

Markov Chain Monte Carlo

Kevin Duh
Mokuyokai 12/10/2009

Problems addressed

- Integration/Expectation

$$E_p[f(x)] = \int_x f(x)p(x)dx$$

- Optimization

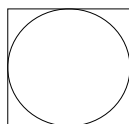
$$\arg \max_x p(x)$$

- Situation suitable for MCMC:

- Direct integration/optimization is hard:
 - X is a large space or f(x) is complex
- P(x) is easily computed (up to normalization constant)

Example: Integration

- Problem: Compute area of circle
 - Suppose we don't know the formula πr^2
- Monte Carlo approach:
 1. Bound the circle with a box
 2. Randomly (uniformly) spread seeds in the box
 3. Count the number of seeds inside the circle



3

Example: Optimization

- Problem: ASR decoding without Viterbi
 - $\arg \max p(x)$ where $p(x)$ is complex acoustic + language model
- Markov Chain Monte Carlo approach:
 1. Start with random hypothesis sentence x_0
 2. Randomly transform x_t into x_{t+1}
 3. If $p(x_t) > p(x_{t+1})$, let $x_{t+1} = x_t$
 4. Repeat steps 2 & 3 until convergence

my hi is name...
hi my is name...
hi is my name...
hi my name is...

4

What we'll cover

- Monte Carlo methods:
 - Rejection sampling
 - Importance sampling
- Markov Chain Monte Carlo (MCMC):
 - Markov Chain review
 - Metropolis-Hastings algorithm
 - Gibbs sampling
- Others: Monte Carlo EM, Slice sampling

5

Monte Carlo Principle

- Approximate density by samples from $p(x)$:

$$p_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x^{(i)}}(x),$$

- Estimate is unbiased and converges to the truth

$$I_N(f) = \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \xrightarrow[N \rightarrow \infty]{a.s.} I(f) = \int_x f(x)p(x)dx.$$

- Advantage over deterministic integration:
 - Samples from high-probability area, so more efficient
 - Question: how to sample from complex $p(x)$?

6

Rejection sampling

- Problem setup:
 - Want to sample $p(x)$, but too hard
 - Instead sample from proposal distribution $q(x)$
 - Requirement: $Mq(x) \geq p(x)$ for all x

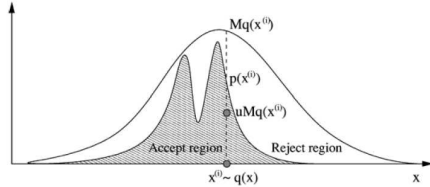


Figure 2. Rejection sampling: Sample a candidate $x^{(i)}$ and a uniform variable u . Accept the candidate sample if $uMq(x^{(i)}) < p(x^{(i)})$, otherwise reject it.

Issues with Rejection sampling

- Difficult to bound $p(x)$ over the whole space with a small M
- If M is too large, most samples will be rejected \rightarrow not efficient

$$\Pr(x \text{ accepted}) = \Pr\left(u < \frac{p(x)}{Mq(x)}\right) = \frac{1}{M}$$

8

Importance Sampling

- Use arbitrary proposal distribution $q(x)$ whose support includes $p(x)$

$$I(f) = \int f(x)p(x)dx = \int f(x) \frac{p(x)}{q(x)} q(x)dx = \int f(x)w(x)q(x)dx$$

↑
Importance weight

- Directly compute expectation (using all samples)

$$\hat{I}_N(f) = \sum_{i=1}^N f(x^{(i)})w(x^{(i)})$$

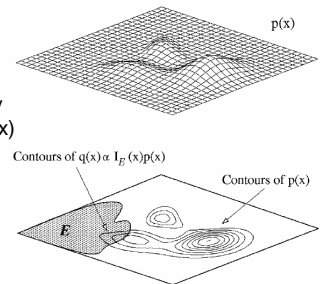
9

What is a good proposal distribution?

- Answer: one that is proportional to $|f(x)|p(x)$

$$q^*(x) = \frac{|f(x)|p(x)}{\int |f(x)|p(x)dx}$$

- Sampling from $q^*(x)$ may be more efficient than $p(x)$
- But this is not always possible



10

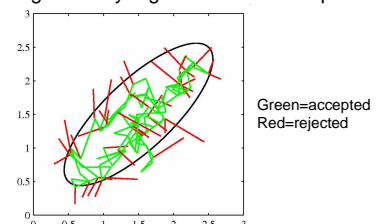
What we'll cover

- Monte Carlo methods:
 - Rejection sampling
 - Importance sampling
- Markov Chain Monte Carlo (MCMC):
 - Markov Chain review
 - Metropolis-Hastings algorithm
 - Gibbs sampling
- Others: Monte Carlo EM, Slice sampling

11

MCMC Motivation

- Monte Carlo methods may not be efficient in high dimensional spaces
- In MCMC, successive samples are correlated via a Markov chain
 - Explores high-density regions of the state-space



12

Markov Chain review

- An invariant $p(x) = p(x)T$ exists if transition matrix T satisfies:
 - Irreducibility: fully-connected
 - Aperiodic: chain not trapped in cycles
- A sufficient (but not necessary) condition is *reversibility*:
 - $p(x_{t+1}) T(x_t|x_{t+1}) = p(x_t) T(x_{t+1}|x_t)$
- In MCMC, we define proposal distribution $T=q(x_{t+1} | x_t)$:
 - After a period of "burn-in" / "mixing", we'll be sampling from the invariant distribution $p(x)$

13

Metropolis-Hastings (MH) algorithm

1. Initialise $x^{(0)}$.
2. For $i = 0$ to $N - 1$
 - Sample $u \sim \mathcal{U}_{[0,1]}$.
 - Sample $x^* \sim q(x^*|x^{(i)})$.
 - If $u < \mathcal{A}(x^{(i)}, x^*) = \min\left\{1, \frac{p(x^*)q(x^{(i)}|x^*)}{p(x^{(i)})q(x^*|x^{(i)})}\right\}$
 - Reject new state?
 - $x^{(i+1)} = x^*$
 - else
 - $x^{(i+1)} = x^{(i)}$

14

MH Examples:
True $p(x)$ is bimodal
Proposal $q(x_{t+1}|x_t) = N(x_t, \sigma^2)$

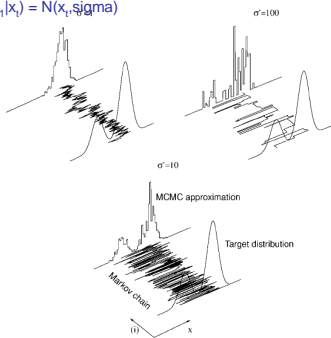


Figure 7. Approximations obtained using the MH algorithm with three Gaussian proposal distributions of different variances.

15

Why does MH work?

The transition kernel for the MH algorithm is

$$K_{MH}(x^{(i+1)} | x^{(i)}) = q(x^{(i+1)} | x^{(i)})\mathcal{A}(x^{(i)}, x^{(i+1)}) + \delta_{x^{(i)}}(x^{(i+1)})r(x^{(i)}),$$

where $r(x^{(i)})$ is the term associated with rejection

$$r(x^{(i)}) = \int_{\mathcal{X}} q(x^* | x^{(i)})(1 - \mathcal{A}(x^{(i)}, x^*)) dx^*.$$

By construction, the reversibility condition is satisfied by K_{MH} with $p(x)$ emerging as the invariant distribution:

$$p(x^{(i)})K_{MH}(x^{(i+1)} | x^{(i)}) = p(x^{(i+1)})K_{MH}(x^{(i)} | x^{(i+1)})$$

The main question is choosing $q()$ so that convergence is fast

- Gibbs sampling, Metropolis method, etc. are instances of MH with different $q()$

16

Gibbs sampler

- Assume a multivariate $p(x)$
- Gibbs sampling uses the conditional probabilities as $q()$
 - Very natural for graphical models
- Acceptance probability $\mathcal{A}()$ turns out to be 1

$$\begin{aligned} \mathcal{A}(x^{(i)}, x^*) &= \min\left\{1, \frac{p(x^*)q(x^{(i)}|x^*)}{p(x^{(i)})q(x^*|x^{(i)})}\right\} \\ &= \min\left\{1, \frac{p(x^*)p(x_j^{(i)}|x_{-j}^{(i)})}{p(x^{(i)})p(x_j^*|x_{-j}^*)}\right\} \\ &= \min\left\{1, \frac{p(x_j^*)}{p(x_j^{(i)})}\right\} = 1. \end{aligned}$$

17

Gibbs sampling pseudocode

1. Initialise $x_{0,1:n}$.
2. For $i = 0$ to $N - 1$
 - Sample $x_1^{(i+1)} \sim p(x_1|x_2^{(i)}, x_3^{(i)}, \dots, x_n^{(i)})$.
 - Sample $x_2^{(i+1)} \sim p(x_2|x_1^{(i+1)}, x_3^{(i)}, \dots, x_n^{(i)})$.
 - \vdots
 - Sample $x_j^{(i+1)} \sim p(x_j|x_1^{(i+1)}, \dots, x_{j-1}^{(i+1)}, x_{j+1}^{(i)}, \dots, x_n^{(i)})$.
 - \vdots
 - Sample $x_n^{(i+1)} \sim p(x_n|x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{n-1}^{(i+1)})$.

18

What we'll cover

- Monte Carlo methods:
 - Rejection sampling
 - Importance sampling
- Markov Chain Monte Carlo (MCMC):
 - Markov Chain review
 - Metropolis-Hastings algorithm
 - Gibbs sampling
- Others: Monte Carlo EM, Slice sampling

19

Monte Carlo EM

1. *E step.* Compute the expected value of the complete log-likelihood function with respect to the distribution of the hidden variables

$$Q(\theta) = \int_{X_h} \log(p(x_h, x_v | \theta)) p(x_h | x_v, \theta^{(old)}) dx_h,$$

where $\theta^{(old)}$ refers to the value of the parameters at the previous time step.

2. *M step.* Perform the following maximisation $\theta^{(new)} = \arg \max_{\theta} Q(\theta)$.

Use Monte Carlo sampling here to:
(1) approximate difficult integrals
(2) get out of local optima

20

Slice sampling

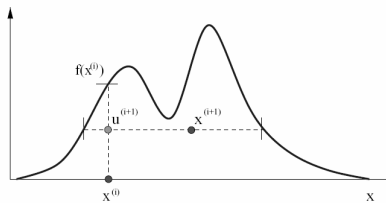
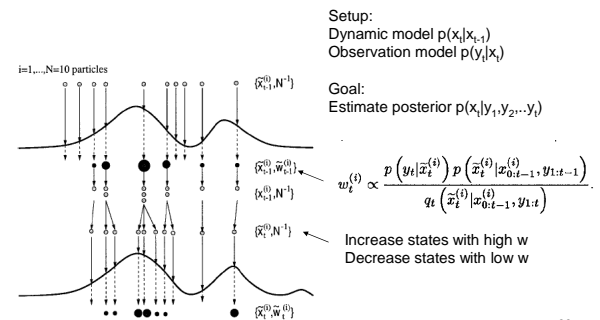


Figure 15. Slice sampling: given a previous sample, we sample a uniform variable $u^{(i+1)}$ between 0 and $f(x^{(i)})$. One then samples $x^{(i+1)}$ in the interval where $f(x) \geq u^{(i+1)}$.

21

Sequential Monte Carlo (Particle Filter)



22

Summary

1. Problems addressed:
 - Difficult integration/optimization where $p(x)$ is easily evaluated for a given x

$$I_N(f) = \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \xrightarrow{N \rightarrow \infty} I(f) = \int_{\mathcal{X}} f(x) p(x) dx.$$

2. All methods use a proposal $q()$ to help sample from $p(x)$
 - Rejection sampling: $Mq(x) \geq p(x)$

3. MCMC differs from Monte Carlo in that successive samples are correlated by $q(x_{t+1} | x_t)$

- Metropolis-Hastings: general $q(x_{t+1} | x_t)$

$$A(x^{(i)}, x^*) = \min \left\{ 1, \frac{p(x^*) q(x^{(i)} | x^*)}{p(x^{(i)}) q(x^* | x^{(i)})} \right\}$$

- Gibbs: $q(x_{t+1} | x_t)$ is conditional probability of multivariate distribution

23

References

- Figures/Equations for these slides come from:
 - Andrieu et. al. "An Intro to MCMC for Machine Learning", Machine Learning 2003
 - Bishop, *Pattern Recognition and Machine Learning* (chapter 11)
- Other useful references:
 - Diaconis, "The MCMC Revolution"
 - Resnik, "Gibbs sampling for the uninitiated"
 - Robert/Casella, *Monte Carlo statistical methods*, 1999
 - Neal, "Probabilistic inference using MCMC methods", 1993

24