Lecture 11 Sampling

University of Amsterdam





- 2 Basic sampling algorithms
 - Transforming samples
 - Rejection Sampling
 - Importance Sampling
- 3 Markov Chain Monte Carlo
 - The Metropolis Algorithm
 - Metropolis-Hastings
 - Gibbs Sampling

4 Wrap-up

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

3



- 2 Basic sampling algorithms
 - Transforming samples
 - Rejection Sampling
 - Importance Sampling
- Markov Chain Monte Carlo
 The Metropolis Algorithm
 Metropolis-Hastings
 Gibbs Sampling

4 Wrap-up

Markov Chain Monte Carlo

(日)、



Recall how graphical models specify the factorisation of a joint probability distribution

- It provides us with algorithms to do marginalisation efficiently
- We are still free to choose the specific form of the distributions
- However choosing the form is tricky: if we do not choose wisely, the distribution may not be computable analytically











Ř







Ř







Sampling

・ロト ・ 雪 ト ・ ヨ ト

э



We're generally not really interested in the posterior distribution for its own sake

• Often, we need the posterior distribution to compute the expectation of some function

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

- Sometimes we want to marginalise out variables
- This is often impossible to do analytically



Introduction 0000000 Sampling Basic sampling algorithms

Markov Chain Monte Carlo

• • • • • • • • • •

Wrap-up 0

UNIVERSITY OF AMSTERDAM

Integrating numerically

We could compute the solution numerically

• Computing the integral by summing approximations

- Becomes intractable very fast as the dimensionality goes up
- How about indefinite integrals (we cannot reach $\infty)$
- Use sampling: if we can draw independent samples $\{z^{(1)} \dots z^{(L)}\}$ of p(z), we can approximate

$$\mathbb{E}[f] = \int f(\mathbf{z}) \, p(\mathbf{z}) \, \mathrm{d}\mathbf{z} \simeq \frac{1}{L} \sum_{l=1}^{L} f(\mathbf{z}^{(l)})$$

Introduction 0000000 Sampling Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up

Integrating numerically

We could compute the solution numerically

- Computing the integral by summing approximations
 - Becomes intractable very fast as the dimensionality goes up
 - How about indefinite integrals (we cannot reach ∞)
- Use sampling: if we can draw independent samples $\{\mathbf{z}^{(1)} \dots \mathbf{z}^{(L)}\}$ of $p(\mathbf{z})$, we can approximate

$$\mathbb{E}[f] = \int f(\mathbf{z}) \, p(\mathbf{z}) \, \mathrm{d}\mathbf{z} \simeq rac{1}{L} \sum_{l=1}^{L} f(\mathbf{z}^{(l)})$$

Intelligent Autonomous Systems

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Introduction 0000000 Sampling Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up O

UNIVERSITY OF AMSTERDAM

Integrating numerically

We could compute the solution numerically

- Computing the integral by summing approximations
 - Becomes intractable very fast as the dimensionality goes up
 - ullet How about indefinite integrals (we cannot reach ∞)
- Use sampling: if we can draw independent samples $\{\mathbf{z}^{(1)} \dots \mathbf{z}^{(L)}\}$ of $p(\mathbf{z})$, we can approximate

$$\mathbb{E}[f] = \int f(\mathbf{z}) \, p(\mathbf{z}) \, \mathrm{d}\mathbf{z} \simeq \frac{1}{L} \sum_{l=1}^{L} f(\mathbf{z}^{(l)})$$

Markov Chain Monte Carlo

イロト 不得 トイヨト イヨト

э



Sampling is therefore useful when:

- We cannot represent the joint distribution tractably but we can evaluate its building blocks
- We cannot compute the expectation of f analytically, but we can evaluate f(x) for given x.

Remember

About Sampling

Samples are easy to manipulate, even if their distribution is not



Markov Chain Monte Carlo



Things to notice:

About Sampling

- The approximation is unbiased: 𝔼[f̂] = 𝔼[f], and var[f̂] = var_{p(z)}[f].
- The approximation is independent of the dimensionality of z. In practice 10-20 independent samples generally suffice for a good approximation.
- If the function we want the expectation of is small where $p(\mathbf{x})$ is large and vice versa, a larger number of samples is needed for good accuracy.
- To get the marginal distribution p(z): sample from the joint distribution p(x, z) and ignore the values of x

Basic sampling algorithms

Markov Chain Monte Carlo

Graphical Models

In the case of Bayesian Networks, we can use *ancestral sampling*:

- Start at root nodes, sample from those
- For each child node, sample from the conditional distribution



However how do we deal with observed values?

- Logic sampling: sample until you reach the observation
- If the sampled value agrees with the observation, keep it
- Otherwise: discard
- Very inefficient and rarely used.

Intelligent Autonomous Systems

Basic sampling algorithms

Markov Chain Monte Carlo

JNIVERSITY OF AMSTERDAM

Graphical Models

In the case of Bayesian Networks, we can use *ancestral sampling*:

- Start at root nodes, sample from those
- For each child node, sample from the conditional distribution



However how do we deal with observed values?

- Logic sampling: sample until you reach the observation
- If the sampled value agrees with the observation, keep it
- Otherwise: discard
- Very inefficient and rarely used.

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



- Basic sampling algorithms
 - Transforming samples
 - Rejection Sampling
 - Importance Sampling
- Markov Chain Monte Carlo
 The Metropolis Algorithm
 Metropolis-Hastings
 Gibbs Sampling

4 Wrap-up

Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up

Transforming samples

Transforming samples

We assume we have a way of generating uniformly distributed pseudo-random numbers z in the range (0,1) (p(z) = 1).

- How can we obtain samples with a different, desired distribution?
- If y = f(z), the distribution of y is given by

$$p(y) = p(z) \left| \frac{dz}{dy} \right|$$



Ň

Basic sampling algorithms

Markov Chain Monte Carlo

(日)、

э

Wrap-up O

Ř

UNIVERSITY OF AMSTERDAM

Transforming samples

Transforming samples



Markov Chain Monte Carlo

(日)、

Ň×

UNIVERSITY OF AMSTERDAM

Transforming samples

Transforming samples

э

The good:

- Fast and exact
- Can easily be extended to multivariate distributions
- Works for the (multivariate) Gaussian distribution
- Building block for more advanced sampling schemes The bad:
 - Only possible for limited, simple distributions



Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up 0

ě

UNIVERSITY OF AMSTERDAM

Rejection Sampling

Rejection sampling

Suppose we know $\tilde{p}(\mathbf{z})$ and want to sample from

$$p(\mathbf{z}) = rac{1}{Z} ilde{p}(\mathbf{z})$$

where $\tilde{p}(\mathbf{z})$ is easy to compute, but Z is unknown.

Example

In a Markov Random Fields, the probability densities are given by the potential functions associated with a clique

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{C} \psi_{C}(\mathbf{x}_{c}),$$

where Z is hard to compute.



Introduction 0000000	Basic sampling algorithms	Markov Chain Monte Carlo 000000000	
Rejection Sampling			
			×



Choose a proposal distribution q(z) and find a constant k such that $kq(z) \ge \tilde{p}(z)$ everywhere

- Sample a value z_0 from q(z)
- Sample a value u uniformly from $[0, kq(z_0)]$
- Reject the pair if $u > \tilde{p}(z_0)$

Intelligent Autonomous Systems

Basic sampling algorithms

Markov Chain Monte Carlo

<ロト <回ト < 注ト < 注ト

æ

Wrap-up O

Rejection Sampling

Rejection sampling



Rejection Sampling

Importance of the proposal distribution

Samples from q(x) are accepted with probability $\tilde{p}(z)/kq(z)$, so that the probability of acceptance is

$$p(ext{accept}) = \int rac{ ilde{p}(z)}{kq(z)} \, q(z) \, \mathrm{d}z = rac{1}{k} \int ilde{p}(z) \, \mathrm{d}z$$

It is therefore important to keep k as small as possible.



- Rejection rate: 1/k
- The optimal $k = (\sigma_q / \sigma_p)^D$ in D dimensions
- If D = 1000 and $\sigma_q = 1.01 \sigma_p$, 1 sample in 20000 is accepted.

Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up O

UNIVERSITY OF AMSTERDAM

Rejection Sampling

Adaptive Rejection Sampling

In practice it is often hard to find a good proposal distribution. However when $\tilde{p}(z)$ is log-concave, we can adapt it on the fly:



- Find the tangent to $\tilde{p}(z)$ at an initial grid of points.
- Sample from the piecewise-linear distribution (this is easy, by transforming uniform samples)
- If the sample is rejected, add the current point to the grid, update q(z).

Basic sampling algorithms

Markov Chain Monte Carlo

・ロト ・ 四 ト ・ ヨ ト ・ ヨ ト

Ξ.

Rejection Sampling

Importance sampling

In order to obtain a good estimate of an expectation, it is not necessary to be able to sample from the distribution

• For example: if we can evaluate p(z) easily, we could discretise the space uniformly and approximate the integrand as

$$\mathbb{E}[f] \simeq \sum_{l=1}^{L} p(\mathbf{z}^{(l)}) f(\mathbf{z}^{(l)})$$

- Problem: Very inaccurate/inefficient, esp. in high dimensional spaces
- Better: Sample from proposal distribution q(z) and reweigh:

$$\mathbb{E}[f] \simeq \frac{1}{L} \sum_{l=1}^{L} \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)})$$

Intelligent Autonomous Systems

ntroduction	Basic sampling algorithms ○○○○○○○○●○○	Markov Chain Monte Carlo 000000000	
mportance Sampling			

Importance Sampling

Often we can only evaluate $p(\mathbf{z})$ up to a normalisation constant. Similarly, it is interesting to be able to sample from $q(\mathbf{z}) = \frac{1}{Z}\tilde{q}(\mathbf{z})$. Then we can define

$$ilde{r}_l = rac{ ilde{p}(\mathbf{z}^{(l)})}{ ilde{q}(\mathbf{z}^{(l)})} ext{ and } w_l = rac{ ilde{r}_l}{\sum_{m=1}^L ilde{r}_m}$$

and compute the expectation as follows

$$\mathbb{E}[f] \simeq \sum_{l=1}^{L} \mathbf{w}_l f(\mathbf{z}^{(l)})$$

(日)、

No samples are thrown away.

Markov Chain Monte Carlo

Importance Sampling

Importance Sampling

pros:

- No need for a scaling constant
- Super-efficient sampling: if q(z) is similar to p(z)f(z), we require fewer samples for a given accuracy than if q(z) is similar to p(z)

cons:

- if p(z)f(z) is strongly varying, few samples will carry most of the weight
- if q(z) is small where p(z)f(z) is large, it is possible for no samples to lie in those regions: the samples may then have low variance while the estimates are severely wrong.

Intelligent Autonomous Systems

Ň

Introduction 0000000	Basic sampling algorithms	Markov Chain Monte Carlo 000000000	
Importance Sampling			
EM Algorithm	า		Ň

Sampling can of course be used for the E-step of the EM algorithm. Recall that we optimise the expectation of the complete log-likelihood:

$$Q(\theta, \theta^{\text{old}}) = \int p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{Z}, \mathbf{X}|\theta) \, \mathrm{d}\mathbf{Z}$$

If we draw from $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$, we get Monte Carlo EM

$$\mathcal{Q}(\boldsymbol{ heta}, \boldsymbol{ heta}^{\mathsf{old}}) \simeq rac{1}{L} \sum_{l=1}^L \ln p(\mathbf{Z}, \mathbf{X} | \boldsymbol{ heta}) \, \mathrm{d} \mathbf{Z}$$

which we can optimise with the usual M-step.

Intelligent Autonomous Systems

・ロト ・母 ト ・ ヨ ト ・ ヨ ・ の へ (?)

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

э.



- Basic sampling algorithms
 Transforming samples
 - Rejection Sampling
 - Importance Sampling
- 3 Markov Chain Monte Carlo
 - The Metropolis Algorithm
 - Metropolis-Hastings
 - Gibbs Sampling

4 Wrap-up

Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up

UNIVERSITY OF AMSTERDAM

The Metropolis Algorithm

Markov Chain Monte Carlo (MCMC)

The techniques we saw until now suffer severe limitations in high-dimensional spaces. The Metropolis algorithm circumvents these by using a conditional proposal distribution:

- Sample candidate sample \mathbf{z}^* from a symmetric proposal distribution $q(\mathbf{z}^* | \mathbf{z}^{(\tau)})$
- Accept the candidate sample with probability

$$A(\mathbf{z}^*, \mathbf{z}^{(au)}) = \min\left(1, rac{ ilde{p}(\mathbf{z}^*)}{ ilde{p}(\mathbf{z}^{(au)})}
ight)$$

- If \mathbf{z}^* is rejected, $\mathbf{z}^{(\tau+1)} = \mathbf{z}^{(\tau)}$.
- This tends to p(z) as $\tau \to \infty$, if $q(z_A|z_B) > 0$ for all z_A , z_B .
- The samples are not independent. Keep every *M*th sample for "independent" samples

Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up

Š

UNIVERSITY OF AMSTERDAM

The Metropolis Algorithm

Example: Metropolis sampling of Gaussian



Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up

The Metropolis Algorithm

Example: Metropolis sampling of Gaussian



Š

Intelligent Autonomous Systems

・ロト・日本・日本・日本・日本・日本

Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Š

The Metropolis Algorithm

Example: Metropolis sampling of Gaussian





Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up

Š

UNIVERSITY OF AMSTERDAM

The Metropolis Algorithm

Example: Metropolis sampling of Gaussian



The Metropolis Algorithm

Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up

Convergence of the Markov Chain

We want the chain to converge to a stationary distribution $p^*(z)$. The distribution $p^*(z)$ is stationary with respect to the chain if

$$p^*(\mathbf{z}) = \sum_{\mathbf{z}'} p^*(\mathbf{z}') T(\mathbf{z}', \mathbf{z})$$

A sufficient condition for this is that the chain exhibits *detailed balance*:

$$p^*(\mathbf{z})T(\mathbf{z},\mathbf{z}') = p^*(\mathbf{z}')T(\mathbf{z}',\mathbf{z})$$

in which case the chain is said to be reversible



Metropolis-Hastings

Metropolis-Hastings

Basic sampling algorithm

Markov Chain Monte Carlo

(日)、

-



If the proposal distribution is not symmetric, we can apply the Metropolis-Hastings algorithm. Accept a candidate sample z^* with probability

$$A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)q(\mathbf{z}^{(\tau)}, \mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)})q(\mathbf{z}^*, \mathbf{z}^{(\tau)})}\right),$$

which exhibits detailed balance with respect to $p(\mathbf{z})$.



Metropolis-Hastings



Proof: Detailed Balance

$$p(\mathbf{z})q(\mathbf{z}|\mathbf{z}')A(\mathbf{z}',\mathbf{z}) = p(\mathbf{z})q(\mathbf{z}|\mathbf{z}')\min(1,\frac{p(\mathbf{z}')q(\mathbf{z}'|\mathbf{z})}{p(\mathbf{z})q(\mathbf{z}|\mathbf{z}')})$$

= min(p(z)q(z|z'), p(z')q(z'|z))
= min(p(z')q(z'|z), p(z)q(z|z'))
= p(z')q(z'|z)A(z,z')

The proposal distribution therefore only affects the efficiency of the sampling, not their end distribution.



Metropolis-Hastings Length scale Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up



UNIVERSITY OF AMSTERDAM

Markov chain methods trade off exploration vs. rejection rate.

- For low rejection rate, proposal distribution scale $\rho \approx \sigma_{\min}$
- Other directions explored by random walk
- Approx. independent samples: (σ_{max}/σ_{min})² steps
- Very different lengths scales ⇒ slow to converge





Metropolis-Hastings Length scale Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up O



Markov chain methods trade off exploration vs. rejection rate.

- For low rejection rate, proposal distribution scale $\rho \approx \sigma_{\min}$
- Other directions explored by random walk
- Approx. independent samples: (σ_{max}/σ_{min})² steps
- Very different lengths scales
 ⇒ slow to converge



Metropolis-Hastings Length scale Basic sampling algorithms

Markov Chain Monte Carlo

Wrap-up O



Markov chain methods trade off *exploration* vs. *rejection rate*.

- For low rejection rate, proposal distribution scale $\rho \approx \sigma_{\min}$
- Other directions explored by random walk
- Approx. independent samples: (σ_{max}/σ_{min})² steps
- Very different lengths scales ⇒ slow to converge









- Other directions explored by random walk
- Approx. independent samples: $(\sigma_{max}/\sigma_{min})^2$ steps
- Very different lengths scales ⇒ slow to converge











 Very different lengths scales ⇒ slow to converge



Introduction 0000000 Gibbs Sampling

Gibbs Sampling

Basic sampling algorithm

Markov Chain Monte Carlo

Wrap-up O



Gibbs sampling is a special case of Metropolis-Hastings sampling, where we sample $p(\mathbf{z}) = p(z_1 \dots z_D)$ by sampling each z_i in turn.

- The proposal distribution is $p(z_i | \mathbf{z}_{\setminus i})$
- All samples are accepted: $q(\mathbf{z}^*|\mathbf{z}) = p(z_i|\mathbf{z}_{\setminus i}), \ \mathbf{z}_{\setminus i}^* = \mathbf{z}_{\setminus i}$ and $p(\mathbf{z}) = p(z_i|\mathbf{z}_{\setminus i})p(\mathbf{z}_{\setminus i})$, so that

$$A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*)q(\mathbf{z}|\mathbf{z}^*)}{p(\mathbf{z})q(\mathbf{z}^*|\mathbf{z})} = \frac{p(z_i^*|\mathbf{z}_{\backslash i}^*)p(\mathbf{z}_{\backslash i}^*)p(z_i|\mathbf{z}_{\backslash i}^*)}{p(z_i|\mathbf{z}_{\backslash i})p(\mathbf{z}_{\backslash i})p(z_i^*|\mathbf{z}_{\backslash i})} = 1$$

Introduction 0000000 Gibbs Sampling Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo



Introduction 0000000 Gibbs Sampling Basic sampling algorithm: 0000000000000 Markov Chain Monte Carlo



Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo

Wrap-up O

UNIVERSITY OF AMSTERDAM

Practical usability of Gibbs sampling

The practical applicability depends on the ease with which we can sample from $p(z_i | \mathbf{z}_{\setminus i})$. For graphical models, the conditional distribution of a node depends on its Markov Blanket:



For directed graphs, the conditional distributions for Gibbs sampling is often complex but log-concave. Adaptive rejection sampling is therefore widely applicable in this setting. One example where MCMC methods perform poorly is for sequential data, such as extensions of the linear dynamical system

- The states are by definition strongly correlated: sampling algorithms converge slowly
- Instead of sampling from the joint distribution, the particle filter will do a forward pass through the data
- Use the process noise as proposal distribution (which is exact)
- Update the weights of the particles according to the probability of the observation given the sample
- Resample, to keep more particles in the high-density areas



イロト 不得 トイヨト イヨト

3



- Basic sampling algorithms
 - Transforming samples
 - Rejection Sampling
 - Importance Sampling
- Markov Chain Monte Carlo
 The Metropolis Algorithm
 Metropolis-Hastings
 - Gibbs Sampling

4 Wrap-up

Wrap-up

Summary

Basic sampling algorithms 0000000000000 Markov Chain Monte Carlo



Today we've seen sampling:

- Why do we do this
- Basic sampling algorithms
- Markov Chain Monte Carlo

(Bishop, p. 523-525) (Bishop, p. 528-534) (Bishop, p. 537-546)

(日)、

э

