

Bayesian Computations – MCMC



Felipe Uribe

Computational Engineering
School of Engineering Sciences
Lappeenranta-Lahti University of Technology (LUT)

Special Course on Inverse Problems
Lappeenranta, FI — January-February, 2024

Overview

MCMC changed our the emphasis from “closed form” solutions to algorithms, expanded our impact to solving “real” applied problems, expanded our impact to improving numerical algorithms using statistical ideas, and led us into a world where “exact” now means “simulated”!!

1. Introduction

Bayesian inverse problems

- All unknowns are represented as random variables with prior densities defined e.g., with respect to the Lebesgue measure (or counting measure if discrete).
- Prior information given probabilistically, i.e., probability density functions; used as regularization.
- We define a Bayesian inverse problem, as the task of characterizing the density

$$\pi(\mathbf{x}) := \pi_{\text{pos}}(\mathbf{x} \mid \mathbf{y}) = \frac{1}{Z} \pi_{\text{like}}(\mathbf{y} \mid \mathbf{x}) \pi_{\text{pr}}(\mathbf{x}). \quad (1)$$

- ▶ $\pi_{\text{pr}}(\mathbf{x})$ is the **prior probability density**.
- ▶ $\pi_{\text{like}}(\mathbf{y} \mid \mathbf{x})$ is the **likelihood function**.
- ▶ $Z = \int_{\mathbb{R}^d} \pi_{\text{like}}(\mathbf{y} \mid \mathbf{x}) \pi_{\text{pr}}(\mathbf{x}) d\mathbf{x}$ is the normalizing constant of the posterior density, called the **model evidence**.

Summarizing posterior inferences

- The *posterior mean* (PM) (or conditional mean) ¹:

$$\mathbf{x}_{\text{PM}} = \mathbb{E}_{\pi_{\text{pos}}}[\mathbf{x}] = \int_{\mathbb{R}^d} \mathbf{x} \pi_{\text{pos}}(\mathbf{x} \mid \mathbf{y}) \, d\mathbf{x}. \quad (2)$$

- Posterior credible sets, posterior quantiles (abusing the notation):

$$\mathbb{P}[\mathbf{x} > \mathbf{a}] = \int_{\mathbf{a}}^{\infty} \pi_{\text{pos}}(\mathbf{x} \mid \mathbf{y}) \, d\mathbf{x}. \quad (3)$$

- Posterior moments and statistics: standard deviation, median, MAD, etc.
- Posterior realizations for direct assessment. Posterior marginals (e.g., via a `plotmatrix`).

¹

Optimality: it minimizes the expected MSE (Bayes risk) for any prior and likelihood, as long as the second order moments exist.

General comments (i)

- Until now we have simply assumed that we can draw random variables, vectors and fields from any desired distribution.
- For some problems, we cannot do this either at all, or in a reasonable amount of time. It is often feasible however to draw dependent samples whose distribution is close to and indeed approaches the desired one.
- In Markov chain Monte Carlo (MCMC) we do this by sampling x_1, x_2, \dots, x_n from a Markov chain constructed so that the distribution of x_i approaches the target distribution.
- The primary method is the Metropolis algorithm, which was named one of the ten most important algorithms of the twentieth century.

General comments (ii)

- We have multiple MCMC algorithms to draw samples from a target posterior distribution.
- As usual, we estimate an expectation $\mu = \mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})\pi(\mathbf{x}) \mathrm{d}\mathbf{x}$:

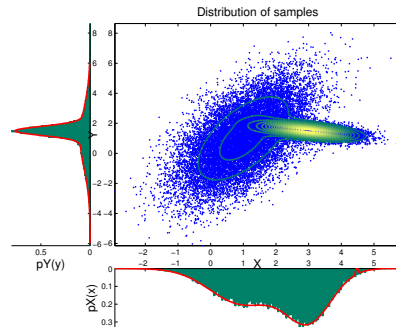
$$\mu \approx \hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i), \quad (4)$$

in the MCMC context, we have two main challenges:

- (i) The draws \mathbf{x}_i have a distribution that approaches π , which is usually not equal to π , hence, the estimate is biased.
- (ii) The \mathbf{x}_i are (in general) statistically dependent, and thus $\hat{\mu}$ is harder to estimate. In extreme cases, the \mathbf{x}_i can get stuck in some subset of their domain and then $\hat{\mu}$ will fail to converge.

MCMC: basic timeline

- **Nicholas Metropolis** (1953) - Metropolis algorithm for generating samples from the Boltzmann distribution.
- **Wilfred K. Hastings** (1970) - Metropolis-Hastings algorithm, the most common MCMC.
- **Julian Besag** - Gibbs sampler (1974), MALA algorithm (1994).
- **Simon Duane et al.** (1987) - Hamiltonian MC.
- **Peter Green** (1995) - Reversible jump MCMC.
- **Joris Bierkens** (2015) - non-reversible MCMC.



This lecture...

- The lecture is based on multiple references. However, we mostly follow Chapters 11 and 12 of the book by **Art Owen**², which is freely available online.

² A. B. Owen. *Monte Carlo theory, methods and examples*. artowen.su.domains/mc/, 2018.

2. Markov chains

Markov chains

- Our solution to hard sampling problems will be to run a Markov chain for a long time so that the values of the chain have a distribution which approaches our target density π which is defined on \mathcal{X} .
- In plain terms, a Markov chain is a “memoryless” stochastic process: to know a future state, we just need to know the current state. This is called the *Markov property*.
- Formally, a discrete time **Markov chain** $X = \{X_n\}$ with a discrete state space \mathcal{X} satisfies the Markov property:

$$\mathbb{P}[X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0] = \mathbb{P}[X_{n+1} = j \mid X_n = i]; \quad (5)$$

“past and future are conditionally independent given the present”.

Markov chains

- In practice, we work with **time-homogeneous** Markov chains for which the *transition probabilities* $\mathbb{P}[X_{n+1} = j \mid X_n = i] = q_{ij}$ are independent of n . In this case, we can arrange the transition probabilities as a (transition) matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$; $\{q_{ij}\}$ (rows sum to one).
- Suppose that at step n , X_n has distribution $\mathbf{s} \in \mathbb{R}^{1 \times m}$ (row probability vector); think of this as a PMF. Here, the element $s_i = \mathbb{P}[X_n = i]$ denotes the probability that the chain is in state i at step n .
- A natural question is to what is the probability that the chain has state value j at step $n+1$? Using the total probability theorem,

$$\mathbb{P}[X_{n+1} = j] = \sum_i \mathbb{P}[X_{n+1} = j \mid X_n = i] \mathbb{P}[X_n = i] = \sum_i s_i q_{ij} \quad j\text{th entry of } \mathbf{sQ}; \quad (6)$$

successive iteration of this equation describes the evolution of the chain, i.e., \mathbf{sQ} is the distribution at $n+1$: “one step in the future, right multiply by \mathbf{Q} ”.

Markov chains

Some properties:

- We say that a Markov chain is **irreducible** (connected), if it is possible to get from anywhere to anywhere. Otherwise, it is called *reducible*.
- We say that a state is **recurrent**, if starting there, the chain has probability 1 of returning to that state. Otherwise, it is called *transient*.
- We say that a Markov chain is **periodic**, if the states start repeating themselves with a given period. Otherwise, it is called *aperiodic*.
- We say that s is **stationary** for a Markov chain with transition matrix Q , if $sQ = s$.

Markov chains

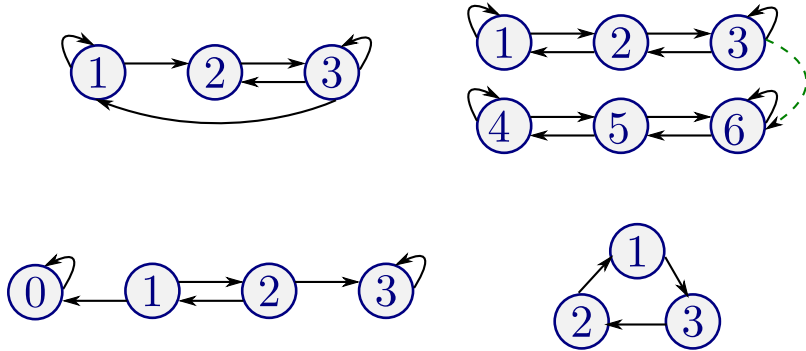


Figure: Four different Markov chains. How can they be classified?

Markov chains

- For an irreducible Markov chain (with finitely many states), we have the following³:
 - (i) The stationary distribution s exists.
 - (ii) The stationary distribution is unique.
 - (iii) If the Markov chain is further aperiodic (Q^m is strictly positive for some m), then $\mathbb{P}[X_n = i] \rightarrow s_i$ as $n \rightarrow \infty$.
- We normally work with Markov chains that are “easy to deal with”. This class of chains is called *time reversible*.
- A Markov chain with transition matrix Q is **reversible**, if there is a s such that

$$s_i q_{ij} = s_j q_{ji} \quad \forall i, j \quad (\text{detailed balance condition}). \quad (7)$$

- If the Markov chain is reversible with respect to s , then s is the stationary distribution.

³

see Theorem 1.8.3 in [12]

Markov chains (example)

- Suppose the state space are (Rain, Sunny, Cloudy) and weather follows a Markov process. Hence, the probability of tomorrow's weather simply depends on today's weather, and not any other previous days (example from [17]).
- The probability transitions are given by

$$\mathbf{Q} = \begin{bmatrix} 0.5 & 0.25 & 0.25 \\ 0.5 & 0 & 0.5 \\ 0.25 & 0.25 & 0.5 \end{bmatrix}. \quad (8)$$

- Suppose today is sunny (and then rainy). What is the expected weather two days from now? Seven days?

$$\text{If } s_0 = [0, 1, 0] \quad s_2 = s_0 \mathbf{Q}^2 = [0.375, 0.25, 0.375] \quad s_7 = s_0 \mathbf{Q}^7 = [0.4, 0.2, 0.4]$$

$$\text{If } s_0 = [1, 0, 0] \quad s_7 = s_0 \mathbf{Q}^7 = [0.4, 0.2, 0.4].$$

Note that after a sufficient amount of time, the expected weather is independent of the starting value, i.e., the chain has reached a stationary distribution.

Markov chains

- The basic idea of discrete-state Markov chains can be generalized to a continuous state Markov process by having a probability kernel $P(\mathbf{x}, \mathbf{x}')$ that satisfies

$$\int_{\mathcal{X}} P(\mathbf{x}, \mathbf{x}') d\mathbf{x}' = 1, \quad \mathbf{x}, \mathbf{x}' \in \mathcal{X}. \quad (9)$$

- Now, the reversibility (detailed balance) condition is:

$$\int_A \pi(\mathbf{x}) P(\mathbf{x}, \mathbf{x}') d\mathbf{x} = \int_B \pi(\mathbf{x}') P(\mathbf{x}', \mathbf{x}) d\mathbf{x}', \quad (10)$$

where $\mathbf{x} \in A$, $\mathbf{x}' \in B$, and $A, B \subset \mathcal{X}$, $P(\cdot, \cdot)$ denotes a Markov transition kernel.

- Different approaches exist to generate kernels P that ensure eq. (10), and hence setting π as the stationary distribution.

3. MCMC – Metropolis–Hastings

Metropolis–Hastings

- If we make a good selection for the transition kernel, it could asymptotically converge to a target distribution independently of where we started from.
- More importantly, we can use the realizations of the Markov chain in Monte Carlo estimators i.e., we can average across the path. However, even if X_n were exact draws, they are not independent anymore.
- In Metropolis–Hastings (MH), the condition eq. (10) is satisfied by separating the transition into two stages: the proposal and the acceptance/rejection steps.
- These are represented respectively by the **proposal distribution** $q(x | y)$ which accounts for the conditional probability of x given the proposed state y , and the **acceptance probability** $\alpha(x, y)$ which expresses the conditional probability of accepting the proposed state y .

Metropolis–Hastings

- The transition kernel can be written as:

$$P(\mathbf{x} \mid \mathbf{y}) = \underbrace{\int_B q(\mathbf{x} \mid \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y}) \, d\mathbf{y}}_{\text{acceptance}} + \underbrace{\mathbb{1}(\mathbf{x} \in B) \int_{\mathcal{X}} q(\mathbf{x} \mid \mathbf{y}) (1 - \alpha(\mathbf{x}, \mathbf{y})) \, d\mathbf{y}}_{\text{rejection}}. \quad (11)$$

- Substituting eq. (11) into eq. (10), yields

$$\int_A \pi(\mathbf{x}) \int_B q(\mathbf{x} \mid \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} \, d\mathbf{x} = \int_B \pi(\mathbf{y}) \int_A q(\mathbf{y} \mid \mathbf{x}) \alpha(\mathbf{y}, \mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}, \quad (12)$$

- Which can be written in compact form as

$$\int_{A \times B} \pi(\mathbf{x}) q(\mathbf{x} \mid \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} \, d\mathbf{x} = \int_{A \times B} \pi(\mathbf{y}) q(\mathbf{y} \mid \mathbf{x}) \alpha(\mathbf{y}, \mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}. \quad (13)$$

Metropolis–Hastings

- The equality eq. (13) holds if

$$\pi(\mathbf{x})q(\mathbf{x} | \mathbf{y})\alpha(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})\alpha(\mathbf{y}, \mathbf{x}) \quad (14a)$$

$$\frac{q(\mathbf{x} | \mathbf{y})\alpha(\mathbf{x}, \mathbf{y})}{q(\mathbf{y} | \mathbf{x})\alpha(\mathbf{y}, \mathbf{x})} = \frac{\pi(\mathbf{y})}{\pi(\mathbf{x})} \quad (14b)$$

$$\frac{\alpha(\mathbf{x}, \mathbf{y})}{\alpha(\mathbf{y}, \mathbf{x})} = \frac{\pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{x})q(\mathbf{x} | \mathbf{y})}. \quad (14c)$$

- An acceptance mechanism that fulfills the condition in eq. (14c) is proposed by [9] (based on [10], which uses a symmetric proposal distribution)

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left(1, \frac{\pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{x})q(\mathbf{x} | \mathbf{y})} \right); \quad (15)$$

an **unnormalized** π can be used as the normalizing constants cancel out in the ratio.

Metropolis–Hastings

Algorithm 1: Metropolis–Hastings

```
1 Initialize  $\mathbf{x}_0$ ;  
2 for  $i = 1$  to  $n$  do  
3   Sample  $\mathbf{x}^* \sim q(\cdot \mid \mathbf{x}_{i-1})$ ;  
4   Compute the acceptance probability  
                                     
$$\alpha = \min \left( 1, \frac{\pi(\mathbf{x}^*)q(\mathbf{x}_{i-1} \mid \mathbf{x}^*)}{\pi(\mathbf{x}_{i-1})q(\mathbf{x}^* \mid \mathbf{x}_{i-1})} \right)$$
  
   Sample  $u \sim \mathcal{U}(0, 1)$ ;  
5   if  $u \leq \alpha$  then  
6     | Accept  $\mathbf{x}_i = \mathbf{x}^*$   
7   else  
8     | Reject  $\mathbf{x}_i = \mathbf{x}_{i-1}$   
9   end  
10 end
```

Metropolis–Hastings (burn-in and thinning)

- One of the difficulties with MH is the potentially strong dependence in the Markov chain states (samples). A common practice is ignore the first $n_b < n$ generated points and estimate μ by

$$\hat{\mu} = \frac{1}{n - n_b} \sum_{i=n_b+1}^n f(\mathbf{x}_i). \quad (16)$$

- This practice is called **burn-in** or warm-up. The distribution of \mathbf{X}_i usually only approaches π as i increases and so the first few observations might be very unrepresentative. Including them could bias the answer.
- We could think of the burn-in period as one way of finding a good starting point for the simulation. Some authors recommend $n_b = 0.5n$.
- We can also “thin” the chain by selecting every other n_t th sample (discarding those in between). This process is also called *subsampling*.

Metropolis–Hastings (general comments)

- The MH acceptance probability is a default choice for $\alpha(\mathbf{x}, \mathbf{y})$ that provides detailed balance. With this default in hand we can then search for good proposal distributions $q(\mathbf{x}, \mathbf{y})$ to suit any given problem.
- It is recommended to formulate the algorithm using log densities, for stability.
- The performance of MH deteriorates with increasing dimension of the parameter space. In this case, the acceptance probability at each step becomes very small as the dimension increases, resulting in slow convergence rates and a large number of repeated samples.

Barker MCMC *

- Recall the Metropolis–Hastings acceptance probability that holds the reversibility condition:

$$\alpha_{\text{MH}}(\mathbf{x}, \mathbf{y}) = \min \left(1, \frac{\pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{x})q(\mathbf{x} | \mathbf{y})} \right). \quad (17)$$

- Several other choices are possible. One alternative proposed by Barker⁴ is the acceptance probability

$$\alpha_{\text{B}}(\mathbf{x}, \mathbf{y}) = \frac{\pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{x})q(\mathbf{x} | \mathbf{y}) + \pi(\mathbf{y})q(\mathbf{y} | \mathbf{x})}. \quad (18)$$

- In general α_{MH} is preferred over α_{B} because, for the same choice of proposal $q(\mathbf{y} | \mathbf{x})$, α_{MH} will result in Markov chains that produce ergodic averages with smallest asymptotic variance⁵. In particular, α_{MH} will maximize the probability of moving from \mathbf{x} to \mathbf{y} .

⁴ A. A. Barker. “Monte Carlo calculations of the radial distribution functions for a proton-electron plasma”. In: *Australian Journal of Physics* 18 (1965), pp. 119–1347.

⁵ L. Tierney. “A note on Metropolis-Hastings kernels for general state spaces”. In: *Annals of applied probability* (1998), pp. 1–9.

Independence sampler

- Perhaps the simplest proposal mechanism is to take iid proposals from some distribution q that does not depend on the present location \mathbf{x} . Then $q(\mathbf{x}^* | \mathbf{x}) = q(\mathbf{y})$.
- The MH proposal for this independence sampler, simplifies to

$$\alpha(\mathbf{x}, \mathbf{x}^*) = \min \left(1, \frac{\pi(\mathbf{x}^*)q(\mathbf{x})}{\pi(\mathbf{x})q(\mathbf{x}^*)} \right); \quad (19)$$

the independence sampler is also called the Metropolized independence sampler.

- The independence sampler is closely related to importance sampling. It is generally safer to have slightly heavy tails in an independence sampler proposal q instead of light tail.

Random walk Metropolis

- A **random walk** is a process where the increments $\mathbf{z}_i = \mathbf{x}_i - \mathbf{x}_{i-1}$ are iid.
- In **random walk Metropolis** (RWM) the proposals take the form $\mathbf{x}_i^* = \mathbf{x}_i + \mathbf{z}_i$, where \mathbf{z}_i are iid random vectors.
- Suppose that $\mathbf{z} \sim q$. Then $q(\mathbf{x}^* | \mathbf{x})$ can be written $q(\mathbf{x}^* - \mathbf{x})$, where we now use q to also represent the probability density of \mathbf{z} .
- We will focus on random walks where the distribution q is symmetric (e.g., a Gaussian) and

$$\alpha(\mathbf{x}, \mathbf{x}^*) = \min \left(1, \frac{\pi(\mathbf{x}^*)}{\pi(\mathbf{x})} \right), \quad (20)$$

because the proposal ratio cancels.

Random walk Metropolis

- It is common in MCMC that we have to tune our proposals. For this we have to specify the so-called **proposal scale** β . Under a Gaussian proposal, this is equal to a standard deviation.
- If the acceptance rate is very small, we can infer that β is too large and then try a smaller value. Conversely, a very high acceptance rate suggests that we should raise β . Under a famous result from [7], it is optimal to tune β so that about 23.4% of proposals are accepted.
- In [7], Gaussian targets and Gaussian proposal were considered, so they were able to study how the proposal variance should depend on the dimension d . They **recommend** $\beta = 2.4$ for $d = 1$, and $\beta = 2.38/\sqrt{d}$ for large d .
- The acceptance rate when using the optimal β is about 44% for $d = 1$ and decreases rapidly to a limiting value of about 23.4% as $d \rightarrow \infty$. In general, as long as the rejection rate is between 15% and 40% the efficiency is close to that of the optimal β .

Random walk Metropolis (example)

- Obtain samples from the probability density,

$$\pi(x) \propto 0.3 \exp(-0.2x^2) + 0.7 \exp(-0.2(x - 10)^2).$$

- We use a Gaussian proposal $\mathcal{N}(y; x, \beta^2)$ (centered at x), and the initial state is defined as $x_0 = 5$.
- To show the effect of the proposal scale, simulations for different values of the standard deviation, $\beta = 0.1$, $\beta = 1$, $\beta = 14$ and $\beta = 50$ are performed.
- We run $n = 5 \times 10^3$ samples without burn-in.

Random walk Metropolis (example)

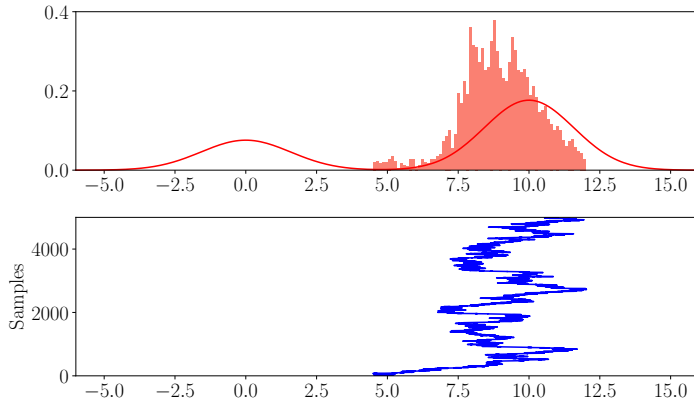


Figure: Bimodal target: scale parameter $\beta = 0.1$. Acceptance probability 0.98. The plot at the bottom is called a **trace plot**.

Random walk Metropolis (example)

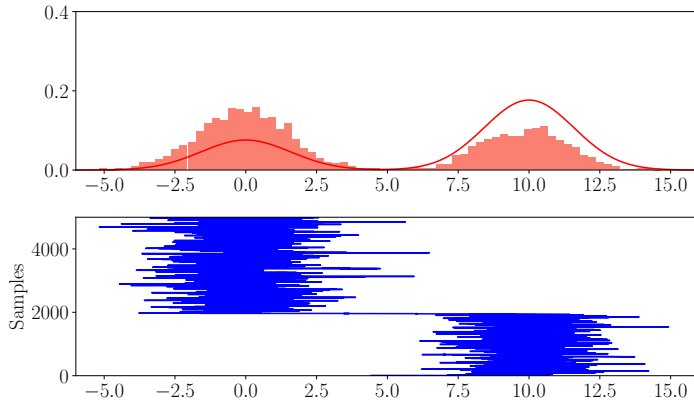


Figure: Bimodal target: scale parameter $\beta = 0.5$. Acceptance probability 0.80.

Random walk Metropolis (example)

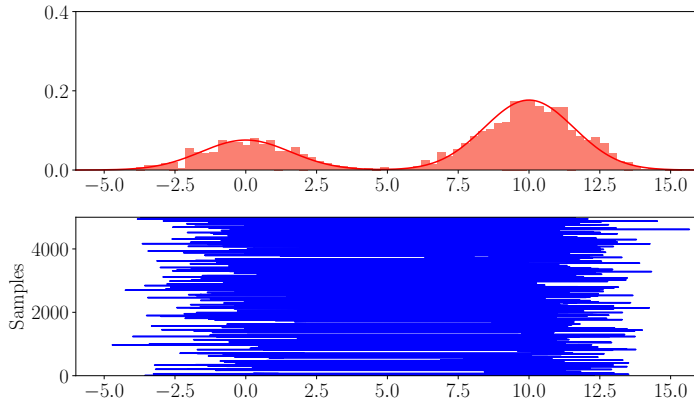


Figure: Bimodal target: scale parameter $\beta = 14$. Acceptance probability 0.23.

Random walk Metropolis (example)

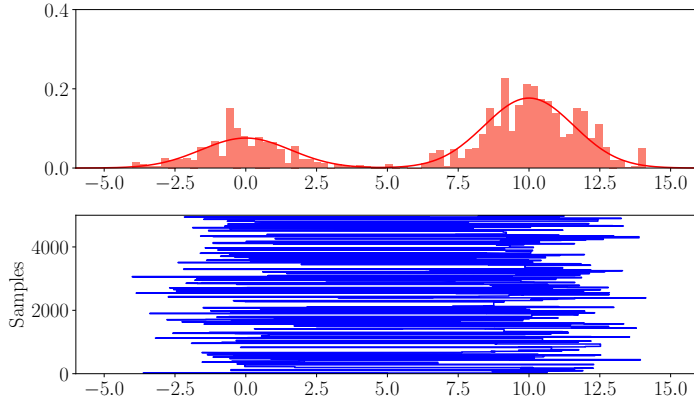


Figure: Bimodal target: scale parameter $\beta = 50$. Acceptance probability 0.08.

Adaptive RWM (I)

- We have to specify a suitable proposal scaling for our MCMC algorithm. This requires some parameter study. We can, however, employ different techniques that perform adaptation of the proposal scale such that the acceptance probability of the sampler is optimal.
- Learn a better proposal $q(\mathbf{x}^* | \mathbf{x})$ from past samples: (i) an appropriate proposal scale, and (ii) an appropriate proposal orientation and anisotropy; this is essential in problems with strong correlation in π .
- Adaptive Metropolis scheme of [8]
 - (i) Covariance matrix at step i

$$\Sigma_i^* = \beta^2 \widehat{\text{Cov}}(\mathbf{x}_1, \dots, \mathbf{x}_i) + \beta^2 \epsilon \mathbf{I}_d, \quad (21)$$

where $\epsilon > 0$ and $\beta^2 = 2.4^2/d$.

- (ii) Proposals are Gaussians at \mathbf{x}_i . Use a covariance Σ_0 for the first i_0 steps, then use Σ_i^* .
- (iii) Chain is not Markov. Nonetheless, one can prove that the chain converges to π .

Adaption with global adaptive scaling (II)

- **Vanishing adaptation** ensures that the chain depends less and less on recently visited states of the chain. This sets controlled MCMC algorithms that produce samples asymptotically distributed according to π .
- Adaptive Metropolis scheme of [1]. At step i update the proposal scale as

$$\log(\beta_{i+1}) = \log(\beta_i) + \gamma_{i+1}(\alpha(x_i, x^*) - \alpha^*),$$

where $\{\gamma_i\} \subset (0, \infty)$ is a sequence of stepsizes which ensures that the variations of the chain states $\{x_i\}$ vanish, α^* is the desired acceptance rate. Note that this is just a standard *Robbins–Monro* recursion.

- The standard approach consists of choosing the sequence $\{\gamma_i\}$ deterministic and non-increasing, but it is also possible to choose $\{\gamma_i\}$ random. A standard choice is to make $\gamma_i = C/i^c$, for $c \in (1/(1 + \beta), 1]$ and constant C .

4. MCMC – diagnostics

MCMC error and diagnostics

- Two of the hardest problems in MCMC are:
 - (i) deciding whether the distribution of x_i has nearly converged to π , and
 - (ii) deciding whether the x_i are mixing well.
- While traces are quite useful, they are only one way diagnostics. When they show us that the chain is not sampling π well, we can believe it. Unfortunately a trace can look perfectly good even when the sample is poor.
- The autocorrelation function has the same problem. If it shows slowly decaying autocorrelations, we know there is a problem, but if they decay rapidly we might still have missed part of the space.
- One approach to generating diagnostics is to run multiple independently generated Markov chains, starting them in different places.

MCMC error and diagnostics

- Recall from Monte Carlo that the mean estimator had a variance $\mathbb{V}[\hat{\mu}_n] = \sigma^2/n$, but the samples were iid. In MCMC, we have:

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i) \quad \text{and} \quad \mathbb{V}[\hat{\mu}_n] = \frac{\sigma^2}{n} \tau_f, \quad (22)$$

where τ_f is the integrated autocorrelation time (or *inefficiency factor*) for the chain $f(\mathbf{X}_i)$.

- The *integrated autocorrelation time* (IACT) for a (weakly) stationary process is defined as⁶

$$\tau_f = \sum_{k=-\infty}^{\infty} \rho_f(k) = 1 + 2 \sum_{k=1}^{\infty} \rho_k, \quad (23)$$

where ρ_k is the normalized autocorrelation function at lag k .

⁶

this is from an asymptotic formula for the variance of the average of a correlated sequence.

MCMC error and diagnostics

- An unbiased estimator of the IACT is

$$\tau_f \approx \hat{\tau}_f = 1 + 2 \sum_{k=1}^{n-1} \left(1 - \frac{k}{n}\right) \rho_k; \quad (24)$$

at longer lags k , the estimator starts to contain more noise than signal, and summing all the way out to n will result in a very noisy estimate of τ_f . Hence, we typically evaluate this estimator at sample size $m \ll n$.

- We have that the variance of the mean estimator using MCMC samples is:

$$\mathbb{V}[\hat{\mu}_n] \approx \frac{\sigma^2}{n} \hat{\tau}_f. \quad (25)$$

MCMC error and diagnostics

- The inefficiency factor can be used to derive the effective sample size:

$$n_{\text{eff}} = \frac{n}{1 + 2 \sum_{k=1}^{\infty} \rho_k} \approx \frac{n}{1 + 2 \sum_{k=1}^{n-1} \left(1 - \frac{k}{n}\right) \rho_k}. \quad (26)$$

the effective sample size is utilized to compare between the variance estimated via correlated MCMC samples and the ideal case of a variance computed from independent draws. Thus, the aim is to obtain an n_{eff} as close as possible to n .

- Other relevant metric used as an indicator of how fast the MCMC chains are mixing is the *mean square jump* (MSJ) distance, defined as

$$\text{MSJ} = \frac{1}{n} \sum_{i=1}^n \|x_i - x_{i-1}\|_2^2, \quad (27)$$

the larger the magnitude of the MSJ, the better the mixing.

MCMC error and diagnostics⁷

- A common test for convergence is the **Geweke test**. This method splits the chain (after burn-in period) into two parts, the first 10% and the last 50%. If the chain is stationary, the averages of these two parts should be approximately equal.
- Other diagnostics are: Gelman–Rubin diagnostic (which compares multiple chains), the Raftery–Lewis diagnostic (which is a method to find a proper burn-in), computing distance between distributions, etc.
- There has long been controversy over whether it is better to simulate one long chain of length n or m chains of length n/m each. There are settings where the single long chain has less bias. Parallel computing changes the picture.

⁷

We recommend ArviZ to analyze and postprocess MCMC results (<https://python.arviz.org/en/stable/>)

MCMC error and diagnostics (RWM example)

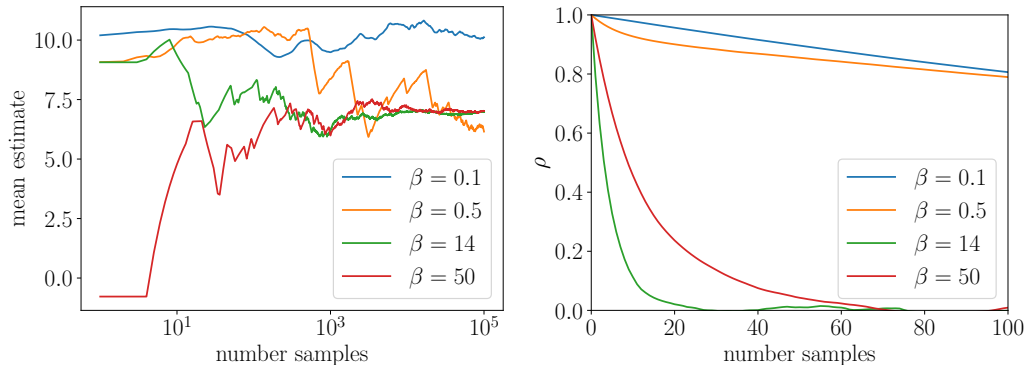


Figure: Cumulative mean and autocorrelation for Markov chain obtained with RWM for proposal different scaling. The n_{eff} is for each case: $[91, 128, 10875, 3578]$, here $n = 10^5$ and $n_b = 0.2n$.

5. MCMC – Gibbs sampler

Gibbs sampler

- We want to sample from $\pi(\mathbf{x})$ where $\mathbf{x} = [x_1, \dots, x_{j-1}, x_j, x_{j+1}, \dots, x_d]$. It might be difficult to construct a good proposal q that changes the whole vector at once.
- However, suppose we can construct good proposals that changes only one (or a few) component(s) of \mathbf{x} . For example, we consider that \mathbf{x} is comprised of x_j and the remaining elements, which we can lump together into \mathbf{x}_{-j} .
- In the Gibbs sampler, we repeatedly sample one component after another from the appropriate conditional distribution. Many models used in statistics and machine learning have simple conditional distributions for which the Gibbs sampler is easy to use.

Gibbs sampler

- We write the **full conditional distribution** of x_j given \mathbf{x}_{-j} as $\pi(x_j \mid \mathbf{x}_{-j})$, with marginal $\pi(\mathbf{x}_{-j})$.
- There are two versions: in the **random scan** Gibbs sampler, the component to update is chosen at random from $1, \dots, d$. In the **systematic scan** Gibbs sampler, the components are updated sequentially.
- We can show directly that sampling component j from its full conditional distribution preserves the stationary distribution π . Suppose that $\mathbf{x} \sim \pi$ and that we replace x_j by a value $z \sim \pi(x_j \mid \mathbf{x}_{-j})$, obtaining the point \mathbf{y} with $y_j = z$ and $y_k = x_k$ for $k \neq j$. Then

$$\pi(\mathbf{y}_{-j})\pi(y_j \mid \mathbf{y}_{-j}) = \pi(\mathbf{y}) \quad (28a)$$

$$\pi(\mathbf{x}_{-j})\pi(z \mid \mathbf{x}_{-j}) = \pi(\mathbf{x}). \quad (28b)$$

Gibbs sampler

- We may also understand the Gibbs sampler by relating it to MH. Suppose that we are at point \mathbf{x} and have decided to modify component j of \mathbf{x} to take the value z . Let \mathbf{y} be the point with $y_j = z$ and $y_k = x_k$ for $k \neq j$. If we use \mathbf{y} as the proposal in MH, then

$$\frac{\pi(\mathbf{y}_{-j})\pi(z \mid \mathbf{x}_{-j})\pi(x_j \mid \mathbf{x}_{-j})}{\pi(\mathbf{x}_{-j})\pi(x_j \mid \mathbf{x}_{-j})\pi(z \mid \mathbf{x}_{-j})} = 1, \quad (29)$$

thus, if we update component j of \mathbf{x} by sampling from its full conditional distribution, then we can view this as a MH proposal that is never rejected.

- Random scan Gibbs has detailed balance, while fixed scan does not, but we can construct a reversible Gibbs sampler with symmetric scan.

Gibbs sampler

- Gibbs sampling provides an alternative generation scheme based on successive generations from the full conditional distributions. The Gibbs sampler described now involves a complete scan over the components.
- Given a state $\mathbf{x}^{(k)} = [x_1^{(k)}, x_2^{(k)}, \dots, x_d^{(k)}]^\top$, samples ($k = 1, \dots, n$) from π are drawn component by component, as follows

$$\begin{aligned}
 x_1^{(k+1)} &\sim \pi \left(x_1 \mid x_2^{(k)}, x_3^{(k)}, \dots, x_d^{(k)} \right), \\
 x_2^{(k+1)} &\sim \pi \left(x_2 \mid x_1^{(k+1)}, x_3^{(k)}, \dots, x_d^{(k)} \right), \\
 &\vdots \\
 x_d^{(k+1)} &\sim \pi \left(x_d \mid x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_{d-1}^{(k+1)} \right).
 \end{aligned} \tag{30}$$

- There are many other possible updating strategies for visiting the components of \mathbf{x} .

Gibbs sampler

Algorithm 2: Random scan Gibbs.

```

1 set initial state  $\mathbf{x}_0$ ;
2 for  $i = 1$  to  $n$  do
3   sample an index  $j$  from  $\mathcal{U}(1, d)$ ;
4   sample component  $z \sim \pi(\cdot \mid \mathbf{x}_{-j}^{(i-1)})$  ;
5   set  $\mathbf{x}^{(i)} = \mathbf{x}^{(i-1)}$  and  $x_j^{(i)} = z$ ;
6 end

```

Algorithm 3: Fixed scan Gibbs.

```

1 set initial state  $\mathbf{x}_0$ ;
2 for  $i = 1$  to  $n$  do
3   set index  $j = ((i - 1) \bmod d) + 1$ ;
4   sample component  $z \sim \pi(\cdot \mid \mathbf{x}_{-j}^{(i-1)})$  ;
5   set  $\mathbf{x}^{(i)} = \mathbf{x}^{(i-1)}$  and  $x_j^{(i)} = z$ ;
6 end

```

Gibbs sampler

- The Gibbs sampler can fail to be irreducible when the space is disconnected.
- For the Gibbs sampler to work properly, it must be possible to reach any point from any other, using only moves parallel to the coordinate axes.
- Main difficulties:
 - (i) If the variables are strongly correlated (negatively or positively) then it may take too long to reach the stationary distribution. In this case, the problem would benefit from a reparametrization in terms of less correlated random variables.
 - (ii) Need to be able to show that the Gibbs sampler Markov chain is ergodic (irreducible and aperiodic). Obvious in many circumstances but sometimes an issue.

Gibbs sampler (example)

- Draw samples from the following (mixed) joint distribution of with $x = \{0, 1, \dots, m\}$ (discrete) and $0 \leq y \leq 1$ (continuous),

$$\pi_{XY}(x, y) \propto \binom{m}{x} y^{x+\alpha-1} (1-y)^{m-x+\beta-1}; \quad (31)$$

where, $m = 16, \alpha = 2, \beta = 4$.

- The conditional densities are:

$$\pi(x | y) = \text{Binomial}(m, y) \quad \text{and} \quad \pi(y | x) = \text{beta}(\alpha + x, n - \beta + x). \quad (32)$$

- The power of the Gibbs sampler is that by computing a sequence of these univariate conditional random variables, we can compute any feature of either marginal distribution.

Gibbs sampler (example)

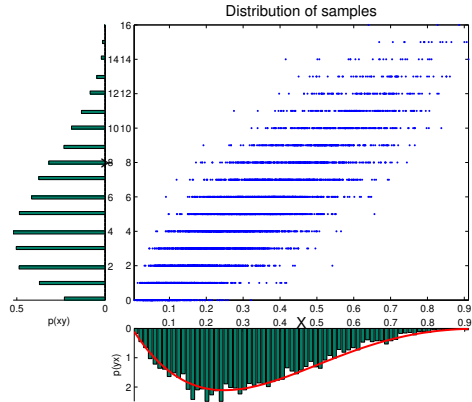


Figure: Samples from the mixed distribution.

Metropolis-within-Gibbs sampler

- Sometimes we can sample most, but not all of the full conditional distributions needed for the Gibbs sampler. For example, suppose that we can sample from $\pi_j | -_j$ for every j from 1 to d , except for one value j^* .
- The hybrid Gibbs or Metropolis-within-Gibbs algorithm uses a MH update in place of the missing j^* th full conditional distribution.
- The sampling schemes inside the Gibbs structure only require the generation of one sample for every Gibbs iteration. This is because the Markov chains associated to the individual conditionals are not required to reach stationarity at each Gibbs iteration. The reasons for this are discussed in [15, p.393].
- Nevertheless, some authors suggest that performing a few extra within-Gibbs iterations can be beneficial to accelerate the convergence of the Gibbs chain and to reduce correlation in the samples.

6. MCMC – preconditioned Crank-Nicolson

Preconditioned Crank–Nicolson (I)

- The performance of the RWM algorithms deteriorates with increasing dimension of the parameter space. The [preconditioned Crank–Nicolson](#) (pCN) method avoids this issue by designing proposal that perform-well in function spaces.
- Given a target posterior, the *preconditioned overdamped Langevin* equation is,

$$\frac{d\mathbf{x}_t}{dt} = -\mathbf{P}\nabla\Psi(\mathbf{x};\mathbf{y}) + \sqrt{2\mathbf{P}}\frac{dW_t}{dt}, \quad (33)$$

where W_t is a standard Brownian motion process, \mathbf{P} is a symmetric and positive semi-definite preconditioner, and the (unnormalized) logarithm of the posterior density Ψ and its gradient can be written as:

$$\Psi(\mathbf{x};\mathbf{y}) = \frac{1}{2} \left\| \Sigma_{\text{pr}}^{-1/2} \mathbf{x} \right\|_2^2 + \Phi(\mathbf{x};\mathbf{y}), \quad \nabla\Psi(\mathbf{x};\mathbf{y}) = \Sigma_{\text{pr}}^{-1} \mathbf{x} + \nabla\Phi(\mathbf{x};\mathbf{y}). \quad (34)$$

Preconditioned Crank–Nicolson (II)

- We can now substitute the gradient in eq. (34) into eq. (33) and solve the resulting stochastic differential equation using a Crank–Nicolson scheme with discretization step Δ [14].
- It can be shown that by making the preconditioner equal to the prior covariance, the following proposal mechanism arises from the Crank–Nicolson discretization [4]:

$$\mathbf{x}^* = \frac{(2 - \Delta)}{(2 + \Delta)} \mathbf{x} + \frac{\sqrt{8\Delta}}{(2 + \Delta)} \boldsymbol{\xi} \quad \longrightarrow \quad \mathbf{x}^* = \sqrt{1 - \beta^2} \mathbf{x} + \beta \boldsymbol{\xi}, \quad (35)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\text{pr}})$ and $\beta = \sqrt{8\Delta} / (2 + \Delta)$.

- A common choice for the scaling parameter is $\beta \in (0, 1]$ for discretization steps $\Delta \in (0, 2]$ (in general $\beta \rightarrow 0$ as $\Delta \rightarrow \infty$).

Preconditioned Crank–Nicolson (III)

- Given a current state $\mathbf{x}^{(i)}$, one has from eq. (35) that the associated proposal distribution is Gaussian with mean vector $\sqrt{1 - \beta^2} \mathbf{x}^{(i)}$ and covariance matrix $\beta^2 \Sigma_{\text{pr}}$.
- The acceptance probability of the pCN algorithm only requires evaluation of the potential function [5]

$$\alpha(\mathbf{x}, \mathbf{x}^*) = \min(1, \exp(\Phi(\mathbf{x}; \mathbf{y}) - \Phi(\mathbf{x}^*; \mathbf{y}))) . \quad (36)$$

- pCN is applicable when the posterior distribution has density with respect to a Gaussian reference measure (prior). However, it can also be extended to more general priors by applying a transformation.

7. MCMC – using gradients

Metropolis-adjusted Langevin algorithm

- In the Metropolis-adjusted Langevin algorithm (**MALA**), new states are proposed using (overdamped) Langevin dynamics, and these are accepted or rejected using the MH algorithm⁸.
- Let π denote a probability density function on \mathbb{R}^d , one from which it is desired to draw an ensemble of iid samples. We consider the overdamped Langevin Itô diffusion

$$\frac{d\mathbf{x}_t}{dt} = \frac{1}{2} \nabla \log \pi(\mathbf{x}) + \frac{dW_t}{dt}, \quad (37)$$

driven by the time derivative of a standard Brownian motion W .

- In the limit, as $t \rightarrow \infty$, the probability distribution $X(t)$ approaches a stationary distribution, which is also invariant under the diffusion. It turns out that this distribution is π .

⁸

If the acceptance/rejection is not applied the resulting approximate method is called **ULA** (unadjusted Langevin algorithm).

Metropolis-adjusted Langevin algorithm

- Approximate sample paths of the Langevin diffusion can be generated by the Euler–Maruyama method with a fixed time step $\varepsilon > 0$. We set \mathbf{x}_0 and then recursively define an approximation to the true solution by

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \frac{\varepsilon^2}{2} \nabla \log \pi(\mathbf{x}_k) + \varepsilon \boldsymbol{\xi}_k, \quad (38)$$

where each $\boldsymbol{\xi}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

- In contrast to the Euler–Maruyama method for simulating the Langevin diffusion. MALA incorporates an additional step. We consider the above update rule as defining a proposal $\tilde{\mathbf{x}}$ for a new state:

$$\tilde{\mathbf{x}}_{k+1} := \mathbf{x}_k + \frac{\varepsilon^2}{2} \nabla \log \pi(\mathbf{x}_k) + \varepsilon \boldsymbol{\xi}_k; \quad (39)$$

this proposal is accepted or rejected according to the MH algorithm.

Metropolis-adjusted Langevin algorithm

- That is, the acceptance probability is:

$$\alpha(\mathbf{x}_k, \tilde{\mathbf{x}}_{k+1}) = \min \left(1, \frac{\pi(\tilde{\mathbf{x}}_{k+1})q(\mathbf{x}_k | \tilde{\mathbf{x}}_{k+1})}{\pi(\mathbf{x}_k)q(\tilde{\mathbf{x}}_{k+1} | \mathbf{x}_k)} \right); \quad (40)$$

where the proposal has the form

$$q(\mathbf{x}' | \mathbf{x}) = \mathcal{N} \left(\mathbf{x}'; \mathbf{x} + \frac{\varepsilon^2}{2} \nabla \log \pi(\mathbf{x}), \varepsilon^2 \mathbf{I}_d \right). \quad (41)$$

- The combined dynamics of the Langevin diffusion and the MH algorithm satisfy the detailed balance conditions necessary for the existence of a unique, invariant, stationary distribution.
- For limited classes of target distributions, the optimal acceptance rate for this algorithm can be shown to be 0.574; this can be used to tune ε .

Hamiltonian MC (intro)



Figure: The mode of the target as a planet and the gradient of the target as its gravitational field. Taken from [3].

- So far we discussed classic MCMC approaches to draw samples from a target distribution $\pi(x)$.
- Hamiltonian Monte Carlo (HMC): use Hamiltonian dynamics to simulate particle trajectories ([6]).
- Define a Hamiltonian function in terms of the target distribution.
- Introduce an auxiliary *momentum* variables, which typically have independent Gaussian distributions.

Hamiltonian MC (intro)

- Hamiltonian dynamics operate on a d -dimensional *position* vector \mathbf{q} , and a d -dimensional *momentum* vector \mathbf{p} , so that the full state space has $2d$ dimensions. The system is described by a function of \mathbf{q} and \mathbf{p} known as the *Hamiltonian* $H(\mathbf{q}, \mathbf{p})$.
- In HMC, one uses Hamiltonian functions that can be written as (closed-system dynamics):

$$H(\mathbf{q}, \mathbf{p}) = \underbrace{U(\mathbf{q})}_{\text{potential energy}} + \underbrace{K(\mathbf{p}, \mathbf{q})}_{\text{kinetic energy}} . \quad (42)$$

- The potential energy is completely determined by the target distribution, indeed $U(\mathbf{q})$ is equal to the logarithm of the target distribution π .
- The kinetic energy is unconstrained and must be specified by the implementation.

Hamiltonian MC

- The Hamiltonian is an energy function for the joint state of 'position-momentum', and so defines a joint distribution for them as follows:

$$\pi(\mathbf{q}, \mathbf{p}) = \frac{1}{Z} \exp \left(-\frac{H(\mathbf{q}, \mathbf{p})}{T} \right) = \frac{1}{Z} \exp(-U(\mathbf{q})) \exp(-K(\mathbf{p})). \quad (43)$$

- There are several ways to set the kinetic energy (density of the auxiliary momentum) [3]:
 - Euclidean–Gaussian kinetic energy: using a fixed covariance \mathbf{M} estimated from the position parameters, $K(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} + \ln(|\mathbf{M}|) + \text{const.}$
 - Riemann–Gaussian kinetic energy: unlike the Euclidean metric, varies as one moves through parameter space, $K(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \mathbf{p}^T \boldsymbol{\Sigma}(\mathbf{q})^{-1} \mathbf{p} + \frac{1}{2} \ln(|\boldsymbol{\Sigma}(\mathbf{q})|) + \text{const.}$
 - Non-Gaussian kinetic energies.

Hamiltonian MC

- Hamilton's equations read as follows:

$$\frac{d\mathbf{q}}{dt} = +\frac{\partial H}{\partial \mathbf{p}} = [\mathbf{M}^{-1}\mathbf{p}] \quad (44a)$$

$$\frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \mathbf{q}} = -\frac{\partial K}{\partial \mathbf{q}} - \frac{\partial U}{\partial \mathbf{q}}, \quad (44b)$$

where $\frac{\partial U}{\partial \mathbf{q}}$ is the gradient of the logarithm of the target density.

- Discretizing Hamilton's equations:
 - (i) Euler's method (**no**).
 - (ii) Modified Euler's method (**a bit better**).
 - (iii) Symplectic integrators: the leapfrog method (**the standard choice**).

Hamiltonian MC (remarks)

Properties

- Hamiltonian dynamics are time-reversible and volume-preserving.
- The dynamics keep the Hamiltonian invariant. A Hamiltonian trajectory will (if simulated exactly) move within a hypersurface of constant probability density.

Each iteration of the HMC algorithm has two steps. Both steps leave the joint distribution of $\pi(q, p)$ invariant (detailed balance) [11].

- In the first step, new values of p are randomly drawn from their Gaussian distribution, independently of the current q .
- In the second step, a Metropolis update is performed, using Hamiltonian dynamics to propose a new state.
- Optimal acceptance rate is 0.65. The step size ε and trajectory length L need to be tuned.

Hamiltonian MC (example)

- More typical behavior of HMC and RWM is illustrated by a 100-dimensional multivariate Gaussian distribution in which the variables are independent, with means of zero, and standard deviations of 0.01, 0.02, ..., 0.99, 1.
- Suppose that we have no knowledge of the details of this distribution, so we will use HMC with the same simple, rotationally symmetric kinetic energy function.
- Consistent with this, we use HMC to this distribution using trajectories with $L = 150$ and with ε randomly selected for each iteration, uniformly from (0.0104, 0.0156); here $n = 10^3$. We compare with a RWM with proposal standard deviation drawn uniformly from (0.0176, 0.0264); with a lag period of 150.
- These are close to optimal settings for both methods. The rejection rate was 0.13 for HMC and 0.75 for RWM.

Hamiltonian MC (example)

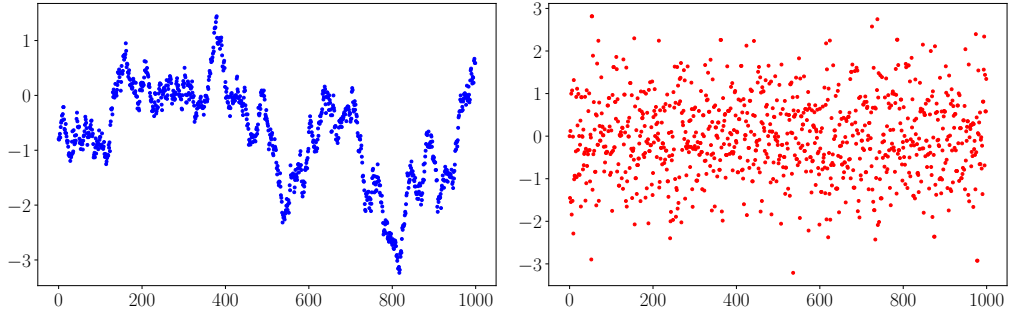


Figure: Values for the variable with largest standard deviation for the 100-dimensional example, from a RWM run and an HMC run.

Hamiltonian MC (example)

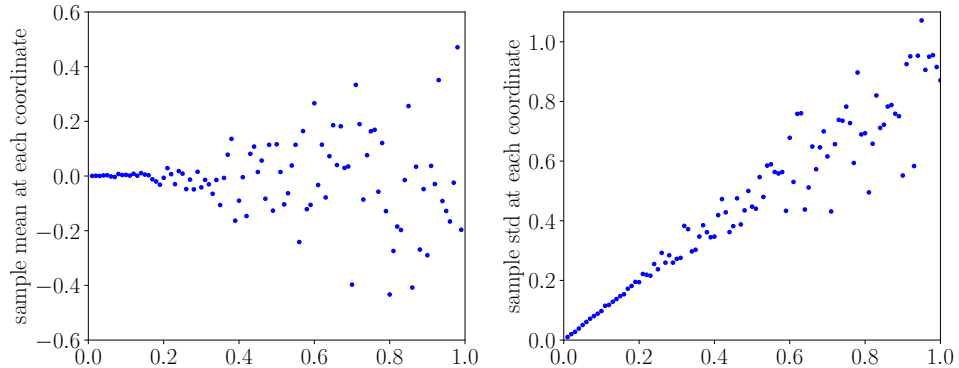


Figure: Estimates of mean and standard deviations for the 100-dimensional example, using RWM.

Hamiltonian MC (example)

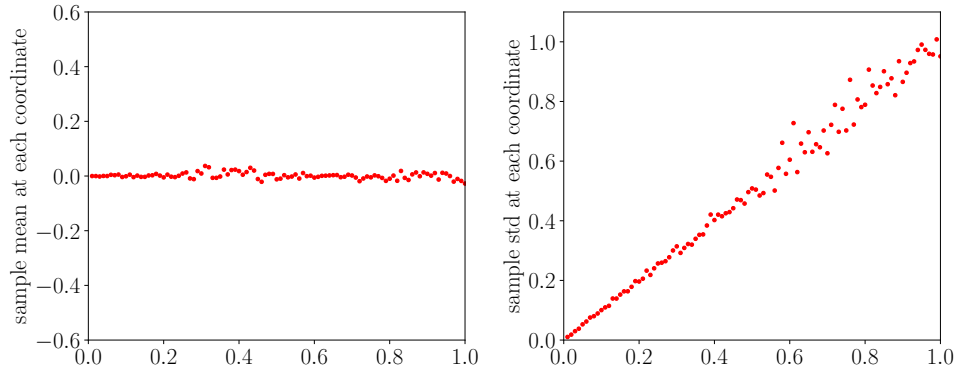


Figure: Estimates of mean and standard deviations for the 100-dimensional example, using HMC.

The no-U-turn sampler (NUTS)

- HMC is an algorithm that avoids the random walk behavior and sensitivity to correlated parameters that plague many MCMC methods by taking a series of steps informed by first-order gradient information.
- However, HMC's performance is highly sensitive to two user-specified parameters: a step size ϵ and a desired number of steps L .
- The No-U-Turn Sampler (NUTS), an extension to HMC that eliminates the need to set a number of steps L , as well as the step-size.
- We simulate in discrete time steps, and to make sure you explore the parameter space properly you simulate steps in one direction and the twice as many in the other direction, turn around again, etc. At some point you want to stop this and a good way of doing that is when you have done a U-turn (i.e., appear to have gone all over the place).

The no-U-turn sampler (NUTS)

- NUTS begins by introducing an auxiliary variable with conditional distribution. After re-sampling from this distribution, NUTS uses the leapfrog integrator to trace out a path forwards and backwards in fictitious time. First running forwards or backwards 1 step, then forwards or backwards 2 steps, then forwards or backwards 4 steps, etc.
- This doubling process implicitly builds a balanced binary tree whose leaf nodes correspond to position-momentum state. The doubling is stopped when the subtrajectory from the leftmost to the rightmost nodes of any balanced subtree of the overall binary tree starts to double back on itself (i.e., the fictional particle starts to make a “U-turn”).
- At this point NUTS stops the simulation and samples from among the set of points computed during the simulation, taking care to preserve detailed balance.

The no-U-turn sampler (NUTS)

- To adapt the step-size, NUTS uses a modified dual averaging algorithm during the burn-in phase.
- The good thing about NUTS is that proposals are made based on the shape of the posterior and they can happen at the other end of the distribution. In contrast, MH makes proposals within a ball, and Gibbs sampling only moves along one (or at least very few) dimensions at a time.

8. Approximation methods

The Laplace approximation (I)

- As an alternative to simulation of integrals, we can also attempt analytic approximations.
- One of the oldest and most useful approximations is the integral Laplace approximation. It is based on the following argument: Suppose that we are interested in evaluating the integral

$$\int_A f(x | y) dx, \quad (45)$$

for a fixed y , and f non-negative and integrable.

- Write

$$f(x | y) = \exp(nh(x | y)),$$

where n is a parameter that can go to infinity.

The Laplace approximation (II)

- Use a Taylor series expansion of $h(x | y)$ about a point x_0 to obtain

$$h(x | y) \approx h(x_0 | y) + (x - x_0)h'(x_0 | y) + \frac{(x - x_0)^2}{2}h''(x_0 | y) + \frac{(x - x_0)^3}{3!}h'''(x_0 | y) + R_n(x)$$

while the remainder $R_n(x)$ satisfies $\lim_{x \rightarrow x_0} R_n(x)/(x - x_0)^3 = 0$.

- Now choose $x_0 = x^*$, the value that satisfies $h'(x^* | y) = 0$ and maximizes $h(x | y)$ for a given value of y . Then, the linear term in the Taylor series is zero and we have the approximation

$$\int_A \exp(nh(x | y))dx = \exp(nh(x^* | y)) \int_A \exp\left(n \left(\frac{(x - x^*)^2}{2}\right) h''(x^* | y)\right) \quad (46a)$$

$$\exp\left(n \left(\frac{(x - x^*)^3}{3!}\right) h'''(x^* | y)\right) dx \quad (46b)$$

which is valid within a neighborhood of x^* .

The Laplace approximation (III)

- The cubic term in the exponent is now expanded in a series around x^* and using the Taylor expansion of the exponential function, we obtain:

$$1 + n \left(\frac{(x - x^*)^3}{3!} \right) h'''(x^*|y) + n^2 \left(\frac{(x - x^*)^6}{2!(3!)^2} \right) [h'''(x^*|y)]^2 \quad (47)$$

and

$$\int_A \exp(nh(x | y)) dx = \exp(nh(x^* | y)) \int_A \exp \left(n \left(\frac{(x - x^*)^2}{2} \right) h''(x^*|y) \right) \quad (48a)$$

$$\left[1 + n \left(\frac{(x - x^*)^3}{3!} \right) h'''(x^*|y) + n^2 \left(\frac{(x - x^*)^6}{2!(3!)^2} \right) [h'''(x^*|y)]^2 + R_n(x) \right] dx \quad (48b)$$

- Excluding R_n , we call the integral approximations of a *first-order*, if it includes only the first term in the right-hand side; of a *second-order*, if it includes the first two terms; and of a *third-order*, if it includes all three terms.

The Laplace approximation (IV)

- We can evaluate these expressions further since the above integrand is the kernel of a Gaussian density with mean x^* and variance $1/(nh''(x^*|y))$. This Gaussian approximation is the so-called [Laplace approximation](#).
- More precisely, letting Φ denote the standard Gaussian CDF, and taking $A = [a, b]$, we can evaluate the integral in the **first-order approximation** to obtain (with $n = 1$)

$$\int_a^b \exp(h(x | y)) dx = \exp(h(x^* | y)) \sqrt{\frac{2\pi}{-h''(x^* | y)}} \quad (49a)$$

$$\left[\Phi \left(\sqrt{-h''(x^* | y)}(b - x^*) \right) - \Phi \left(\sqrt{-h''(x^* | y)}(a - x^*) \right) \right]. \quad (49b)$$

The Laplace approximation (V)

- The Laplace approximation is reasonable in the central region of the density, it becomes quite unacceptable in the tails.
- In problems where Monte Carlo calculations are prohibitive because of computing time, the Laplace approximation can be useful as a guide to the solution of the problem [15].
- Also, the corresponding Taylor series can be used as a proposal density (for MCMC), which is particularly useful in problems where no obvious proposal exists.

The saddlepoint approximation (I)

- The saddlepoint approximation, in contrast to the Laplace approximation, is mainly a technique for approximating a function rather than an integral, although it naturally leads to an integral approximation (Proposed by Henry E. Daniels in 1954).
- Suppose that we are interested in evaluating the integral

$$g(y) = \int_A \pi(x | y) dx, \quad (50)$$

for a range of values of y .

- One interpretation of a saddlepoint approximation is that for each value of y , we do a Laplace approximation centered at x^* (the saddle point⁹).
- One way to derive the saddlepoint approximation is to use an **Edgeworth expansion** (aka Gram–Charlier A series).

⁹ also minimax point; it is a point on the surface of a function where the slopes (derivatives) in orthogonal directions are all zero, but which is not a local extremum of the function.

The saddlepoint approximation (II)

- As a result of a quite detailed derivation, we obtain the approximation to the density of X to be

$$\pi_X(x) = \frac{\sqrt{n}}{\sigma} \varphi\left(\frac{x - \mu}{\sigma/\sqrt{n}}\right) \left\{ 1 + \frac{\kappa}{6\sqrt{n}} \left(\left(\frac{x - \mu}{\sigma/\sqrt{n}}\right)^3 - 3 \left(\frac{x - \mu}{\sigma/\sqrt{n}}\right) \right) + \mathcal{O}(1/n) \right\}. \quad (51)$$

- Ignoring the term within braces produces the usual Gaussian approximation, which is accurate to $\mathcal{O}(1/\sqrt{n})$.
- If we are using eq. (51) for values of x near μ , then the value of the expression in braces is close to zero, and the approximation will then be accurate to $\mathcal{O}(1/n)$. The trick of the saddlepoint approximation is to make this always be the case,
- To do so, we use a family of densities such that, for each x , we can choose a density from the family to cancel the term in braces in eq. (51).

The saddlepoint approximation (III)

- One method of creating such a family is through a technique known as **exponential tilting** (aka or Exponential Change of Measure). The result of the exponential tilt is a family of Edgeworth expansions for $\pi_X(x)$ indexed by a parameter τ , that is

$$\pi_X(x) = \exp(-n[\tau x - K(\tau)]) \frac{\sqrt{n}}{\sigma_\tau} \varphi\left(\frac{x - \mu_\tau}{\sigma_\tau/\sqrt{n}}\right) \quad (52a)$$

$$\left\{ 1 + \frac{\kappa_\tau}{6\sqrt{n}} \left(\left(\frac{x - \mu_\tau}{\sigma_\tau/\sqrt{n}} \right)^3 - 3 \left(\frac{x - \mu_\tau}{\sigma_\tau/\sqrt{n}} \right) \right) + \mathcal{O}(1/n) \right\}. \quad (52b)$$

- As the parameter τ free to choose, for each x we choose $\tau = \tau(x)$ so that the mean satisfies $\mu_\tau = x$. This choice cancels the middle term in the square brackets), thereby improving the order of the approximation.

The saddlepoint approximation (IV)

- If $K(\tau) = \log(\mathbb{E}[\exp \tau X])$ is the cumulant generating function, we can choose τ so that $K'(\tau) = x$, which is the saddlepoint equation.
- Denoting $\sigma_\tau = K''(\tau)$, and $\hat{\tau}$ the value obtained from the saddlepoint equation, we get the saddlepoint approximation (with $n = 1$):

$$\pi_X(x) = \exp(K(\hat{\tau}) - \hat{\tau}x) \frac{1}{\sigma_{\hat{\tau}}} \varphi(0) + (1 + \mathcal{O}(1)) \quad (53a)$$

$$\approx \frac{1}{\sqrt{2\pi\sigma_{\hat{\tau}}}} \exp(K(\hat{\tau}) - \hat{\tau}x). \quad (53b)$$

- The saddlepoint can also be used to approximate the tail area of a distribution [15].
- This better error rates are obtained by renormalizing the approximation so that it integrates to 1.

Final remarks

- Other algorithms we did not cover in class are:
 - (i) Auxiliary variable: slice sampler, simulated annealing, simulated tempering, Hit-and-run.
 - (ii) Sequential algorithms: sequential importance sampling, population Monte Carlo, etc.
 - (iii) Approximate methods: Laplace approximations, approximate Bayesian computation (likelihood-free), variational Bayesian inference, transport maps, Stein variational gradient descent, etc.
- We covered the most common MCMC algorithms used in practice. The particular choice of a method will depend on your application and computational resources.
- Quite often we have to work with approximate methods. Sampling is a computationally intensive task, which complicates application of UQ in general inverse problems.
- UQ keeps growing by the day; new samplers and techniques are addressing complicated inference tasks.
- Check this cool demo ! [Link Chi Feng's MCMC demo.](#)

References

- [1] C. Andrieu et al. "A tutorial on adaptive MCMC". In: *Statistics and Computing* 18 (2008), 343—373.
- [2] A. A. Barker. "Monte Carlo calculations of the radial distribution functions for a proton-electron plasma". In: *Australian Journal of Physics* 18 (1965), pp. 119–1347.
- [3] M. Betancourt. "A conceptual introduction to Hamiltonian Monte Carlo". In: *arXiv:1701.02434v2* (2017), pp. 1–60.
- [4] S. L. Cotter et al. "MCMC methods for functions: modifying old algorithms to make them faster". In: *Statistical Science* 28.3 (2013), pp. 424–446.
- [5] M. Dashti et al. "The Bayesian approach to inverse problems". In: *Handbook of uncertainty quantification*. Ed. by R. Ghanem et al. Springer International Publishing, 2017. Chap. 10, pp. 311–428.
- [6] S. Duane et al. "Hybrid Monte Carlo". In: *Physics Letters B* 195.2 (1987), pp. 216–222.
- [7] A. Gelman et al. "Efficient Metropolis jumping rules". In: *Bayesian Analysis* 5 (1996), 599–607.
- [8] H. Haario et al. "An Adaptive Metropolis Algorithm". In: *Bernoulli* 7.2 (2001), pp. 223–242.
- [9] W. K. Hastings. "Monte Carlo sampling methods using Markov chains and their applications". In: *Biometrika* 57.1 (1970), pp. 97–109.
- [10] N. Metropolis et al. "Equation of state calculations by fast computing machines". In: *Chemical Physics* 21.6 (1953), pp. 1087–1092.
- [11] R. M. Neal. "MCMC using Hamiltonian dynamics". In: *Handbook of Markov chain Monte Carlo*. Ed. by S. Brooks et al. Chapman & Hall/CRC, 2011. Chap. 5, pp. 113–162.
- [12] J. R. Norris. *Markov Chains*. Cambridge University Press, 1998.

- [13] A. B. Owen. *Monte Carlo theory, methods and examples*. artowen.su.domains/mc/, 2018.
- [14] W. H. Press et al. *Numerical recipes in C: the art of scientific computing*. 3rd ed. Cambridge University Press, 2007.
- [15] C. P. Robert et al. *Monte Carlo statistical methods*. 2nd ed. Springer, 2004.
- [16] L. Tierney. "A note on Metropolis-Hastings kernels for general state spaces". In: *Annals of applied probability* (1998), pp. 1–9.
- [17] B. Walsh. *Markov chain Monte Carlo and Gibbs sampling*. Lecture notes for EEB 581, version 26, 2004.

Disclaimer: all figures are either generated by the Author or under Creative Commons licenses (except the Figure in slide 61, which has been referenced)