(x_1, ..., x_d).$$

  which is such that its full conditionals are easy to sample.

- Mixture models, Capture-recapture models, Tobit models, Probit models etc.

# Mixtures of Normals



- Independent data $y_1, ..., y_n$

$$Y_i | \theta \sim \sum_{k=1}^{K} p_k \mathcal{N}\left(\mu_k, \sigma_k^2\right)$$

where $\theta = \left(p_1, ..., p_K, \mu_1, ..., \mu_K, \sigma_1^2, ..., \sigma_K^2\right)$.

# Bayesian Model

- Likelihood function

$$p(y_1,...,y_n|\theta) = \prod_{i=1}^n p(y_i|\theta) = \prod_{i=1}^n \left( \sum_{k=1}^K \frac{p_k}{\sqrt{2\pi\sigma_k^2}} \exp\left( -\frac{(y_i - \mu_k)^2}{2\sigma_k^2} \right) \right).$$

Let's fix $K = 2$, $\sigma_k^2 = 1$ and $p_k = 1/K$ for all $k$.

- Prior model

$$p(\theta) = \prod_{k=1}^K p(\mu_k)$$

where

$$\mu_k \sim \mathcal{N}(\alpha_k, \beta_k).$$

Let us fix $\alpha_k = 0, \beta_k = 1$ for all $k$.

- Not obvious how to sample $p(\mu_1 \mid \mu_2, y_1, \ldots, y_n)$.

# Auxiliary Variables for Mixture Models

- Associate to each $Y_i$ an auxiliary variable $Z_i \in \{1, ..., K\}$ such that

$$\mathbb{P}(Z_i = k | \theta) = p_k \text{ and } Y_i | Z_i = k, \theta \sim \mathcal{N}\left(\mu_k, \sigma_k^2\right)$$

so that

$$p\left(y_i | \theta\right) = \sum_{k=1}^{K} \mathbb{P}(Z_i = k) \mathcal{N}\left(y_i; \mu_k, \sigma_k^2\right)$$

- The extended posterior is given by

$$p\left(\theta, z_1, ..., z_n | y_1, ..., y_n\right) \propto p(\theta) \prod_{i=1}^{n} \mathbb{P}(z_i | \theta) p\left(y_i | z_i, \theta\right).$$

- Gibbs samples alternately

$$\mathbb{P}(z_{1:n} | y_{1:n}, \mu_{1:K})$$
$$p\left(\mu_{1:K} | y_{1:n}, z_{1:n}\right).$$

# Gibbs Sampling for Mixture Model

- We have
$$\mathbb{P}\left(z_{1:n}\middle| y_{1:n},\theta\right) = \prod_{i=1}^{n} \mathbb{P}\left(z_i\middle| y_i,\theta\right)$$

  where

$$\mathbb{P}\left(z_i\middle| y_i,\theta\right) = \frac{\mathbb{P}\left(z_i|\theta\right)p\left(y_i\middle| z_i,\theta\right)}{\sum_{k=1}^{K}\mathbb{P}\left(z_i=k|\theta\right)p\left(y_i\middle| z_i=k,\theta\right)}$$

- Let $n_k = \sum_{i=1}^{n} \mathbf{1}_{\{k\}}\left(z_i\right), n_k\overline{y}_k = \sum_{i=1}^{n} y_i\mathbf{1}_{\{k\}}\left(z_i\right)$ then

$$\mu_k\middle| z_{1:n}, y_{1:n} \sim \mathcal{N}\left(\frac{n_k\overline{y}_k}{1+n_k}, \frac{1}{1+n_k}\right).$$

# Mixtures of Normals



Figure: 200 points sampled from $\frac{1}{2}\mathcal{N}(-2,1) + \frac{1}{2}\mathcal{N}(2,1)$.
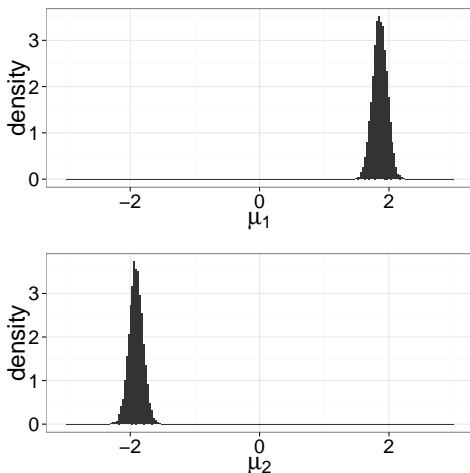
# Mixtures of Normals



Figure: Histogram of the parameters obtained by 10,000 iterations of Gibbs sampling.
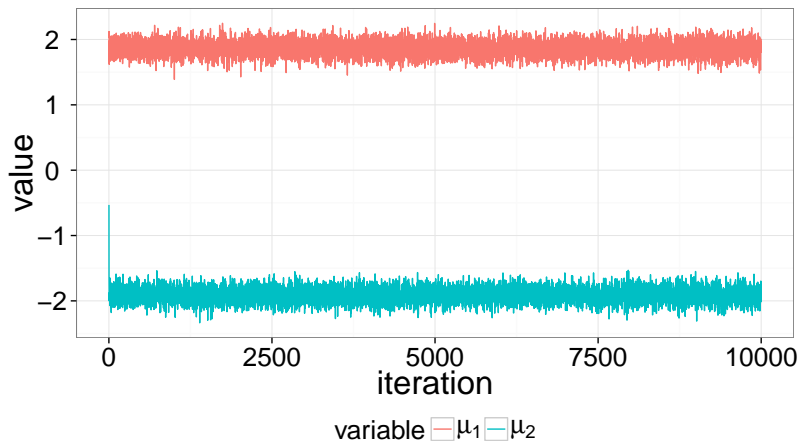
# Mixtures of Normals



Figure: Traceplot of the parameters obtained by 10,000 iterations of Gibbs sampling.

# Gibbs sampling in practice

- Many posterior distributions can be automatically decomposed into conditional distributions by computer programs.

- This is the idea behind BUGS (Bayesian inference Using Gibbs Sampling), JAGS (Just another Gibbs Sampler).

# Gibbs Recap

- Given a target $\pi(x) = \pi(x_1, x_2, ..., x_d)$, Gibbs sampling works by sampling from $\pi_{X_j|X_{-j}}(x_j|x_{-j})$ for $j = 1, ..., d$.

- Sampling exactly from one of these full conditionals might be a hard problem itself.

- Even if it is possible, the Gibbs sampler might converge slowly if components are highly correlated.

- If the components are not highly correlated then Gibbs sampling performs well, even when $d \to \infty$, e.g. with an error increasing "only" polynomially with $d$.

- Metropolis–Hastings algorithm (1953, 1970) is a more general algorithm that can bypass these problems.

- Additionally Gibbs can be recovered as a special case.

# Metropolis–Hastings algorithm

- Target distribution on $\mathbb{X} = \mathbb{R}^d$ of density $\pi(x)$.
- Proposal distribution: for any $x, x' \in \mathbb{X}$, we have $q(x'|x) \geq 0$ and $\int_{\mathbb{X}} q(x'|x) \, dx' = 1$.
- Starting with $X^{(1)}$, for $t = 2, 3, \ldots$
  (a) Sample $X^\star \sim q\left(\cdot | X^{(t-1)}\right)$.
  (b) Compute

$$\alpha\left(X^\star | X^{(t-1)}\right) = \min\left(1, \frac{\pi(X^\star) \, q\left(X^{(t-1)} | X^\star\right)}{\pi(X^{(t-1)}) \, q(X^\star | X^{(t-1)})}\right).$$

  (c) Sample $U \sim \mathcal{U}_{[0,1]}$. If $U \leq \alpha\left(X^\star | X^{(t-1)}\right)$, set $X^{(t)} = X^\star$, otherwise set $X^{(t)} = X^{(t-1)}$.
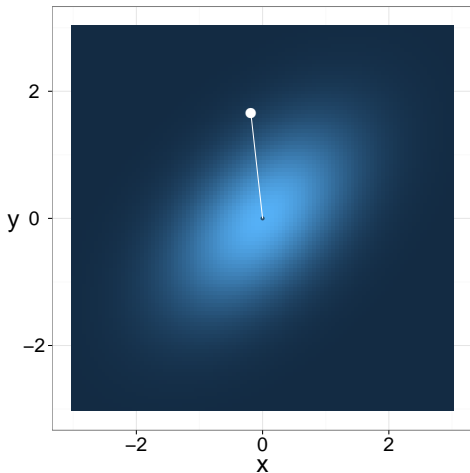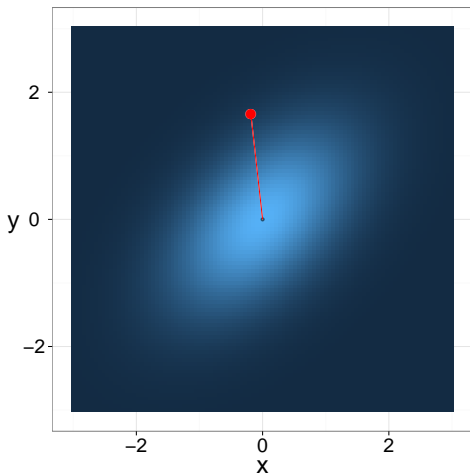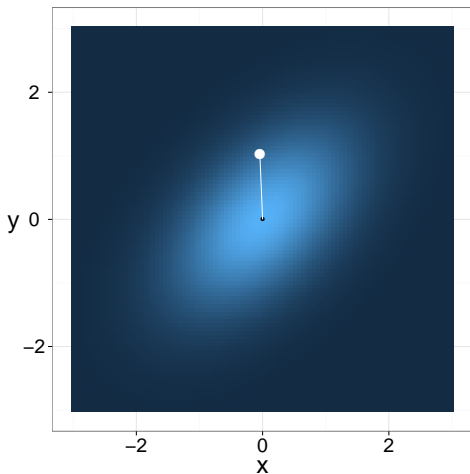
# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
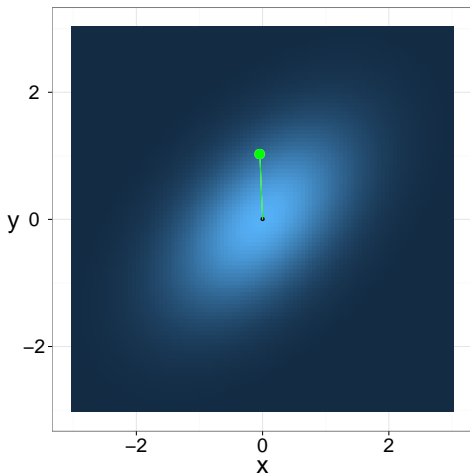


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
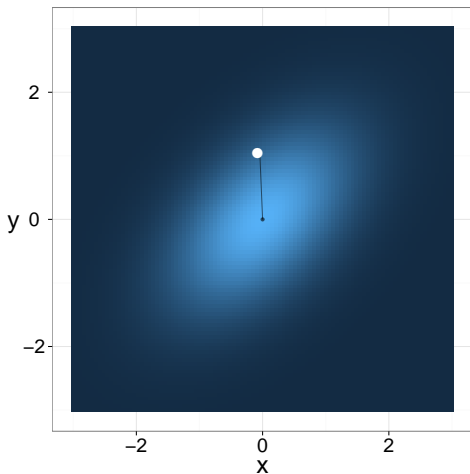


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
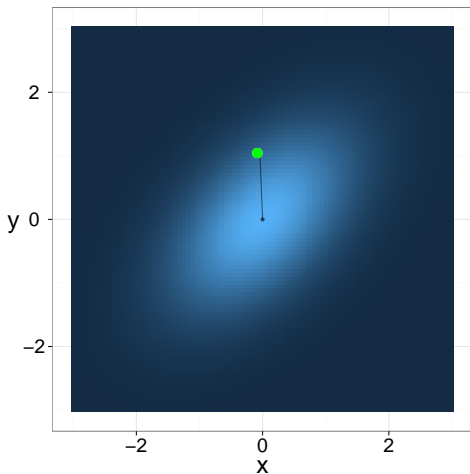


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
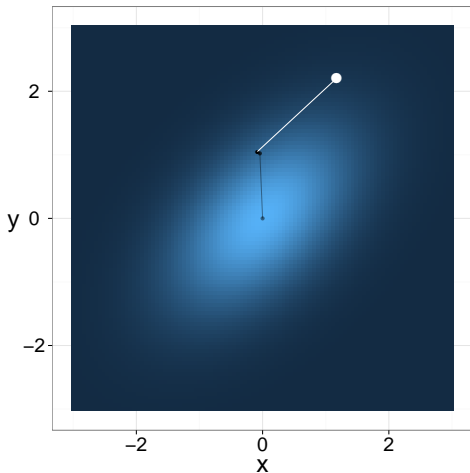


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
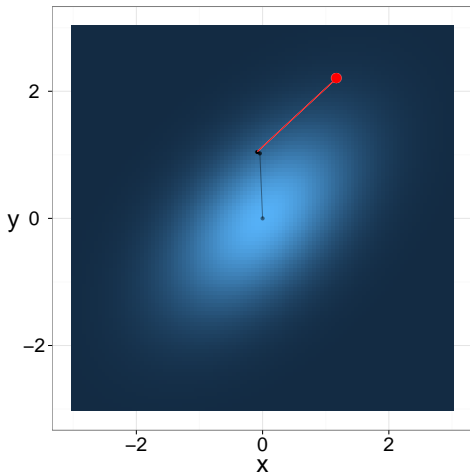


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
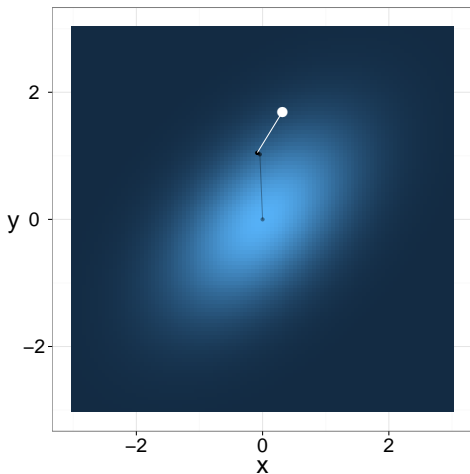


Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm
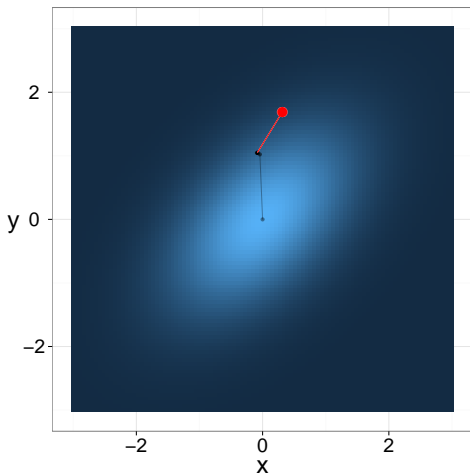


Figure: Metropolis–Hastings on a bivariate Gaussian target.

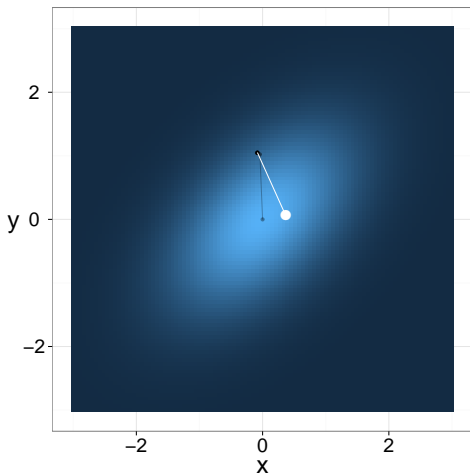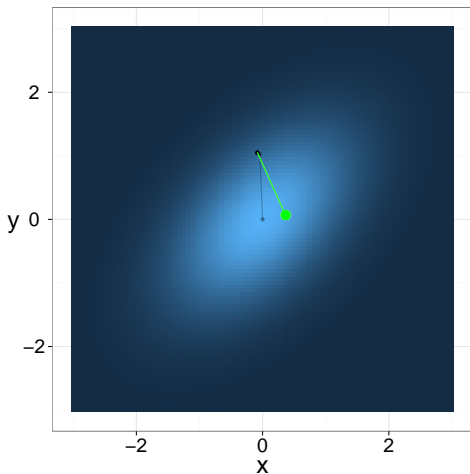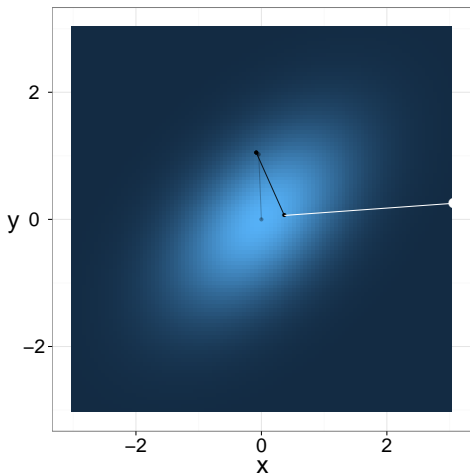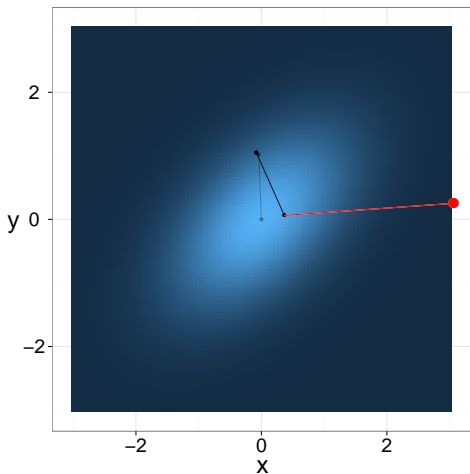# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target.

# Metropolis–Hastings algorithm

- Metropolis–Hastings only requires point-wise evaluations of $\pi(x)$ up to a normalizing constant; indeed if $\widetilde{\pi}(x) \propto \pi(x)$ then

$$\frac{\pi(x^\star)\, q\left(x^{(t-1)}\,\middle|\,x^\star\right)}{\pi(x^{(t-1)})\, q\left(x^\star\,\middle|\,x^{(t-1)}\right)} = \frac{\widetilde{\pi}(x^\star)\, q\left(x^{(t-1)}\,\middle|\,x^\star\right)}{\widetilde{\pi}(x^{(t-1)})\, q\left(x^\star\,\middle|\,x^{(t-1)}\right)}.$$

- At each iteration $t$, a candidate is proposed.

- The **average acceptance probability** from the current state is

$$a\left(x^{(t-1)}\right) := \int_{\mathbb{X}} \alpha\left(x\,\middle|\,x^{(t-1)}\right) q\left(x\,\middle|\,x^{(t-1)}\right) dx$$

  in which case $X^{(t)} = X$, otherwise $X^{(t)} = X^{(t-1)}$.

- This algorithm clearly defines a Markov chain $(X^{(t)})_{t \geq 1}$.

# Transition Kernel and Reversibility

### Lemma

*The kernel of the Metropolis–Hastings algorithm is given by*

$$K(y \mid x) \equiv K(x, y) = \alpha(y \mid x) q(y \mid x) + (1 - a(x)) \delta_x(y).$$

### Proof.

We have

$$K(x, y)$$
$$= \int q(x^\star \mid x) \{ \alpha(x^\star \mid x) \delta_{x^\star}(y) + (1 - \alpha(x^\star \mid x)) \delta_x(y) \} dx^\star$$
$$= q(y \mid x) \alpha(y \mid x) + \left\{ \int q(x^\star \mid x)(1 - \alpha(x^\star \mid x)) dx^\star \right\} \delta_x(y)$$
$$= q(y \mid x) \alpha(y \mid x) + \left\{ 1 - \int q(x^\star \mid x) \alpha(x^\star \mid x) dx^\star \right\} \delta_x(y)$$
$$= q(y \mid x) \alpha(y \mid x) + \left\{ 1 - a(x) \right\} \delta_x(y). \qquad \Box$$

# Reversibility

## Proposition

*The Metropolis–Hastings kernel $K$ is $\pi$–reversible and thus admit $\pi$ as invariant distribution.*

## Proof.

For any $x, y \in \mathbb{X}$, with $x \neq y$

$$
\begin{aligned}
\pi(x)K(x,y) &= \pi(x)q(y \mid x)\alpha(y \mid x) \\
&= \pi(x)q(y \mid x)\left(1 \wedge \frac{\pi(y)q(x \mid y)}{\pi(x)q(y \mid x)}\right) \\
&= \left(\pi(x)q(y \mid x) \wedge \pi(y)q(x \mid y)\right) \\
&= \pi(y)q(x \mid y)\left(\frac{\pi(x)q(y \mid x)}{\pi(y)q(x \mid y)} \wedge 1\right) = \pi(y)K(y,x).
\end{aligned}
$$

If $x = y$, then obviously $\pi(x)K(x,y) = \pi(y)K(y,x)$. $\qquad\square$

# Reducibility and periodicity of Metropolis–Hastings

- Consider the target distribution

$$\pi(x) = \left( \mathcal{U}_{[0,1]}(x) + \mathcal{U}_{[2,3]}(x) \right)/2$$

and the proposal distribution

$$q\left(x^\star \mid x\right) = \mathcal{U}_{(x-\delta, x+\delta)}\left(x^\star\right).$$

- The MH chain is reducible if $\delta \leq 1$: the chain stays either in $[0,1]$ or $[2,3]$.

- Note that the MH chain is aperiodic if it always has a non-zero chance of staying where it is.

# Some results

**Proposition**

If $q(x^\star|x) > 0$ for any $x, x^\star \in supp(\pi)$ then the Metropolis-Hastings chain is *irreducible*, in fact every state can be reached in a single step (strongly irreducible).

Less strict conditions in (Roberts & Rosenthal, 2004).

**Proposition**

If the MH chain is *irreducible* then it is also *Harris recurrent* (see Tierney, 1994).

# LLN for MH

## Theorem

*If the Markov chain generated by the Metropolis–Hastings sampler is $\pi$–irreducible, then we have for any integrable function $\varphi : \mathbb{X} \to \mathbb{R}$:*

$$\lim_{t \to \infty} \frac{1}{t} \sum_{i=1}^{t} \varphi \left( X^{(i)} \right) = \int_{\mathbb{X}} \varphi (x) \, \pi (x) \, dx$$

*for every starting value $X^{(1)}$.*

# Random Walk Metropolis–Hastings

- In the Metropolis–Hastings, pick $q(x^\star \mid x) = g(x^\star - x)$ with $g$ being a *symmetric* distribution, thus

$$X^\star = X + \varepsilon, \quad \varepsilon \sim g;$$

  e.g. $g$ is a zero-mean multivariate normal or t-student.

- Acceptance probability becomes

$$\alpha(x^\star \mid x) = \min\left(1, \frac{\pi(x^\star)}{\pi(x)}\right).$$

- We accept...
  - a move to a more probable state with probability $1$;
  - a move to a less probable state with probability

$$\pi(x^\star)/\pi(x) \le 1.$$

# Independent Metropolis–Hastings

- **Independent proposal**: a proposal distribution $q(x^\star \mid x)$ which does not depend on $x$.

  - Acceptance probability becomes

    $$\alpha(x^\star \mid x) = \min\left(1, \frac{\pi(x^\star)q(x)}{\pi(x)q(x^\star)}\right).$$

  - For instance, multivariate normal or t-student distribution.

- If $\pi(x)/q(x) < M$ for all $x$ and some $M < \infty$, then the chain is **uniformly ergodic**.

- The acceptance probability at stationarity is at least $1/M$ (Lemma 7.9 of Robert & Casella).

- On the other hand, if such an $M$ does not exist, the chain is not even geometrically ergodic!

# Choosing a good proposal distribution

- Goal: design a Markov chain with small correlation $\rho\left(X^{(t-1)}, X^{(t)}\right)$ between subsequent values (why?).

- Two sources of correlation:
  - between the current state $X^{(t-1)}$ and proposed value $X \sim q\left(\cdot \mid X^{(t-1)}\right)$,
  - correlation induced if $X^{(t)} = X^{(t-1)}$, if proposal is rejected.

- Trade-off: there is a compromise between
  - proposing large moves,
  - obtaining a decent acceptance probability.

- For multivariate distributions: covariance of proposal should reflect the covariance structure of the target.

# Choice of proposal

- Target distribution, we want to sample from

$$\pi(x) = \mathcal{N}\left(x; \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}\right).$$

- We use a random walk Metropolis—Hastings algorithm with

$$g(\varepsilon) = \mathcal{N}\left(\varepsilon; 0, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right).$$

- What is the optimal choice of $\sigma^2$?
- We consider three choices: $\sigma^2 = 0.1^2, 1, 10^2$.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 0.1^2$, the acceptance rate is $\approx 94\%$.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 0.1^2$, the acceptance rate is $\approx 94\%$.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 1$, the acceptance rate is $\approx 52\%$.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 1$, the acceptance rate is $\approx 52\%$.

# Metropolis–Hastings algorithm
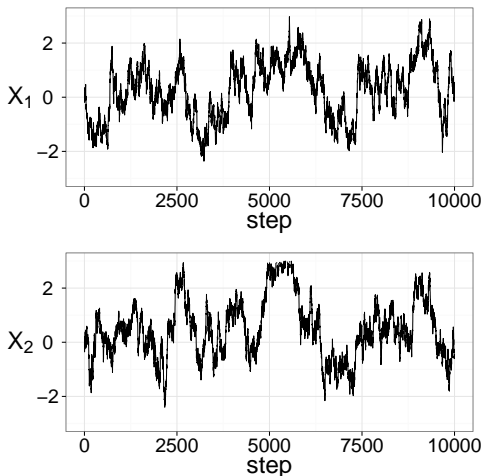


Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 10$, the acceptance rate is $\approx 1.5\%$.

# Metropolis–Hastings algorithm



Figure: Metropolis–Hastings on a bivariate Gaussian target. With $\sigma^2 = 10$, the acceptance rate is $\approx 1.5\%$.

# Choice of proposal

- Aim at some intermediate acceptance ratio: 20%? 40%? Some hints come from the literature on "optimal scaling".

- Literature suggest tuning to get .234...

- Maximize the expected square jumping distance:

$$\mathbb{E}\left[||X_{t+1} - X_t||^2\right]$$

- In multivariate cases, try to mimick the covariance structure of the target distribution.

  Cooking recipe: run the algorithm for $T$ iterations, check some criterion, tune the proposal distribution accordingly, run the algorithm for $T$ iterations again . . .

  "Constructing a chain that mixes well is somewhat of an art."

  *All of Statistics*, L. Wasserman.

# The adaptive MCMC approach

- One can make the transition kernel $K$ adaptive, i.e. use $K_t$ at iteration $t$ and choose $K_t$ using the past sample $(X_1, \ldots, X_{t-1})$.

- The Markov chain is not homogeneous anymore: the mathematical study of the algorithm is much more complicated.

- Adaptation can be counterproductive in some cases (see Atchadé & Rosenthal, 2005)!

- Adaptive Gibbs samplers also exist.

⚠ Extreme care is needed when designing adaptive algorithms: it's easy to make an algorithm with the wrong invariant distribution.

## Sophisticated Proposals

- "Langevin" proposal relies on

$$X^\star = X^{(t-1)} + \frac{\sigma}{2} \nabla \log \pi \left( X^{(t-1)} \right) + \sigma W$$

where $W \sim \mathcal{N}(0, I_d)$, so the Metropolis-Hastings acceptance ratio is

$$\frac{\pi(X^\star) q(X^{(t-1)} \mid X^\star)}{\pi(X^{(t-1)}) q(X^\star \mid X^{(t-1)})}$$

$$= \frac{\pi(X^\star)}{\pi(X^{(t-1)})} \frac{\mathcal{N}(X^{(t-1)}; X^\star + \frac{\sigma}{2}.\nabla \log \pi \left( X^\star \right); \sigma^2)}{\mathcal{N}(X^\star; X^{(t-1)} + \frac{\sigma}{2}.\nabla \log \pi \left( X^{(t-1)} \right); \sigma^2)}.$$

- Possibility to use higher order derivatives:

$$X^\star = X^{(t-1)} + \frac{\sigma}{2} \left[ \nabla^2 \log \pi \left( X^{(t-1)} \right) \right]^{-1} \nabla \log \pi \left( X^{(t-1)} \right) + \sigma W.$$

# Sophisticated Proposals

- We can use

$$q(X^\star | X^{(t-1)}) = g(X^\star; \varphi(X^{(t-1)}))$$

where $g$ is a distribution on $\mathbb{X}$ of parameters $\varphi(X^{(t-1)})$ and $\varphi$ is a deterministic mapping

$$\frac{\pi(X^\star) q(X^{(t-1)} | X^\star)}{\pi(X^{(t-1)}) q(X^\star | X^{(t-1)})} = \frac{\pi(X^\star) g(X^{(t-1)}; \varphi(X^\star))}{\pi(X^{(t-1)}) g(X^\star; \varphi(X^{(t-1)}))}.$$

- For instance, use heuristics borrowed from optimization techniques.

# Sophisticated Proposals

The following link shows a comparison of

- adaptive Metropolis-Hastings,
- Gibbs sampling,
- No U-Turn Sampler (e.g. Hamiltonian MCMC)
  on a simple linear model.

  twiecki.github.io/blog/2014/01/02/visualizing-mcmc/

# Sophisticated Proposals

- Assume you want to sample from a target $\pi$ with $\mathrm{supp}(\pi) \subset \mathbb{R}^+$, e.g. the posterior distribution of a variance/scale parameter.

- Any proposed move, e.g. using a normal random walk, to $\mathbb{R}^-$ is a waste of time.

- Given $X^{(t-1)}$, propose $X^\star = \exp(\log X^{(t-1)} + \varepsilon)$ with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. What is the acceptance probability then?

$$\alpha(X^\star \mid X^{(t-1)}) = \min\left(1, \frac{\pi(X^\star)}{\pi(X^{(t-1)})} \frac{q(X^{(t-1)} \mid X^\star)}{q(X^\star \mid X^{(t-1)})}\right)$$
$$= \min\left(1, \frac{\pi(X^\star)}{\pi(X^{(t-1)})} \frac{X^\star}{X^{(t-1)}}\right).$$

Why?

$$\frac{q(y|x)}{q(x \mid y)} = \frac{\frac{1}{y\sigma\sqrt{2\pi}} \exp\left[-\frac{(\log y - \log x)^2}{2\sigma^2}\right]}{\frac{1}{x\sigma\sqrt{2\pi}} \exp\left[-\frac{(\log x - \log y)^2}{2\sigma^2}\right]} = \frac{x}{y}.$$

## Random Proposals

- Assume you want to use $q_{\sigma^2}(X^\star|X^{(t-1)}) = \mathcal{N}(X; X^{(t-1)}, \sigma^2)$ but you don't know how to pick $\sigma^2$. You decide to pick a random $\sigma^{2,\star}$ from a distribution $f(\sigma^2)$:

$$\sigma^{2,\star} \sim f(\sigma^{2,\star}), \ X^\star|\sigma^{2,\star} \sim q_{\sigma^{2,\star}}(\cdot|X^{(t-1)})$$

so that

$$q(X^\star|X^{(t-1)}) = \int q_{\sigma^{2,\star}}(X^\star|X^{(t-1)}) f(\sigma^{2,\star}) d\sigma^{2,\star}.$$

- Perhaps $q(X^\star|X^{(t-1)})$ cannot be evaluated, e.g. the above integral is intractable. Hence the acceptance probability

$$\min\{1, \frac{\pi(X^\star)q(X^{(t-1)}|X^\star)}{\pi(X^{(t-1)})q(X^\star|X^{(t-1)})}\}$$

cannot be computed.

# Random Proposals

- Instead you decide to accept your proposal with probability

$$\alpha_t = \min \left\{ 1, \frac{\pi\left(X^\star\right) q_{\sigma^{2,(t-1)}}\left(X^{(t-1)}\middle| X^\star\right)}{\pi\left(X^{(t-1)}\right) q_{\sigma^{2,\star}}\left(X^\star | X^{(t-1)}\right)} \right\}$$

  where $\sigma^{2,(t-1)}$ corresponds to parameter of the last accepted proposal.

- With probability $\alpha_t$, set $\sigma^{2,(t)} = \sigma^{2,\star}$, $X^{(t)} = X^\star$, otherwise $\sigma^{2,(t)} = \sigma^{2,(t-1)}$, $X^{(t)} = X^{(t-1)}$.

- **Question**: Is it valid? If so, why?

## Random Proposals

- Consider the extended target

$$\widetilde{\pi}\left(x, \sigma^2\right) := \pi(x) f\left(\sigma^2\right).$$

- Previous algorithm is a Metropolis-Hastings of target $\widetilde{\pi}(x, \sigma^2)$ and proposal

$$q(y, \tau^2 | x, \sigma^2) = f(\tau^2) q_{\tau^2}(y | x)$$

- Indeed, we have

$$\frac{\widetilde{\pi}(y, \tau^2)}{\widetilde{\pi}(x, \sigma^2)} \frac{q(x, \sigma^2 | y, \tau^2)}{q(y, \tau^2 | x, \sigma^2)}$$

$$= \frac{\pi(y) f(\tau^2)}{\pi(x) f(\sigma^2)} \frac{f(\sigma^2) q_{\sigma^2}(x | y)}{f(\tau^2) q_{\tau^2}(y | x)} = \frac{\pi(y)}{\pi(x)} \frac{q_{\sigma^2}(x | y)}{q_{\tau^2}(y | x)}$$

- **Remark**: we just need to be able to sample from $f(\cdot)$, not to evaluate it.

# Using multiple proposals

- Consider a target of density $\pi(x)$ where $x \in \mathbb{X}$.
- To sample from $\pi$, you might want to use various proposals for Metropolis-Hastings $q_1(x'|x)$, $q_2(x'|x),...,q_p(x'|x)$.
- One way to achieve this is to build a proposal

$$q(x'|x) = \sum_{j=1}^{p} \beta_j q_j(x'|x), \ \beta_j > 0, \sum_{j=1}^{p} \beta_j = 1,$$

and Metropolis-Hastings requires evaluating

$$\alpha\left(X^{\star}|X^{(t-1)}\right) = \min\left(1, \frac{\pi(X^{\star})q\left(X^{(t-1)}|X^{\star}\right)}{\pi(X^{(t-1)})q(X^{\star}|X^{(t-1)})}\right),$$

and thus evaluating $q_j\left(X^{\star}|X^{(t-1)}\right)$ for $j = 1,...,p$.

# Motivating Example

- Let

$$q\left(x'\mid x\right) = \beta_1 \mathcal{N}\left(x'; x, \Sigma\right) + \left(1 - \beta_1\right) \mathcal{N}\left(x'; \mu(x), \Sigma\right)$$

  where $\mu\colon \mathbb{X} \to \mathbb{X}$ is a clever but computationally expensive deterministic optimisation algorithm.

- Using $\beta_1 \approx 1$ will make most proposed points come from the cheaper proposal distribution $\mathcal{N}\left(x'; x, \Sigma\right)$...

- ... but you won't save time as $\mu\left(X^{(t-1)}\right)$ needs to be evaluated at every step.

# Composing kernels

- How to use different proposals to sample from $\pi$ without evaluating all the densities at each step?

- What about combining different Metropolis-Hastings updates $K_j$ using proposal $q_j$ instead? i.e.

$$K_j \left( x, x' \right) = \alpha_j \left( x' \,|\, x \right) q_j \left( x' \,|\, x \right) + \left( 1 - a_j \left( x \right) \right) \delta_x \left( x' \right)$$

where

$$\alpha_j(x'|x) = \min\left( 1, \frac{\pi(x')q_j(x|x')}{\pi(x)q_j(x'|x)} \right)$$
$$a_j(x) = \int \alpha_j(x'|x) q_j(x'|x)\, dx'.$$

# Composing kernels

Generally speaking, assume

- $p$ possible updates characterised by kernels $K_j(\cdot,\cdot)$,

- each kernel $K_j$ is $\pi$-invariant.

  Two possibilities of combining the $p$ MCMC updates:

- **Cycle**: perform the MCMC updates in a deterministic order.

- **Mixture**: Pick an MCMC update at random.

# Cycle of MCMC updates

- Starting with $X^{(1)}$ iterate for $t = 2, 3, ...$
- (a) Set $Z^{(t,0)} := X^{(t-1)}$.
- (b) For $j = 1, ..., p$, sample $Z^{(t,j)} \sim K_j\left(Z^{(t,j-1)}, \cdot\right)$.
- (c) Set $X^{(t)} := Z^{(t,p)}$.
- Full cycle transition kernel is

$$K\left(x^{(t-1)}, x^{(t)}\right) = \int \cdots \int K_1\left(x^{(t-1)}, z^{(t,1)}\right) K_2\left(z^{(t,1)}, z^{(t,2)}\right)$$
$$\cdots K_p\left(z^{(t,p-1)}, x^{(t)}\right) dz^{(t,1)} \cdots dz^{(t,p-1)}.$$

- $K$ is $\pi$-invariant.

# Mixture of MCMC updates

- Starting with $X^{(1)}$ iterate for $t = 2, 3, \ldots$
- (a) Sample $J$ from $\{1, \ldots, p\}$ with $\mathbb{P}(J = k) = \beta_k$.
- (b) Sample $X^{(t)} \sim K_J\left(X^{(t-1)}, \cdot\right)$.
- Corresponding transition kernel is

$$K\left(x^{(t-1)}, x^{(t)}\right) = \sum_{j=1}^{p} \beta_j K_j\left(x^{(t-1)}, x^{(t)}\right).$$

- $K$ is $\pi$-invariant.
- The algorithm is *different* from using a mixture proposal

$$q\left(x' \mid x\right) = \sum_{j=1}^{p} \beta_j q_j\left(x' \mid x\right).$$

# Metropolis-Hastings Design for Multivariate Targets

- If $\dim(\mathbb{X})$ is large, it might be very difficult to design a "good" proposal $q(x'|x)$.

- As in Gibbs sampling, we might want to partition $x$ into $x = (x_1,...,x_d)$ and denote $x_{-j} := x \setminus \{x_j\}$.

- We propose "local" proposals where only $x_j$ is updated

$$q_j(x'|x) = \underbrace{q_j\left(x'_j\big|x\right)}_{\text{propose new component } j} \underbrace{\delta_{x_{-j}}\left(x'_{-j}\right)}_{\text{keep other components fixed}} \quad .$$

# Metropolis-Hastings Design for Multivariate Targets

- This yields

$$\alpha_j(x,x') = \min\left(1, \frac{\pi(x'_{-j},x'_j)q_j(x_j|x_{-j},x'_j)}{\pi(x_{-j},x_j)q_j(x'_j|x_{-j},x_j)}\underbrace{\frac{\delta_{x'_{-j}}(x_{-j})}{\delta_{x_{-j}}(x'_{-j})}}_{=1}\right)$$

$$= \min\left(1, \frac{\pi(x_{-j},x'_j)q_j(x_j|x_{-j},x'_j)}{\pi(x_{-j},x_j)q_j(x'_j|x_{-j},x_j)}\right)$$

$$= \min\left(1, \frac{\pi_{X_j|X_{-j}}(x'_j|x_{-j})q_j(x_j|x_{-j},x'_j)}{\pi_{X_j|X_{-j}}(x_j|x_{-j})q_j(x'_j|x_{-j},x_j)}\right).$$

# One-at-a-time MH (cycle/systematic scan)

Starting with $X^{(1)}$ iterate for $t = 2, 3, ...$
For $j = 1, ..., d$,

- Sample $X^\star \sim q_j(\cdot | X_1^{(t)}, ..., X_{j-1}^{(t)}, X_j^{(t-1)}, ..., X_d^{(t-1)})$.

- Compute

$$
\begin{aligned}
\alpha_j \;=\; \min\Bigg( & 1, \frac{\pi_{X_j | X_{-j}}\left(X_j^\star \mid X_1^{(t)} \dots X_{j-1}^{(t)}, X_{j+1}^{(t-1)} \dots X_d^{(t-1)}\right)}{\pi_{X_j | X_{-j}}\left(X_j^{(t-1)} \mid X_1^{(t)} \dots X_{j-1}^{(t)}, X_{j+1}^{(t-1)} \dots X_d^{(t-1)}\right)} \\
& \times \frac{q_j\left(X_j^{(t-1)} \Big| X_1^{(t)} \dots X_{j-1}^{(t)}, X_j^\star, X_{j+1}^{(t-1)} \dots X_d^{(t-1)}\right)}{q_j\left(X_j^\star \Big| X_1^{(t)} \dots X_{j-1}^{(t)}, X_j^{(t-1),}, X_{j+1}^{(t-1)} \dots X_d^{(t-1)}\right)} \Bigg).
\end{aligned}
$$

- With probability $\alpha_j$, set $X^{(t)} = X^\star$, otherwise set $X^{(t)} = X^{(t-1)}$.

# One-at-a-time MH (mixture/random scan)

Starting with $X^{(1)}$ iterate for $t = 2, 3, ...$

- Sample $J$ from $\{1, ..., d\}$ with $\mathbb{P}(J = k) = \beta_k$.
- Sample $X^\star \sim q_J\left(\cdot \,|\, X_1^{(t)}, ..., X_d^{(t-1)}\right)$.
- Compute

$$
\begin{aligned}
\alpha_J = \min\Bigg( 1, &\frac{\pi_{X_J|X_{-J}}\left(X_J^\star \,|\, X_1^{(t-1)} \dots X_{J-1}^{(t-1)}, X_{J+1}^{(t-1)} \dots\right)}{\pi_{X_J|X_{-J}}\left(X_J^{(t-1)} \,|\, X_1^{(t-1)} \dots X_{J-1}^{(t-1)}, X_{J+1}^{(t-1)} \dots\right)} \\
&\times \frac{q_J\left(X_J^{(t-1)} \,\Big|\, X_1^{(t-1)} \dots X_{J-1}^{(t-1)}, X_J^\star, X_{J+1}^{(t-1)} \dots X_d^{(t-1)}\right)}{q_J\left(X_J^\star \,\Big|\, X_1^{(t-1)} \dots X_{J-1}^{(t-1)}, X_J^{(t-1)}, X_{J+1}^{(t-1)} \dots X_d^{(t-1)}\right)} \Bigg).
\end{aligned}
$$

- With probability $\alpha_J$ set $X^{(t)} = X^\star$, otherwise $X^{(t)} = X^{(t-1)}$.

# Gibbs Sampler as a Metropolis-Hastings algorithm

## Proposition

*The systematic Gibbs sampler is a cycle of one-at-a time MH whereas the random scan Gibbs sampler is a mixture of one-at-a time MH where*

$$q_j\left(x_j'\middle| x\right) = \pi_{X_j|X_{-j}}\left(x_j'\middle| x_{-j}\right).$$

## Proof.

It follows from

$$\frac{\pi\left(x_{-j}, x_j'\right)}{\pi\left(x_{-j}, x_j\right)} \frac{q_j\left(x_j\middle| x_{-j}, x_j'\right)}{q_j\left(x_j'\middle| x_{-j}, x_j\right)}$$

$$= \frac{\pi\left(x_{-j}\right)\pi_{X_j|X_{-j}}\left(x_j'\middle| x_{-j}\right)}{\pi\left(x_{-j}\right)\pi_{X_j|X_{-j}}\left(x_j\middle| x_{-j}\right)} \frac{\pi_{X_j|X_{-j}}\left(x_j\middle| x_{-j}\right)}{\pi_{X_j|X_{-j}}\left(x_j'\middle| x_{-j}\right)} = 1.$$

# This is not a Gibbs sampler

Consider a case where $d = 2$. From $X_1^{(t-1)}, X_2^{(t-1)}$ at time $t-1$:

- Sample $X_1^\star \sim \pi(X_1 \mid X_2^{(t-1)})$, then $X_2^\star \sim \pi(X_2 \mid X_1^\star)$. The proposal is then $X^\star = (X_1^\star, X_2^\star)$.

- Compute

$$\alpha_t = \min\left(1, \frac{\pi(X_1^\star, X_2^\star)}{\pi(X_1^{(t-1)}, X_2^{(t-1)})} \frac{q(X^{(t-1)} \mid X^\star)}{q(X^\star \mid X^{(t-1)})}\right)$$

- Accept $X^\star$ or not based on $\alpha_t$, where here

$$\alpha_t \neq 1$$

!!