Lecture 9: Markov Chain Monte Carlo

Jacob Steinhardt

September 23, 2021

- HW2 is due next week
 - Intro to HW OH on Friday; HW party and extra OH next week
- Limit one slip day for HW2, so that we can release solutions before the exam
- Midterm 1 coming up on 10/5
 - Covers everything through Lecture 12 (next Thurs), Lab 4 (next week), HW2, and Discussion 5 (next week)
 - Will release more info + practice tests soon



- Rejection sampling
- Markov chain review



- Rejection sampling
- Markov chain review

This time: Markov chain Monte Carlo

- Gibbs sampling
- Metropolis-Hastings

Running example: person *i*, gender $z_i \in \{0, 1\}$, height $x_i \in \mathbb{R}$, hyperparameters $\theta = (\alpha, \beta, \pi)$

Running example: person *i*, gender $z_i \in \{0, 1\}$, height $x_i \in \mathbb{R}$, hyperparameters $\theta = (\alpha, \beta, \pi)$

How to do inference in this model?

- Method 1: place prior on θ , sample $p(\theta, z \mid x)$
- Method 2: maximize $\log p(x \mid \theta) = \log (\sum_{z} p(x, z \mid \theta))$
 - "half-Bayesian" (not this class)

Running example: person *i*, gender $z_i \in \{0, 1\}$, height $x_i \in \mathbb{R}$, hyperparameters $\theta = (\alpha, \beta, \pi)$

How to do inference in this model?

- Method 1: place prior on θ , sample $p(\theta, z \mid x)$
- Method 2: maximize $\log p(x \mid \theta) = \log (\sum_{z} p(x, z \mid \theta))$
 - "half-Bayesian" (not this class)

How many possibilities for *z*? Height/gender example:

Running example: person *i*, gender $z_i \in \{0, 1\}$, height $x_i \in \mathbb{R}$, hyperparameters $\theta = (\alpha, \beta, \pi)$

How to do inference in this model?

- Method 1: place prior on θ , sample $p(\theta, z \mid x)$
- Method 2: maximize $\log p(x \mid \theta) = \log (\sum_z p(x, z \mid \theta))$
 - "half-Bayesian" (not this class)

How many possibilities for *z*? Height/gender example:

- 100 people, genders *z*₁,...,*z*₁₀₀
- $2^{100} \approx 10^{30}$ possibilities

Running example: person *i*, gender $z_i \in \{0, 1\}$, height $x_i \in \mathbb{R}$, hyperparameters $\theta = (\alpha, \beta, \pi)$

How to do inference in this model?

- Method 1: place prior on θ , sample $p(\theta, z \mid x)$
- Method 2: maximize $\log p(x \mid \theta) = \log (\sum_{z} p(x, z \mid \theta))$
 - "half-Bayesian" (not this class)

How many possibilities for *z*? Height/gender example:

- 100 people, genders *z*₁,...,*z*₁₀₀
- $2^{100} \approx 10^{30}$ possibilities

This is why we need good sampling algorithms! (General term: "approximate inference") • Have an arbitrary distribution $p(x_1, \ldots, x_n)$ that we want to sample from

- Have an arbitrary distribution $p(x_1, ..., x_n)$ that we want to sample from
- Current tool: rejection sampling
 - Proposal distribution $q(x_1, \ldots, x_n)$ for all x_i at once
 - Issue: too slow (typically exponentially small acceptance rate in n)
 - E.g. even if x_i are independent, and $q(x_i)/p(x_i) \le 1.1$, need 1.1^n tries ($\approx 2.5 \cdot 10^{41}$ for n = 1000)

- Have an arbitrary distribution $p(x_1, ..., x_n)$ that we want to sample from
- Current tool: rejection sampling
 - Proposal distribution $q(x_1, \ldots, x_n)$ for all x_i at once
 - Issue: too slow (typically exponentially small acceptance rate in n)
 - E.g. even if x_i are independent, and $q(x_i)/p(x_i) \le 1.1$, need 1.1^{*n*} tries ($\approx 2.5 \cdot 10^{41}$ for n = 1000)
- Idea behind Gibbs sampling: change one variable at a time (Markov chain)

Algorithm:

- Initialize (x_1, \ldots, x_n) arbitrarily
- Repeat:
 - Pick *i* (randomly or sequentially)
 - Re-sample x_i from $p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ (often denote $p(x_i | x_{-i})$)

Algorithm:

- Initialize (x_1, \ldots, x_n) arbitrarily
- Repeat:
 - Pick *i* (randomly or sequentially)
 - Re-sample x_i from $p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ (often denote $p(x_i | x_{-i})$)

Defines a Markov chain, and can prove that the stationary distribution is $p(x_1, ..., x_n)$ (!!).



























Recall hierarchical models (e.g. height and gender example)



Recall hierarchical models (e.g. height and gender example)



Suppose we want to do Gibbs sampling for this model

Recall hierarchical models (e.g. height and gender example)



Suppose we want to do Gibbs sampling for this model

• Sample
$$z_i$$
: $p(z_i | x_i, \theta) \propto \underbrace{p(z_i | \theta)}_{\text{prior}} \underbrace{p(x_i | z_i)}_{\text{likelihood}}$

Recall hierarchical models (e.g. height and gender example)



Suppose we want to do Gibbs sampling for this model

• Sample
$$z_i$$
: $p(z_i | x_i, \theta) \propto \underbrace{p(z_i | \theta)}_{\text{prior}} \underbrace{p(x_i | z_i)}_{\text{likelihood}}$

• Sample θ (e.g. μ_0 for height/gender model):

$$p(\mu_0 \mid z_{1:n}, x_{1:n}) \propto \underbrace{p(\mu_0)}_{\text{prior}} \cdot \underbrace{\prod_{i:z_i=0} \exp(-(x_i - \mu_0)^2 / 2\sigma^2)}_{\text{likelihood}}$$

Assuming chain is ergodic, just need to show stationary distribution is preserved.

Suppose $x \sim p$ and x' is obtained from x by Gibbs sampling update. Want to show that x' is also distributed according to p.

If index *i* is updated, then $x' = (x_1, ..., x_{i-1}, x'_i, x_{i+1}, ...)$, where $x'_i \sim p(x_i \mid x_1, ..., x_{i-1}, x_{i+1}, ...)$.

Indices $\neq i$ distributed according to *p*, and $x'_i \mid x'_{-i}$ is as well, so x' follows *p*.

Suppose that $x_1, x_2 \in \{0, 1\}$ with following probability table:

	0	1
0	0.5	0.0
1	0.0	0.5

What will Gibbs sampling do?

- Repeatedly sample from $p(x_i | x_{-i})$
- Creates Markov chain whose stationary distribution is $p(x_1,...,x_n)$
- Flexible: conditional $p(x_i | x_{-i})$ one-dimensional, easy to sample from
- Don't need to "get lucky" with graphical model structure
- Extensions, e.g. block Gibbs sampling

• Gibbs sampling: one possible Markov chain

- Gibbs sampling: one possible Markov chain
- Is there a more general strategy?

- Gibbs sampling: one possible Markov chain
- Is there a more general strategy?
- Yes! Combine with idea of rejection sampling

- Gibbs sampling: one possible Markov chain
- Is there a more general strategy?
- Yes! Combine with idea of rejection sampling
- Given any "proposed Markov chain" $q(x^{\text{new}} | x^{\text{old}})$, will combine with an accept/reject step to create new Markov chain with the correct stationary distribution

Given *x*^{old}:

 x^{new})

- Sample *x*^{new} from *q*
- With probability

, accept (replace x^{old} with

Given *x*^{old}:

• Sample *x*^{new} from *q*



Given *x*^{old}:

- Sample *x*^{new} from *q*
- With probability x^{new})

 $rac{p(x^{
m new})}{p(x^{
m old})} rac{q(x^{
m old}|x^{
m new})}{q(x^{
m new}|x^{
m old})}$

, accept (replace x^{old} with

Given x^{old} :

• Sample *x*^{new} from *q*

• With probability
$$\min\left(1, \frac{p(x^{\text{new}})}{p(x^{\text{old}})} \frac{q(x^{\text{old}}|x^{\text{new}})}{q(x^{\text{new}}|x^{\text{old}})}\right)$$
, accept (replace x^{old} with x^{new})

Given x^{old} :

- Sample *x*^{new} from *q*
- With probability $\left[\min\left(1, \frac{p(x^{\text{new}})}{p(x^{\text{old}})} \frac{q(x^{\text{old}}|x^{\text{new}})}{q(x^{\text{new}}|x^{\text{old}})} \right) \right]$, accept (replace x^{old} with x^{new})
- Otherwise, reject (keep x^{old})

Gibbs sampling: special choice of q where we always accept!

Can show that if an ergodic Markov chain satisfies $\bar{p}(x)A(x' \mid x) = \bar{p}(x')A(x \mid x')$ for all x, x', then it has stationary distribution \bar{p} .

This condition is called **detailed balance**.

Metropolis-Hastings sets probabilities so that detailed balance holds.

Performance of MCMC algorithms governed by **mixing time**: how long it takes to get close to stationary distribution.

Mixing time can vary dramatically, from nearly linear to exponential in number of variables.

[mixing time examples: on board]