

## Additional topics on Monte Carlo

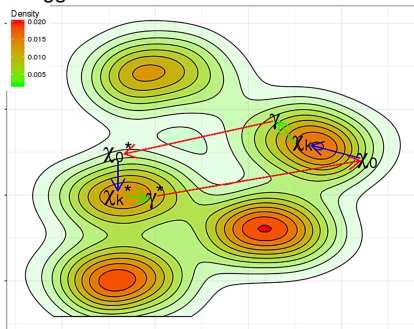
Geir Storvik

Geilo Winter school 2023

## Outline

- 1 Mode jumping MCMC
- 2 Particle MCMC
- 3 Reversible jump MCMC
- 4 Non-reversible MCMC
- 5 Continuous time Markov processes
- 6 Additional topics in MCMC
- 7 Practice

- MCMC: Work reasonable well for **unimodal** distributions
  - Struggle more with multimodal distributions



- Several possible approaches:
  - Simulated tempering: Use  $\pi_k(\mathbf{x}) = \pi_k(\mathbf{x})^{T_k}$ ,  $T_k \leq 1$ , move between different "models"
  - SMC: Similar sequence
- Here: **Mode jumping MCMC** Tjelmeland and Hegstad (2001)

## Mode jumping MCMC

- Aim: Allow for **large** changes
- Main problem: Large move in space will typically result in low density value
  - M-H: Very low acceptance rate
- Main idea
  - 1 Make a large change  $\mathbf{x} \rightarrow \mathbf{x}_0^*$
  - 2 Perform a local optimization  $\mathbf{x}_0^* \rightarrow \mathbf{x}_k^*$ 
    - Possibly through  $k$  steps of some optimization routine
  - 3 Small perturbation:  $\mathbf{x}_k^* \rightarrow \mathbf{x}^*$
  - 4 Accept  $\mathbf{x}^*$  through an M-H step
- M-H: Detailed balance require possibility for moving backwards as well
  - The small perturbation in step 3 makes this possible

## Algorithm

---

### Algorithm 1 MJMCMC step from current state $\mathbf{x}$

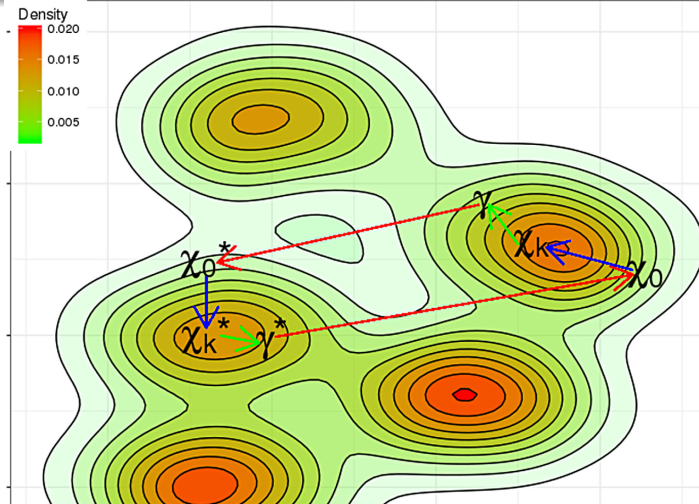
---

- 1: Generate  $\mathbf{x}_0^* = \mathbf{x}^* + \varepsilon^*$ ,  $\varepsilon^* \sim N(\mathbf{0}, \sigma_L^2 \mathbf{R})$ ,  $\sigma_L$  large
- 2: Optimize  $\mathbf{x}_0^* \rightarrow \mathbf{x}_k^*$
- 3: Small perturbation:  $\mathbf{x}^* \sim g_S(\mathbf{x}^* | \mathbf{x}_k^*)$
- 4: Generate  $\mathbf{x}_0 = \mathbf{x}^* - \varepsilon^*$
- 5: Optimize  $\mathbf{x}_0 \rightarrow \mathbf{x}_k$
- 6: Calculate

$$r = \frac{\pi(\mathbf{x}^*) q_r(\mathbf{x} | \mathbf{x}_k)}{\pi(\mathbf{x}) q_r(\mathbf{x}^* | \mathbf{x}_k)}$$

- 7: Accept  $\mathbf{x}^*$  with probability  $\min\{1, r\}$
-

## MJMCMC - graphical illustration



## MCMCMC for model selection

- Consider a model

$$y_i \sim f(y_i; \eta_i, \phi)$$

$$\eta_i = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j z_{i,j}$$

$$\gamma_j \sim$$

Bern( $q$ )

$$\beta_j | \gamma_j = 1 \sim N(0, \sigma_\beta^2)$$

- Aim:  $p(\boldsymbol{\gamma} | \mathbf{y})$ .
- $2^p$  possible models, in addition unknown  $\beta_j$ 's
- Possible:  $p(\boldsymbol{\gamma}, \boldsymbol{\beta} | \mathbf{y})$  through Reversible jump MCMC
- Pseudo-Marginal MCMC:

- 1 Generate proposal  $\boldsymbol{\gamma}^*$  from  $\boldsymbol{\gamma}$
- 2 Accept  $\boldsymbol{\gamma}^*$  with probability  $\min\{1, r\}$  where

$$r = \frac{p(\boldsymbol{\gamma}^* | \mathbf{y}) g(\boldsymbol{\gamma} | \boldsymbol{\gamma}^*)}{p(\boldsymbol{\gamma} | \mathbf{y}) g(\boldsymbol{\gamma}^* | \boldsymbol{\gamma})}$$

- Hubin and Storvik (2018): Linear models
- Hubin et al. (2021): Neural network type models

# Particle MCMC

- Andrieu et al. (2010)
- Ideal MCMC ( $p(\theta|\mathbf{y}) \propto p(\theta)L(\theta)$ ):
  - 1 Sample  $\theta^* \sim g(\theta^*|\theta)$
  - 2 Calculate M-H ratio  $r = \frac{p(\theta^*)L(\theta^*)g(\theta|\theta^*)}{p(\theta)L(\theta)g(\theta^*|\theta)}$
  - 3 Accept  $\theta^*$  with prob  $\min\{1, r\}$



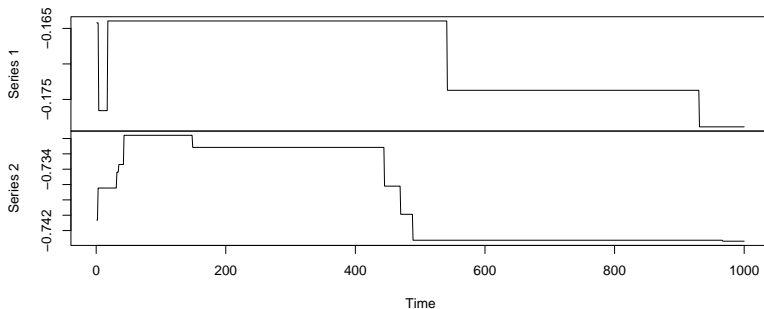
# Particle MCMC

- **Andrieu et al. (2010)**
- Ideal MCMC ( $p(\theta|\mathbf{y}) \propto p(\theta)L(\theta)$ ):
  - 1 Sample  $\theta^* \sim g(\theta^*|\theta)$
  - 2 Calculate M-H ratio  $r = \frac{p(\theta^*)L(\theta^*)g(\theta|\theta^*)}{p(\theta)L(\theta)g(\theta^*|\theta)}$
  - 3 Accept  $\theta^*$  with prob  $\min\{1, r\}$
- Pseudo-Marginal algorithm:
  - 1 Sample  $\theta^* \sim g(\theta^*|\theta)$
  - 2 Calculate  $\hat{L}(\theta^*)$
  - 3 Calculate M-H ratio  $\hat{r} = \frac{\pi(\theta^*)p(\theta|\theta^*)}{\pi(\theta)p(\theta^*|\theta)}$
  - 4 Accept  $\theta^*$  with prob  $\min\{1, \hat{r}\}$
- **Particle MCMC**: Use SMC to calculate  $\hat{L}(\theta^*)$

# Lemmings - AR(2) process

Now:  $x_t = a_1 x_{t-1} + a_2 x_{t-2} + \varepsilon_t$

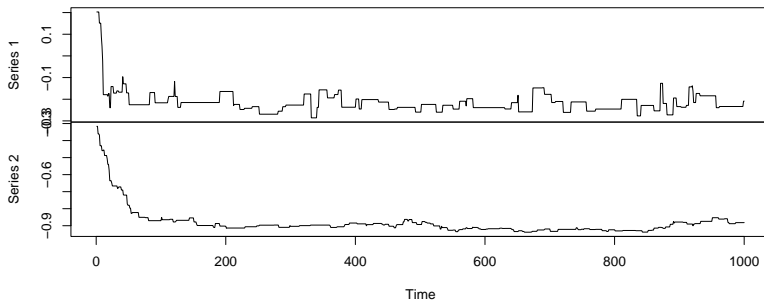
a.M



# Lemmings - AR(2) process

Now:  $x_t = a_1 x_{t-1} + a_2 x_{t-2} + \varepsilon_t$

a.M



# Reversible jump MCMC

- Examples of changing dimensions
  - $Y_i = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij} + \varepsilon_i$
  - Neural networks with some weights put to zero.
- **Reversible Jump MCMC**
  - Assume several **models**  $\mathcal{M}_1, \dots, \mathcal{M}_K$
  - Corresponding parameters  $\theta_1, \dots, \theta_K$  **of different dimensions!**
  - Aim: Simulate  $\mathbf{x} = (\mathcal{M}, \theta_{\mathcal{M}})$
  - RJMCMC: **Green (1995)**
  - RJMCMC: M-H method for moving between spaces of different dimensions
  - Main challenge: When changing  $\mathcal{M} \rightarrow \mathcal{M}^*$ , how to propose  $\theta_{\mathcal{M}^*}$ ?

# Reversible jump MCMC

- Examples of changing dimensions
  - $Y_i = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij} + \varepsilon_i$
  - Neural networks with some weights put to zero.
- **Reversible Jump MCMC**
  - Assume several **models**  $\mathcal{M}_1, \dots, \mathcal{M}_K$
  - Corresponding parameters  $\theta_1, \dots, \theta_K$  **of different dimensions!**
  - Aim: Simulate  $\mathbf{x} = (\mathcal{M}, \theta_{\mathcal{M}})$
  - RJMCMC: **Green (1995)**
  - RJMCMC: M-H method for moving between spaces of different dimensions
  - Main challenge: When changing  $\mathcal{M} \rightarrow \mathcal{M}^*$ , how to propose  $\theta_{\mathcal{M}^*}$ ?

## Changing dimensions

- Assume several **models**  $\mathcal{M}_1, \dots, \mathcal{M}_K$
- Corresponding parameters  $\theta_1, \dots, \theta_K$  **of different dimensions!**
- Aim: Simulate  $\mathbf{x} = (\mathcal{M}, \theta_{\mathcal{M}})$
- RJMCMC: M-H method for moving between spaces of different dimensions
- Main challenges:
  - When changing dimensions, how to compare densities on different spaces?
  - When changing  $\mathcal{M} \rightarrow \mathcal{M}^*$ , how to propose  $\theta_{\mathcal{M}^*}$ ?

# Reversible jump MCMC

- **Green (1995)**: Include auxiliary variables to match dimensions.
- Consider change  $(\mathcal{M}_1, \theta_1)$  to  $(\mathcal{M}_2, \theta_2)$  with  $|\theta_1| < |\theta_2|$ 
  - $j(1 \rightarrow 2)$  probability for moving from  $\mathcal{M}_1$  to  $\mathcal{M}_2$

- Algorithm

- 1 Generate  $\mathbf{u}_1$  such that  $|\theta_1| + |\mathbf{u}_1| = |\theta_2|$
- 2 Propose  $\theta_2 = \theta_2(\theta_1, \mathbf{u}_1)$  (bijective)
- 3 Calculate acceptance ratio

$$r = \frac{\pi(\mathcal{M}_2, \theta_2)q(2 \rightarrow 1)}{\pi(\mathcal{M}_1, \theta_1)q(1 \rightarrow 2)q(\mathbf{u}_1)} \left| \frac{\partial(\theta_2)}{\partial(\theta_1, \mathbf{u}_1)} \right|$$

- 4 Accept with probability  $\min\{1, r\}$ .
- Use  $1/r$  for opposite move
  - More general settings possible

# Logistic regression

- Assume model

$$Y_i \sim \text{Binom}(p_i) \qquad \text{logit}(p_i) = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij}$$

$$\Pr(\gamma_j = 1) = q \qquad \beta_j | \gamma_j = 1 \sim N(0, \sigma_\beta^2)$$

- Assume  $\gamma_j = 0$ , want to change to  $\gamma_j^* = 1$
- Generate  $u_1 \sim g_j()$
- Put

$$\beta_k^* = \begin{cases} \beta_k & k \neq j; \\ u_1 & k = j. \end{cases}$$

- Accept with probability  $\min\{1, r\}$  where

$$r = \frac{\pi(\beta^*, \gamma^*)}{\pi(\beta, \gamma) g_j(\beta_j^*)} \times 1$$

- Script `Log_reg_RJ.R`



## Non-reversible MCMC

- Main criterion ( $\pi$ -invariance)

$$\pi(\mathbf{x}^*) = \int_{\mathbf{x}} \pi(\mathbf{x}) P(\mathbf{x}^* | \mathbf{x}) d\mathbf{x}$$

- **Sufficient** criterion for stationarity

$$\pi(\mathbf{x}) P(\mathbf{x}^* | \mathbf{x}) = \pi(\mathbf{x}^*) P(\mathbf{x} | \mathbf{x}^*) \quad \text{Detailed balance}$$

- Results in a reversible MCMC (moving backwards is similar to moving forwards)

## Non-reversible MCMC

- Main criterion ( $\pi$ -invariance)

$$\pi(\mathbf{x}^*) = \int_{\mathbf{x}} \pi(\mathbf{x}) P(\mathbf{x}^* | \mathbf{x}) d\mathbf{x}$$

- **Sufficient** criterion for stationarity

$$\pi(\mathbf{x}) P(\mathbf{x}^* | \mathbf{x}) = \pi(\mathbf{x}^*) P(\mathbf{x} | \mathbf{x}^*) \quad \text{Detailed balance}$$

- Results in a reversible MCMC (moving backwards is similar to moving forwards)
- Assume now  $x \in \mathcal{Z}$
- Introduce  $v \in \{-1, 1\}$  and consider extended distribution  $\bar{\pi}(x, v) = 0.5\pi(x)I(v \in \{-1, 1\})$ .
- Define Markov chain

$$P(x^*, v | x, w) = \alpha(x, v)I(x^* = x + v, w = v) + (1 - \alpha(x, v))I(x^* = x, w = -v)$$

with  $\alpha(x, v) = \min\{1, \pi(x + v)/\pi(x)\}$

## The zig-zag process

- Continuous-time Markov process
- Can use sub-sampling with an **exact approximate scheme**
- Can be **super-efficient** when combined with control-covariate ideas
- References: **Bierkens et al. (2019)** (and references therein)
- Main idea:
  - Move all components linearly in a given direction:  $x_i(t) = x_i^k + z_i^k t$
  - Change direction of  $z_i^k$  at random (continuous) time points

## Algorithm

Let  $(T^0, \mathbf{x}^0, \theta^0) = (0, \xi, \theta)$

**for**  $k = 1, 2, \dots$  **do**

Let  $\xi^k(t) \equiv \mathbf{x}^{k-1} + \theta^{k-1}t, t \geq 0$

For  $i = 1, \dots, p$ , let  $\tau_i^k$  be distributed according to

$$\Pr(\tau_i^k \geq t) = \exp\left(-\int_0^t \lambda_i(\xi^k(s), \mathbf{z}^{k-1})ds\right)$$

Let  $i_0 \equiv \arg \min_{i \in \{1, \dots, p\}} \tau_i^k$

Let  $T^k \equiv T^{k-1} + \tau_{i_0}^k$

Let  $\mathbf{x}^k \equiv \xi^k(T^k)$

Let

$$z_i^k = \begin{cases} z_i^{k-1} & \text{if } i \neq i_0 \\ -z_i^{k-1} & \text{if } i = i_0 \end{cases}$$

**end for**

# Trajectories

- Piecewise deterministic trajectories  $(\mathbf{x}(t), \theta(t))$ :

$$(\mathbf{x}(t), \theta(t)) = (\mathbf{x}^k + \mathbf{z}^k(t - T^k), \mathbf{z}^k) \quad \text{for } t \in [T^k, T^{k+1}), k = 0, 1, 2, \dots$$

- Monte Carlo estimate:

$$\begin{aligned}\hat{\mu}_i &= \frac{1}{T} \int_0^T x_i(t) dt \\ &= \frac{1}{T} \sum_{k=0}^K \int_{T^k}^{T^{k+1}} [x_i^k + z_i^k(t - T^k)] dt \\ &= x_i^k(T^{k+1} - T^k) + 0.5 z_i^k(T^{k+1} - T^k)^2\end{aligned}$$

## What does it converge to?

- Distribution depending on the functions  $\lambda_i(\xi, \mathbf{z})$ .
- Assume  $\theta_i \in \{-1, 1\}$
- Assume a Bayesian setting:

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

- Define

$$\Psi(\mathbf{x}) = -\log p(\mathbf{x}) - \log p(\mathbf{y}|\mathbf{x})$$

$$\lambda_i(\mathbf{x}, \theta) = (\theta_i \partial_i \Psi(\mathbf{x}))^+ + \gamma_i(\mathbf{x}, \theta)$$

where  $\gamma_i$  is non-negative and  $\gamma_i(\mathbf{x}, \theta) = \gamma_i(\mathbf{x}, \theta_{-i})$  with  $\theta_{-i}$  is equal to  $\theta$  except for the  $i$ th component which is flipped.

- Then (under some regularity conditions)
  - The Zig-Zag process has  $p(\mathbf{x}|\mathbf{y})$  as invariant distribution
  - The process is ergodic:

$$\lim_{t \rightarrow \infty} \int_0^t f(\mathbf{x}(s)) ds = \int f(\mathbf{x}) \pi(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

## Additional topics in MCMC

- **Adaptive MCMC**: Automatic tuning of proposal distributions
  - Main challenge: Specifying proposal based on history of chain **breaks down the Markov property**
  - Solution: Reduce the amount of tuning as the number of iterations increases

## Additional topics in MCMC

- **Adaptive MCMC**: Automatic tuning of proposal distributions
  - Main challenge: Specifying proposal based on history of chain **breaks down the Markov property**
  - Solution: Reduce the amount of tuning as the number of iterations increases
- **Simulated tempering**
  - Define  $f^i(\mathbf{x}) \propto \pi(\mathbf{x})^{1/\tau_i}$ ,  $1 = \tau_1 < \tau_2 < \dots < \tau_m$
  - Simulate  $(\mathbf{x}, I)$ , where  $I$  changes distribution
  - Easier to move around when  $\tau_i > 1$
  - Keep samples for which  $I = 1$



## Additional topics in MCMC

- **Adaptive MCMC**: Automatic tuning of proposal distributions
  - Main challenge: Specifying proposal based on history of chain **breaks down the Markov property**
  - Solution: Reduce the amount of tuning as the number of iterations increases
- **Simulated tempering**
  - Define  $f^i(\mathbf{x}) \propto \pi(\mathbf{x})^{1/\tau_i}$ ,  $1 = \tau_1 < \tau_2 < \dots < \tau_m$
  - Simulate  $(\mathbf{x}, I)$ , where  $I$  changes distribution
  - Easier to move around when  $\tau_i > 1$
  - Keep samples for which  $I = 1$
- **Multiple-Try M-H**
  - Generate  $k$  proposals  $\mathbf{x}_1^*, \dots, \mathbf{x}_k^*$  from  $g(\cdot | \mathbf{x}^{(t)})$
  - Select  $\mathbf{x}_j^*$  with probability  $w(\mathbf{x}^{(t)}, \mathbf{x}_j^*) = \pi(\mathbf{x}^{(t)})g(\mathbf{x}_j^* | \mathbf{x}^{(t)})\lambda(\mathbf{x}^{(t)}, \mathbf{x}_j^*)$ ,  $\lambda$  symmetric
  - Sample  $\mathbf{x}_1^{**}, \dots, \mathbf{x}_{k-1}^{**}$  from  $g(\cdot | \mathbf{x}_j^*)$ , put  $\mathbf{x}_k^{**} = \mathbf{x}^{(t)}$
  - Use **Generalized M-H ratio**

$$R_g = \frac{\sum_{i=1}^k w(\mathbf{x}^{(t)}, \mathbf{x}_i^*)}{\sum_{i=1}^k w(\mathbf{x}_j^*, \mathbf{x}_i^{**})}$$

## Implement your SMC algorithm

- Consider the model for lemmings data
- Model

$$\mathbf{y}_t \sim \text{Binom} \left( 1, \frac{\exp(\mathbf{x}_t)}{1 + \exp(\mathbf{x}_t)} \right)$$

$$\mathbf{x}_t = a\mathbf{x}_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$a \sim \text{Uniform}[0, 1]$$

Use  $a = 0.5, \sigma = 1$

- Some data are missing: How to handle this?

## Implement your MCMC algorithm

- Consider the model for lemmings data
- Model

$$\mathbf{y}_t \sim \text{Binom} \left( 1, \frac{\exp(\mathbf{x}_t)}{1 + \exp(\mathbf{x}_t)} \right)$$

$$\mathbf{x}_t = a\mathbf{x}_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N(0, \sigma^2)$$

$$a \sim \text{Uniform}[0, 1]$$

Use  $a = 0.5, \sigma = 1$

- Note: if only changing  $\mathbf{x}_t \rightarrow \mathbf{x}_t^*$ :

$$\frac{\pi(\mathbf{x}^*)}{\pi(\mathbf{x})} = \frac{p(\mathbf{x}_t^*|\mathbf{x}_{t-1})p(\mathbf{x}_{t+1}|\mathbf{x}_t^*)p(\mathbf{y}_t|\mathbf{x}_t^*)}{p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t+1}|\mathbf{x}_t)p(\mathbf{y}_t|\mathbf{x}_t)}$$

## References

- C. Andrieu, A. Doucet, and R. Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, 2010.
- J. Bierkens, P. Fearnhead, and G. Roberts. The zig-zag process and super-efficient sampling for bayesian analysis of big data. *The Annals of Statistics*, 47(3):1288–1320, 2019.
- P. J. Green. Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika*, 82(4):711–732, 1995.
- A. Hubin and G. Storvik. Mode jumping mcmc for bayesian variable selection in glmm. *Computational Statistics & Data Analysis*, 127:281–297, 2018.
- A. Hubin, G. Storvik, and F. Frommlet. Flexible bayesian nonlinear model configuration. *Journal of Artificial Intelligence Research*, 72:901–942, 2021.
- H. Tjelmeland and B. K. Hegstad. Mode jumping proposals in mcmc. *Scandinavian journal of statistics*, 28(1):205–223, 2001.