

# Markov chain Monte Carlo

Geir Storvik

Geilo Winter school 2023



UiO : Universitetet i Oslo



# Outline

- 1 Theoretical properties
- 2 Hamiltonian MC
- 3 Pseudo-marginal MCMC
- 4 Reversible jump MCMC
- 5 Additional topics

## Theoretical properties



# Convergence issues of MCMC

- **Theoretical properties:**

$$\mathbf{x}^{(t)} \xrightarrow{\mathcal{D}} \pi(\mathbf{x}), \quad \text{as } t \rightarrow \infty$$

$$\hat{\theta}_1 = \frac{1}{L} \sum_{t=1}^L h(\mathbf{x}^{(t)}) \rightarrow E^p[h(\mathbf{x})] \quad \text{as } L \rightarrow \infty$$

- **Note:** We also have

$$\hat{\theta}_2 = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)}) \rightarrow E^p[h(\mathbf{x})] \quad \text{as } L \rightarrow \infty$$

- **Advantage:** Remove those variables with distribution very different from  $\pi(\mathbf{x})$
- **Disadvantage:** Need more samples
- **Question:** How to specify  $D$  and  $L$ ?
  - $D$ : Large enough so that  $\mathbf{x}^{(t)} \approx \pi(\mathbf{x})$  for  $t > D$  (bias small)
  - $L$ : Large enough so that  $\text{Var}[\hat{\theta}_2]$  is small enough

# Central limit theorems

- Under additional requirements (see e.g. [Robert and Casella \(1999\)](#), ch 6)) we have

$$\sqrt{L}(\hat{\theta} - \theta) \approx N(0, \gamma^2)$$

- Essential requirement:

$$\gamma^2 = \text{Var}_{\pi}[h(X_0)] + 2 \sum_{k=1}^{\infty} \text{Cov}_{\pi}[h(X_0), h(X_k)] < \infty$$

- Interpretation: Correlations should decay fast enough

# Mixing

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \text{Var}[h(\mathbf{x}^{(t)})] + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \right]$$

# Mixing

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \text{Var}[h(\mathbf{x}^{(t)})] + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \right]$$

Assume  $D$  large, so "converged":

$$\text{Var}[h(\mathbf{x}^{(t)})] \approx \sigma_h^2, \quad \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \approx \sigma_h^2 \rho(t-s)$$



# Mixing

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \text{Var}[h(\mathbf{x}^{(t)})] + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \right]$$

Assume  $D$  large, so "converged":

$$\text{Var}[h(\mathbf{x}^{(t)})] \approx \sigma_h^2, \quad \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \approx \sigma_h^2 \rho(t-s)$$

gives

$$\begin{aligned} \text{Var}[\hat{\theta}] &\approx \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \sigma_h^2 + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \sigma_h^2 \rho(t-s) \right] \\ &= \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{L-1} \frac{L-k}{L} \rho(k) \right] \end{aligned}$$

# Mixing

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \text{Var}[h(\mathbf{x}^{(t)})] + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \right]$$

Assume  $D$  large, so "converged":

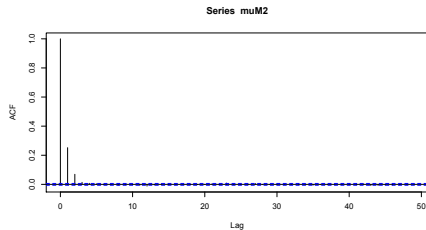
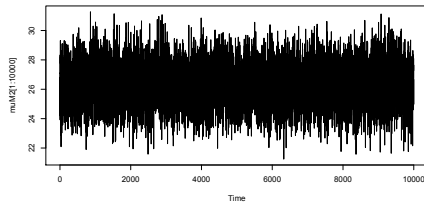
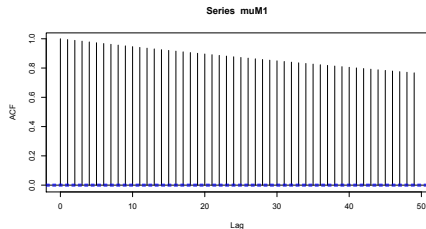
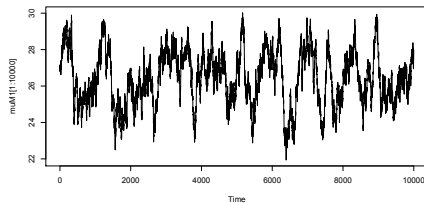
$$\text{Var}[h(\mathbf{x}^{(t)})] \approx \sigma_h^2, \quad \text{Cov}[h(\mathbf{x}^{(s)}), h(\mathbf{x}^{(t)})] \approx \sigma_h^2 \rho(t-s)$$

gives

$$\begin{aligned} \text{Var}[\hat{\theta}] &\approx \frac{1}{L^2} \left[ \sum_{t=D+1}^{D+L} \sigma_h^2 + 2 \sum_{s=D+1}^{D+L-1} \sum_{t=s+1}^{D+L} \sigma_h^2 \rho(t-s) \right] \\ &= \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{L-1} \frac{L-k}{L} \rho(k) \right] \end{aligned}$$

- Good mixing:**  $\rho(k)$  decreases fast with  $k$ !

# Example



# How to assess convergence?

- Graphical diagnostics:
  - Sample paths:
    - Plot  $h(\mathbf{x}^{(t)})$  as function of  $t$
    - Useful with **different**  $h(\cdot)$  functions!

# How to assess convergence?

- Graphical diagnostics:

- Sample paths:

- Plot  $h(\mathbf{x}^{(t)})$  as function of  $t$
    - Useful with **different**  $h(\cdot)$  functions!

- Cusum diagnostics

- Plot  $\sum_{i=1}^t [h(\mathbf{x}^{(i)}) - \hat{\theta}_n]$  versus  $t$
    - Wiggly and small excursions from 0: Indicate chain is mixing well

# How to assess convergence?

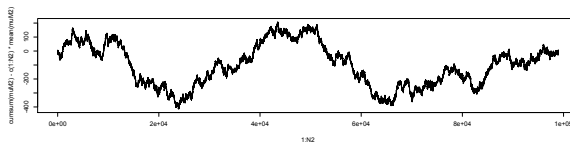
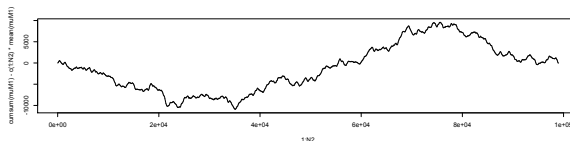
- Graphical diagnostics:

- Sample paths:

- Plot  $h(\mathbf{x}^{(t)})$  as function of  $t$
    - Useful with **different**  $h(\cdot)$  functions!

- Cusum diagnostics

- Plot  $\sum_{i=1}^t [h(\mathbf{x}^{(i)}) - \hat{\theta}_n]$  versus  $t$
    - Wiggly and small excursions from 0: Indicate chain is mixing well



# The Gelman-Rubin diagnostic

- Motivated from **analysis of variance**
- Assume  $J$  chains run in parallel
- $j$ th chain:  $x_j^{(D+1)}, \dots, x_j^{(D+L)}$  (first  $D$  discarded)

# The Gelman-Rubin diagnostic

- Motivated from **analysis of variance**
- Assume  $J$  chains run in parallel
- $j$ th chain:  $x_j^{(D+1)}, \dots, x_j^{(D+L)}$  (first  $D$  discarded)
- Define

$$\bar{x}_j = \frac{1}{L} \sum_{t=D+1}^{D+L} x_j^{(t)}$$

$$\bar{x} = \frac{1}{J} \sum_{j=1}^J \bar{x}_j$$

$$B = \frac{L}{J-1} \sum_{j=1}^J (\bar{x}_j - \bar{x})^2$$

$$W = \frac{1}{J} \sum_{j=1}^J s_j^2$$

$$s_j^2 = \frac{1}{L-1} \sum_{t=D+1}^{D+L} (x_j^{(t)} - \bar{x}_j)^2$$



# The Gelman-Rubin diagnostic

- Motivated from **analysis of variance**
- Assume  $J$  chains run in parallel
- $j$ th chain:  $x_j^{(D+1)}, \dots, x_j^{(D+L)}$  (first  $D$  discarded)
- Define

$$\bar{x}_j = \frac{1}{L} \sum_{t=D+1}^{D+L} x_j^{(t)}$$

$$\bar{x} = \frac{1}{J} \sum_{j=1}^J \bar{x}_j$$

$$B = \frac{L}{J-1} \sum_{j=1}^J (\bar{x}_j - \bar{x})^2$$

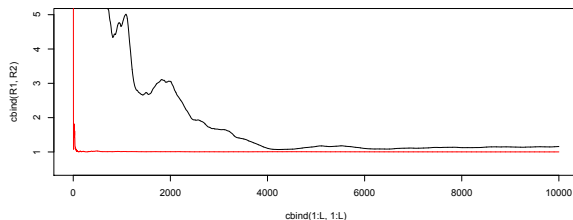
$$W = \frac{1}{J} \sum_{j=1}^J s_j^2$$

$$s_j^2 = \frac{1}{L-1} \sum_{t=D+1}^{D+L} (x_j^{(t)} - \bar{x}_j)^2$$

- If converged, both  $B$  and  $W$  estimates  $\sigma^2 = \text{Var}_f[X]$
- Diagnostic:  $R = \frac{\frac{L-1}{L} W + \frac{1}{L} B}{W}$
- "Rule":  $\sqrt{R} < 1.1$  indicate  $D$  and  $L$  are sufficient

# Example

- $D = 100, L = 1000$ :  $\sqrt{R_1} = 1.588, \sqrt{R_2} = 1.002$ ,
- $D = 1000, L = 1000$ :  $\sqrt{R_1} = 1.700, \sqrt{R_2} = 1.004$ ,
- $D = 1000, L = 10000$ :  $\sqrt{R_1} = 1.049, \sqrt{R_2} = 1.0008$







# M-H: Choice of proposal distribution

- Independence chain:
  - $g(\cdot) \approx p(\cdot)$
  - High acceptance rate
  - Tail properties most important:  $f/g$  should be bounded

# M-H: Choice of proposal distribution

- Independence chain:
  - $g(\cdot) \approx p(\cdot)$
  - High acceptance rate
  - Tail properties most important:  $f/g$  should be bounded
- Random walk proposal
  - Tune variance so that acceptance rate is between 25 and 50%

# Effective sample size for MCMC

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{L-1} \frac{L-k}{L} \rho(k) \right] \xrightarrow{L \rightarrow \infty} \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{\infty} \rho(k) \right]$$

# Effective sample size for MCMC

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{L-1} \frac{L-k}{L} \rho(k) \right] \xrightarrow{L \rightarrow \infty} \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{\infty} \rho(k) \right]$$

- If independent samples:

$$\text{Var}[\hat{\theta}] = \frac{\sigma_h^2}{L}$$



# Effective sample size for MCMC

- For  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

$$\text{Var}[\hat{\theta}] = \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{L-1} \frac{L-k}{L} \rho(k) \right] \xrightarrow{L \rightarrow \infty} \frac{\sigma_h^2}{L} \left[ 1 + 2 \sum_{k=1}^{\infty} \rho(k) \right]$$

- If independent samples:

$$\text{Var}[\hat{\theta}] = \frac{\sigma_h^2}{L}$$

- Effective sample size:  $\frac{L}{1 + 2 \sum_{k=1}^{\infty} \rho(k)}$
- Use empirical estimates  $\hat{\rho}(k)$
- Usual to truncate the summation when  $\hat{\rho}(k) < 0.1$ .

# Number of chains

- Assume possible to perform  $N$  iterations
  - One long chain of length  $N$ , or
  - $J$  parallel chains, each of length  $N/J$ ?

# Number of chains

- Assume possible to perform  $N$  iterations
  - One long chain of length  $N$ , or
  - $J$  parallel chains, each of length  $N/J$ ?
- **Burnin:**
  - One long chain: Only need to discard  $D$  samples
  - Parallel chains: Need to discard  $J \cdot D$  samples

# Number of chains

- Assume possible to perform  $N$  iterations
  - One long chain of length  $N$ , or
  - $J$  parallel chains, each of length  $N/J$ ?
- **Burnin:**
  - One long chain: Only need to discard  $D$  samples
  - Parallel chains: Need to discard  $J \cdot D$  samples
- **Check of convergence**
  - Easier with many parallel chains

# Number of chains

- Assume possible to perform  $N$  iterations
  - One long chain of length  $N$ , or
  - $J$  parallel chains, each of length  $N/J$ ?
- **Burnin:**
  - One long chain: Only need to discard  $D$  samples
  - Parallel chains: Need to discard  $J \cdot D$  samples
- **Check of convergence**
  - Easier with many parallel chains
- **Efficiency**
  - Parallel chains give more independent samples

# Number of chains

- Assume possible to perform  $N$  iterations
  - One long chain of length  $N$ , or
  - $J$  parallel chains, each of length  $N/J$ ?
- **Burnin:**
  - One long chain: Only need to discard  $D$  samples
  - Parallel chains: Need to discard  $J \cdot D$  samples
- **Check of convergence**
  - Easier with many parallel chains
- **Efficiency**
  - Parallel chains give more independent samples
- **Computational issues**
  - Possible to utilize multiple cores with parallel chains

# Example

- Number of positive tests of Covid-19 Nov 1 - Nov 7 2021 (county level):

Day	1	2	3	4	5	6	7	8	9	10
Cases	845	331	76	47	1105	126	186	58	156	258
Population	693494	479892	265238	241235	1241165	371385	419396	307231	636531	468702

- Model:

$$y_j \sim \text{Poisson}(N_j \theta_j)$$

$$\theta_j \sim \text{Gamma}(a, b)$$

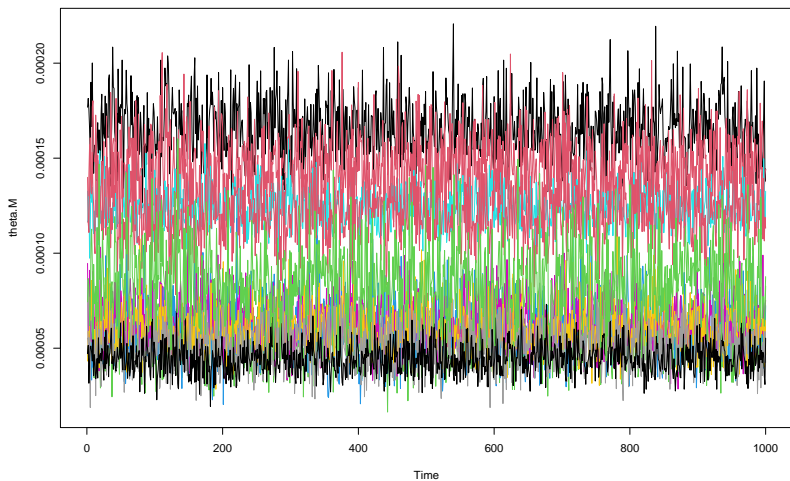
$$p(b|a) \sim \text{Gamma}(\alpha, \beta) \quad a \text{ assumed fixed}$$

- Can show

$$p(\theta_{1:11}|a, b, \mathbf{y}) = \prod_{j=1}^{11} \text{Gamma}(a + y_j, b + N_j)$$

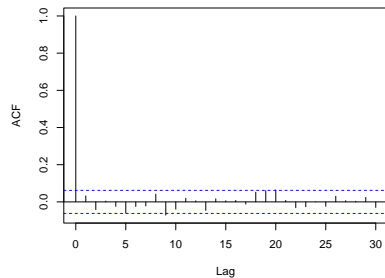
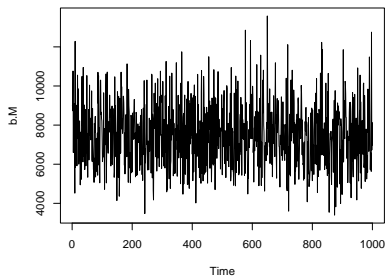
$$p(b|a, \theta, \mathbf{y}) = \text{Gamma}(\alpha + 11a, \beta + \sum_{j=1}^{11} 1\theta_j)$$

# Res-covid





# Res-covid



# Data uncertainty and Monte Carlo uncertainty

- **Parameter:**  $\theta = E^p[h(\mathbf{x})]$
- **Estimator:**  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :

# Data uncertainty and Monte Carlo uncertainty

- **Parameter:**  $\theta = E^p[h(\mathbf{x})]$
- **Estimator:**  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :
- **Two types of uncertainty**
  - Variability in  $h(\mathbf{x})$ :  $\sigma_h^2 = \text{Var}^p[h(\mathbf{x})]$ 
    - Estimator:  $\hat{\sigma}_h^2 = \frac{1}{L} \sum_{t=D+1}^{D+L} [h(\mathbf{x}^{(t)}) - \hat{\theta}]^2$

# Data uncertainty and Monte Carlo uncertainty

- **Parameter:**  $\theta = E^p[h(\mathbf{x})]$
- **Estimator:**  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :
- **Two types of uncertainty**
  - Variability in  $h(\mathbf{x})$ :  $\sigma_h^2 = \text{Var}^p[h(\mathbf{x})]$ 
    - Estimator:  $\hat{\sigma}_h^2 = \frac{1}{L} \sum_{t=D+1}^{D+L} [h(\mathbf{x}^{(t)}) - \hat{\theta}]^2$
  - MC variability in  $\hat{\theta}$ :
    - Estimator: Divide data into **batches** of size  $b = \lfloor L^{1/a} \rfloor$ , make estimates  $\hat{\theta}$  within each batch and variance from these

# Data uncertainty and Monte Carlo uncertainty

- **Parameter:**  $\theta = E^p[h(\mathbf{x})]$
- **Estimator:**  $\hat{\theta} = \frac{1}{L} \sum_{t=D+1}^{D+L} h(\mathbf{x}^{(t)})$ :
- **Two types of uncertainty**
  - Variability in  $h(\mathbf{x})$ :  $\sigma_h^2 = \text{Var}^p[h(\mathbf{x})]$ 
    - Estimator:  $\hat{\sigma}_h^2 = \frac{1}{L} \sum_{t=D+1}^{D+L} [h(\mathbf{x}^{(t)}) - \hat{\theta}]^2$
  - MC variability in  $\hat{\theta}$ :
    - Estimator: Divide data into **batches** of size  $b = \lfloor L^{1/a} \rfloor$ , make estimates  $\hat{\theta}$  within each batch and variance from these
- **Recommendation:** Specify  $L$  so that MC variability is less than 5% of variability in  $h(\mathbf{x})$ .

# Uncertainty in Bayesian setting

- Interest in  $\theta \sim p(\theta|\mathbf{y})$
- Monte Carlo estimate of mean:  $\hat{\mu}_\theta = \frac{1}{L} \sum_{i=1}^L \theta^i$
- Can show:

$$E[\hat{\mu}_\theta] = \mu_\theta = E[\theta|\mathbf{y}]$$

$$E[(\theta - \hat{\mu}_\theta)^2] = \left(1 + \frac{1}{L}\right) E[(\theta - \mu_\theta)^2]$$

## Hamiltonian MC

# Hamiltonian MC

- Common trick in Monte Carlo: Introduce **auxiliary variables**
- Hamiltonian MC (**Neal et al., 2011**):

$$\pi(\mathbf{q}) \propto \exp(-U(\mathbf{q}))$$

Distribution of interest

$$\pi(\mathbf{q}, \mathbf{p}) \propto \exp(-U(\mathbf{q}) - 0.5\mathbf{p}^T\mathbf{p})$$

Extended distribution

$$= \exp(-H(\mathbf{q}, \mathbf{p}))$$

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + 0.5\mathbf{p}^T\mathbf{p}$$



# Hamiltonian MC

- Common trick in Monte Carlo: Introduce **auxiliary variables**
- Hamiltonian MC (**Neal et al., 2011**):

$$\pi(\mathbf{q}) \propto \exp(-U(\mathbf{q}))$$

Distribution of interest

$$\pi(\mathbf{q}, \mathbf{p}) \propto \exp(-U(\mathbf{q}) - 0.5\mathbf{p}^T \mathbf{p})$$

Extended distribution

$$= \exp(-H(\mathbf{q}, \mathbf{p}))$$

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + 0.5\mathbf{p}^T \mathbf{p}$$

- Note
  - $\mathbf{q}$  and  $\mathbf{p}$  are **independent**
  - $\mathbf{p} \sim N(\mathbf{0}, I)$ .
  - Usually  $\dim(\mathbf{p}) = \dim(\mathbf{q})$

# Hamiltonian MC

- Common trick in Monte Carlo: Introduce **auxiliary variables**
- Hamiltonian MC (Neal et al., 2011):

$$\pi(\mathbf{q}) \propto \exp(-U(\mathbf{q}))$$

Distribution of interest

$$\pi(\mathbf{q}, \mathbf{p}) \propto \exp(-U(\mathbf{q}) - 0.5\mathbf{p}^T \mathbf{p})$$

Extended distribution

$$= \exp(-H(\mathbf{q}, \mathbf{p}))$$

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + 0.5\mathbf{p}^T \mathbf{p}$$

- Note
  - $\mathbf{q}$  and  $\mathbf{p}$  are **independent**
  - $\mathbf{p} \sim N(\mathbf{0}, I)$ .
  - Usually  $\dim(\mathbf{p}) = \dim(\mathbf{q})$
- Algorithm ( $\mathbf{q}$ ) current value
  - 1 Simulate  $\mathbf{p} \sim N(\mathbf{0}, I)$
  - 2 Generate  $(\mathbf{q}^*, \mathbf{p}^*)$  such that  $H(\mathbf{q}^*, \mathbf{p}^*) \approx H(\mathbf{q}, \mathbf{p})$
  - 3 Accept  $(\mathbf{q}^*, \mathbf{p}^*)$  by a Metropolis-Hastings step

# Hamiltonian MC

- Common trick in Monte Carlo: Introduce **auxiliary variables**
- Hamiltonian MC (Neal et al., 2011):

$$\pi(\mathbf{q}) \propto \exp(-U(\mathbf{q}))$$

Distribution of interest

$$\pi(\mathbf{q}, \mathbf{p}) \propto \exp(-U(\mathbf{q}) - 0.5\mathbf{p}^T \mathbf{p})$$

Extended distribution

$$= \exp(-H(\mathbf{q}, \mathbf{p}))$$

$$H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + 0.5\mathbf{p}^T \mathbf{p}$$

- Note
  - $\mathbf{q}$  and  $\mathbf{p}$  are **independent**
  - $\mathbf{p} \sim N(\mathbf{0}, I)$ .
  - Usually  $\dim(\mathbf{p}) = \dim(\mathbf{q})$
- Algorithm ( $\mathbf{q}$ ) current value
  - 1 Simulate  $\mathbf{p} \sim N(\mathbf{0}, I)$
  - 2 Generate  $(\mathbf{q}^*, \mathbf{p}^*)$  such that  $H(\mathbf{q}^*, \mathbf{p}^*) \approx H(\mathbf{q}, \mathbf{p})$
  - 3 Accept  $(\mathbf{q}^*, \mathbf{p}^*)$  by a Metropolis-Hastings step
- Step 1 is a Gibbs sampling step!
- Main challenge: Generate  $(\mathbf{q}^*, \mathbf{p}^*)$

# Hamiltonian dynamics

- Consider  $(\mathbf{q}, \mathbf{p})$  as a time-process  $(\mathbf{q}(t), \mathbf{p}(t))$
- **Hamiltonian dynamics**: Change through

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = - \frac{\partial H}{\partial q_i}$$

# Hamiltonian dynamics

- Consider  $(\mathbf{q}, \mathbf{p})$  as a time-process  $(\mathbf{q}(t), \mathbf{p}(t))$
- Hamiltonian dynamics**: Change through

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = - \frac{\partial H}{\partial q_i}$$

This gives

$$\begin{aligned} \frac{dH}{dt} &= \sum_{i=1}^d \left[ \frac{\partial H}{\partial q_i} \frac{dq_i}{dt} + \frac{\partial H}{\partial p_i} \frac{dp_i}{dt} \right] \\ &= \sum_{i=1}^d \left[ \frac{\partial H}{\partial q_i} \frac{\partial H}{\partial p_i} - \frac{\partial H}{\partial p_i} \frac{\partial H}{\partial q_i} \right] = 0 \end{aligned}$$

# Hamiltonian dynamics

- Consider  $(\mathbf{q}, \mathbf{p})$  as a time-process  $(\mathbf{q}(t), \mathbf{p}(t))$
- Hamiltonian dynamics**: Change through

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = - \frac{\partial H}{\partial q_i}$$

This gives

$$\begin{aligned} \frac{dH}{dt} &= \sum_{i=1}^d \left[ \frac{\partial H}{\partial q_i} \frac{dq_i}{dt} + \frac{\partial H}{\partial p_i} \frac{dp_i}{dt} \right] \\ &= \sum_{i=1}^d \left[ \frac{\partial H}{\partial q_i} \frac{\partial H}{\partial p_i} - \frac{\partial H}{\partial p_i} \frac{\partial H}{\partial q_i} \right] = 0 \end{aligned}$$

- If we can change  $(\mathbf{q}, \mathbf{p})$  exactly by the Hamiltonian dynamics,  $H$  will not change!
- In practice, only possible to make numerical approximations

# Hamiltonian dynamics - Eulers method

- Assume

$$\begin{aligned} p_i(t + \varepsilon) &= p_i(t) + \varepsilon \frac{dp_i}{dt}(t) \\ &= p_i(t) - \varepsilon \frac{\partial U}{\partial q_i}(q_i(t)) \end{aligned}$$

$$\begin{aligned} q_i(t + \varepsilon) &= q_i(t) + \varepsilon \frac{dq_i}{dt}(t) \\ &= q_i(t) + \varepsilon p_i(t) \end{aligned}$$

- Note: **Derivatives** of  $U(\mathbf{q})$  are used.
- However, not very exact.

# Hamiltonian dynamics - Eulers method

- Assume

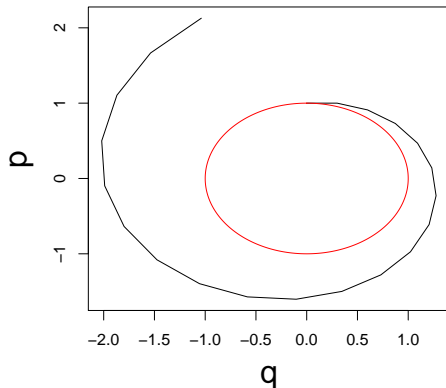
$$p_i(t + \varepsilon) = p_i(t) + \varepsilon \frac{dp_i}{dt}(t)$$

$$= p_i(t) - \varepsilon \frac{\partial U}{\partial q_i}(q_i(t))$$

$$q_i(t + \varepsilon) = q_i(t) + \varepsilon \frac{dq_i}{dt}(t)$$

$$= q_i(t) + \varepsilon p_i(t)$$

- Note: **Derivatives** of  $U(\mathbf{q})$  are used.
- However, not very exact.





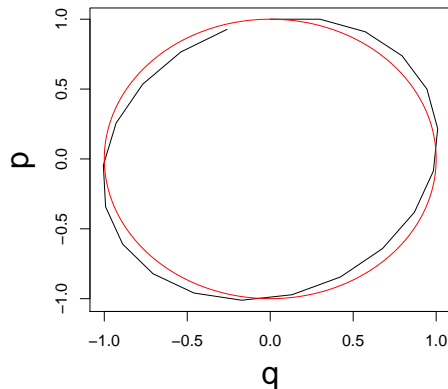
# Hamiltonian dynamics - the modified Eulers method

- Assume

$$p_i(t + \varepsilon) = p_i(t) - \varepsilon \frac{\partial U}{\partial q_i}(q(t))$$

$$q_i(t + \varepsilon) = q_i(t) + \varepsilon p_i(t + \varepsilon)$$

- Better than Eulers method.



# Hamiltonian dynamics - the Leapfrog method

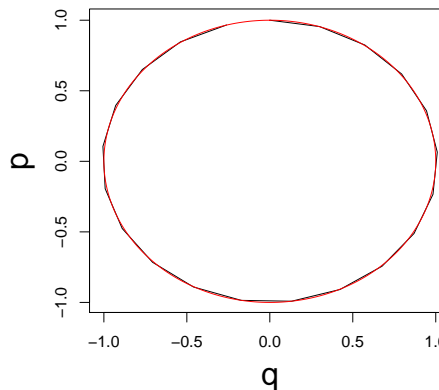
- Assume

$$p_i(t + \frac{\varepsilon}{2}) = p_i(t) - \frac{\varepsilon}{2} \frac{\partial U}{\partial q_i}(q(t))$$

$$q_i(t + \varepsilon) = q_i(t) + \varepsilon p_i(t + \frac{\varepsilon}{2})$$

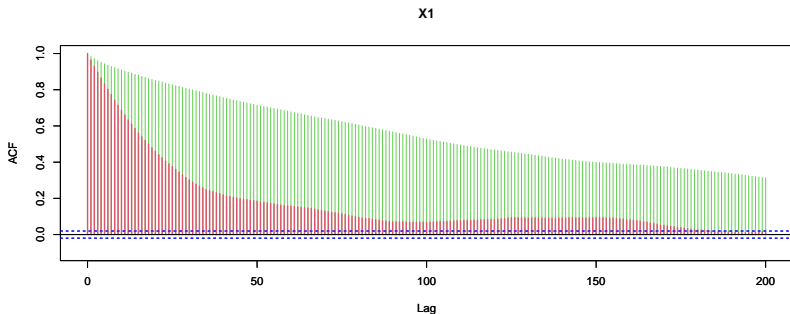
$$p_i(t + \varepsilon) = p_i(t + \frac{\varepsilon}{2}) - \frac{\varepsilon}{2} \frac{\partial U}{\partial q_i}(q(t + \varepsilon))$$

- Quite exact!
- Idea: Use this  $L$  steps



## Example - 2-dimensional Gaussian

- Assume  $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ ,  $\Sigma = \begin{pmatrix} 1 & 0.95 \\ 0.95 & 1 \end{pmatrix}$
- $H(\mathbf{x}, \mathbf{p}) = 0.5\mathbf{x}^T \Sigma^{-1} \mathbf{x} + 0.5\mathbf{p}^T \mathbf{p}$
- Use  $L = 5$  leapfrog steps, with stepsize  $\varepsilon = 0.1$
- Leapfrog\_Gauss2.R



# Example - mixture Gaussians

- Assume

$$\pi(x) = pN(x; \mu_1, \sigma_1^2) + (1 - p)N(x; \mu_2, \sigma_2^2)$$

- $H(\mathbf{x}, \mathbf{p}) = -\log(\pi(x) + 0.5\mathbf{p}^T \mathbf{p})$
- Use  $L = 5$  leapfrog steps, with stepsize  $\varepsilon = 0.1$
- `Leapfrog_mixture.R`

## Pseudo-marginal MCMC

# Pseudo-marginal MCMC - Andrieu and Roberts (2009)

- Assume a Hierarchical model

$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}; \theta)$$

Observations

$$\mathbf{x} \sim p(\mathbf{x}|\theta)$$

Latent process

$$\theta \sim p(\theta)$$

Prior

- Main interest in  $\theta$ ,

$$\pi(\theta) = p(\theta|\mathbf{y}) = \int_{\mathbf{x}} p(\theta, \mathbf{x}|\mathbf{y}) d\mathbf{x}$$

- More general setting: Interest in only a small subsets of the unknowns
- Possible: Construct MCMC algorithm for  $p(\mathbf{x}, \theta|\mathbf{y})$
- More efficient: MCMC algorithm directly for  $p(\theta|\mathbf{y})$ 
  - M-H: Require calculation of  $\frac{\pi(\theta^*)}{\pi(\theta)}$ , difficult!
- Idea of pseudo-marginal MCMC: Replace  $\pi(\theta^*)$  by an (unbiased) estimate!

# Pseudo-marginal MCMC - cont

- Algorithm, current values  $\theta_t, \hat{\pi}(\theta_t)$ :

- 1 Sample  $\theta^* \sim g(\theta^*|\theta)$
- 2 Construct estimate  $\hat{\pi}(\theta^*)$
- 3 Accept  $\theta^*$  with probability

$$\hat{R}(\theta^*) = \min \left\{ 1, \frac{\hat{\pi}(\theta^*)g(\theta|\theta^*)}{\hat{\pi}(\theta)g(\theta^*|\theta)} \right\}$$

- Property: If

- the M-H algorithm with **exact**  $\pi(\theta)$  converges "properly"
- the estimates  $\hat{\pi}(\theta)$  are unbiased

then the approximate M-H algorithm will also converge "properly"

- but convergence is slower!
- See **Andrieu and Roberts (2009)** for more details

# How to estimate $\pi(\theta)$ ?

- We have

$$\begin{aligned}
 \pi(\theta) &= \int_{\mathbf{x}} \pi(\mathbf{x}, \theta) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \frac{\pi(\mathbf{x}, \theta)}{g_{\theta}(\mathbf{x})} g_{\theta}(\mathbf{x}) d\mathbf{x} \\
 &\approx \frac{1}{N} \frac{\pi(\mathbf{x}^{(i)}, \theta)}{g_{\theta}(\mathbf{x}^{(i)})} \quad \mathbf{x}^{(i)} \sim g_{\theta}(\mathbf{x})
 \end{aligned}$$

- Unknown normalization constant cancels out in ratio.
- Alternatives possible



# Why do pseudo-marginal MCMC work?

- Consider extended distribution

$$\tilde{\pi}(\theta, \mathbf{x}_{1:N}) = \frac{1}{N} \sum_{k=1}^N \pi(\theta, \mathbf{x}_k) \prod_{i \neq k} g_{\theta}(\mathbf{x}_i)$$

- We have

$$\int_{\mathbf{x}_{1:N}} \tilde{\pi}(\theta, \mathbf{x}_{1:N}) d\mathbf{x}_{1:N} = \pi(\theta)$$

- M-H algorithm (given  $(\theta, \mathbf{x}_{1:N})$ ):

- 1 Generate  $\theta^* \sim g(\theta^*|\theta)$
- 2 Generate  $\mathbf{x}_{1:N}^* \sim \prod_{k=1}^N g_{\theta^*}(\mathbf{x}_k^*)$
- 3 Calculate acceptance ratio

$$r = \frac{\tilde{\pi}(\theta^*, \mathbf{x}_{1:N}^*)}{\tilde{\pi}(\theta, \mathbf{x}_{1:N})} \times \frac{g(\theta^*|\theta) \prod_{k=1}^N g_{\theta^*}(\mathbf{x}_k^*)}{g(\theta^*|\theta) \prod_{k=1}^N g_{\theta^*}(\mathbf{x}_k^*)}$$

# Logist regression

- Assume

$$y_i \sim \text{Binom} \left( 1, \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right)$$

$$\eta_i = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij}$$

- Different models depending on  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$
- Marginal distributions:

$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\gamma}) &= \int_{\boldsymbol{\beta}_\gamma} p(\mathbf{y}|\boldsymbol{\beta}_\gamma) p(\boldsymbol{\beta}_\gamma|\boldsymbol{\gamma}) d\boldsymbol{\beta}_\gamma \\ &= \int_{\boldsymbol{\beta}_\gamma} p(\mathbf{y}|\boldsymbol{\beta}_\gamma) \frac{p(\boldsymbol{\beta}_\gamma|\boldsymbol{\gamma})}{g(\boldsymbol{\beta}_\gamma|\boldsymbol{\gamma})} g(\boldsymbol{\beta}_\gamma|\boldsymbol{\gamma}) d\boldsymbol{\beta}_\gamma \end{aligned}$$

- Main challenge: Choice of  $g()$
- Script: `Logreg_Pseudo.R`

## Reversible jump MCMC

# 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  $\theta_2 = \theta_2(\theta_1, \mathbf{u}_1)$
- 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

# Logist regression

- Assume

$$y_i \sim \text{Binom} \left( 1, \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right)$$

$$\eta_t = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij}$$

- Different models depending on  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$
- RJMCMC with changing **one**  $\gamma_j$  at a time:  $\gamma_j \rightarrow 1 - \gamma_j$

# RJ - logistic regression

- Adding component:  $\gamma_j = 0 \rightarrow \gamma_j = 1$
- Increase dimension by 1.
- $u_1$ : Simulation of  $\beta_j$
- Jacobian=1
- Script: `Logreg_RJ.R`

Additional topics



# Additional topics in MCMC

- Mode jumping MCMC
- Reversible jump MCMC
- Non-reversible MCMC
- Subsampling MCMC
- Continuous-time Markov processes
- **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

- Mode jumping MCMC
- Reversible jump MCMC
- Non-reversible MCMC
- Subsampling MCMC
- Continuous-time Markov processes
- **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

- Mode jumping MCMC
- Reversible jump MCMC
- Non-reversible MCMC
- Subsampling MCMC
- Continuous-time Markov processes
- **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**

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

- C. Andrieu and G. O. Roberts. The pseudo-marginal approach for efficient monte carlo computations. *The Annals of Statistics*, 37(2):697–725, 2009.
- P. J. Green. Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika*, 82(4):711–732, 1995.
- R. M. Neal et al. MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2(11):2, 2011.
- C. P. Robert and G. Casella. *Monte Carlo statistical methods*, volume 2. Springer, 1999.