

Markov Chain Monte Carlo algorithms

Modèles probabilistes pour l'apprentissage

Théo Moins¹

¹Statify, Inria Grenoble Rhône-Alpes

November 21, 2022

Table of Contents

1. Introduction
2. Metropolis–Hastings
3. Introduction to Hamiltonian Monte Carlo
4. MCMC Convergence diagnostics

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Model with parameters $\boldsymbol{\theta} := \theta_1, \dots, \theta_D$: $p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{i=1}^n p(x_i \mid \boldsymbol{\theta})$

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Model with parameters $\boldsymbol{\theta} := \theta_1, \dots, \theta_D$: $p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{i=1}^n p(x_i \mid \boldsymbol{\theta})$

Prior uncertainty on $\boldsymbol{\theta}$: $p(\boldsymbol{\theta})$

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Model with parameters $\boldsymbol{\theta} := \theta_1, \dots, \theta_D$: $p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{i=1}^n p(x_i \mid \boldsymbol{\theta})$

Prior uncertainty on $\boldsymbol{\theta}$: $p(\boldsymbol{\theta})$

\implies **Bayesian update:**

$$p(\boldsymbol{\theta} \mid \mathbf{x}) \propto p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})$$

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Model with parameters $\boldsymbol{\theta} := \theta_1, \dots, \theta_D$: $p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{i=1}^n p(x_i \mid \boldsymbol{\theta})$

Prior uncertainty on $\boldsymbol{\theta}$: $p(\boldsymbol{\theta})$

\implies **Bayesian update:**

$$p(\boldsymbol{\theta} \mid \mathbf{x}) \propto p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})$$

What's next?

Reminder: Bayesian paradigm

Data: $\mathbf{x} := x_1, \dots, x_n$

Model with parameters $\boldsymbol{\theta} := \theta_1, \dots, \theta_D$: $p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{i=1}^n p(x_i \mid \boldsymbol{\theta})$

Prior uncertainty on $\boldsymbol{\theta}$: $p(\boldsymbol{\theta})$

\implies **Bayesian update:**

$$p(\boldsymbol{\theta} \mid \mathbf{x}) \propto p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})$$

What's next? All computations reduce to posterior means of quantity of interest $f(\boldsymbol{\theta})$:

$$\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\boldsymbol{\theta})] = \int f(\boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathbf{x})d\boldsymbol{\theta}$$

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$,

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$,
 \implies computation of an integral
(highly dimensional, non-explicit, etc.)

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$,
 \implies computation of an integral
(highly dimensional, non-explicit, etc.)

MCMC (Markov Chain Monte Carlo):

Monte Carlo	Markov Chain
$\mathbb{E}[f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta^{(i)})$	$\theta^{(i+1)} \mid \theta^{(i)} \sim P(\theta^{(i)}, \cdot)$

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$,
 \implies computation of an integral
(highly dimensional, non-explicit, etc.)

MCMC (Markov Chain Monte Carlo):

Monte Carlo	Markov Chain
$\mathbb{E}[f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta^{(i)})$	$\theta^{(i+1)} \mid \theta^{(i)} \sim P(\theta^{(i)}, \cdot)$

With a correct choice of $P(\cdot, \cdot)$ one can often prove that

$$\frac{1}{N} \sum_{i=1}^N f(\theta^{(i)}) \rightarrow \mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)], \quad \text{when } N \rightarrow +\infty.$$

Table of Contents

1. Introduction
2. Metropolis–Hastings
3. Introduction to Hamiltonian Monte Carlo
4. MCMC Convergence diagnostics

Metropolis–Hastings

```
1 Initialize  $\theta^{(0)}$ ;  
2 for  $i \leftarrow 1$  to  $N$  do  
3   | Sample a candidate  $\theta^* \sim q(\theta^* | \theta^{(i-1)})$ ;  
4   | Compute an acceptance ratio  $\alpha(\theta^* | \theta^{(i-1)}) \in [0, 1]$ ;  
5   | Accept the candidate ( $\theta^{(i)} = \theta^*$ ) with probability  $\alpha(\theta^* | \theta^{(i-1)})$ ;  
6   | Otherwise reject ( $\theta^{(i)} = \theta^{(i-1)}$ );  
7 end
```

Algorithm 1: Metropolis-Hastings Algorithm

Illustration: <https://chi-feng.github.io/mcmc-demo/app.html>

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

- The ratio $\frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})}$ favours the candidates with the largest posterior density,

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

- The ratio $\frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})}$ favours the candidates with the largest posterior density,
- The ratio $\frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})}$ can be seen as a correction term,

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

- The ratio $\frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})}$ favours the candidates with the largest posterior density,
- The ratio $\frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})}$ can be seen as a correction term,
- $\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)})$ is only an **acceptance probability**!

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

- The ratio $\frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})}$ favours the candidates with the largest posterior density,
- The ratio $\frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})}$ can be seen as a correction term,
- $\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)})$ is only an **acceptance probability**!

Special cases:

- **Metropolis:** $q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*) = q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})$ (symmetrical proposal).

Metropolis–Hastings

Acceptation ratio for the general framework:

$$\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = \min \left(1, \frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})} \frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})} \right)$$

- The ratio $\frac{p(\boldsymbol{\theta}^* \mid \mathbf{x})}{p(\boldsymbol{\theta}^{(i)} \mid \mathbf{x})}$ favours the candidates with the largest posterior density,
- The ratio $\frac{q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*)}{q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})}$ can be seen as a correction term,
- $\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)})$ is only an **acceptance probability**!

Special cases:

- **Metropolis:** $q(\boldsymbol{\theta}^{(i)} \mid \boldsymbol{\theta}^*) = q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i)})$ (symmetrical proposal).
- **Gibbs:** $q(\boldsymbol{\theta}_k^* \mid \boldsymbol{\theta}^{(i)}) = p(\boldsymbol{\theta}_k^* \mid \boldsymbol{\theta}_1^{(i)}, \dots, \boldsymbol{\theta}_D^{(i)}, \mathbf{x})$ (full conditional, $\alpha(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(i-1)}) = 1$)

Metropolis–Hastings

```
1 Initialize  $\theta^{(0)}$ ;  
2 for  $i \leftarrow 1$  to  $N$  do  
3   | Sample  $\theta^* \sim q(\theta^* | \theta^{(i-1)})$ ;  
4   | Sample  $u \sim \mathcal{U}[0, 1]$ ;  
5   | if  $u < \alpha(\theta^* | \theta^{(i-1)})$  then  
6   |   |  $\theta^{(i)} = \theta^*$                                 /* Accept */  
7   | else  
8   |   |  $\theta^{(i)} = \theta^{(i-1)}$                         /* Reject */  
9   | end  
10 end
```

Algorithm 2: Metropolis–Hastings Algorithm

Performance of Random Walk Metropolis

To conclude:

1. **Proposal distribution** $q(\theta^* | \theta^{(i-1)})$: favour large volumes . . .
2. **Acceptance ratio** $\alpha(\theta^* | \theta^{(i-1)})$: favour large densities.

\implies The combination makes the selection toward the typical set.

Drawback: tuning is extremely hard in high dimension!

To sum up:

MCMC algorithms

- Iterative sampling of the probability density,

To sum up:

MCMC algorithms

- Iterative sampling of the probability density,
- Next sample depends on the precedent sample,

To sum up:

MCMC algorithms

- Iterative sampling of the probability density,
- Next sample depends on the precedent sample,
- Based on an Acceptation-Rejection Rule,

To sum up:

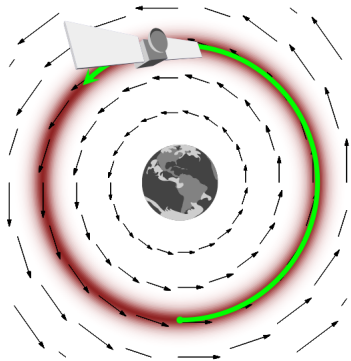
MCMC algorithms

- Iterative sampling of the probability density,
- Next sample depends on the precedent sample,
- Based on an Acceptation-Rejection Rule,
- Examples: Metropolis-Hastings (MH), Gibbs Sampler, Hamiltonian Monte Carlo (HMC) etc.

Table of Contents

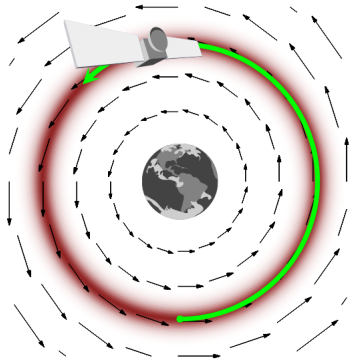
1. Introduction
2. Metropolis–Hastings
3. Introduction to Hamiltonian Monte Carlo
4. MCMC Convergence diagnostics

Introduction to HMC



From [Betancourt \(2017\)](#)

Introduction to HMC

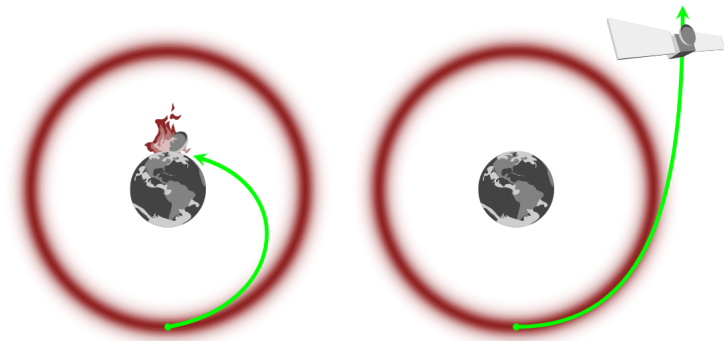


- Probabilistic system \rightarrow Physical system
- Mode of the distribution \rightarrow Massive planet
- Gradient of the density \rightarrow Gravitational field

From [Betancourt \(2017\)](#)

Introduction to HMC

Conservative exploration:



From [Betancourt \(2017\)](#)

Introduction to HMC

For each dimension k , we add a momentum ξ_k to the position θ_k ($2D$ parameters).

Introduction to HMC

For each dimension k , we add a momentum ξ_k to the position θ_k ($2D$ parameters).

$$\implies p(\theta, \xi \mid \mathbf{x}) = p(\xi \mid \theta, \mathbf{x})p(\theta \mid \mathbf{x})$$

Introduction to HMC

For each dimension k , we add a momentum ξ_k to the position θ_k ($2D$ parameters).

$$\implies p(\boldsymbol{\theta}, \boldsymbol{\xi} \mid \mathbf{x}) = p(\boldsymbol{\xi} \mid \boldsymbol{\theta}, \mathbf{x})p(\boldsymbol{\theta} \mid \mathbf{x})$$

Hamiltonian definition:

$$\begin{aligned} H(\boldsymbol{\theta}, \boldsymbol{\xi}) &= -\log p(\boldsymbol{\xi}, \boldsymbol{\theta} \mid \mathbf{x}) \\ &= -\log p(\boldsymbol{\xi} \mid \boldsymbol{\theta}, \mathbf{x}) - \log p(\boldsymbol{\theta} \mid \mathbf{x}) \\ &= K + V, \quad \text{with } \begin{cases} K : \text{kinetic energy} \\ V : \text{potential energy} \end{cases} \end{aligned}$$

Introduction to HMC

Hamilton's equation of motion:

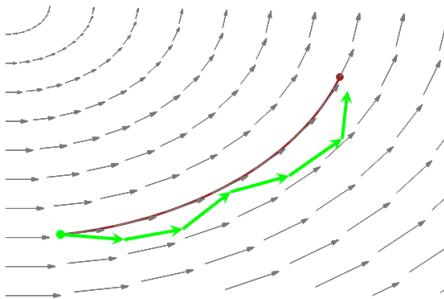
$$\frac{d\theta_k}{dt} = \frac{\partial H}{\partial \xi_k}, \quad \frac{d\xi_k}{dt} = \frac{\partial H}{\partial \theta_k}. \quad (1)$$

Introduction to HMC

Hamilton's equation of motion:

$$\frac{d\theta_k}{dt} = \frac{\partial H}{\partial \xi_k}, \quad \frac{d\xi_k}{dt} = -\frac{\partial H}{\partial \theta_k}. \quad (1)$$

Solving (1) leads to obtain the position $\phi_t(\boldsymbol{\theta}, \boldsymbol{\xi}) \in \mathbb{R}^{2D}$ across time t .



From [Betancourt \(2017\)](#)

Introduction to HMC

```
1 Initialize  $\theta^{(0)}$ ;  
2 for  $i \leftarrow 1$  to  $N$  do  
3   | Sample a momentum  $\xi^{(i-1)} \sim q(\xi^{(i-1)} \mid \theta^{(i-1)})$ ;  
4   | Sample an amount of time  $t \sim \mathcal{U}[0, T]$ ;  
5   | Numerically solve the Hamilton's equation  $(\theta_t, \xi_t) \leftarrow \phi_t(\theta^{(i-1)}, \xi^{(i-1)})$ ;  
6   | Accept the candidate  $(\theta^{(i)} = \theta_t)$  with probability  $\alpha(\theta_t, -\xi_t \mid \theta^{(i-1)}, \xi^{(i-1)})$ ;  
7   | Otherwise reject  $(\theta^{(i)} = \theta^{(i-1)})$ ;  
8 end
```

Algorithm 3: HMC Algorithm

Introduction to HMC

```
1 Initialize  $\theta^{(0)}$ ;
2 for  $i \leftarrow 1$  to  $N$  do
3   Sample  $\xi^* \sim \mathcal{N}(0, M)$ ;
4   Let  $\theta_0 = \theta^{(i-1)}$  and  $\xi_0 = \xi^*$ ;
5   for  $l \leftarrow 1$  to  $L$  do
6      $\xi_{l-1/2} = \xi_{l-1} - \epsilon \nabla_{\theta} V(\theta_0)/2$ ;
7      $\theta_l = \theta_{l-1} + \epsilon \xi_{l-1}$ ;
8      $\xi_l = \xi_{l-1/2} - \epsilon \nabla_{\theta} V(\theta_{Leaps})/2$ ;
9   end
10  Sample  $u \sim \mathcal{U}[0, 1]$ ;
11  if  $u < \alpha(\theta_L, -\xi_L \mid \theta_0, \xi_0)$  then
12     $\theta^{(i)} = \theta_L$ 
13  else
14     $\theta^{(i)} = \theta_0$ 
15  end
```

/* Accept */

/* Reject */

Introduction to HMC - Summary

Hamiltonian Monte Carlo

- Algorithm based on Hamiltonian Mechanics
- Time discretization simulated by the Leap Frog Algorithm
- Acceptation rule based on the discretization error in simulating Hamiltonian mechanics
- A lot of hyperparameters!

References:

Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.

Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical science*, 7(4):457–472.

Table of Contents

1. Introduction
2. Metropolis–Hastings
3. Introduction to Hamiltonian Monte Carlo
4. MCMC Convergence diagnostics

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$.

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$.

MCMC (Markov Chain Monte Carlo):

Monte Carlo

$$\mathbb{E}[f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta^{(i)})$$

Markov Chain

$$\theta^{(i+1)} \mid \theta^{(i)} \sim P(\theta^{(i)}, \cdot)$$

Reminder: MCMC

Bayesian inference on $\theta \sim p(\theta \mid \mathbf{x}) \implies$ computation of $\mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)] = \int f(\theta)p(\theta \mid \mathbf{x})d\theta$.

MCMC (Markov Chain Monte Carlo):

Monte Carlo	Markov Chain
$\mathbb{E}[f(\theta)] \approx \frac{1}{N} \sum_{i=1}^N f(\theta^{(i)})$	$\theta^{(i+1)} \mid \theta^{(i)} \sim P(\theta^{(i)}, \cdot)$

With a correct choice of $P(\cdot, \cdot)$ one can often prove that

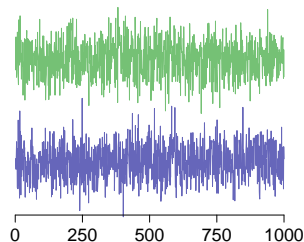
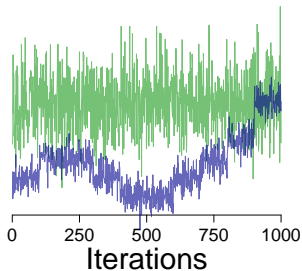
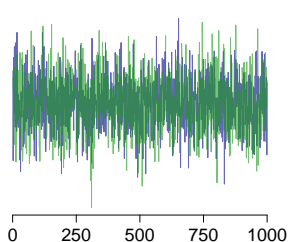
$$\frac{1}{N} \sum_{i=1}^N f(\theta^{(i)}) \rightarrow \mathbb{E}_{p(\cdot \mid \mathbf{x})}[f(\theta)], \quad \text{when } N \rightarrow +\infty.$$

\hookrightarrow How to choose N ?

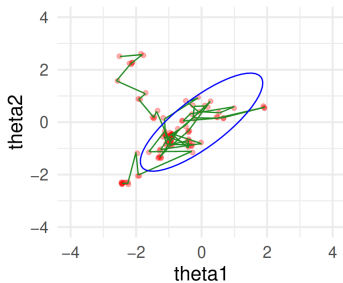
Has the chain(s) converged?

Two kind of convergence issues: mixing and stationarity (univariate case)

Simulations

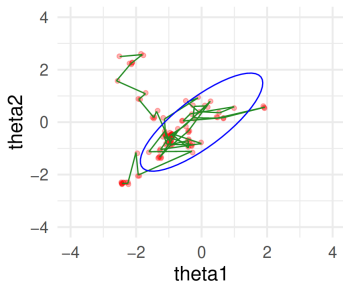


Has the chain(s) converged?



From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

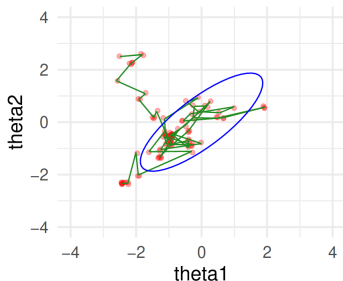
Has the chain(s) converged?



From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

Warm-up: discard the early iterations (usually the first half),

Has the chain(s) converged?

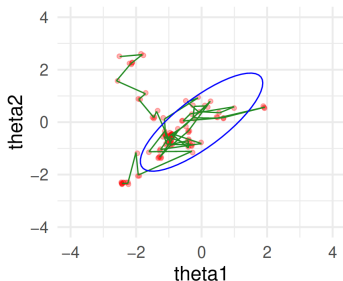


From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

Warm-up: discard the early iterations (usually the first half),

Thinning: keep every k th iterations from each sequence,

Has the chain(s) converged?



From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

Warm-up: discard the early iterations (usually the first half),

Thinning: keep every k th iterations from each sequence,

Multiple chains: reduce the influence of the starting point.

Autocorrelation

The dependency decreases the amount of effective information contained in each element of the chain.

Autocorrelation

The dependency decreases the amount of effective information contained in each element of the chain.

↪ Increase of the variance on the target parameters.

Autocorrelation

The dependency decreases the amount of effective information contained in each element of the chain.

↔ Increase of the variance on the target parameters.

Autocorrelation function at lag- t : correlation between elements of the sequence distant from t steps.

$$\text{ACF}_t(\theta) = \frac{\text{ACov}_t(\theta)}{\text{Var}(\theta)} = \frac{\frac{1}{S-t} \sum_{s=1}^{S-t} (\theta^{(s)} - \bar{\theta})(\theta^{(s+t)} - \bar{\theta})}{\frac{1}{S-1} \sum_{s=1}^S (\theta^{(s)} - \bar{\theta})^2}$$

Autocorrelation

The dependency decreases the amount of effective information contained in each element of the chain.

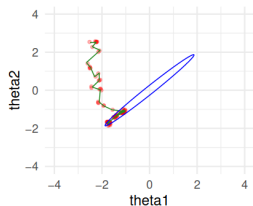
↪ Increase of the variance on the target parameters.

Autocorrelation function at lag- t : correlation between elements of the sequence distant from t steps.

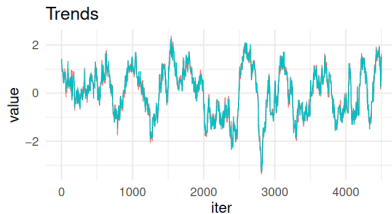
$$\text{ACF}_t(\theta) = \frac{\text{ACov}_t(\theta)}{\text{Var}(\theta)} = \frac{\frac{1}{S-t} \sum_{s=1}^{S-t} (\theta^{(s)} - \bar{\theta})(\theta^{(s+t)} - \bar{\theta})}{\frac{1}{S-1} \sum_{s=1}^S (\theta^{(s)} - \bar{\theta})^2}$$

Ideal case: uncorrelated draws $\implies \text{ACF}_t(\theta) = 0 \quad \forall t > 1$

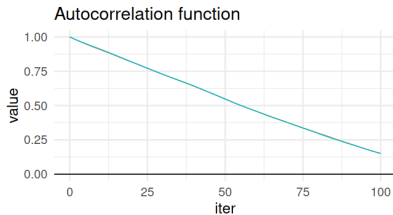
Autocorrelation



• Draws — Steps of the sampler — 90% HP

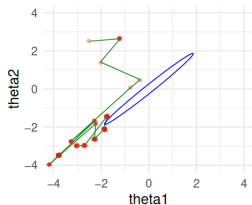


— theta1 — theta2

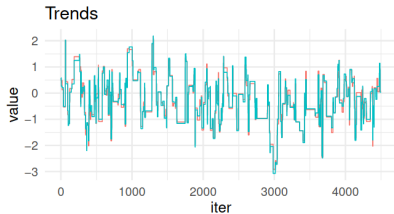


From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

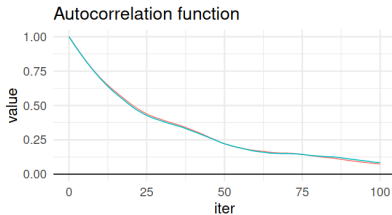
Autocorrelation



• Draws — Steps of the sampler — 90% HP



— theta1 — theta2



From Aki Vehtari : https://avehtari.github.io/BDA_course_Aalto/

Effective Sample Size (ESS)

In the independent case ($M = 1$ chain):

$$N.\text{Var}(\bar{\theta}) = N.\mathbb{E} [(\bar{\theta} - \theta_0)^2] \xrightarrow[n \rightarrow +\infty]{} \text{Var}(\theta \mid \mathbf{x})$$

Effective Sample Size (ESS)

In the independent case ($M = 1$ chain):

$$N.\text{Var}(\bar{\theta}) = N.\mathbb{E} [(\bar{\theta} - \theta_0)^2] \xrightarrow{n \rightarrow +\infty} \text{Var}(\theta \mid \mathbf{x})$$

Otherwise:

$$N.\text{Var}(\bar{\theta}) \xrightarrow{n \rightarrow +\infty} \left(1 + 2 \sum_{t=1}^{\infty} \text{ACF}_t(\theta) \right) \text{Var}(\theta \mid \mathbf{x}),$$

Effective Sample Size (ESS)

In the independent case ($M = 1$ chain):

$$N.\text{Var}(\bar{\theta}) = N.\mathbb{E} [(\bar{\theta} - \theta_0)^2] \xrightarrow{n \rightarrow +\infty} \text{Var}(\theta \mid \mathbf{x})$$

Otherwise:

$$N.\text{Var}(\bar{\theta}) \xrightarrow{n \rightarrow +\infty} \left(1 + 2 \sum_{t=1}^{\infty} \text{ACF}_t(\theta) \right) \text{Var}(\theta \mid \mathbf{x}),$$

ESS: Equivalent number of independent draws

$$\text{ESS} = \frac{N}{1 + 2 \sum_{t=1}^{\infty} \text{ACF}_t(\theta)}$$

\hookrightarrow Number of samples to obtain the same variance in the i.i.d case.

\hat{R} (aka potential scale reduction factor)

Introduced by [Gelman and Rubin \(1992\)](#).

Consider m chains of size n , with $\theta^{(i,j)}$ denoting the i th draw from chain j .

Comparison of the **between-variance** B and the **within-variance** W of the chains:

$$\hat{R} = \sqrt{\frac{\hat{W} + \hat{B}}{\hat{W}}}$$

\hat{R} (aka potential scale reduction factor)

Introduced by [Gelman and Rubin \(1992\)](#).

Consider m chains of size n , with $\theta^{(i,j)}$ denoting the i th draw from chain j .

Comparison of the **between-variance** B and the **within-variance** W of the chains:

$$\hat{R} = \sqrt{\frac{\hat{W} + \hat{B}}{\hat{W}}}$$

$$\text{Between var : } \hat{B} = \frac{1}{m-1} \sum_{j=1}^m (\bar{\theta}^{(\cdot,j)} - \bar{\theta}^{(\cdot,\cdot)})^2, \quad \text{where } \bar{\theta}^{(\cdot,j)} = \frac{1}{n} \sum_{i=1}^n \theta^{(i,j)}, \quad \bar{\theta}^{(\cdot,\cdot)} = \frac{1}{m} \sum_{j=1}^m \bar{\theta}^{(\cdot,j)},$$

$$\text{Within var : } \hat{W} = \frac{1}{m} \sum_{j=1}^m s_j^2, \quad \text{where } s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\theta^{(i,j)} - \bar{\theta}^{(\cdot,j)})^2.$$

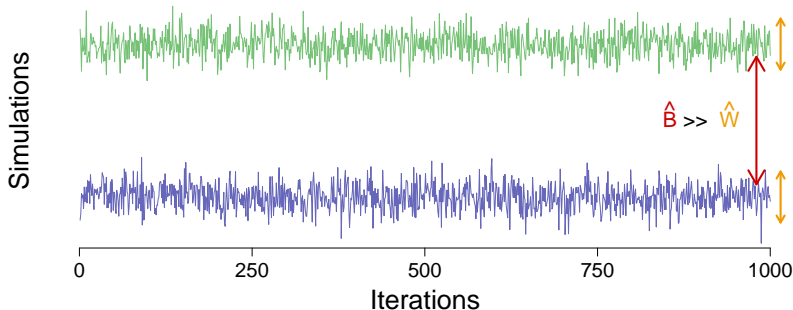
\hat{R} (aka potential scale reduction factor)

Introduced by [Gelman and Rubin \(1992\)](#).

Consider m chains of size n , with $\theta^{(i,j)}$ denoting the i th draw from chain j .

Comparison of the **between-variance** B and the **within-variance** W of the chains:

$$\hat{R} = \sqrt{\frac{\hat{W} + \hat{B}}{\hat{W}}}$$



Limitations of \hat{R}

The main limitations are:

- It is not robust to certain types of non-convergence.
 \hat{R} and potentially also rank- \hat{R}

Limitations of \hat{R}

The main limitations are:

- It is not robust to certain types of non-convergence.
 \hat{R} and potentially also rank- \hat{R}
- It suffers from a lack of interpretability.
What is R associated to \hat{R} ?

Limitations of \hat{R}

The main limitations are:

- It is not robust to certain types of non-convergence.
 \hat{R} and potentially also rank- \hat{R}
- It suffers from a lack of interpretability.
What is R associated to \hat{R} ?
- It must be compared to an arbitrary chosen threshold.
 $\hat{R} \geq 1.1$? 1.01 ?

Limitations of \hat{R}

The main limitations are:

- It is not robust to certain types of non-convergence.
 \hat{R} and potentially also rank- \hat{R}
- It suffers from a lack of interpretability.
What is R associated to \hat{R} ?
- It must be compared to an arbitrary chosen threshold.
 $\hat{R} \geq 1.1$? 1.01 ?
- It is associated with a univariate parameter.
How to manage multiple parameters?

Limitations of \hat{R}

The main limitations are:

- It is not robust to certain types of non-convergence.
 \hat{R} and potentially also rank- \hat{R}
- It suffers from a lack of interpretability.
What is R associated to \hat{R} ?
- It must be compared to an arbitrary chosen threshold.
 $\hat{R} \geq 1.1$? 1.01?
- It is associated with a univariate parameter.
How to manage multiple parameters?

BDA book recommends to use \hat{R} and ESS in the following way:

$$\begin{aligned}\hat{R} \in [1, 1.01] &\implies \text{"Chains are mixing well"}. \\ \text{ESS} > 400 &\implies \text{"Enough data for estimation"}.\end{aligned}$$

Examples

Example 1: Bayesian logistic regression

$$\boldsymbol{\beta} \sim \mathcal{N}(0, 0.35^2 \cdot \mathbf{I}_4), \quad y_j \sim \text{Bernoulli} \left(\frac{1}{1 + e^{-\mathbf{x}_j^\top \boldsymbol{\beta}}} \right).$$

Examples

Example 1: Bayesian logistic regression

$$\beta \sim \mathcal{N}(0, 0.35^2 \cdot \mathbf{I}_4), \quad y_j \sim \text{Bernoulli} \left(\frac{1}{1 + e^{-\mathbf{x}_j^\top \beta}} \right).$$

Example 2: Hierarchical model (8 Schools)

↪ Test the effectiveness of coaching courses.

y_j : coaching effect for the school j

$$\mu \sim \mathcal{N}(0, 5), \quad \tau \sim \mathcal{N}(0, 10),$$
$$\theta_j \sim \mathcal{N}(\mu, \tau), \quad y_j \sim \mathcal{N}(\theta_j, \sigma_j^2)$$

Model "in between" the separate model and the joint model.

School	Estimated treatment effect, y_j	Standard error of effect estimate, σ_j
A	28	15
B	8	10
C	-3	16
D	7	11
E	-1	9
F	1	11
G	18	10
H	12	18

From BDA book