# Teaching Bayes: The Essential Parts
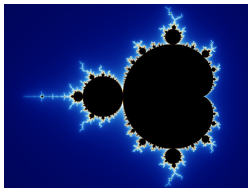
Rebecca C. Steorts
ISBA World Meeting 2016

Lecture 2: Intro to Gibbs Sampling

# Intro to Markov chain Monte Carlo (MCMC)

Goal: sample from $f(x)$, or approximate $E_f[h(X)]$.

Function $f(x)$ is very complicated and hard to sample from.

How to deal with this?

1. What's a simple way?
2. What are two other ways?
3. What happens in high dimensions?

# High dimensional spaces

- In low dimensions, importance and rejection sampling work pretty well.
- But in high dimensions, a proposal $g(x)$ that worked in 2-D, often doesn't mean that it will work in any dimension.
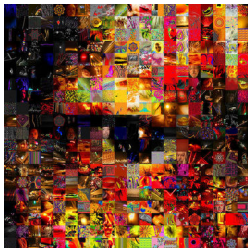- Why? It's hard to capture high dimensional spaces!



Figure 1: A high dimensional space (many images).

We turn to Markov chain Monte Carlo (MCMC).

# Intution

Imagine that we have a complicated function $f$ below and it's high probability regions are represented in green.
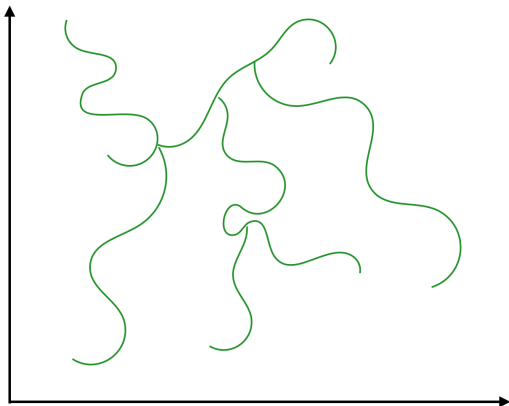


Figure 2: Example of a Markov chain
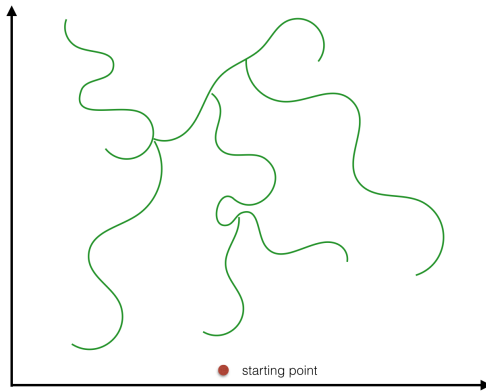
# Intution



Figure 3: Example of a Markov chain and red starting point
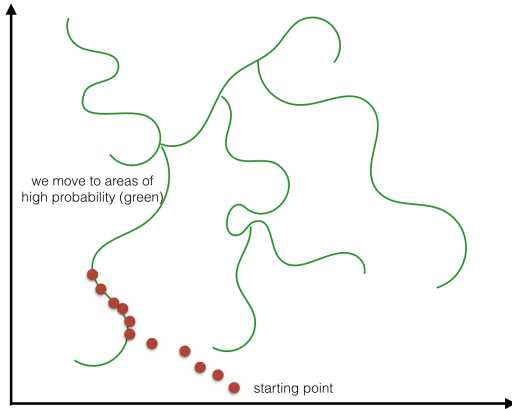
# Intution



Figure 4: Example of a Markov chain and moving from the starting point to a high probability region.

# What is Markov Chain Monte Carlo

- Markov Chain – where we go next only depends on our last state (the Markov property).
- Monte Carlo – just simulating data.

# Why MCMC?

(a) the region of high probability tends to be "connected"
  ▶ That is, we can get from one point to another without going through a low-probability region, and
(b) we tend to be interested in the expectations of functions that are relatively smooth and have lots of "symmetries"
  ▶ That is, one only needs to evaluate them at a small number of representative points in order to get the general picture.

# Advantages/Disadvantages of MCMC:

Advantages:

- ▶ applicable even when we can't directly draw samples
- ▶ works for complicated distributions in high-dimensional spaces, even when we don't know where the regions of high probability are
- ▶ relatively easy to implement
- ▶ fairly reliable

Disadvantages:

- ▶ slower than simple Monte Carlo or importance sampling (i.e., requires more samples for the same level of accuracy)
- ▶ can be very difficult to assess accuracy and evaluate convergence, even empirically

# Two-stage Gibbs sampler

- Suppose $p(x, y)$ is a p.d.f. or p.m.f. that is difficult to sample from directly.
- Suppose, though, that we *can* easily sample from the conditional distributions $p(x|y)$ and $p(y|x)$.
- The Gibbs sampler proceeds as follows:
    1. set $x$ and $y$ to some initial starting values
    2. then sample $x|y$, then sample $y|x$, then $x|y$, and so on.

# Two-stage Gibbs sampler

0. Set $(x_0, y_0)$ to some starting value.

1. Sample $x_1 \sim p(x|y_0)$, that is, from the conditional distribution $X \mid Y = y_0$.
   Current state: $(x_1, y_0)$
   Sample $y_1 \sim p(y|x_1)$, that is, from the conditional distribution $Y \mid X = x_1$.
   Current state: $(x_1, y_1)$

2. Sample $x_2 \sim p(x|y_1)$, that is, from the conditional distribution $X \mid Y = y_1$.
   Current state: $(x_2, y_1)$
   Sample $y_2 \sim p(y|x_2)$, that is, from the conditional distribution $Y \mid X = x_2$.
   Current state: $(x_2, y_2)$
   $\vdots$

Repeat iterations 1 and 2, M times.

This procedure defines a sequence of pairs of random variables

$$(X_0, Y_0), (X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \ldots$$

# Markov chain and dependence

$$(X_0, Y_0), (X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \ldots$$

satisfies the property of being a Markov chain.

The conditional distribution of $(X_i, Y_i)$ given all of the previous pairs depends only on $(X_{i-1}, Y_{i-1})$

$(X_0, Y_0), (X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \ldots$ are not iid samples (Think about why).

# Ideal Properties of MCMC

- $(x_0, y_0)$ chosen to be in a region of high probability under $p(x, y)$, but often this is not so easy.
- We run the chain for M iterations and discard the first $B$ samples $(X_1, Y_1), \ldots, (X_B, Y_B)$. This is called *burn-in*.
- Typically: if you run the chain long enough, the choice of $B$ doesn't matter.
- Roughly speaking, the performance of an MCMC algorithm—that is, how quickly the sample averages $\frac{1}{N} \sum_{i=1}^{N} h(X_i, Y_i)$ converge—is referred to as the *mixing rate*.
- An algorithm with good performance is said to "have good mixing", or "mix well".

## Toy Example

Suppose we want to sample from the bivariate distribution:

$$p(x, y) \propto e^{-xy} \mathbb{1}(x, y \in (0, c))$$

where $c > 0$, and $(0, c)$ denotes the (open) interval between $0$ and $c$. (This example is due to Casella & George, 1992.)

# Toy Example

- The Gibbs sampling approach is to alternately sample from $p(x|y)$ and $p(y|x)$.
- Note $p(x, y)$ is symmetric with respect to $x$ and $y$.
- Hence, only need to derive one of these and then we can get the other one by just swapping $x$ and $y$.
- Let's look at $p(x|y)$.

# Toy Example

$$p(x, y) \propto e^{-xy} \mathbb{1}(x, y \in (0, c))$$

$$p(x|y) \underset{x}{\propto} p(x, y) \underset{x}{\propto} e^{-xy} \mathbb{1}(0 < x < c) \underset{x}{\propto} \text{Exp}(x|y) \mathbb{1}(x < c).^1$$

▶ $p(x|y)$ is a *truncated* version of the $\text{Exp}(y)$ distribution
▶ It is the same as taking $X \sim \text{Exp}(y)$ and conditioning on it being less than $c$, i.e., $X \mid X < c$.
▶ Let's refer to this as the $\text{TExp}(y, (0, c))$ distribution.

---

[1]Under $\propto$, we write the random variable $(x)$ for clarity.

# Toy Example

An easy way to generate a sample from $Z \sim \text{TExp}(\theta, (0, c))$, is:

1. Sample $U \sim \text{Uniform}(0, F(c|\theta))$ where

$$F(x|\theta) = 1 - e^{-\theta x}$$

   is the $\text{Exp}(\theta)$ c.d.f.

2. Set $Z = F^{-1}(U|\theta)$ where

$$F^{-1}(u|\theta) = -(1/\theta) \log(1 - u)$$

   is the inverse c.d.f. for $u \in (0, 1)$.

Verify the last step on your own.

Let's apply Gibbs sampling, denoting $S = (0, c)$.

0. Initialize $x_0, y_0 \in S$.
1. Sample $x_1 \sim \text{TExp}(y_0, S)$, then sample $y_1 \sim \text{TExp}(x_1, S)$.
2. Sample $x_2 \sim \text{TExp}(y_1, S)$, then sample $y_2 \sim \text{TExp}(x_2, S)$.
   ⋮
$N$. Sample $x_N \sim \text{TExp}(y_{N-1}, S)$, sample $y_N \sim \text{TExp}(x_N, S)$.

Figure 5 demonstrates the algorithm, with $c = 2$ and initial point $(x_0, y_0) = (1, 1)$.
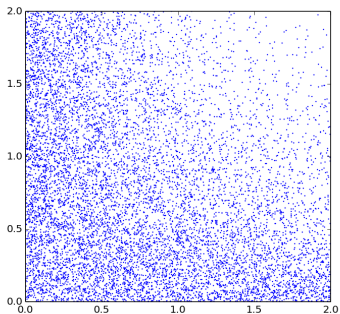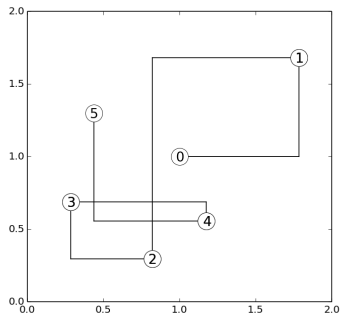
Figure 5: (Left) Schematic representation of the first 5 Gibbs sampling iterations/sweeps/scans. (Right) Scatterplot of samples from $10^4$ Gibbs sampling iterations.

# Pareto example

Distributions of sizes and frequencies often tend to follow a "power law" distribution.

- ▶ wealth of individuals
- ▶ size of oil reserves
- ▶ size of cities
- ▶ word frequency
- ▶ returns on stocks

# Power Law Distribution

The Pareto distribution with shape $\alpha > 0$ and scale $c > 0$ has p.d.f.

$$\mathrm{Pareto}(x|\alpha, c) = \frac{\alpha c^{\alpha}}{x^{\alpha+1}} \mathbb{1}(x > c) \propto \frac{1}{x^{\alpha+1}} \mathbb{1}(x > c).$$

This is referred to as a power law distribution, because the p.d.f. is proportional to $x$ raised to a power. Notice that $c$ is a lower bound on the observed values. In this example, we'll see how Gibbs sampling can be used to perform inference for $\alpha$ and $c$.

| Rank | City | Population |
|------|------|------------|
| 1 | Charlotte | 731424 |
| 2 | Raleigh | 403892 |
| 3 | Greensboro | 269666 |
| 4 | Durham | 228330 |
| 5 | Winston-Salem | 229618 |
| 6 | Fayetteville | 200564 |
| 7 | Cary | 135234 |
| 8 | Wilmington | 106476 |
| 9 | High Point | 104371 |
| 10 | Greenville | 84554 |
| 11 | Asheville | 85712 |
| 12 | Concord | 79066 |
| ⋮ | ⋮ | ⋮ |
| 44 | Havelock | 20735 |
| 45 | Carrboro | 19582 |
| 46 | Shelby | 20323 |
| 47 | Clemmons | 18627 |
| 48 | Lexington | 18931 |
| 49 | Elizabeth City | 18683 |
| 50 | Boone | 17122 |

# Parameter Interpretations

- $\alpha$ tells us the scaling relationship between the size of cities and their probability of occurring.
  - Let $\alpha = 1$.
  - Density looks like $1/x^{\alpha+1} = 1/x^2$.
  - Cities with 10,000–20,000 inhabitants occur roughly $10^{\alpha+1} = 100$ times as frequently as cities with 100,000–110,000 inhabitants.
- $c$ represents the cutoff point—any cities smaller than this were not included in the dataset.

To keep things as simple as possible, let's use an (improper) default prior:

$$p(\alpha, c) \propto \mathbb{1}(\alpha, c > 0).$$

Recall from Module 4:

- An *improper/default prior* is a nonnegative function of the parameters which integrates to infinity.
- Often (but not always!) the resulting "posterior" will be proper.
- It is important that the "posterior" be proper, since otherwise the whole Bayesian framework breaks down.

Recall

$$p(x|\alpha, c) = \frac{\alpha c^{\alpha}}{x^{\alpha+1}} \mathbb{1}(x > c) \tag{1}$$

$$\mathbb{1}(\alpha, c > 0) \tag{2}$$

Let's derive the posterior:

$$p(\alpha, c|x_{1:n}) \underset{\alpha,c}{\overset{\text{def}}{\propto}} p(x_{1:n}|\alpha, c)p(\alpha, c)$$

$$\underset{\alpha,c}{\propto} \mathbb{1}(\alpha, c > 0) \prod_{i=1}^{n} \frac{\alpha c^{\alpha}}{x_i^{\alpha+1}} \mathbb{1}(x_i > c)$$

$$= \frac{\alpha^n c^{n\alpha}}{(\prod x_i)^{\alpha+1}} \mathbb{1}(c < x_*) \mathbb{1}(\alpha, c > 0) \tag{3}$$

where $x_* = \min\{x_1, \ldots, x_n\}$.

As a joint distribution on $(\alpha, c)$,

- ▶ this does not seem to have a recognizable form,
- ▶ and it is not clear how we might sample from it directly.

Let's try Gibbs sampling!

To use Gibbs, we need to be able to sample $\alpha | c, x_{1:n}$ and $c | \alpha, x_{1:n}$.

By Equation 3, we find that

$$\begin{aligned}
p(\alpha | c, x_{1:n}) &\underset{\alpha}{\propto} p(\alpha, c | x_{1:n}) \underset{\alpha}{\propto} \frac{\alpha^n c^{n\alpha}}{(\prod x_i)^\alpha} \mathbb{1}(\alpha > 0) \\
&= \alpha^n \exp\left(-\alpha(\sum \log x_i - n \log c)\right) \mathbb{1}(\alpha > 0) \\
&\underset{\alpha}{\propto} \text{Gamma}\left(\alpha \mid n + 1, \sum \log x_i - n \log c\right),
\end{aligned}$$

and

$$p(c | \alpha, x_{1:n}) \underset{c}{\propto} p(\alpha, c | x_{1:n}) \underset{c}{\propto} c^{n\alpha} \mathbb{1}(0 < c < x_*),$$

which we will define to be $\text{Mono}(\alpha, x_*)$

# Defining the Mono distribution

For $a > 0$ and $b > 0$, define the distribution $\mathrm{Mono}(a, b)$ (for monomial) with p.d.f.

$$\mathrm{Mono}(x|a, b) \propto x^{a-1}\mathbb{1}(0 < x < b).$$

Since $\int_0^b x^{a-1}dx = b^a/a$, we have

$$\mathrm{Mono}(x|a, b) = \frac{a}{b^a}x^{a-1}\mathbb{1}(0 < x < b),$$

and for $0 < x < b$, the c.d.f. is

$$F(x|a, b) = \int_0^x \mathrm{Mono}(y|a, b)dy = \frac{a}{b^a}\frac{x^a}{a} = \frac{x^a}{b^a}.$$

To use the inverse c.d.f. technique, we solve for the inverse of $F$ on $0 < x < b$: Let $u = \frac{x^a}{b^a}$ and solve for $x$.

$$u = \frac{x^a}{b^a} \tag{4}$$

$$b^a u = x^a \tag{5}$$

$$b u^{1/a} = x \tag{6}$$

Can sample from $\mathrm{Mono}(a, b)$ by drawing $U \sim \mathrm{Uniform}(0, 1)$ and setting $X = bU^{1/a}$.[2]

---

[2]It turns out that this is an inverse of the Pareto distribution, in the sense that if $X \sim \mathrm{Pareto}(\alpha, c)$ then $1/X \sim \mathrm{Mono}(\alpha, 1/c)$.

So, in order to use the Gibbs sampling algorithm to sample from the posterior $p(\alpha, c | x_{1:n})$, we initialize $\alpha$ and $c$, and then alternately update them by sampling:

$$\alpha | c, x_{1:n} \sim \text{Gamma}\left(n+1, \sum \log x_i - n \log c\right)$$
$$c | \alpha, x_{1:n} \sim \text{Mono}(n\alpha + 1, x_*).$$

# Ways of visualizing results

**Traceplots**. A traceplot simply shows the sequence of samples, for instance $\alpha_1, \ldots, \alpha_N$, or $c_1, \ldots, c_N$. Traceplots are a simple but very useful way to visualize how the sampler is behaving.
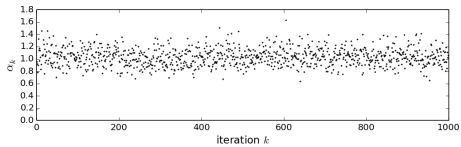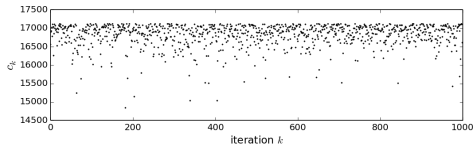
Figure 6: Traceplot of $\alpha$



Figure 7: Traceplot of c.

**Estimated density**. We are primarily interested in the posterior on $\alpha$, since it tells us the scaling relationship between the size of cities and their probability of occurring.

By making a histogram of the samples $\alpha_1, \ldots, \alpha_N$, we can estimate the posterior density $p(\alpha|x_{1:n})$.

The two vertical lines indicate the lower $\ell$ and upper $u$ boundaries of an (approximate) 90% credible interval $[\ell, u]$—that is, an interval that contains 90% of the posterior probability:

$$\mathbb{P}\big(\boldsymbol{\alpha} \in [\ell, u]\big|x_{1:n}\big) = 0.9.$$
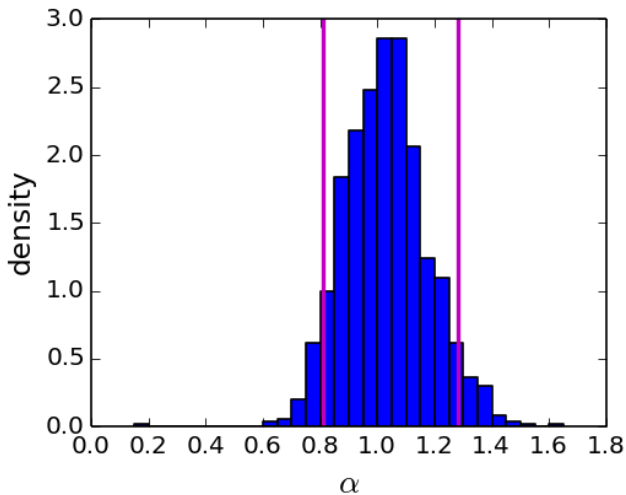
Figure 8: Estimated density of $\alpha|x_{1:n}$ with $\approx 90$ percent credible intervals.

**Running averages**. Panel (d) shows the running average $\frac{1}{k} \sum_{i=1}^{k} \alpha_i$ for $k = 1, \ldots, N$.

In addition to traceplots, running averages such as this are a useful heuristic for visually assessing the convergence of the Markov chain.

The running average shown in this example still seems to be meandering about a bit, suggesting that the sampler needs to be run longer (but this would depend on the level of accuracy desired).
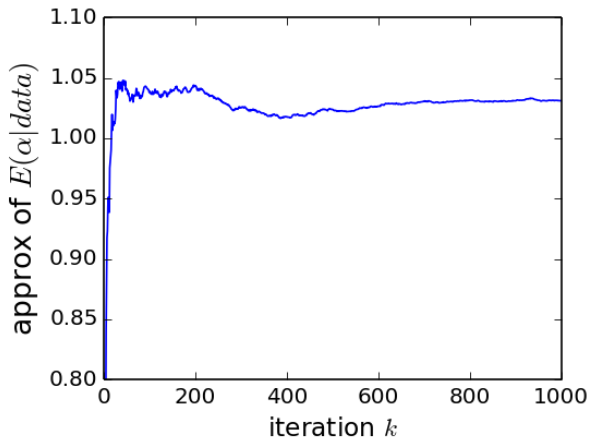
Figure 9: Running average plot

# Survival function

A survival function is defined to be

$$S(x) = \mathbb{P}(X > x) = 1 - \mathbb{P}(X \le x).$$

Power law distributions are often displayed by plotting their survival function $S(x)$, on a log-log plot.

Why? $S(x) = (c/x)^\alpha$ for the $\mathrm{Pareto}(\alpha, c)$ distribution and on a log-log plot this appears as a line with slope $-\alpha$.

The posterior survival function (or more precisely, the posterior predictive survival function), is $S(x|x_{1:n}) = \mathbb{P}(X_{n+1} > x \mid x_{1:n})$.

Figure 10(e) shows an empirical estimate of the survival function (based on the empirical c.d.f., $\hat{F}(x) = \frac{1}{n}\sum_{i=1}^{n} \mathbb{1}(x \geq x_i)$) along with the posterior survival function, approximated by

$$S(x|x_{1:n}) = \mathbb{P}(X_{n+1} > x \mid x_{1:n}) \tag{7}$$

$$= \int \mathbb{P}(X_{n+1} > x \mid \alpha, c) p(\alpha, c|x_{1:n}) d\alpha dc \tag{8}$$

$$\approx \frac{1}{N}\sum_{i=1}^{N} \mathbb{P}(X_{n+1} > x \mid \alpha_i, c_i) = \frac{1}{N}\sum_{i=1}^{N} (c_i/x)^{\alpha_i}. \tag{9}$$

This is computed for each $x$ in a grid of values.
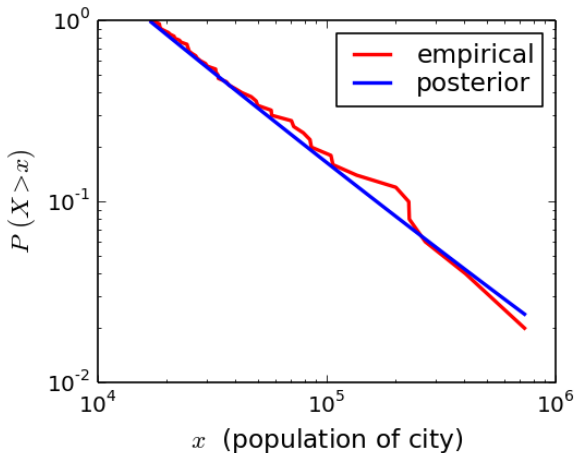
[Think about why each line is true on your own].

Figure 10: Empirical vs posterior survival function

# Questions you should be able to answer!

- When should we use MCMC in a Bayesian setting?
- When would we use an MCMC over Importance sampling and Rejection sampling?
- What is a Gibbs sampler?
- What are simple diagnostics of MCMC?
- Are we guaranteed convergence of the Markov chain emprically?
- What do are diagnostics really tell us?