

Lecture 16: Sampling

Recommended reading:

Bishop: Chapter 11

MacKay: Chapter 29

Duncan Gillies and Marc Deisenroth

Department of Computing
Imperial College London

February 22–24, 2016

Monte Carlo Methods—Motivation

- ▶ Monte Carlo methods are computational techniques that make use of **random numbers**
- ▶ Two typical problems:
 1. Problem 1: Generate samples $\{\mathbf{x}^{(s)}\}$ from a given probability distribution $p(\mathbf{x})$

Monte Carlo Methods—Motivation

- ▶ Monte Carlo methods are computational techniques that make use of **random numbers**
- ▶ Two typical problems:
 1. Problem 1: Generate samples $\{\mathbf{x}^{(s)}\}$ from a given probability distribution $p(\mathbf{x})$
 2. Problem 2: Estimate expectations of functions under that distribution:

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

Monte Carlo Methods—Motivation

- ▶ Monte Carlo methods are computational techniques that make use of **random numbers**
- ▶ Two typical problems:
 1. Problem 1: Generate samples $\{\mathbf{x}^{(s)}\}$ from a given probability distribution $p(\mathbf{x})$
 2. Problem 2: Estimate expectations of functions under that distribution:

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

▶▶ Example: Means/variances of distributions, marginal likelihood

Complication: Integral cannot be evaluated analytically

Monte Carlo Estimation

- ▶ Statistical sampling can be applied to compute **expectations**

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim p(\mathbf{x})\end{aligned}$$

Monte Carlo Estimation

- ▶ Statistical sampling can be applied to compute **expectations**

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim p(\mathbf{x})\end{aligned}$$

- ▶ Example: Making predictions (e.g., Bayesian linear regression with a training set $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ at test input \mathbf{x}_*)

$$\begin{aligned}p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) &= \int p(\mathbf{y}_*|\boldsymbol{\theta}, \mathbf{x}_*)p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \\ &\approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}_*|\boldsymbol{\theta}^{(s)}, \mathbf{x}_*), \quad \boldsymbol{\theta}^{(s)} \sim p(\boldsymbol{\theta}|\mathcal{D})\end{aligned}$$

Monte Carlo Estimation

- ▶ Statistical sampling can be applied to compute **expectations**

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim p(\mathbf{x})\end{aligned}$$

- ▶ Example: Making predictions (e.g., Bayesian linear regression with a training set $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ at test input \mathbf{x}_*)

$$\begin{aligned}p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) &= \int p(\mathbf{y}_*|\boldsymbol{\theta}, \mathbf{x}_*)p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \\ &\approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}_*|\boldsymbol{\theta}^{(s)}, \mathbf{x}_*), \quad \boldsymbol{\theta}^{(s)} \sim p(\boldsymbol{\theta}|\mathcal{D})\end{aligned}$$

- ▶ If we can sample from $p(\mathbf{x})$ (or $p(\boldsymbol{\theta})$) we can approximate these integrals

Properties of Monte Carlo Sampling

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim p(\mathbf{x})\end{aligned}$$

- ▶ Estimator is **unbiased**
- ▶ **Variance shrinks** $\propto 1/S$, regardless of the dimensionality of \mathbf{x}

Alternatives to Monte Carlo

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

To evaluate these expectations we can use other methods than Monte Carlo:

- ▶ Numerical integration (low-dimensional problems)
- ▶ Deterministic approximations, e.g., **Variational Bayes**, **Expectation Propagation**

Back to Monte Carlo Estimation

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim p(\mathbf{x})\end{aligned}$$

- ▶ How do we get these samples?
- ▶▶ Need to solve Problem 1
 - ▶ Sampling from simple distributions
 - ▶ Sampling from complicated distributions

Important Example

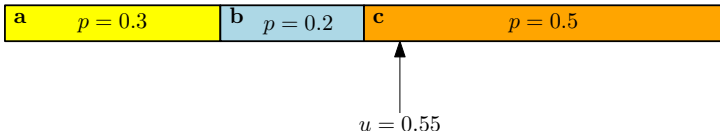
- ▶ By specifying the model, we know the prior $p(\boldsymbol{\theta})$ and the likelihood $p(\mathcal{D}|\boldsymbol{\theta})$
- ▶ The **unnormalized posterior** is

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

and there is often no hope to compute the normalization constant

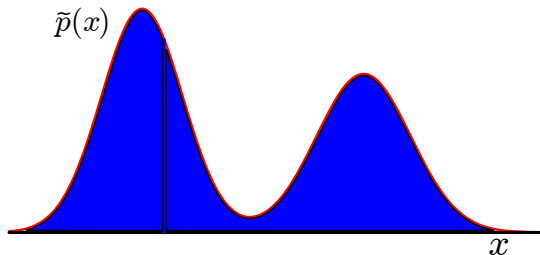
- ▶ Samples are a good way to characterize this posterior (important for model comparison, Bayesian predictions, ...)

Sampling Discrete Values



- ▶ $u \sim \mathcal{U}[0, 1]$, where \mathcal{U} is the uniform distribution
- ▶ $u = 0.55 \Rightarrow x = c$

Continuous Variables

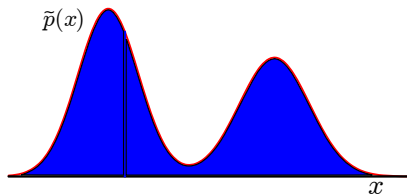


More complicated.

Geometrically, sample uniformly from the area under the curve

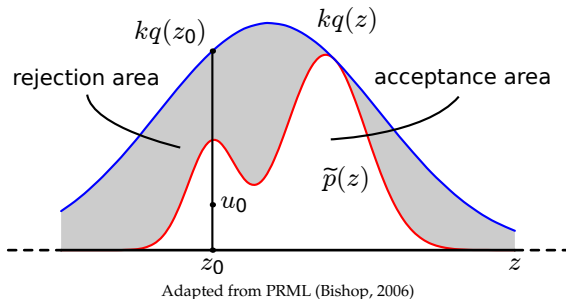
Rejection Sampling

Rejection Sampling: Setting



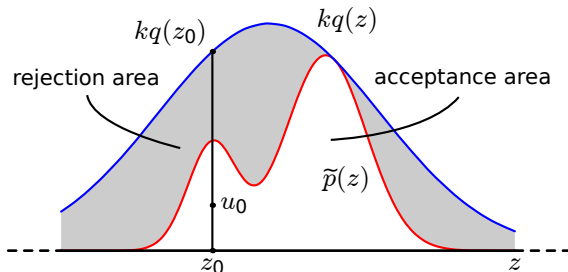
- ▶ Assume sampling from $p(z)$ is difficult
- ▶ Evaluating $\tilde{p}(z) = Zp(z)$ is easy (and Z may be unknown)
- ▶ Find a simpler distribution (**proposal distribution**) $q(z)$ from which we can easily draw samples (e.g., Gaussian)
- ▶ Find an upper bound $kq(z) \geq \tilde{p}(z)$

Algorithm



1. Generate $z_0 \sim q(z)$
2. Generate $u_0 \sim \mathcal{U}[0, kq(z_0)]$
3. If $u_0 > \tilde{p}(z_0)$, reject the sample. Otherwise, retain z_0

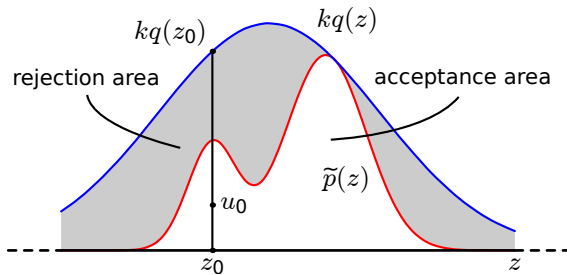
Properties



Adapted from PRML (Bishop, 2006)

- ▶ Accepted pairs (z, u) are uniformly distributed under the curve of $\tilde{p}(z)$
- ▶ Probability density of the z -coordinates of accepted points must be proportional to $\tilde{p}(z)$
- ▶ Samples are independent samples from $p(z)$

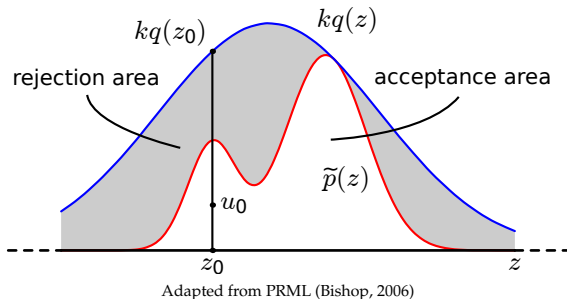
Shortcomings



Adapted from PRML (Bishop, 2006)

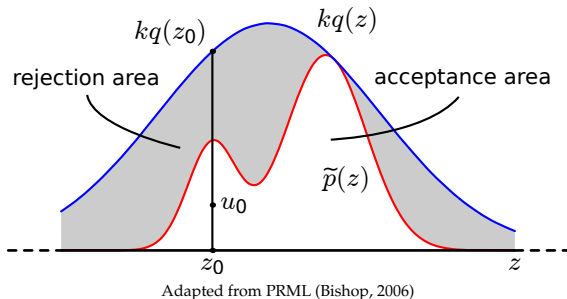
- ▶ Finding k is tricky

Shortcomings



- ▶ Finding k is tricky
- ▶ In high dimensions the factor k is probably huge

Shortcomings



- ▶ Finding k is tricky
- ▶ In high dimensions the factor k is probably huge
- ▶ **Low acceptance rate**

Importance Sampling

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\mathbb{E}_p[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\begin{aligned}\mathbb{E}_p[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int f(\mathbf{x})p(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x}\end{aligned}$$

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\begin{aligned}\mathbb{E}_p[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int f(\mathbf{x})p(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x} = \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x}\end{aligned}$$

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\begin{aligned}\mathbb{E}_p[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int f(\mathbf{x})p(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x} = \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \\ &= \mathbb{E}_q\left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}\right]\end{aligned}$$

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\begin{aligned}\mathbb{E}_p[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int f(\mathbf{x})p(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x} = \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \\ &= \mathbb{E}_q\left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}\right]\end{aligned}$$

If we choose q in a way that we can easily sample from it, we can approximate this last expectation by Monte Carlo:

$$E_q\left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}\right] \approx \frac{1}{S}\sum_{s=1}^S f(\mathbf{x}^{(s)})\frac{p(\mathbf{x}^{(s)})}{q(\mathbf{x}^{(s)})}, \quad \mathbf{x}^{(s)} \sim q(\mathbf{x})$$

Importance Sampling

Key idea: Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution q (**proposal distribution**):

$$\begin{aligned}\mathbb{E}_p[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int f(\mathbf{x})p(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x} = \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \\ &= \mathbb{E}_q\left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}\right]\end{aligned}$$

If we choose q in a way that we can easily sample from it, we can approximate this last expectation by Monte Carlo:

$$E_q\left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}\right] \approx \frac{1}{S}\sum_{s=1}^S f(\mathbf{x}^{(s)})\frac{p(\mathbf{x}^{(s)})}{q(\mathbf{x}^{(s)})} = \frac{1}{S}\sum_{s=1}^S w_s f(\mathbf{x}^{(s)}), \quad \mathbf{x}^{(s)} \sim q(\mathbf{x})$$

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p
- ▶ Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
 - ▶▶ **Degeneracy** (see also Particle Filtering)

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p
- ▶ Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
 - ▶▶ **Degeneracy** (see also Particle Filtering)
- ▶ **Many draws** from proposal density q required, especially in high dimensions

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p
- ▶ Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
 - ▶▶ **Degeneracy** (see also Particle Filtering)
- ▶ **Many draws** from proposal density q required, especially in high dimensions
- ▶ Requires to be able to evaluate true p . Generalization exists for \tilde{p} . This generalization is biased (but consistent).

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p
- ▶ Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
 - ▶▶ **Degeneracy** (see also Particle Filtering)
- ▶ **Many draws** from proposal density q required, especially in high dimensions
- ▶ Requires to be able to evaluate true p . Generalization exists for \tilde{p} . This generalization is biased (but consistent).
- ▶ Does not scale to interesting problems

Properties

- ▶ Unbiased if $q > 0$ where $p > 0$ and if we can evaluate p
- ▶ Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
 - ▶▶ **Degeneracy** (see also Particle Filtering)
- ▶ **Many draws** from proposal density q required, especially in high dimensions
- ▶ Requires to be able to evaluate true p . Generalization exists for \tilde{p} . This generalization is biased (but consistent).
- ▶ Does not scale to interesting problems
- ▶▶ Different approach to sample from complicated (high-dimensional) distributions

Markov Chains

Objective

Generate samples from an unknown target distribution.

Markov Chains

Key idea: Instead of independent samples, use a proposal density q that depends on the state $\mathbf{x}^{(t)}$

Markov Chains

Key idea: Instead of independent samples, use a proposal density q that depends on the state $\mathbf{x}^{(t)}$

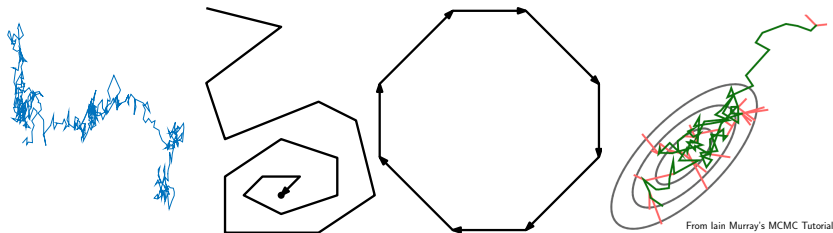
- ▶ **Markov property:** $p(\mathbf{x}^{(t+1)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)})$ only depends on the previous setting/state of the chain
- ▶ T is called a **transition operator**
- ▶ Example: $T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \sigma^2 \mathbf{I})$

Markov Chains

Key idea: Instead of independent samples, use a proposal density q that depends on the state $\mathbf{x}^{(t)}$

- ▶ **Markov property:** $p(\mathbf{x}^{(t+1)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)})$ only depends on the previous setting/state of the chain
- ▶ T is called a **transition operator**
- ▶ Example: $T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \sigma^2 \mathbf{I})$
- ▶ Samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ form a **Markov chain**
- ▶ Samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ are **no longer independent**, but **unbiased**
 - ▶▶ We can still average them

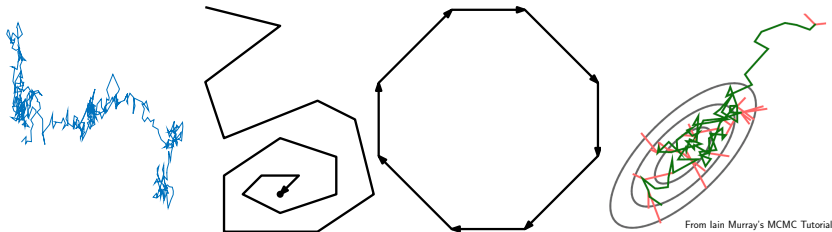
Behavior of Markov Chains



Four different behaviors of Markov chains:

- ▶ Diverge (e.g., random walk diffusion where $\mathbf{x}^{(t+1)} \sim \mathcal{N}(\mathbf{x}^{(t)}, \mathbf{I})$)

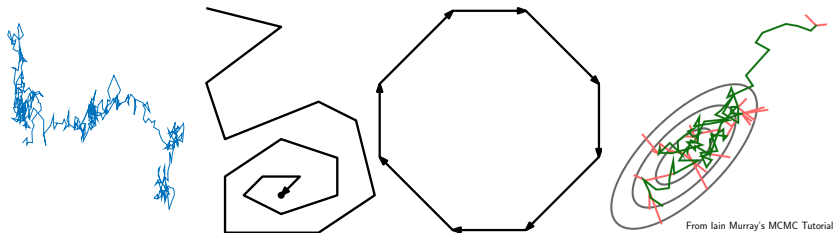
Behavior of Markov Chains



Four different behaviors of Markov chains:

- ▶ Diverge (e.g., random walk diffusion where $\mathbf{x}^{(t+1)} \sim \mathcal{N}(\mathbf{x}^{(t)}, \mathbf{I})$)
- ▶ Converge to an absorbing state

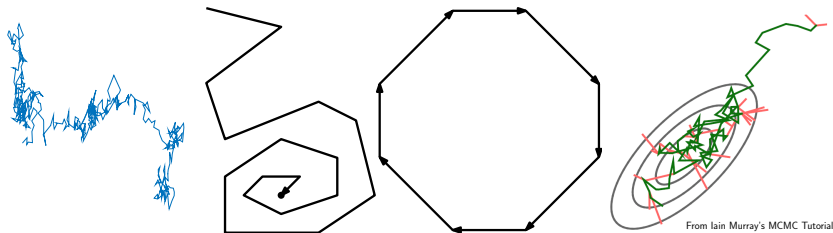
Behavior of Markov Chains



Four different behaviors of Markov chains:

- ▶ Diverge (e.g., random walk diffusion where $\mathbf{x}^{(t+1)} \sim \mathcal{N}(\mathbf{x}^{(t)}, \mathbf{I})$)
- ▶ Converge to an absorbing state
- ▶ Converge to a (deterministic) limit cycle

Behavior of Markov Chains



Four different behaviors of Markov chains:

- ▶ Diverge (e.g., random walk diffusion where $x^{(t+1)} \sim \mathcal{N}(x^{(t)}, I)$)
- ▶ Converge to an absorbing state
- ▶ Converge to a (deterministic) limit cycle
- ▶ Converge to an equilibrium distribution p^* : Markov chain remains in a region, bouncing around in a random way

Converging to an Equilibrium Distribution

- ▶ Remember objective: Explore/sample parameters that may have generated our data (generate samples from posterior)
 - ▶▶ Bouncing around in an equilibrium distribution is a good thing

Converging to an Equilibrium Distribution

- ▶ Remember objective: Explore/sample parameters that may have generated our data (generate samples from posterior)
 - ▶▶ Bouncing around in an equilibrium distribution is a good thing
- ▶ Design the Markov chain such that the equilibrium distribution is the desired posterior
- ▶ We know the equilibrium distribution (the one we want to sample from)
 - ▶▶ Generate a Markov chain that converges to that equilibrium distribution (independent of start state)
- ▶ Although successive samples are dependent we can effectively generate independent samples by running the Markov chain long enough: Discard most of the samples, retain only every M th sample

Conditions for Converging to an Equilibrium Distribution

Markov chain conditions:

- ▶ **Invariance/Stationarity:** If you run the chain for a long time and you are in the equilibrium distribution, you stay in equilibrium if you take another step.
 - ▶ Self-consistency property
- ▶ **Ergodicity:** Any state can be reached from any state.
 - ▶ Equilibrium distribution is the same no matter where we start

Conditions for Converging to an Equilibrium Distribution

Markov chain conditions:

- ▶ **Invariance/Stationarity:** If you run the chain for a long time and you are in the equilibrium distribution, you stay in equilibrium if you take another step.
 - ▶ Self-consistency property
- ▶ **Ergodicity:** Any state can be reached from any state.
 - ▶ Equilibrium distribution is the same no matter where we start

Property

Ergodic Markov chains only have one equilibrium distribution

Conditions for Converging to an Equilibrium Distribution

Markov chain conditions:

- ▶ **Invariance/Stationarity:** If you run the chain for a long time and you are in the equilibrium distribution, you stay in equilibrium if you take another step.
 - ▶ Self-consistency property
- ▶ **Ergodicity:** Any state can be reached from any state.
 - ▶ Equilibrium distribution is the same no matter where we start

Property

Ergodic Markov chains only have one equilibrium distribution

- ▶ Use ergodic and stationary Markov chains to generate samples from the equilibrium distribution

Invariance and Detailed Balance

- ▶ Invariance: Each step leaves the distribution p^* invariant (we stay in p^*):

$$p^*(\mathbf{x}') = \sum_x T(\mathbf{x}'|\mathbf{x})p^*(\mathbf{x})$$

Invariance and Detailed Balance

- ▶ Invariance: Each step leaves the distribution p^* invariant (we stay in p^*):

$$p^*(\mathbf{x}') = \sum_x T(\mathbf{x}'|\mathbf{x})p^*(\mathbf{x})$$

Once we sample from p^* , the transition operator will not change this, i.e., we do not fall back to some funny distribution $p \neq p^*$

Invariance and Detailed Balance

- ▶ Invariance: Each step leaves the distribution p^* invariant (we stay in p^*):

$$p^*(\mathbf{x}') = \sum_x T(\mathbf{x}'|\mathbf{x})p^*(\mathbf{x})$$

Once we sample from p^* , the transition operator will not change this, i.e., we do not fall back to some funny distribution $p \neq p^*$

- ▶ **Sufficient condition** for p^* being invariant:

Detailed balance:

$$p^*(\mathbf{x})T(\mathbf{x}|\mathbf{x}') = p^*(\mathbf{x}')T(\mathbf{x}'|\mathbf{x})$$

- ▶▶ Also ensures that the Markov chain is reversible.

Metropolis-Hastings

Metropolis-Hastings

- ▶ Assume that $\tilde{p} = Zp$ can be evaluated easily (in practice: $\log \tilde{p}$)
- ▶ Proposal density $q(\mathbf{x}'|\mathbf{x}^{(t)})$ depends on last sample $\mathbf{x}^{(t)}$.
Example: Gaussian centered at $\mathbf{x}^{(t)}$

Metropolis-Hastings

- ▶ Assume that $\tilde{p} = Zp$ can be evaluated easily (in practice: $\log \tilde{p}$)
- ▶ Proposal density $q(\mathbf{x}'|\mathbf{x}^{(t)})$ depends on last sample $\mathbf{x}^{(t)}$.
Example: Gaussian centered at $\mathbf{x}^{(t)}$

Metropolis-Hastings Algorithm

1. Generate $\mathbf{x}' \sim q(\mathbf{x}'|\mathbf{x}^{(t)})$

2. If

$$\frac{q(\mathbf{x}^{(t)}|\mathbf{x}')\tilde{p}(\mathbf{x}')}{q(\mathbf{x}'|\mathbf{x}^{(t)})\tilde{p}(\mathbf{x}^{(t)})} \geq u, \quad u \sim U[0, 1]$$

accept the sample $\mathbf{x}^{(t+1)} = \mathbf{x}'$. Otherwise set $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)}$.

Metropolis-Hastings

- ▶ Assume that $\tilde{p} = Zp$ can be evaluated easily (in practice: $\log \tilde{p}$)
- ▶ Proposal density $q(\mathbf{x}'|\mathbf{x}^{(t)})$ depends on last sample $\mathbf{x}^{(t)}$.
Example: Gaussian centered at $\mathbf{x}^{(t)}$

Metropolis-Hastings Algorithm

1. Generate $\mathbf{x}' \sim q(\mathbf{x}'|\mathbf{x}^{(t)})$

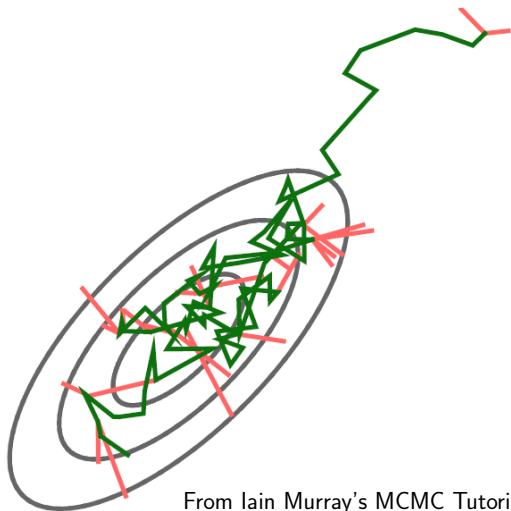
2. If

$$\frac{q(\mathbf{x}^{(t)}|\mathbf{x}')\tilde{p}(\mathbf{x}')}{q(\mathbf{x}'|\mathbf{x}^{(t)})\tilde{p}(\mathbf{x}^{(t)})} \geq u, \quad u \sim U[0, 1]$$

accept the sample $\mathbf{x}^{(t+1)} = \mathbf{x}'$. Otherwise set $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)}$.

- ▶ If proposal distribution is symmetric: [Metropolis Algorithm](#) (Metropolis et al., 1953); Otherwise [Metropolis-Hastings Algorithm](#) (Hastings, 1970)

Example



Step-Size Demo

- ▶ Explore $p(x) = \mathcal{N}(x | 0, 1)$ for different step sizes σ .
- ▶ We can only evaluate $\log \tilde{p}(x) = -x^2/2$
- ▶ Proposal distribution q : Gaussian $\mathcal{N}(x^{(t+1)} | x^{(t)}, \sigma^2)$ centered at the current state for various step sizes σ
- ▶ Expect to explore the space between $-2, 2$.

Step-Size Demo: Discussion

- ▶ Acceptance rate depends on the step size of the proposal distribution
 - ▶▶ Exploration parameter

Step-Size Demo: Discussion

- ▶ Acceptance rate depends on the step size of the proposal distribution
 - ▶▶ Exploration parameter
- ▶ If we do not reject enough, the method does not work.

Step-Size Demo: Discussion

- ▶ Acceptance rate depends on the step size of the proposal distribution
 - ▶▶ Exploration parameter
- ▶ If we do not reject enough, the method does not work.
- ▶ In rejection sampling we do not like rejections, but in MH rejections tell you where the target distribution is.

Step-Size Demo: Discussion

- ▶ Acceptance rate depends on the step size of the proposal distribution
 - ▶▶ Exploration parameter
- ▶ If we do not reject enough, the method does not work.
- ▶ In rejection sampling we do not like rejections, but in MH rejections tell you where the target distribution is.
- ▶ Theoretical results: in 1D 44%, in higher dimensions about 25% acceptance rate for good mixing properties

Step-Size Demo: Discussion

- ▶ Acceptance rate depends on the step size of the proposal distribution
 - ▶▶ Exploration parameter
- ▶ If we do not reject enough, the method does not work.
- ▶ In rejection sampling we do not like rejections, but in MH rejections tell you where the target distribution is.
- ▶ Theoretical results: in 1D 44%, in higher dimensions about 25% acceptance rate for good mixing properties
- ▶ Tune the step size

Properties

- ▶ Unlike rejection sampling, the previous sample is used to reset the chain (if a sample was discarded)
- ▶ If $q > 0$, we will end up in the equilibrium distribution:
$$p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p^*(\mathbf{x})$$
- ▶ Explore the state space by random walk
 - ▶▶ May take a while in high dimensions
- ▶ No further catastrophic problems in high dimensions

Gibbs Sampling

Gibbs Sampling

- ▶ Assumption: $p(\mathbf{x})$ is too complicated to draw samples from directly, but its conditionals $p(x_i | \mathbf{x}_{\setminus i})$ are tractable to work with
- ▶ Example:

$$y_i \sim \mathcal{N}(\mu, \tau^{-1}), \quad \mu \sim \mathcal{N}(0, 1), \quad \tau \sim \text{Gamma}(2, 1)$$

Gibbs Sampling

- ▶ Assumption: $p(\mathbf{x})$ is too complicated to draw samples from directly, but its conditionals $p(x_i|x_{\setminus i})$ are tractable to work with
- ▶ Example:

$$y_i \sim \mathcal{N}(\mu, \tau^{-1}), \quad \mu \sim \mathcal{N}(0, 1), \quad \tau \sim \text{Gamma}(2, 1)$$

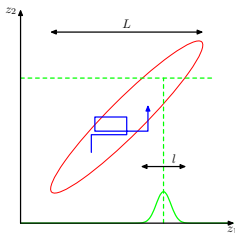
Then

$$\begin{aligned} p(\mathbf{y}, \mu, \tau) &= \prod_{i=1}^n p(y_i | \mu, \tau) p(\mu) p(\tau) \\ &\propto \tau^{n/2} \exp\left(-\frac{\tau}{2} \sum_i (y_i - \mu)^2\right) \exp\left(-\frac{1}{2}\mu^2\right) \tau \exp(-\tau) \end{aligned}$$

$$p(\mu | \tau) = \mathcal{N}\left(\frac{\tau \sum_i y_i}{1 + n\tau}, (1 + n\tau)^{-1}\right)$$

$$p(\tau | \mu) = \text{Gamma}\left(2 + \frac{n}{2}, 1 + \frac{1}{2} \sum (y_i - \mu)^2\right)$$

Algorithm

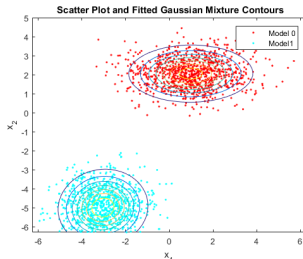


From PRML (Bishop, 2006)

Assuming n parameters x_1, \dots, x_n , Gibbs sampling samples individual variables conditioned on all others:

1. $x_1^{(t+1)} \sim p(x_1 | x_2^{(t)}, \dots, x_n^{(t)})$
2. $x_2^{(t+1)} \sim p(x_2 | x_1^{(t+1)}, x_3^{(t)}, \dots, x_n^{(t)})$
3. \vdots
4. $x_n^{(t+1)} \sim p(x_n | x_1^{(t+1)}, \dots, x_{n-1}^{(t+1)})$

Gibbs Sampling: Ergodicity



- ▶ $p(x)$ is invariant
- ▶ **Ergodicity**: Sufficient to show that all conditionals are greater than 0.
 - ▶▶ Then any point in x -space can be reached from any other point (potentially with low probability) in a finite number of steps involving one update of each of the component variables.

Properties

- ▶ Gibbs is Metropolis-Hastings with acceptance probability 1:
Sequence of proposal distributions q is defined in terms of conditional distributions of the joint $p(\mathbf{x})$
 - ▶▶ Converge to equilibrium distribution: $p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p(\mathbf{x})$
 - ▶▶ Exploration by random walk behavior can be slow

¹<http://mc-stan.org/>

²<http://www.mrc-bsu.cam.ac.uk/software/bugs/>

³<http://mcmc-jags.sourceforge.net/>

Properties

- ▶ Gibbs is Metropolis-Hastings with acceptance probability 1: Sequence of proposal distributions q is defined in terms of conditional distributions of the joint $p(\mathbf{x})$
 - ▶ **Converge** to equilibrium distribution: $p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p(\mathbf{x})$
 - ▶ Exploration by random walk behavior can be slow
- ▶ **No adjustable parameters** (e.g., step size)

¹<http://mc-stan.org/>

²<http://www.mrc-bsu.cam.ac.uk/software/bugs/>

³<http://mcmc-jags.sourceforge.net/>

Properties

- ▶ Gibbs is Metropolis-Hastings with acceptance probability 1:
Sequence of proposal distributions q is defined in terms of conditional distributions of the joint $p(\mathbf{x})$
 - ▶ **Converge** to equilibrium distribution: $p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p(\mathbf{x})$
 - ▶ Exploration by random walk behavior can be slow
- ▶ **No adjustable parameters** (e.g., step size)
- ▶ Applicability depends on how easy it is to draw samples from the conditionals

¹<http://mc-stan.org/>

²<http://www.mrc-bsu.cam.ac.uk/software/bugs/>

³<http://mcmc-jags.sourceforge.net/>

Properties

- ▶ Gibbs is Metropolis-Hastings with acceptance probability 1:
Sequence of proposal distributions q is defined in terms of conditional distributions of the joint $p(\mathbf{x})$
 - ▶ **Converge** to equilibrium distribution: $p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p(\mathbf{x})$
 - ▶ Exploration by random walk behavior can be slow
- ▶ **No adjustable parameters** (e.g., step size)
- ▶ Applicability depends on how easy it is to draw samples from the conditionals
- ▶ May not work well if the **variables are correlated**

¹<http://mc-stan.org/>

²<http://www.mrc-bsu.cam.ac.uk/software/bugs/>

³<http://mcmc-jags.sourceforge.net/>

Properties

- ▶ Gibbs is Metropolis-Hastings with acceptance probability 1:
Sequence of proposal distributions q is defined in terms of conditional distributions of the joint $p(\mathbf{x})$
 - ▶ **Converge** to equilibrium distribution: $p^{(t)}(\mathbf{x}) \xrightarrow{t \rightarrow \infty} p(\mathbf{x})$
 - ▶ Exploration by random walk behavior can be slow
- ▶ **No adjustable parameters** (e.g., step size)
- ▶ Applicability depends on how easy it is to draw samples from the conditionals
- ▶ May not work well if the **variables are correlated**
- ▶ **Statistical software** derives the conditionals of the model, and it works out how to do the updates: STAN¹, WinBUGS², JAGS³

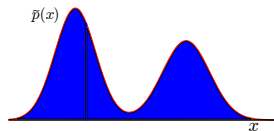
¹<http://mc-stan.org/>

²<http://www.mrc-bsu.cam.ac.uk/software/bugs/>

³<http://mcmc-jags.sourceforge.net/>

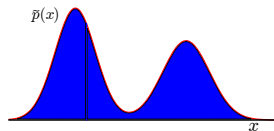
Slice Sampling

Key Idea behind Slice Sampling



- ▶ **Idea:** Sample point (random walk) uniformly under the curve $\tilde{p}(x)$

Key Idea behind Slice Sampling

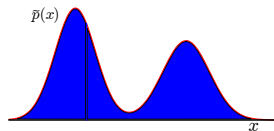


- ▶ **Idea:** Sample point (random walk) uniformly under the curve $\tilde{p}(x)$

- ▶ Introduce additional variable u , define joint $\hat{p}(x, u)$:

$$\hat{p}(x, u) = \begin{cases} 1/Z_p & \text{if } 0 \leq u \leq \tilde{p}(x) \\ 0 & \text{otherwise} \end{cases}, \quad Z_p = \int \tilde{p}(x) dx$$

Key Idea behind Slice Sampling



- ▶ **Idea:** Sample point (random walk) uniformly under the curve $\tilde{p}(x)$

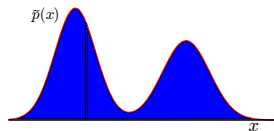
- ▶ Introduce additional variable u , define joint $\hat{p}(x, u)$:

$$\hat{p}(x, u) = \begin{cases} 1/Z_p & \text{if } 0 \leq u \leq \tilde{p}(x) \\ 0 & \text{otherwise} \end{cases}, \quad Z_p = \int \tilde{p}(x) dx$$

- ▶ The marginal distribution over x is then

$$\int \hat{p}(x, u) du = \int_0^{\tilde{p}(x)} 1/Z_p du = \tilde{p}(x)/Z_p = p(x)$$

Key Idea behind Slice Sampling



- ▶ **Idea:** Sample point (random walk) uniformly under the curve $\tilde{p}(x)$

- ▶ Introduce additional variable u , define joint $\hat{p}(x, u)$:

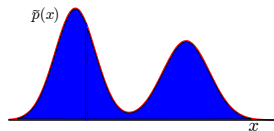
$$\hat{p}(x, u) = \begin{cases} 1/Z_p & \text{if } 0 \leq u \leq \tilde{p}(x) \\ 0 & \text{otherwise} \end{cases}, \quad Z_p = \int \tilde{p}(x) dx$$

- ▶ The marginal distribution over x is then

$$\int \hat{p}(x, u) du = \int_0^{\tilde{p}(x)} 1/Z_p du = \tilde{p}(x)/Z_p = p(x)$$

- ▶▶ Obtain samples from unknown $p(x)$ by sampling from $\hat{p}(x, u)$ and then ignore u values

Key Idea behind Slice Sampling



- ▶ **Idea:** Sample point (random walk) uniformly under the curve $\tilde{p}(x)$

- ▶ Introduce additional variable u , define joint $\hat{p}(x, u)$:

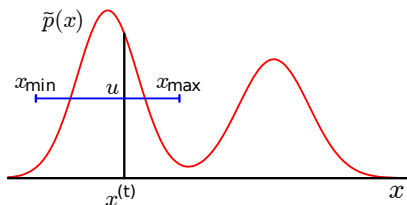
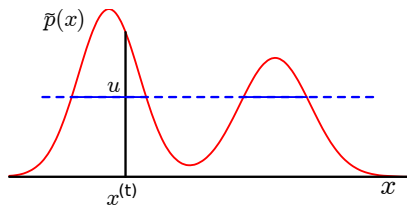
$$\hat{p}(x, u) = \begin{cases} 1/Z_p & \text{if } 0 \leq u \leq \tilde{p}(x) \\ 0 & \text{otherwise} \end{cases}, \quad Z_p = \int \tilde{p}(x) dx$$

- ▶ The marginal distribution over x is then

$$\int \hat{p}(x, u) du = \int_0^{\tilde{p}(x)} 1/Z_p du = \tilde{p}(x)/Z_p = p(x)$$

- ▶ Obtain samples from unknown $p(x)$ by sampling from $\hat{p}(x, u)$ and then ignore u values
- ▶ Gibbs sampling: **Update one variable at a time**

Slice Sampling Algorithm

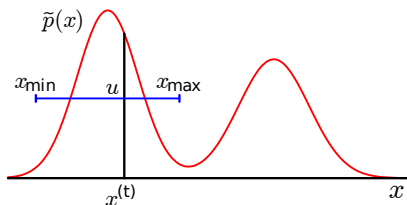
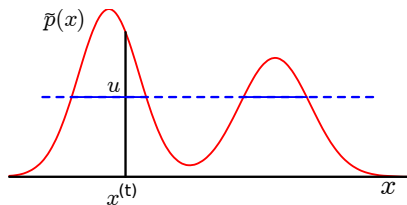


Adapted from PRML (Bishop, 2006)

► Repeat the following steps:

1. Draw $u|x^{(t)} \sim \mathcal{U}[0, \tilde{p}(x)]$
2. Draw $x^{(t+1)}|u \sim \mathcal{U}[\{x : \tilde{p}(x) > u\}]$ ►► slice

Slice Sampling Algorithm



Adapted from PRML (Bishop, 2006)

- ▶ Repeat the following steps:
 1. Draw $u|x^{(t)} \sim \mathcal{U}[0, \tilde{p}(x)]$
 2. Draw $x^{(t+1)}|u \sim \mathcal{U}[\{x : \tilde{p}(x) > u\}]$ ▶▶ slice
- ▶ In practice, we sample $x^{(t+1)}|u$ uniformly from an interval $[x_{\min}, x_{\max}]$ around $x^{(t)}$.
- ▶ The interval is found adaptively (see Neal (2003) for details)

Relation to other Sampling Methods

Similar to:

- ▶ **Metropolis:** Just need to be able to evaluate $\tilde{p}(x)$
More robust to the choice of parameters (e.g., step size is automatically adapted)
- ▶ **Gibbs:** 1-dimensional transitions in state space
No longer required that we can easily sample from 1-D conditionals
- ▶ **Rejection:** Asymptotically draw samples from the volume under the curve described by \tilde{p}
No upper-bounding of \tilde{p} required

Properties

- ▶ Slice sampling can be applied to multivariate distributions by repeatedly sampling each variable in turn (similar to Gibbs sampling). ▶▶ See (Neal, 2003; Murray et al., 2010) for more details
- ▶ This requires to compute a function that is proportional to $p(x_i | \mathbf{x}_{\setminus i})$ for all variables x_i .

Properties

- ▶ Slice sampling can be applied to multivariate distributions by repeatedly sampling each variable in turn (similar to Gibbs sampling). ▶ See (Neal, 2003; Murray et al., 2010) for more details
- ▶ This requires to compute a function that is proportional to $p(x_i | \mathbf{x}_{\setminus i})$ for all variables x_i .
- ▶ No rejections
- ▶ Adaptive step sizes
- ▶ Easy to implement
- ▶ Broadly applicable

Discussion MCMC

- ▶ **Asymptotic guarantee to converge** to the equilibrium distribution for any kind of model
- ▶ **General-purpose method** to draw samples in any kind of probabilistic model ▶▶ **Probabilistic Programming**
- ▶ **Convergence difficult to assess**
- ▶ **Long chains required in high dimensions**
- ▶ **Choice of proposal distribution is hard**
- ▶ **Need to store all samples** (subsequent computations require to work with these samples)

References I

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
- [2] D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, The Edinburgh Building, Cambridge CB2 2RU, UK, 2003.