Markov Chain Monte Carlo methods

Tomas McKelvey and Lennart Svensson

Signal Processing Group Department of Signals and Systems Chalmers University of Technology, Sweden

November 26, 2012



Today's learning objectives

After this lecture you should be able to

- Explain what Monte Carlo methods are and why they are important.
- Describe how importance sampling works.
- Use the Metropolis-Hastings algorithm to sample from a general distribution.
- Explain what Gibbs sampling is and when it is useful.

Decisions vs integrals Deterministic solutions

Bayesian Inference

- Data: a realization x
- Parameters, latent variables: $\theta = [\theta_1, \theta_2, \dots, \theta_p]$
- Likelihood: $L(x|\theta)$
- Inference based on the joint posterior

$$\pi(\theta|x) = \frac{L(x|\theta)\pi_0(\theta)}{\int L(x|\theta)\pi_0(\theta) \, d\theta}$$
$$\propto L(x|\theta)\pi_0(\theta)$$
Posterior \propto Likelihood \times Prior

Decisions vs integrals Deterministic solutions

Making decisions

Decisions are made by **minimizing posterior expected loss**. Two examples:

Istimation:

$$\hat{ heta} = \mathbb{E}\{ heta | x\} = \int heta \pi(heta | x) \, d heta$$

2 Model selection among two candidates m_1 and m_2 :

$$\int f(x|\theta_1,m_1)\pi(\theta_1|m_1) d\theta_1 \stackrel{m_1}{\underset{m_2}{\overset{>}{\leq}}} \int f(x|\theta_2,m_2)\pi(\theta_2|m_2) d\theta_2$$

Decisions vs integrals Deterministic solutions

Decisions by solving integrals

• In general, decisions are often made by computing integrals

$$I = \int g(\theta) \pi(\theta \big| x) \, d\theta$$

Today's problem formulation

Given: a distribution $p(\theta)$, known up to a normalization constant, and function, $g(\theta)$. **Objective:** find

$$I = \int g(\theta) p(\theta) \, d\theta.$$

Analytical solutions

- If conjugate prior, posterior p(θ) = π(θ|x) has nice analytical expression
 Examples: Beta-Binomial, Gauss-Gauss, Gamma-Poisson, etc.
- For some functions $g(\theta)$ the integral then has a closed form expression \Rightarrow we are done!
- Limitations: rarely ever possible in problems of practical interest.

More common: conjugate priors are used but solutions hard due to nuisance parameters

Deterministic approximations

• There are many techniques to approximate such integrals (or to perform inference) deterministically.

• Some examples:

- quadrature methods (related to unscented transform)
- Laplace approximation (assumes unimodality)
- message passing (uses conditional independence)
- variational Bayes (approximate independence)
- Very important and useful techniques!
- Weakness: limited accuracy. In particular when approximations do not match model structure.

Monte Carlo methods

Method

• Generate independent and identically distributed (iid) samples $\theta_1, \ldots, \theta_N$ from $p(\theta)$. Approximate

$$I = \int g(\theta) p(\theta) \, d\theta pprox rac{1}{N} \sum_{i=1}^N g(heta_i).$$



Monte Carlo methods – properties

•
$$\hat{I} = \frac{1}{N} \sum_{i=1}^{N} g(\theta_i)$$
 is an *unbiased estimate* of *I*

$$\mathbb{E}\left\{\frac{1}{N}\sum_{i=1}^{N}g(\theta_i)\right\}=I$$

Error covariance is

$$\mathsf{Cov}\{\hat{I}\} = \mathbb{E}\left\{(\hat{I} - I)(\hat{I} - I)^{\mathsf{T}}\right\} = \dots$$
$$= ?$$

 \rightarrow ..., independently of dimensionality of θ! ● Difficulty: how can we generate $θ_1, ..., θ_N$?

Monte Carlo methods – properties

•
$$\hat{l} = \frac{1}{N} \sum_{i=1}^{N} g(\theta_i)$$
 is an *unbiased estimate* of *l*

$$\mathbb{E}\left\{\frac{1}{N}\sum_{i=1}^{N}g(\theta_i)\right\}=I$$

Error covariance is

$$Cov{\hat{l}} = \mathbb{E}\left\{ (\hat{l} - l)(\hat{l} - l)^T \right\} = \dots$$
$$= \frac{1}{N}Cov\{g(\theta)\}$$

 \rightsquigarrow vanishes as 1/N, independently of dimensionality of θ !

• **Difficulty:** how can we generate $\theta_1, \ldots, \theta_N$?

Monte Carlo methods – feasible?

- Usually p(θ) is not a simple distribution, like Gaussian or Beta
 ⇒ we cannot use built-in random generators directly!
- In fact, $p(\theta)$ is often the posterior

 $p(\theta) \propto \pi(\theta) f(x|\theta)$

with unknown normalization.

• Is it still possible to use the Monte Carlo idea?

Importance sampling – basic idea

Importance sampling

• Generate $\theta_1, \ldots, \theta_N$ from a proposal distribution $q(\theta)$ and use approximation

$$I = \int g(\theta) \frac{p(\theta)}{q(\theta)} q(\theta) \, d\theta \approx \frac{1}{N} \sum_{i=1}^{N} g(\theta_i) \frac{p(\theta_i)}{q(\theta_i)}$$

• Often, weights $\tilde{w}_i = \frac{p(\theta_i)}{q(\theta_i)N}$ contain unknown normalization, therefore replaced by \tilde{w}_i .

$$w_i = \frac{w_i}{\sum_{n=1}^N \tilde{w}_n}$$

Now

CHAL

$$I \approx \sum_{i=1}^{N} w_i g(\theta_i)$$

with no need to know normalization of $p(\theta)$.

S Chalmers University of Technology

Monte Carlo methods Importance sampling

Importance sampling – an illustration

• An illustration of samples and weights



• Note that some samples are given zero weight.

Importance sampling – remarks

- Importance sampling "works" as long as support of $q(\theta)$ includes support of $p(\theta)$.
- Proposal should also be easy to generate samples from and similar to $p(\theta)$.
- Generates independent samples that can be used for Monte Carlo integration.
- Works very well if proposal selected carefully.
- Normally suffers from curse of dimensionality

 proposal mismatch grows quickly and makes most weights zero!

The MCMC concept

- Can Markov chains be used to perform Monte Carlo sampling?
- Suppose that we wish to generate samples from



• Samples can be obtained using a carefully designed Markov



CHALMERS Chalmers University of Technology

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Markov chains and kernels

- A Markov chain is a sequence of random variables, $\theta_1, \theta_2, \ldots$, that can be thought of as evolving over time.
- A key property is the conditional independence (Markov property):

$$f(\theta_m | \theta_{m-1}, \theta_{m-2}, \dots, \theta_1) = f(\theta_m | \theta_{m-1})$$

• The distribution $f(\theta_m | \theta_{m-1})$ is called the transition kernel

$$\begin{split} & \mathcal{K}(\theta_m \big| \theta_{m-1}) = f(\theta_m \big| \theta_{m-1}). \\ \Rightarrow \begin{cases} 1 = \int \mathcal{K}(\theta_m \big| \theta_{m-1}) \, d\theta_m \\ f(\theta_m) = \int \mathcal{K}(\theta_m \big| \theta_{m-1}) f(\theta_{m-1}) \, d\theta_{m-1} \end{cases} \end{split}$$

Stationary distributions

• Many Markov chains converge to a **stationary distribution** over time

 $\rightsquigarrow \theta_{m-1}$ and θ_m (almost) identically distributed for large m.

• $\pi(\theta)$ is a stationary distribution if

$$\pi(\theta_m) = \int \mathcal{K}(\theta_m | \theta_{m-1}) \pi(\theta_{m-1}) \, d\theta_{m-1}$$

• Example: compute $\pi_1 = \pi(\theta = 1)$ and $\pi_2 = \pi(\theta = 2)$.



Stationary distributions

• Many Markov chains converge to a **stationary distribution** over time

 $\rightsquigarrow \theta_{m-1}$ and θ_m (almost) identically distributed for large m.

• $\pi(\theta)$ is a stationary distribution if

$$\pi(\theta_m) = \int \mathcal{K}(\theta_m | \theta_{m-1}) \pi(\theta_{m-1}) \, d\theta_{m-1}$$

• Example: compute $\pi_1 = \pi(\theta = 1)$ and $\pi_2 = \pi(\theta = 2)$.

$$0.2 \\ \theta = 1 \\ 0.4 \\ \theta = 2 \\ 0.8 \\ P = \begin{pmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{pmatrix}$$

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Reversible chains

• A Markov chain is reversible if $\exists \pi(\theta)$:

$$\pi(\theta') \mathcal{K}(\theta | \theta') = \pi(\theta) \mathcal{K}(\theta' | \theta)$$
(1)

 \rightsquigarrow probability of beeing at θ and passing to θ' and from beeing at θ' and passing to θ are the same.

• Example: verify that the above chain is reversible.

$$p = pP$$

• Note: (1) is called detailed balance condition and implies that $\pi(\theta)$ is a stationary distribution!

Proof: integrate both sides with respect to θ

$$\int \pi(heta') \mathsf{K}(hetaigert heta') \, \mathsf{d} heta = \int \pi(heta) \mathsf{K}(heta'igert heta) \, \mathsf{d} heta$$

where lhs is simply $\pi(\theta')$. The condition for stationarity!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Reversible chains

• A Markov chain is reversible if $\exists \pi(\theta)$:

$$\pi(\theta') \mathcal{K}(\theta | \theta') = \pi(\theta) \mathcal{K}(\theta' | \theta)$$
(1)

 \rightsquigarrow probability of beeing at θ and passing to θ' and from beeing at θ' and passing to θ are the same.

• Example: verify that the above chain is reversible.

$$P = \begin{pmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{pmatrix} \qquad \qquad p = pP$$

• Note: (1) is called detailed balance condition and implies that $\pi(\theta)$ is a stationary distribution!

Proof: integrate both sides with respect to θ

$$\int \pi(heta') {\cal K}(hetaigert heta') \, d heta = \int \pi(heta) {\cal K}(heta'igert heta) \, d heta$$

where lhs is simply $\pi(\theta')$. The condition for stationarity!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Reversible chains

• A Markov chain is reversible if $\exists \pi(\theta)$:

$$\pi(\theta') \mathcal{K}(\theta | \theta') = \pi(\theta) \mathcal{K}(\theta' | \theta)$$
(1)

 \rightsquigarrow probability of beeing at θ and passing to θ' and from beeing at θ' and passing to θ are the same.

• Example: verify that the above chain is reversible.

$$P = \begin{pmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{pmatrix}$$
 $P^{\infty} = \begin{pmatrix} 1/3 & 2/3 \\ 1/3 & 2/3 \end{pmatrix}$ $p = pP$

• Note: (1) is called detailed balance condition and implies that $\pi(\theta)$ is a stationary distribution!

Proof: integrate both sides with respect to θ

$$\int \pi(heta') \mathsf{K}(hetaigert heta') \, \mathsf{d} heta = \int \pi(heta) \mathsf{K}(heta'igert heta) \, \mathsf{d} heta$$

where lhs is simply $\pi(\theta')$. The condition for stationarity!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Reversible chains

• A Markov chain is reversible if $\exists \pi(\theta)$:

$$\pi(\theta') \mathcal{K}(\theta | \theta') = \pi(\theta) \mathcal{K}(\theta' | \theta)$$
(1)

 \rightsquigarrow probability of beeing at θ and passing to θ' and from beeing at θ' and passing to θ are the same.

• Example: verify that the above chain is reversible.

$$P = \begin{pmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{pmatrix} P^{\infty} = \begin{pmatrix} 1/3 & 2/3 \\ 1/3 & 2/3 \end{pmatrix} p = pP$$

1/3 * 0.4 = 2/3 * 0.2

• Note: (1) is called detailed balance condition and implies that $\pi(\theta)$ is a stationary distribution!

Proof: integrate both sides with respect to θ

$$\int \pi(heta') \mathsf{K}(hetaigert heta') \, \mathsf{d} heta = \int \pi(heta) \mathsf{K}(heta'igert heta) \, \mathsf{d} heta$$

where lhs is simply $\pi(\theta')$. The condition for stationarity!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Metropolis-Hastings (MH) algorithm

- An MCMC algorithm which is flexible and very simple to use.
- A technique to design a reversible Markov chain with stationary distribution $\pi(\theta) = p(\theta)$.
- MH has two components
 - **Proposal distribution** $q(\theta'|\theta)$: suggests a move from θ to θ' . Often simply $\theta' = \theta + n$ where $n \sim \mathcal{N}(0, I)$.
 - Acceptance probability a(θ' | θ): the proposed state is accepted with probability a(θ' | θ)
 → a(θ' | θ) is selected to ensure that the detailed balance conditioned is satisfied for π(θ) = p(θ).

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Metropolis-Hastings (MH) algorithm

Summary of the MH algorithm

- Start with an arbitrary initial value θ_0 .
- Update from θ_m to θ_{m+1} $(m = 0, 1, \dots, N)$ by
 - Generate $\xi \sim q(\xi | \theta_m)$.

2 Take

$$heta_{m+1} = egin{cases} \xi & ext{with probability } a(\xi ig heta_m) \ heta_m & ext{otherwise}. \end{cases}$$

• Given a realization of the chain we do Monte Carlo sampling:

$$\frac{1}{N-N_b+1}\sum_{m=N_b}^N g(\theta_m)$$

where N_b is when we hope that the chain has reaches its stationary distribution (burn-in time)

CHALMERS

Chalmers University of Technology

Markov chains Metropolis-Hastings algorithm Gibbs sampling

The acceptance probability

• We wish to select $a(\theta'|\theta)$ such that

$$p(heta')K(hetaig| heta')=p(heta)K(heta'ig| heta)$$

 \rightsquigarrow ensures that $p(\theta)$ is a stationary distribution.

• What is the transition kernel, $K(\theta'|\theta)$ (when $\theta' \neq \theta$)?

$$K(\theta'|\theta) = ?$$

• How can we select $a(\theta'|\theta)$ to satisfy the above detailed balance condition?

Markov chains Metropolis-Hastings algorithm Gibbs sampling

The acceptance probability

• We wish to select $a(\theta'|\theta)$ such that

$$p(heta')K(hetaig| heta')=p(heta)K(heta'ig| heta)$$

 \rightsquigarrow ensures that $p(\theta)$ is a stationary distribution.

• What is the transition kernel, $K(\theta'|\theta)$ (when $\theta' \neq \theta$)?

$$K(\theta'|\theta) = q(\theta'|\theta)a(\theta'|\theta).$$

• How can we select $a(\theta'|\theta)$ to satisfy the above detailed balance condition?

Markov chains Metropolis-Hastings algorithm Gibbs sampling

The acceptance probability

• We wish to select $a(\theta'|\theta)$ such that

$$p(heta')K(hetaig| heta')=p(heta)K(heta'ig| heta)$$

 \rightsquigarrow ensures that $p(\theta)$ is a stationary distribution.

• What is the transition kernel, $K(\theta'|\theta)$ (when $\theta' \neq \theta$)?

$$K(\theta'|\theta) = q(\theta'|\theta)a(\theta'|\theta).$$

• How can we select $a(\theta'|\theta)$ to satisfy the above detailed balance condition?

$$p(heta')q(heta| heta')a(heta| heta') = p(heta)q(heta'| heta)a(heta'| heta) \ \Leftrightarrow rac{a(heta| heta')}{a(heta'| heta)} = rac{p(heta)q(heta'| heta)}{p(heta')q(heta| heta')}$$

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Metropolis-Hastings (MH) algorithm

Summary of the MH algorithm

- Start with an arbitrary initial value θ_0 .
- Update from $heta_m$ to $heta_{m+1}$ $(m=0,1,\ldots,N)$ by
 - **1** Generate $\xi \sim q(\xi | \theta_m)$.

2 Take

$$\theta_{m+1} = \begin{cases} \xi & \text{ with probability } \min\left\{1, \frac{p(\xi)q(\theta_m|\xi)}{p(\theta_m)q(\xi|\theta_m)}\right\}\\ \theta_m & \text{ otherwise.} \end{cases}$$

• Note: algorithm involves point-wise evaluation of $p(\xi)/p(\theta_m)$ \rightarrow possible even when the normalization constant is unknown!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

The proposal distribution $q(\theta'|\theta)$

- We want to select q(θ|θ') such that the chain "explores the space quickly".
 If it does, it is said to "mix" well.
- The most common choice for q is the random walk proposal

$$q(\theta' | \theta) = f(||\theta' - \theta||)$$

$$\Leftrightarrow \xi = \theta_m + v,$$

where v is a symmetric random variable.

• For this choice, the acceptance probability simplifies to

$$a(\xi | heta_m) = \min\left\{1, rac{p(\xi)}{p(heta_m)}
ight\}$$

 \rightsquigarrow always accept proposals that move upwards.

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Metropolis-Hastings – summary

• Pros:

- The MH algorithm can be applied to virtually any problem.
- We can collect all our samples from a single run of the Markov chain.
- The freedom in the choice of proposal gives us the possibility to design a chain that mixes quickly.
- Cons:
 - We need to select the proposal density, $q(\theta|\theta')$.
 - A poor choice will give us samples that are highly correlated and do not represent the full distribution.
 - It is also difficult (impossible?) to know when the chain has reached the stationary distribution. Use convergence diagnostics!

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Gibbs sampling – general algorithm

- Gibbs sampling is an MCMC algorithm that make use of the model structure.
- Idea: sample one variable (dimension) at a time conditioned on all the others.

The Gibbs sampler

CHAL

- Start with arbitrary initial vector $\theta_0 = [\theta_0(1), \theta_0(2), \dots, \theta_0(d)]^T$.
- For $m = 0, 1, \dots, N$ generate

1)
$$\theta_{m+1}(1) \sim p(\theta(1)|\theta_m(2),\ldots,\theta_m(d))$$

2) $\theta_{m+1}(2) \sim p(\theta(2)|\theta_{m+1}(1),\theta_m(3),\ldots,\theta_m(d))$

d)
$$\theta_{m+1}(d) \sim p(\theta(d)|\theta_{m+1}(1),\ldots,\theta_{m+1}(d-1))$$

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Gibbs sampling - toy example

• The Gibbs sampler for the model

$$\begin{bmatrix} \theta(1) \\ \theta(2) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right)$$

is

1)
$$\theta_{m+1}(1)|\theta_m(2) \sim \mathcal{N}(\rho\theta_m(2), 1-\rho^2)$$

2) $\theta_{m+1}(2)|\theta_{m+1}(1) \sim \mathcal{N}(\rho\theta_{m+1}(1), 1-\rho^2)$

– For $\rho = 0.7$ it can look like this:



Markov chains Metropolis-Hastings algorithm Gibbs sampling

Gibbs sampling vs Hierarchical models

 $\phi_{m+1} \sim \pi(\phi | \theta_{m+1})$

• Gibbs sampling is particularly useful for Hierachical models.



Figure: A very small hierachical model.

- While generating ϕ we can ignore x due to conditional independence.
 - \rightsquigarrow this benefit is substantially larger in bigger networks.

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Gibbs sampling – summary

Pros:

- Few design choices makes it simple to use.
- Utilizes the model structure.
- Generates high dimensional variables using a sequence of low dimensional simulations.
- Cons:
 - Mixes poorly if variables are highly correlated.
 - Requires knowledge of $p(\theta)$.
 - Only applicable if conditionals are easy to simulate.

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Other sampling methods

There are many other sampling methods that we could not cover today:

- Rejection sampling a classical method to generate independent samples from a distribution.
- Slice sampler a more general type of Gibbs sampler.
- Population Monte Carlo a combination of MH and particle filters.
- Hamiltonian (or Hybrid) Monte Carlo introduces additional latent variables that enables large steps in the state space.
- Adaptive MCMC methods to adaptively improve the proposal distribution based on what we learn from the Markov chain.
- Reversible Jump MCMC an extension to situations where the dimensionality of θ is unknown.

Markov chains Metropolis-Hastings algorithm Gibbs sampling

Today's learning objectives

After this lecture you should be able to

- Explain what Monte Carlo methods are and why they are important.
- Describe how importance sampling works.
- Use the Metropolis-Hastings algorithm to sample from a general distribution.
- Explain what Gibbs sampling is and when it is useful.

Markov chains Metropolis-Hastings algorithm Gibbs sampling

References

Some further readings:

- Gelfand, A.E. and Smith, A.F.M., "Sampling Based Approaches to Calculating Marginal Densities", Journal of the American Statistical Association, Vol. 85, pp. 398-409, 1990.
- Casella G. and George E.I. "Explaining the Gibbs Sampler", The American Statistician, Vol. 46 No. 3, pp. 167-174, 1992
- Hastings W.K. "Monte Carlo Sampling Methods Using Markov Chains and Their Applications", *Biometrika*, Vol. 57, No. 1, pp. 97-109, 1970
- Gilks W.R., Richardson S. and Spiegelhalter D.J., Markov Chain Monte Carlo in Practice. Chapman & Hall/CRC, 1996.