# Lecture 7: Brief reminder

- Problem: Sample from $\pi(x)$, $x \in S$.
- MCMC idea:
  - ▶ Construct Markov chain with $\pi(x)$ as limiting distribution.
  - ▶ Simulate the Markov chain for a long time so that it has time to converge.
  - ▶ Most MCMC samplers are based on reversible Markov chains $\Rightarrow$ Their convergence is proved by checking the detailed balance equation.

# Review: Metropolis-Hastings construction

- $$P(y \mid x) = \begin{cases} Q(y \mid x)\alpha(y \mid x), & y \neq x \\ 1 - \sum_{z \neq x} Q(z \mid x)\alpha(z \mid x), & y = x \end{cases}$$

- $$\alpha(y \mid x) = \min\left\{1, \frac{\pi(y)}{\pi(x)} \cdot \frac{Q(x \mid y)}{Q(y \mid x)}\right\}$$

# Review: Metropolis-Hastings algorithm

1: Init $x_0 \sim g(x_0)$
2: **for** $i = 1, 2, \ldots$ **do**
3:  Generate a proposal $y \sim Q(y|x_{i-1})$
4:  $u \sim U(0, 1)$
5:  **if** $u < \min\left(1, \dfrac{\pi(y)}{\pi(x_{i-1})} \times \underbrace{\dfrac{Q(x_{i-1}|y)}{Q(y|x_{i-1})}}_{\text{Proposal ratio}}\right)$ **then**

$\underbrace{\phantom{u < \min\left(1, \dfrac{\pi(y)}{\pi(x_{i-1})} \times \dfrac{Q(x_{i-1}|y)}{Q(y|x_{i-1})}\right)}}_{\text{Acceptance probability } \alpha}$

6:   $x_i \leftarrow y$
7:  **else**
8:   $x_i \leftarrow x_{i-1}$
9:  **end if**
10: **end for**

# What about

- Irreducible: Must be checked in each case. Must choose $Q(y \mid x)$ so that this is ok.

# What about

- Irreducible: Must be checked in each case. Must choose $Q(y \mid x)$ so that this is ok.

- Aperiodic: Sufficient that $P(x \mid x) > 0$ for one $x \in S$, so sufficient that $\alpha(y \mid x) < 1$ for one pair $y, x \in S$.

# What about

- Irreducible: Must be checked in each case. Must choose $Q(y \mid x)$ so that this is ok.

- Aperiodic: Sufficient that $P(x \mid x) > 0$ for one $x \in S$, so sufficient that $\alpha(y \mid x) < 1$ for one pair $y, x \in S$.

- Positive recurrent: for finite $S$, irreducibility is sufficient. More difficult in general, but if Markov chain is not recurrent we will see this as drift in the simulations. (In practice usually no problem).

# Metropolis-Hastings algorithm

Elements of the problem:

- Target distribution $\pi(x)$: Given by the problem

- Proposal distribution $Q(y|x)$: Chosen by the user

- Acceptance probability $\alpha(y|x)$: Derived in order to fullfill the detailed balance condition.

# Remarks on the Metropolis–Hastings algorithm

- Under some regularity conditions it can be shown that the Metropolis-Hasting algorithm converges to the target distribution regardless of the specific choice of $Q(y|x)$.

For more comments and details see: Chib, S. and Greenberg, E. (1995), *Understanding the Metropolis-Hastings algorithm, The American Statistician, 49: 327–335*

# Remarks on the Metropolis-Hastings algorithm

- Under some regularity conditions it can be shown that the Metropolis-Hasting algorithm converges to the target distribution regardless of the specific choice of $Q(y|x)$.

- However, the speed of convergence and the dependence between the successive samples depends strongly on the proposal distribution.

For more comments and details see: Chib, S. and Greenberg, E. (1995), *Understanding the Metropolis-Hastings algorithm, The American Statistician, 49: 327–335*

# Remarks on the Metropolis-Hastings algorithm

- Under some regularity conditions it can be shown that the Metropolis-Hasting algorithm converges to the target distribution regardless of the specific choice of $Q(y|x)$.

- However, the speed of convergence and the dependence between the successive samples depends strongly on the proposal distribution.

- Since we only need to compute the ratio $\pi(y)/\pi(x)$, the proportionality constant is irrelevant.

For more comments and details see: Chib, S. and Greenberg, E. (1995), Understanding the Metropolis-Hastings algorithm, The American Statistician, 49: 327–335

# Remarks on the Metropolis-Hastings algorithm

- Under some regularity conditions it can be shown that the Metropolis-Hasting algorithm converges to the target distribution regardless of the specific choice of $Q(y|x)$.

- However, the speed of convergence and the dependence between the successive samples depends strongly on the proposal distribution.

- Since we only need to compute the ratio $\pi(y)/\pi(x)$, the proportionality constant is irrelevant.

- Similarly, we only care about $Q(.)$ up to a constant.

For more comments and details see: Chib, S. and Greenberg, E. (1995), *Understanding the Metropolis-Hastings algorithm, The American Statistician, 49: 327–335*

# Remarks on the Metropolis-Hastings algorithm

- Under some regularity conditions it can be shown that the Metropolis-Hasting algorithm converges to the target distribution regardless of the specific choice of $Q(y|x)$.

- However, the speed of convergence and the dependence between the successive samples depends strongly on the proposal distribution.

- Since we only need to compute the ratio $\pi(y)/\pi(x)$, the proportionality constant is irrelevant.

- Similarly, we only care about $Q(.)$ up to a constant.

- Often it is advantageous to calculate the acceptance probability on log-scale, which makes the computations more stable.

For more comments and details see: Chib, S. and Greenberg, E. (1995), Understanding the Metropolis-Hastings algorithm, The American Statistician, 49: 327–335

# Special cases of the Metropolis-Hastings algorithm

Depending on the choice of $Q(y|x_{i-1})$ different special cases result. In particular, two classes are important

- The independence proposal
- The Metropolis algorithm
  - Random walk proposals

# Independence proposal

- The proposal distribution does not depend on the current value $x_{i-1}$

$$Q(x|x_{i-1}) = Q(x).$$

- $Q(x)$ is an approximation to $\pi(x)$
  $\Rightarrow$ Acceptance rate should be close to 1.

- The sampler is closer to rejection sampler. However, here if we reject, then we retain the sample.

Experience:

- Performance is either very good or very bad, usually very bad.

- The tails of the proposal distribution should be at least as heavy as the tails of the target distribution.
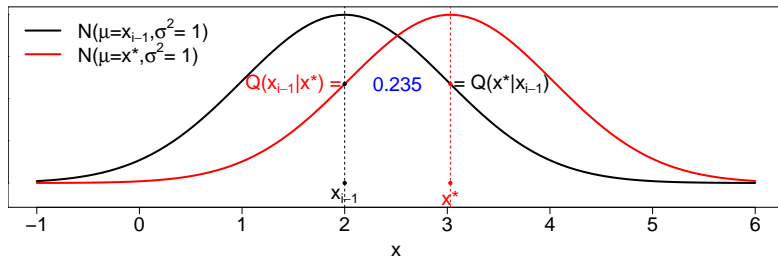
# The Metropolis algorithm

The proposal density is symmetric around the current value, that means

$$Q(x_{i-1}|y) = Q(y|x_{i-1}).$$

Hence,

$$\alpha = \min\left(1, \frac{\pi(y)}{\pi(x_{i-1})} \times \frac{Q(x_{i-1}|y)}{Q(y|x_{i-1})}\right) = \min\left(1, \frac{\pi(y)}{\pi(x_{i-1})}\right)$$

A particular case is the random walk proposal, defined as the current value $x_{i-1}$ plus a random variate of a 0-centred symmetric distribution.

# Examples for random walks proposal

Assume $x$ is scalar.

Then all proposal kernels, which <span style="color:red">add a random variable generated from a zero-symmetrical distribution to the current value $x_{i-1}$</span>, are random walk proposals. For example:

$$y \sim \mathcal{N}(x_{i-1}, \sigma^2)$$

$$y \sim t_\nu(x_{i-1}, \sigma^2)$$

$$y \sim \mathcal{U}(x_{i-1} - d, x_{i-1} + d)$$

See R-code `demo_mcmcRW_2D.R`.

# Efficiency of the Metropolis-Hastings algorithm

The efficiency and performance of the Metropolis-Hastings algorithm depends crucially on the relative frequency of acceptance.

# Efficiency of the Metropolis-Hastings algorithm

The efficiency and performance of the Metropolis-Hastings algorithm depends crucially on the relative frequency of acceptance.

An acceptance rate of one is not always good. Consider the random walk proposal:

- Too large acceptance rate $\Rightarrow$ Slow exploration of the target density.
- Too small acceptance rate $\Rightarrow$ Large moves are proposed, but rarely accepted.

# Efficiency of the Metropolis-Hastings algorithm

The efficiency and performance of the Metropolis-Hastings algorithm depends crucially on the relative frequency of acceptance.

An acceptance rate of one is not always good. Consider the random walk proposal:

- Too large acceptance rate $\Rightarrow$ Slow exploration of the target density.
- Too small acceptance rate $\Rightarrow$ Large moves are proposed, but rarely accepted.

Tuning the acceptance rate:

- For random walk proposals, acceptance rates between 20% and 50% are typically recommended. They can be achieved by changing the variance of the proposal distribution.
- For independence proposals a high acceptance rate is desired, which means that the proposal density is close to the target density.

# Example: Random walk proposal

Exploration of a standard Gaussian distribution ($\mathcal{N}(0,1)$) using a random walk Metropolis algorithm. As proposal assume a Gaussian distribution with variance $\sigma^2$, where.

- $\sigma = 0.24$
- $\sigma = 2.4$
- $\sigma = 24$

See R-code `demo_mcmcRW.R`.

# Example of Rao (1973)

The vector $\boldsymbol{y} = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$ is multinomial distributed with probabilities

$$\left\{ \frac{1}{2} + \frac{\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4} \right\}$$

We would like to simulate from the posterior distribution (assuming a uniform prior)

$$f(\theta|\boldsymbol{y}) \propto (2+\theta)^{y_1}(1-\theta)^{y_2+y_3}\theta^{y_4}.$$

using MCMC and compare two proposal kernels:

1. independence proposal

2. random walk proposal

See R-code `demo_mcmcRao.R`.

# Rao: Independence proposal

$$\theta^\star \sim \mathcal{N}(\text{Mod}(\theta|\boldsymbol{y}), F^2 \times I_p^{-1}), \tag{1}$$

where $\text{Mod}(\theta|\text{data})$ denotes the posterior mode, $I_p$ the negative curvature of the log posterior at the mode, and $F$ a factor to blow up the standard deviation.

# Rao: Random walk proposal

$$\theta^\star \sim \mathsf{U}(\theta^{(k)} - d, \theta^{(k)} + d),$$

where $\theta^{(k)}$ denotes the current state of the Markov chain and $d = \sqrt{12}/2 \cdot 0.1$.

# Comments on the Metropolis-Hasting algorithm

- A trivial special case results when

$$Q(x^\star | x_{i-1}) = \pi(x^\star),$$

  That means, we propose realisations from the target distribution. Then $\alpha = 1$ and all proposals are accepted.

- The advantage of the MH-algorithm is that arbitrary proposal kernels can be used. The algorithm will always converge to the target distribution.

- However, the speed of convergence and the dependence between the successive samples depends strongly on the proposal distribution.

## Numerical Note

How to compute

$$\alpha(y|x) = \min\left\{1, \frac{\pi(y)}{\pi(x)} \frac{Q(x|y)}{Q(y|x)}\right\}$$

Naive strategy: Compute $\pi(y)$, $\pi(x)$, $Q(y|x)$, $Q(x|y)$. Then compute the ratio.

# Numerical Note

How to compute

$$\alpha(y|x) = \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \frac{Q(x|y)}{Q(y|x)} \right\}$$

Naive strategy: Compute $\pi(y)$, $\pi(x)$, $Q(y|x)$, $Q(x|y)$. Then compute the ratio.

Solution:

- Simplify the expression as much as possible
- Compute everything in log-scale

# Remarks on Gibbs sampling

- High dimensional updates of $x$ can be boiled down to scalar updates.

- Visiting schedule: Various approaches exist (and can be justified) to ordering the variables in the sampling loop. One approach is random sweeps: variables are chosen at random to resample.

- Gibbs sampling assumes that it is easy to sample from the full-conditional distribution. This is sometimes not so easy. Alternatively, a Metropolis-Hastings proposal can be used for the $j$-th component, i.e. Metropolis-within-Gibbs $\Rightarrow$ Hybrid Gibbs sampler.

# Remarks on Gibbs sampling

- Blocking or grouping is possible, that means not all elements of $x$ are treated individually. Might be useful when elements of $x$ are correlated.

- Care must be taken when improper prior are used, which may lead to an improper posterior distribution. Impropriety implies that there does not exist a joint density to which the full-conditional distributions correspond.

# Example : Conjugate gamma-Poisson hierarchical model

Example from George et al. (1993) regarding the analysis of 10 power plants.

- $y_i$ number of failures of pump $i$
- $t_i$ length of operation time of pump $i$ (in kilo hours)

Model:
$$y_i \mid \lambda_i \sim \mathsf{Po}(\lambda_i t_i)$$

Conjugate prior for $\lambda_i$:
$$\lambda_i \mid \alpha, \beta \sim \mathsf{G}(\alpha, \beta)$$

Hyper-prior on $\alpha$ and $\beta$:

$$\alpha \sim \mathsf{Exp}(1.0) \qquad\qquad \beta \sim \mathsf{G}(0.1, 10.0)$$

# Conjugate gamma-Poisson hierarchical model (II)

The posterior of the 12 parameters $(\alpha, \beta, \lambda_1, \ldots, \lambda_{10})$ given $y_1, \ldots, y_{10}$ is proportional to

$$\pi(\alpha, \beta, \lambda_1, \ldots, \lambda_{10} \mid y_1, \ldots, y_{10}) \propto \pi(\alpha)\pi(\beta) \prod_{i=1}^{10} [\pi(\lambda_i \mid \alpha, \beta)\pi(y_i \mid \lambda_i)]$$

$$\propto e^{-\alpha}\beta^{0.1-1}e^{-10\beta} \left\{ \prod_{i=1}^{10} \exp(-\lambda_i t_i)\lambda_i^{y_i} \right\} \left\{ \prod_{i=1}^{10} \exp(-\beta\lambda_i)\lambda_i^{\alpha-1} \right\} \left[ \frac{\beta^\alpha}{\Gamma(\alpha)} \right]^{10}.$$

This posterior is not of closed form.

# Conjugate gamma-Poisson hierarchical model (II)

The posterior of the 12 parameters $(\alpha, \beta, \lambda_1, \ldots, \lambda_{10})$ given $y_1, \ldots, y_{10}$ is proportional to

$$\pi(\alpha, \beta, \lambda_1, \ldots, \lambda_{10} \mid y_1, \ldots, y_{10}) \propto \pi(\alpha)\pi(\beta) \prod_{i=1}^{10} [\pi(\lambda_i \mid \alpha, \beta)\pi(y_i \mid \lambda_i)]$$

$$\propto e^{-\alpha} \beta^{0.1-1} e^{-10\beta} \left\{ \prod_{i=1}^{10} \exp(-\lambda_i t_i)\lambda_i^{y_i} \right\} \left\{ \prod_{i=1}^{10} \exp(-\beta\lambda_i)\lambda_i^{\alpha-1} \right\} \left[ \frac{\beta^\alpha}{\Gamma(\alpha)} \right]^{10}.$$

This posterior is not of closed form.

What are the full conditional distributions?