

MATH-414 – Stochastic simulation

Lecture 12: Metropolis-Hastings algorithms

Prof. Fabio Nobile

Outline

Metropolis-Hastings algorithms

Convergence diagnostics

Independence sampler

Idea: take proposal density $q(x, y) = g(y)$ independent of the current state x , where $g : \mathcal{X} \rightarrow \mathbb{R}_+$ is a probability density function on \mathcal{X} that dominates f (i.e. $g(x) = 0 \Rightarrow f(x) = 0$)

Algorithm: Independence sampler Metropolis-Hastings

Given: $X_0 \sim \lambda$, $\text{supp}(\lambda) \subset \text{supp}(f)$

1 **for** $n = 0, 1, \dots$, **do**

2 Generate $Y_{n+1} \sim g$

3 Compute $\alpha(X_n, Y_{n+1}) = \min \left\{ \frac{f(Y_{n+1})}{f(X_n)} \frac{g(X_n)}{g(Y_{n+1})}, 1 \right\}$

4 Generate $U \sim \mathcal{U}(0, 1)$ and set

$$X_{n+1} = Y_{n+1}, \text{ if } U \leq \alpha(X_n, Y_{n+1}), \quad X_{n+1} = X_n, \text{ otherwise}$$

5 **end**

Similar to Acceptance-Rejection sampling but:

- ▶ Whenever the proposal is rejected, the current state is repeated in the chain, contrary to AR (\rightsquigarrow induces correlation in the sequence)
- ▶ No need to estimate the constant $C = \sup_{x \in \mathcal{X}} g(x)/f(x)$

Convergence of Independence sampler

Lemma

Let $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ be a Markov transition kernel with invariant measure π . If there exists $\epsilon \in (0, 1)$ and a probability measure ν on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ such that

$$P(x, A) \geq \epsilon \nu(A), \quad \forall x \in \mathcal{X}, A \in \mathcal{B}(\mathcal{X}), \quad (1)$$

then

$$\|\pi^{n,\lambda} - \pi\|_{TV} \leq 2(1 - \epsilon)^n. \quad (2)$$

- ▶ The condition (1) is called *uniform minorizing condition*.
- ▶ An exponential convergence of the type (2) is called *geometric ergodicity*

Proof

Consider two coupled chains $\{X_n\} \sim \text{Markov}(\lambda, P)$ and $\{Y_n\} \sim \text{Markov}(\pi, P)$ constructed using the following algorithm. (Rk. $\{Y_n\}$ is at stationarity)

```
1 Let  $X_0 \sim \lambda, Y_0 \sim \pi$ 
2 for  $n = 0, 1, \dots$ , do
3   Draw  $Z_n \sim \text{Be}(\epsilon), \mathbb{P}(Z_n = 1) = \epsilon, \mathbb{P}(Z_n = 0) = 1 - \epsilon$ 
4   if  $Z_n = 1$  then
5     draw  $W \sim \nu$  and set  $X_{n+1} = Y_{n+1} = W$ 
6   else
7     draw  $X_{n+1} \sim \frac{P(X_n, \cdot) - \epsilon \nu(\cdot)}{1 - \epsilon}$  and  $Y_{n+1} \sim \frac{P(Y_n, \cdot) - \epsilon \nu(\cdot)}{1 - \epsilon}$  independently
8   end
9 end
```

- ▶ It is easy to verify that $\{X_n\} \sim \text{Markov}(\lambda, P)$ and $\{Y_n\} \sim \text{Markov}(\pi, P)$.
- ▶ Let $T = \inf\{n \geq 0 : Z_n = 1\}$, which satisfies $\mathbb{P}(T \geq n) = (1 - \epsilon)^n$.
- ▶ After T , the two chains have the same distribution $X_n \sim Y_n, n > T$.

$$\begin{aligned} \|\pi^{n, \lambda} - \pi\|_{\text{TV}} &= 2 \sup_{A \in \mathcal{B}(\mathcal{X})} |\pi^{n, \lambda}(A) - \pi(A)| = 2 \sup_{A \in \mathcal{B}(\mathcal{X})} |\mathbb{P}(X_n \in A) - \mathbb{P}(Y_n \in A)| \\ &= 2 \sup_A |\mathbb{P}(X_n \in A, T < n) + \mathbb{P}(X_n \in A, T \geq n) - \mathbb{P}(Y_n \in A, T < n) - \mathbb{P}(Y_n \in A, T \geq n)| \\ &= 2 \sup_A |\mathbb{P}(X_n \in A, T \geq n) - \mathbb{P}(Y_n \in A, T \geq n)| \\ &= 2 \sup_A |\mathbb{P}(X_n \in A, Y_n \notin A, T \geq n) - \mathbb{P}(X_n \notin A, Y_n \in A, T \geq n)| \leq 2\mathbb{P}(T \geq n). \end{aligned}$$

Convergence of Independence sampler

Theorem

If there exists $M < +\infty$ such that $f(x) \leq Mg(x)$ for all $x \in \mathcal{X}$, then the chain generated by the independence sampler is uniformly ergodic and

$$\|\pi^{n,\lambda} - \pi\|_{TV} \leq 2 \left(1 - \frac{C}{M}\right)^n, \quad \text{for any } \lambda, \text{ with } C = \int_{\mathcal{X}} f(x) dx.$$

Proof: If f is not normalized, let $\tilde{f} = f/C$, $C = \int_{\mathcal{X}} f$. Notice that

$$\alpha(x, y)q(x, y) = g(y) \min \left\{ \frac{f(y)g(x)}{f(x)g(y)}, 1 \right\} = f(y) \min \left\{ \frac{g(x)}{f(x)}, \frac{g(y)}{f(y)} \right\} \geq \frac{1}{M} f(y).$$

It follows that for any $A \in \mathcal{B}(\mathcal{X})$,

$$P(x, A) = \int_A \alpha(x, y)q(x, y)dy + (1 - \alpha^*(x))\mathbb{1}_A(x) \geq \frac{1}{M} \int_A f(y) dy \geq \frac{C}{M} \pi(A)$$

and the result follows from Lemma 1.

Random walk Metropolis (RWM)

Idea: perform only local moves with proposed increment distributions identical and symmetric, i.e. $q(x, y) = q(|y - x|)$.

Typical case $q(x, \cdot) = N(x, \sigma^2 I_{d \times d})$.

This algorithm leads to geometric ergodicity under the following (sufficient) conditions (see [Jarner-Hansen, 2000])

- ▶ f has super-exponential tails, i.e. it is positive, continuous and satisfies

$$\lim_{|x| \rightarrow \infty} \frac{x}{|x|} \cdot \nabla \log f(x) = -\infty$$

- ▶ f satisfies

$$\limsup_{|x| \rightarrow \infty} \frac{x}{|x|} \cdot \frac{\nabla f(x)}{|\nabla f(x)|} < 0$$

- ▶ q is bounded away from zero in some region around zero:

$$\exists \delta_q, \epsilon_q > 0 \text{ s.t. } q(x) \geq \epsilon_q, \text{ for } |x| \leq \delta_q$$

One variable at a time

- ▶ Suppose $\mathcal{X} = \mathcal{X}^{(1)} \times \dots \times \mathcal{X}^{(d)}$ and $x \in \mathcal{X}$ has components $x = (x^{(1)}, \dots, x^{(d)})$; Notation: $x^{(-i)} = (x^{(1)}, \dots, x^{(i-1)}, x^{(i+1)}, \dots, x^{(d)})$.
- ▶ Consider a family of proposal transition densities $q^{(i)} : \mathcal{X} \times \mathcal{X}^{(i)} \rightarrow \mathbb{R}_+$

Idea: update one component at the time either chosen randomly or by performing a systematic sweep over the components.

Algorithm: One variable at a time MH with **random selection**.

- 1 Generate $X_0 \sim \lambda$
- 2 **for** $n = 0, 1, \dots$ **do**
 - 3 Draw index $i_n \sim \beta$ (p.m.f on $\{1, \dots, d\}$). Set $x = X_n^{(i_n)}$
 - 4 Draw $y \sim q_{i_n}(X_n, \cdot)$ and set $Y_{n+1} = (y, X_n^{(-i_n)})$
 - 5 Compute $\alpha_{i_n}(X_n, Y_{n+1}) = \min \left\{ \frac{f(Y_{n+1})}{f(X_n)} \frac{q_{i_n}(Y_{n+1}, X)}{q_{i_n}(X_n, y)}, 1 \right\}$
 - 6 Set $X_{n+1} = \begin{cases} Y_{n+1} & \text{with prob. } \alpha_{i_n}(X_n, Y_{n+1}) \\ X_n & \text{otherwise} \end{cases}$
- 7 **end**

One variable at a time

Algorithm: One variable at a time MH with **systematic sweep**.

```
1 Generate  $X_0 \sim \lambda$ 
2 for  $n = 0, 1, \dots$  do
3   Set  $Y_{n+1,0} = X_n$ 
4   for  $i = 1, \dots, d$  do
5     Draw  $y \sim q_i(X_n, \cdot)$  and set  $\tilde{Y} = (y, Y_{n+1,i-1}^{(-i)})$ 
6     Set  $Y_{n+1,i} = \begin{cases} \tilde{Y}, & \text{with prob. } \alpha_i(Y_{n+1,i-1}, \tilde{Y}) \\ Y_{n+1,i-1}, & \text{otherwise} \end{cases}$ 
7   end
8    $X_{n+1} = Y_{n+1,d}$ 
9 end
```

Detailed balance for one variable at a time

Let p_i be the transition density of the algorithm when the index i is selected:

$$\begin{aligned} p_i(x, y) &= p_i((x^{(i)}, x^{(-i)}), (y^{(i)}, y^{(-i)})) \\ &= \left(\alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) + (1 - \alpha_i^*(x)) \delta_{x^{(i)}}(y^{(i)}) \right) \delta_{x^{(-i)}}(y^{(-i)}) \\ &= \underbrace{\alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)})}_{p_i^1(x, y)} + \underbrace{(1 - \alpha_i^*(x)) \delta_{x^{(i)}}(y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)})}_{p_i^2(x, y)} \end{aligned}$$

with $\alpha_i^*(x) = \int_{\mathcal{X}} \alpha_i(x, y^{(i)}) q_i(x, y^{(i)}) dy^{(i)}$

Denote $P_i(x, A) = \int_A p_i(x, y) dy$ the corresponding Markov transition kernel and \mathcal{P}_i the associated Markov transition operator.

Lemma

Let $P_i : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ is in detailed balance with f .

Proof of the Lemma

$$\int_A P_i(x, B) f(x) dx = \underbrace{\int_A \int_B p_i^1(x, y) f(x) dy dx}_{T_1} + \underbrace{\int_A \int_B p_i^2(x, y) dy dx}_{T_2}$$

$$\begin{aligned} T_1 &= \int_A \int_B \min \left\{ 1, \frac{f(y^{(i)}, x^{(-i)}) q_i((y^{(i)}, x^{(-i)}), x^{(i)})}{f(x) q_i(x, y^{(i)})} \right\} q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(x) dx dy \\ &= \int_A \int_B \min \left\{ 1, \frac{f(y) q_i(y, x^{(i)})}{f(x) q_i(x, y^{(i)})} \right\} q_i(x, y^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(x) dy dx \\ &= \int_A \int_B \min \left\{ \frac{f(x) q_i(x, y^{(i)})}{f(y) q_i(y, x^{(i)})}, 1 \right\} q_i(y, x^{(i)}) \delta_{x^{(-i)}}(y^{(-i)}) f(y) dy dx \\ &= \int_A \int_B \min \left\{ \frac{f(x^{(i)}, y^{(-i)}) q_i((x^{(i)}, y^{(-i)}), y^{(i)})}{f(y) q_i(y, x^{(i)})}, 1 \right\} q_i(y, x^{(i)}) \delta_{y^{(-i)}}(x^{(-i)}) f(y) dy dx \\ &= \int_A \int_B \alpha_i(y, x^{(i)}) q_i(y, x^{(i)}) \delta_{y^{(-i)}}(x^{(-i)}) f(y) dx dy = \int_B \int_A p_i^1(y, x) f(y) dx dy \end{aligned}$$

$$T_2 = \int_A \int_B (1 - \alpha_i^*(x)) \delta_x(y) f(x) dy dx = \int_B \int_A p_i^2(y, x) f(y) dx dy$$

Detailed balance for one variable at a time

Random selection algorithm. Markov Transition kernel:

$$\begin{aligned} P^{\text{rand}}(x, A) &= \mathbb{P}(X_{n+1} \in A \mid X_n = x) \\ &= \sum_{i=1}^d \mathbb{P}(X_{n+1} \in A \mid X_n = x, i_n = i) \mathbb{P}(i_n = i) = \frac{1}{d} \sum_{i=1}^d P_i(x, A) \end{aligned}$$

Associated operator: $\mathcal{P}^{\text{rand}} = \frac{1}{d} \sum_{i=1}^d \mathcal{P}_i$

(P_i, f) in detailed balance $\implies (P^{\text{rand}}, f)$ in detailed balance (hence $\mathcal{P}^{\text{rand}}$ reversible and with invariant distribution f)

Detailed balance for one variable at a time

Systematic sweep algorithm. Markov Transition kernel:

$$\begin{aligned} P^{\text{sweep}}(x, A) &= \mathbb{P}(X_{n+1} \in A \mid X_n = x) \\ &= \int_A \int_{\mathcal{X}^{d-1}} P_d(y_{d-1}, dy_d) \cdots P_2(y_1, dy_2) P_1(x, dy_1) \end{aligned}$$

Associated operator: $\mathcal{P}^{\text{sweep}} = \mathcal{P}_1 \cdots \mathcal{P}_d$

f invariant for $P_i \implies f$ invariant for $\mathcal{P}^{\text{sweep}}$. However $\mathcal{P}^{\text{sweep}}$ is not reversible!

Markov transition operator of the reversed chain $\widehat{\mathcal{P}}^{\text{sweep}} = \mathcal{P}_d \cdots \mathcal{P}_1$

Possible fix to recover reversibility (if important)

$$\mathcal{P}^{\text{sweep}} = \mathcal{P}_1 \cdots \mathcal{P}_d \cdots \mathcal{P}_1$$

forward sweep $i = 1, \dots, d$ followed by backward sweep $i = d - 1, \dots, 1$.

Gibbs sampler

Consider a **one variable at a time** sampler using the **conditional distributions** as proposal densities $q_i(x, \cdot) = f_{X^{(i)} | X^{(-i)}}(\cdot | x^{(-i)})$

Given $x = (x^{(i)}, x^{(-i)})$ and $y = (y^{(i)}, x^{(-i)})$, the acceptance rate is

$$\begin{aligned}\alpha_i(x, y) &= \min \left\{ \frac{f(y) f_{X^{(i)} | X^{(-i)}}(x^{(i)} | x^{(-i)})}{f(x) f_{X^{(i)} | X^{(-i)}}(y^{(i)} | x^{(-i)})}, 1 \right\} \\ &= \min \left\{ \frac{f(y) f(x) / f_{X^{(-i)}}(x^{(-i)})}{f(x) f(y) / f_{X^{(-i)}}(x^{(-i)})}, 1 \right\} = 1\end{aligned}$$

Hence, in Gibbs sampler all the moves are accepted, provided one is able to generate exactly from the conditional distributions $f_{X^{(i)} | X^{(-i)}}(\cdot | x^{(-i)})$.

Algorithm: Gibbs with random sweep.

- 1 Generate $X_0 \sim \lambda$
- 2 **for** $n = 0, 1, \dots$ **do**
- 3 Draw i_n from a pmf β on $\{1, \dots, d\}$
- 4 Generate $y^{(i_n)} \sim f_{X^{(i_n)} | X^{(-i_n)}}(\cdot | X_n^{(-i_n)})$
- 5 Set $X_{n+1} = (y^{(i_n)}, X_n^{(-i_n)})$
- 6 **end**

Metropolis Adjusted Langevin Algorithm (MALA)

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}_+$ be our target probability density and consider the following Stochastic Differential Equation (Langevin dynamics)

$$dX_t = \nabla \log f(X_t) + \sqrt{2} dW_t, \quad t > 0, \quad X_0 \sim \lambda \quad (3)$$

with W_t a standard Wiener process and λ a probability measure on \mathbb{R}^d .

Let us denote by $\rho(x, t)$ the probability density function of X_t (provided it exists):

$$\int_A \rho(x, t) dx = \mathbb{P}_\lambda(X_t \in A)$$

Under quite general conditions on f , one has $\lim_{t \rightarrow \infty} \rho(x, t) = f(x)$, i.e. the distribution of X_t converges to f and f is an invariant probability density function for (3) (time continuous Markov chain)

Problem: usually, exact solutions of (3) are not available

Metropolis Adjusted Langevin Algorithm (MALA)

Remedy: use numerical discretization, e.g. Euler-Maruyama method

$$X_{n+1} = X_n + \Delta t \nabla \log f(X_n) + \sqrt{2\Delta t} \xi_n, \quad \xi_n \sim N(0, I) \quad (4)$$

However, the discrete time Markov chain $\{X_n\}_n$ will not have anymore f as invariant distribution due to the numerical discretization error

Idea: use (4) as a proposal distribution within a Metropolis-Hastings Algorithm

Algorithm: Metropolis Adjusted Langevin Algorithm (MALA).

- 1 Generate $X_0 \sim \lambda$
 - 2 **for** $n = 0, 1, \dots$ **do**
 - 3 Generate $Y \sim N(X_n + \Delta t \nabla \log f(X_n), 2\Delta t I)$
 - 4 Compute $\alpha(X_n, Y) = \min \left\{ 1, \frac{f(Y)}{f(X_n)} \frac{\exp(-\|X_n - Y - \Delta t \nabla \log f(Y)\|^2 / 2\Delta t)}{\exp(-\|Y - X_n - \Delta t \nabla \log f(X_n)\|^2 / 2\Delta t)} \right\}$
 - 5 Set $X_{n+1} = \begin{cases} Y & \text{with prob. } \alpha(X_n, Y) \\ X_n & \text{otherwise} \end{cases}$
 - 6 **end**
-

- ▶ Similar to a RWM; but proposal is not symmetric and uses gradients information

Ergodic estimator

- ▶ Let $\{X_n\}_n \sim \text{Markov}(\lambda, P)$ on $\mathcal{X} \subset \mathbb{R}^d$, with unique invariant distribution π .
- ▶ We assume moreover that $\{X_n\}_n$ is **geometrically ergodic**, i.e. there exist $\gamma > 0$ and $h : \mathcal{X} \rightarrow \mathbb{R}_+$ s.t.

$$\|\pi^{n,\lambda} - \pi\|_{\text{TV}} \leq \lambda(h)e^{-\gamma n}, \quad \lambda(h) = \int_{\mathcal{X}} h(x)d\lambda(x)$$

Recall that for $\mu, \nu \in \mathcal{M}_1(\mathcal{X})$

$$\|\mu - \nu\|_{\text{TV}} = 2 \sup_{A \in \mathcal{B}(\mathcal{X})} |\mu(A) - \nu(A)| = \sup_{\phi \in L^\infty(\mathcal{X})} \frac{|\int_{\mathcal{X}} \phi(x)d\mu(x) - \int_{\mathcal{X}} \phi(x)d\nu(x)|}{\|\phi\|_{L^\infty(\mathcal{X})}}$$

- ▶ Given a π -integrable function $\psi : \mathcal{X} \rightarrow \mathbb{R}$, we estimate $\mu = \mathbb{E}_\pi[\psi]$ by the ergodic estimator

$$\hat{\mu}_{N,b}^{\text{MCMC}} = \frac{1}{N} \sum_{i=1}^N \psi(X_{i+b})$$

- ▶ **Question:** how to monitor the approximation error $|\hat{\mu}_{N,b}^{\text{MCMC}} - \mu|$

Bias

If the chain is at stationarity ($\lambda = \pi$), then $\hat{\mu}_{N,b}^{MCMC}$ is unbiased. Indeed, $X_n \sim f, \forall n$ and

$$\mathbb{E}_{\pi}[\hat{\mu}_{N,b}^{MCMC}] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\pi}[\psi(X_{i+b})] = \mu$$

If, instead, the chain is not at stationarity ($\lambda \neq \pi$), the estimator $\hat{\mu}_{N,b}^{MCMC}$ is biased ! However

$$\begin{aligned} |\mathbb{E}_{\lambda}[\hat{\mu}_{N,b}^{MCMC} - \mu]| &= \left| \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\lambda}[\psi(X_{i+b}) - \mu] \right| \\ &\leq \frac{1}{N} \sum_{i=1}^N \left| \int_{\mathcal{X}} \psi(y) d\pi^{i+b,\lambda}(y) - \int_{\mathcal{X}} \psi(y) d\pi(y) \right| \\ &\leq \frac{1}{N} \sum_{i=1}^N \|\psi\|_{L^{\infty}(\mathcal{X})} \|\pi^{i+b,\lambda} - \pi\|_{\text{TV}} \\ &\leq \frac{1}{N} \|\psi\|_{L^{\infty}(\mathcal{X})} \lambda(h) \sum_{i=b+1}^N e^{-\gamma i} \leq \frac{e^{-\gamma b}}{N} \frac{\|\psi\|_{L^{\infty}(\mathcal{X})} \lambda(h)}{1 - e^{-\gamma}} \end{aligned}$$

Bias

- ▶ The Bias decays as $O(\frac{1}{N})$, faster than the standard deviation (which is $O(\frac{1}{\sqrt{N}})$)
- ▶ Moreover, it decays as $O(e^{-\gamma b})$ and can be dramatically reduced by increasing the *burn-in* b .
- ▶ Reasonable values of b can be guessed from a trace-plot of the chain $\{\psi(X_n)\}_n$ (smallest time after which the chain looks at stationarity)

Asymptotic variance

Assume that a sufficient burn-in period has been removed and the chain is essentially at stationarity. Then, the following result on the *asymptotic variance* holds

Lemma

Let $\{X_n\} \sim \text{Markov}(\pi, P)$ with π invariant for P , and denote

$$c(k) = \text{Cov}_\pi(\psi(X_0), \psi(X_k)) = \text{Cov}_\pi(\psi(X_j), \psi(X_{j+k})).$$

Then

$$\mathbb{V}\text{ar}_\pi[\hat{\mu}_{N,b}^{MCMC}] = \frac{\sigma_{MCMC,N}^2}{N}, \quad \text{with } \sigma_{MCMC,N}^2 = c(0) + 2 \sum_{\ell=1}^{N-1} \left(1 - \frac{\ell}{N}\right) c(\ell).$$

Moreover, if $\sum_{k=0}^{\infty} |c(k)| < +\infty$, then

$$\lim_{N \rightarrow \infty} N \mathbb{V}\text{ar}[\hat{\mu}_{N,b}^{MCMC}] = \sigma_{MCMC}^2$$

with $\sigma_{MCMC}^2 = c(0) + 2 \sum_{k=1}^{\infty} c(k)$.

Proof

$$\begin{aligned}\text{Var}_\pi[\hat{\mu}_{N,b}^{MCMC}] &= \mathbb{E}_\pi \left[\left(\frac{1}{N} \sum_{j=1}^N \psi(X_{j+b}) - \mu \right)^2 \right] \\ &= \frac{1}{N^2} \sum_{j=1}^N \sum_{k=1}^N \mathbb{E}_\pi [(\psi(X_{j+b}) - \mu)(\psi(X_{k+b}) - \mu)] \\ &= \frac{1}{N^2} \left[\sum_{j=1}^N \underbrace{\text{Var}_\pi[\psi(X_{j+b})]}_{c(0)} + 2 \sum_{j=1}^{N-1} \sum_{k=j+1}^N \underbrace{\text{Cov}_\pi(\psi(X_{j+b}), \psi(X_{k+b}))}_{c(k-j)} \right] \\ &= \frac{c(0)}{N} + \frac{2}{N^2} \sum_{j=1}^{N-1} \sum_{\ell=1}^{N-j} c(\ell) \\ &= \frac{c(0)}{N} + \frac{2}{N} \sum_{\ell=1}^{N-1} \frac{N-\ell}{N} c(\ell) \\ &= \frac{1}{N} \left(c(0) + 2 \sum_{\ell=1}^{N-1} \left(1 - \frac{\ell}{N}\right) c(\ell) \right).\end{aligned}$$

If now $\sum_{\ell=0}^{\infty} |c(\ell)| < +\infty$, it follows that $\lim_{N \rightarrow \infty} N \text{Var}_\pi[\hat{\mu}_{N,b}^{MCMC}] = \sigma_{MCMC}^2$.

Asymptotic variance

- ▶ The quantity

$$\sigma_{MCMC}^2 = c(0) + 2 \sum_{k=1}^{\infty} c(k)$$

is called **time-average variance constant** (TAVC) or **asymptotic variance**

- ▶ If $\{X_n\}_n$ were iid and distributed as f (pure Monte Carlo sampling) then the variance of the Monte Carlo estimator would be $\text{Var} [\hat{\mu}_N^{MC}] = \frac{c(0)}{N}$.
- ▶ Given N , we call **effective sample size** (ESS) the sample size that a Monte Carlo estimator would use to achieve the same variance as the MCMC one:

$$\text{Var} [\hat{\mu}_{N,b}^{MCMC}] \rightarrow \frac{\sigma_{MCMC}^2}{N} = \frac{c(0)}{ESS} \quad \implies \quad ESS = N \frac{c(0)}{\sigma_{MCMC}^2}$$

- ▶ For reversible, geometrically ergodic, Markov chains, a CLT holds

$$\sqrt{N}(\hat{\mu}_{N,b}^{MCMC} - \mu) \xrightarrow{d} N(0, \sigma_{MCMC}^2)$$

Estimating the asymptotic variance – covariance method

Given a path $\{X_n\}_n$ and a burn-in time, we can estimate the covariances

$$\hat{c}(k) = \frac{1}{N-k-1} \sum_{j=1}^{N-k} (\psi(X_{j+b}) - \hat{\mu}_{N,b}^{MCMC})(\psi(X_{j+b+k}) - \hat{\mu}_{N,b}^{MCMC})$$

and

$$\hat{\sigma}_{MCMC}^2 = \hat{c}(0) + 2 \sum_{k=1}^{N-2} \hat{c}(k).$$

However, the last terms in the sum are very unstable. Better estimator

$$\hat{\sigma}_M^2 = \hat{c}(0) + 2 \sum_{k=1}^M \hat{c}(k), \quad \text{with } M = 2 \min\{k : \hat{c}(2k) + \hat{c}(2k+1) < 0\}.$$

(valid for reversible Markov Chains)

Estimating the asymptotic variance - batch means

An alternative idea to estimate σ_{MCMC}^2 is to split the sequence $\{X_n\}_{n=b+1}^{N+b}$ into M blocks of size $T = N/M$

Then we can build M different sample averages

$$\hat{\mu}^{(i)} = \frac{1}{T} \sum_{j=(i-1)T+b+1}^{iT+b} \psi(X_j), \quad \text{and} \quad \hat{\mu}_{N,b}^{MCMC} = \frac{1}{M} \sum_{i=1}^M \hat{\mu}^{(i)}.$$

If T is sufficiently large (larger than the relaxation time), the M blocks are nearly independent so that

$$\text{Var} [\hat{\mu}_{N,b}^{MCMC}] \approx \frac{\sigma_{MCMC}^2}{N} \approx \frac{\text{Var} [\hat{\mu}^{(1)}]}{M}$$

and $\text{Var} [\hat{\mu}^{(1)}]$ can be estimated by a sample variance estimator

$$\text{Var} [\hat{\mu}^{(1)}] \approx \hat{\sigma}_{\hat{\mu}^{(1)}}^2 = \frac{1}{M-1} \sum_{i=1}^M \left(\hat{\mu}^{(i)} - \hat{\mu}_{N,b}^{MCMC} \right)^2.$$

Finally, an estimator for σ_{MCMC}^2 is

$$\hat{\sigma}_{MCMC}^2 = \frac{N}{M} \hat{\sigma}_{\hat{\mu}^{(1)}}^2 = \frac{T}{M-1} \sum_{i=1}^M \left(\hat{\mu}^{(i)} - \hat{\mu}_{N,b}^{MCMC} \right)^2.$$