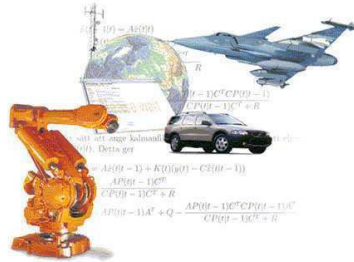


Machine Learning, Lecture 10 MCMC and sampling methods



Thomas Schön

Division of Automatic Control
Linköping University
Linköping, Sweden.

Email: schon@isy.liu.se,
Phone: 013 - 281373,
Office: House B, Entrance 27.

Outline of lecture 10

2(48)

1. About the exam
2. Summary of lecture 9
3. Motivation for Monte Carlo methods
4. Basic sampling methods
 - Transformation methods
 - Rejection sampling
 - Importance sampling
5. Markov Chain Monte Carlo (MCMC)
 - General properties
 - Metropolis Hastings sampler
 - Gibbs sampler

(Chapter 11)

About the exam (I/II)

3(48)

- If you have followed the course and completed the exercises you will not be surprised when you see the exam.
- You will learn new things during the exam.

Practicalities:

- Time frame: 3 days (72h), somewhere in the time frame March 17, 2011 - April 26, 2013 (April 26, 2013 is the last day to start the exam).
- Within 72 hours after you have collected the exam, you put your solutions in an envelope (seal it) and hand it in to the coordinator at the Division of Automatic Control, Ninna Stensgård.

About the exam (II/II)

4(48)

As usual the **graduate exam honor code** applies. This means,

- The course books, other books and MATLAB are all allowed aids.
- Internet services such as email, web browsers and other communication with the surrounding world concerning the exam is NOT allowed.
- You are NOT allowed to actively search for the solutions in books, papers, the Internet or anywhere else.
- You are NOT allowed to talk to others (save for the responsible teachers) about the exam at all.
- You are NOT allowed to look at exams from earlier version of the course.
- If anything is unclear concerning what is allowed and not, just ask me.

We worked with three different graphs:

1. **Directed acyclic graphs:** a.k.a. Bayesian networks
2. **Undirected graphs:** a.k.a. Markov random fields
3. **Factor graphs:** Used for inference

Used **d-separation** as a means to check conditional independence among random variables.

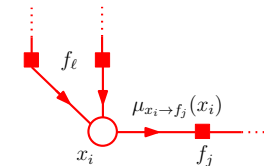
Introduced the Hammersley-Clifford theorem to find the PDF for an undirected graph.

Inference in graphical models amounts to computing the posterior distribution of one or more of the nodes that are not observed. To do this we have to compute marginals.

The inference algorithm is expressed in terms of a **message passing** algorithm, where local messages are propagated around the graph. Two interconnected types of messages are considered:

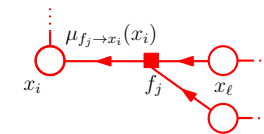
- Messages from variable nodes to factor nodes

$$\mu_{x_i \rightarrow f_j}(x_i) = \prod_{f_\ell \in \text{ne}(x_i) \setminus f_j} \mu_{f_\ell \rightarrow x_i}(x_i)$$



- Messages from factor nodes to variable nodes

$$\mu_{f_j \rightarrow x_i}(x_i) = \sum_{x_\ell \in \text{ne}(f_j) \setminus x_i} f_j \prod_{x_\ell \in \text{ne}(f_j) \setminus x_i} \mu_{x_\ell \rightarrow f_j}(x_\ell)$$



Sum-Product Algorithm

- Calculate messages from variable nodes to factor nodes

$$\mu_{x_i \rightarrow f_j}(x_i) = \prod_{f_\ell \in \text{ne}(x_i) \setminus f_j} \mu_{f_\ell \rightarrow x_i}(x_i)$$

- Calculate messages from factor nodes to variable nodes

$$\mu_{f_j \rightarrow x_i}(x_i) = \sum_{x_\ell \in \text{ne}(f_j) \setminus x_i} f_j \prod_{x_\ell \in \text{ne}(f_j) \setminus x_i} \mu_{x_\ell \rightarrow f_j}(x_\ell)$$

- Iterate messages until convergence. (Different iteration schemes can be designed.)
- After convergence, the marginals are calculated as

$$p(x_i) \propto \prod_{f_\ell \in \text{ne}(x_i)} \mu_{f_\ell \rightarrow x_i}(x_i)$$

In solving inference problems we are sooner or later typically faced with various **integration problems**, which tend to live in high dimensional spaces.

This holds for both Maximum likelihood and Bayesian approaches.

To be concrete, we have the following general problems

1. Expectation
2. Marginalization (includes normalization)

An expected value often provides an interesting (and interpretable) point estimate.

Computing an expectation amounts to solving the following integral

$$E[g(z)] = \int_{\mathcal{Z}} g(z)p(z)dz,$$

for some function $g : \mathcal{Z} \rightarrow \mathbb{R}^{n_g}$.

Example: Computing a point estimate ($g(x_t) = x_t$).

If we are interested in the properties of a stochastic variable z_1 and have access to the PDF $p(z_1, z_2 | y_{1:T})$, then we can marginalize out the variable z_2 , resulting in $p(z_1 | y_{1:T})$.

$$p(z_1 | y_{1:T}) = \int_{\mathcal{Z}_2} p(z_1, z_2 | y_{1:T}) dz_2$$

Examples: Normalization $p(y_{1:T}) = \int p(y_{1:T} | z)p(z)dz$ (used in e.g., empirical Bayes). As another example (in using the EM algorithm for nonlinear ML identification) we need the two-step smoothing densities $p(x_{t:t+1} | y_{1:T})$, whereas several smoothing algorithms provides the entire joint smoothing density $p(x_{1:T} | y_{1:T})$.

Many of the models we are currently interested in do **not** allow for closed form expressions. We are forced to approximations. Broadly speaking there are two classes,

1. **Deterministic analytical approximations:** Either approximate the model or restrict the solution to belong to an analytically tractable form. Examples, Laplace approximation, variational Bayes (VB), expectation propagation (EP).
2. **Stochastic approximations:** Keep the model and approximate the solution without imposing any restrictions other than the computational resources available.

In this lecture we are concerned with stochastic approximations (deterministic approximations are already covered).

Monte Carlo methods provides **computational solutions**, where the obtained accuracy is limited only by our computational resources.

Monte Carlo methods respects the model and the general solution. The approximation does not impose any restricting assumptions on the model or the solution.

(Very) restrictive assumption: Assume that we have N samples $\{z^i\}_{i=1}^N$ from the target density $\pi(z)$,

$$\hat{\pi}(z) = \sum_{i=1}^N \frac{1}{N} \delta_{z^i}(z)$$

Allows for the following approximation of the integral,

$$E[g(z)] = \int g(z)\pi(z)dz \approx \int g(z) \sum_{i=1}^N \frac{1}{N} \delta_{z^i}(z) dz = \frac{1}{N} \sum_{i=1}^N g(z^i)$$

$$" \int + \delta \rightarrow \sum "$$

The integral

$$I(g(z)) \triangleq E_{\pi(z)} [g(z)] = \int g(z)\pi(z)dz.$$

is approximated by

$$\hat{I}^M(g(z)) = \frac{1}{M} \sum_{i=1}^M g(z^i).$$

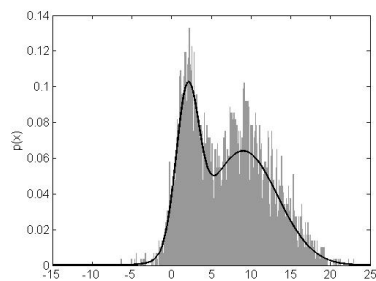
The strong law of large numbers tells us that

$$\hat{I}^M(g(z)) \xrightarrow{\text{a.s.}} I(g(z)), \quad M \rightarrow \infty,$$

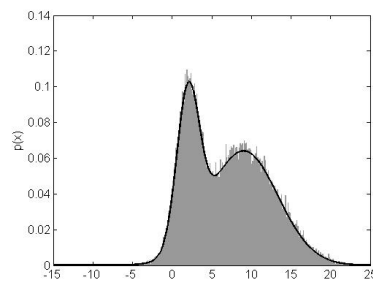
and the central limit theorem state that

$$\frac{\sqrt{M} \left(\hat{I}^M(g(z)) - I(g(z)) \right)}{\sigma_g} \xrightarrow{d} \mathcal{N}(0, 1), \quad M \rightarrow \infty.$$

$$\pi(z) = 0.3\mathcal{N}(z | 2, 2) + 0.7\mathcal{N}(z | 9, 19)$$



5 000 samples



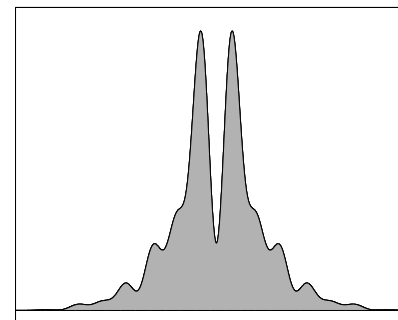
50 000 samples

Obvious problem: In general we are **not** able to directly sample from the density we are interested in.

Rejection sampling is a Monte Carlo method that produce i.i.d. samples from a target distribution

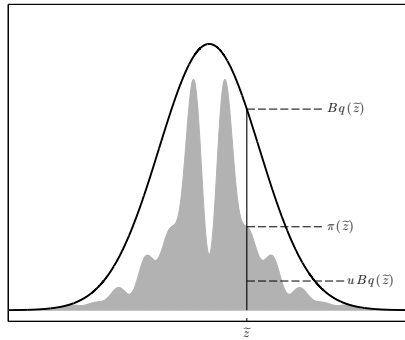
$$\pi(z) = \frac{\tilde{\pi}(z)}{C_\pi},$$

where $\tilde{\pi}(z)$ can be evaluated and C_π is a normalization constant.



Key idea: Generate random numbers uniformly from the area under the graph of the target distribution $\pi(z)$.

Just as hard as the original problem, but what if...



Generate a sample \tilde{z} from a proposal distribution $q(z)$ and a sample $u \sim \mathcal{U}[0, 1]$.

The sample \tilde{z} is then an i.i.d. sample from the target if

$$u \leq \frac{\tilde{\pi}(\tilde{z})}{Bq(\tilde{z})}$$



Assumptions:

1. It is easy to sample from $q(z)$.
2. There exists a constant B such that $\pi(z) \leq Bq(z), \forall z \in \mathcal{Z}$.
3. The support of $q(z)$ includes the support of $\pi(z)$, i.e., $q(z) > 0$ when $\pi(z) > 0$.

Algorithm 1 Rejection sampling (RS)

1. Sample $\tilde{z} \sim q(z)$.
2. Sample $u \sim \mathcal{U}[0, 1]$.
3. If $u \leq \frac{\tilde{\pi}(\tilde{z})}{Bq(\tilde{z})}$ accept \tilde{z} as a sample from $\pi(z)$ and go to 1.
4. Otherwise, reject \tilde{z} and go to 1.

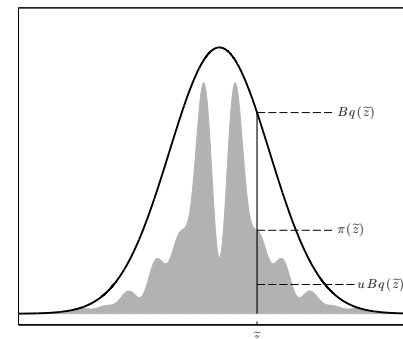


- The procedure can be used with multivariate densities in the same way.
- The rejection rate depends on B , choose B as small as possible, while still satisfying $\pi(z) \leq Bq(z), \forall z \in \mathcal{Z}$.
- Choosing a good proposal distribution $q(z)$ is very important.
- Rejection sampling is used to construct fast **particle smoothers** via backward simulation.



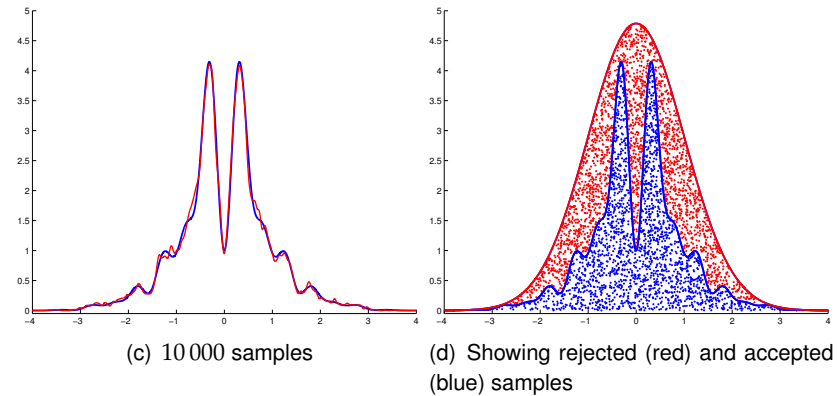
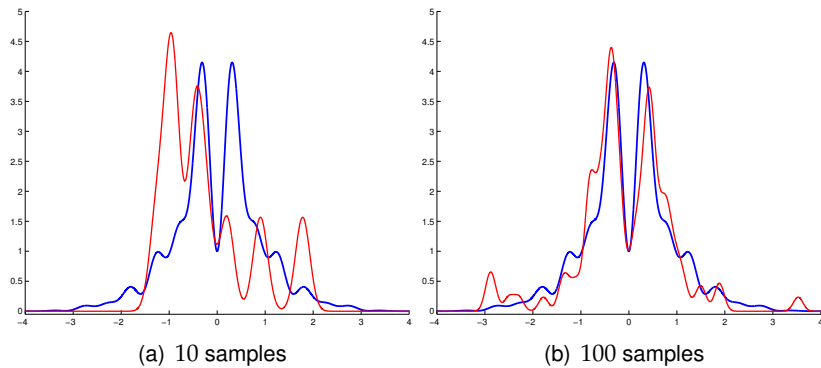
Task: Generate M i.i.d. samples from the following distribution,

$$\pi(z) = \frac{1}{C\pi} e^{-\frac{1}{2}z^2} \left(\sin(6z)^2 + 3 \cos(z)^2 \sin(4z)^2 + 1 \right),$$



Solution: Use rejection sampling where $q(z) = \mathcal{N}(z | 0, 1)$ and $B = 12$.





Algorithm 2 Importance sampler (IS)

1. Sample $z^i \sim q(z)$.
2. Compute the weights $\tilde{w}^i = \tilde{\pi}(z^i) / q(z^i)$.
3. Normalize the weights $w^i = \tilde{w}^i / \sum_{j=1}^N \tilde{w}^j$.

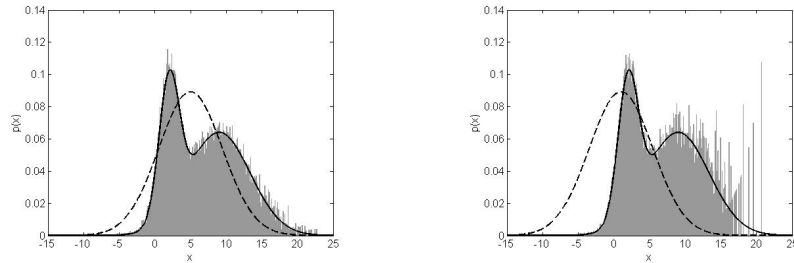
Each step is carried out for $i = 1, \dots, M$.

IS does not provide samples from the target density, but the samples $\{z^i\}_{i=1}^M$ together with the normalized weights $\{w^i\}_{i=1}^M$ provides an **empirical approximation** of the target density,

$$\hat{\pi}(z) = \sum_{i=1}^M w^i \delta_{z^i}(z).$$

When this approximation is inserted into $I(g(z)) = \int g(z)\pi(z)dz$ the resulting estimate is

$$\hat{I}^M(g(z)) = \sum_{i=1}^M w^i g(z^i).$$



$q_1(x) = \mathcal{N}(5, 20)$ (dashed curve) $q_2(x) = \mathcal{N}(1, 20)$ (dashed curve)

50 000 samples used in both simulations.

Lesson learned: It is important to be careful in selecting the importance density.

Let us revisit the same problem (scalar LGSS) used in illustrating the EM and the VB algorithms,

$$\begin{aligned} x_{t+1} &= \theta x_t + v_t, \\ y_t &= \frac{1}{2}x_t + e_t, \end{aligned} \quad \begin{pmatrix} v_t \\ e_t \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix} \right).$$

$p(x_1) = \mathcal{N}(x_1 | 0, 0.1)$. The true parameter value for θ is given by $\theta^* = 0.9$. We use $p(\theta) = \mathcal{N}(\theta | \mu_\theta, \sigma_\theta^2)$ as prior distribution for θ .

The identification problem is now to determine the parameter θ on the basis of the observations $y_{1:T}$ and the above model, using the IS algorithm. The result will be an estimate of the posterior distribution $p(\theta | y_{1:T})$.

The importance sampler will target

$$\pi(\theta) = p(\theta | y_{1:T}) = \frac{p(y_{1:T} | \theta)p(\theta)}{p(y_{1:T})} \propto p(y_{1:T} | \theta)p(\theta).$$

Choose the proposal distribution to be the same as the prior,

$$q(\theta) = \mathcal{N}(\theta | \mu_\theta, \sigma_\theta^2).$$

The importance weights are then computed according to

$$\tilde{w}^i = \frac{\tilde{\pi}(\theta^i)}{\tilde{q}(\theta^i)} = p(y_{1:T} | \theta^i), \quad i = 1, \dots, M,$$

i.e., the likelihood.

The Kalman filter straightforwardly allows us to evaluate the importance weights $\tilde{w}^i = p(y_{1:T} | \theta^i)$,

$$\begin{aligned} p(y_{1:T} | \theta^i) &= \prod_{t=1}^T p(y_t | y_{1:t-1}, \theta^i) = \prod_{t=1}^T \mathcal{N}(y_t | \hat{y}_{t|t-1}(\theta^i), S_{t|t-1}(\theta^i)), \\ \hat{y}_{t|t-1}(\theta^i) &= 0.5\hat{x}_{t|t-1}(\theta^i), \\ S_{t|t-1}(\theta^i) &= 0.5^2 P_{t|t-1}(\theta^i) + 0.1, \end{aligned}$$

where $\hat{x}_{t|t-1}(\theta^i)$ and $P_{t|t-1}(\theta^i)$ are provided by the Kalman filter.

Algorithm 3 Importance sampler targeting $p(\theta | y_{1:T})$

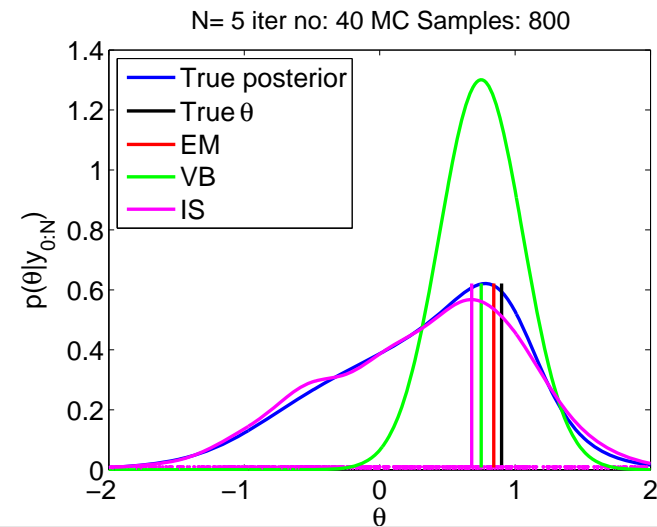
1. Generate M i.i.d. samples from the proposal distribution,

$$\theta^i \sim \mathcal{N}(\theta | \mu_{\theta}, \sigma_{\theta}^2).$$

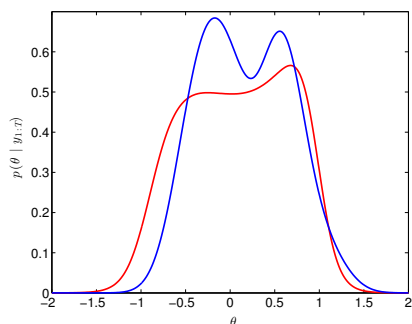
2. Compute importance weights $\tilde{w}^i = p(y_{1:T} | \theta^i)$.

3. Normalize the importance weights $w^i = \frac{\tilde{w}^i}{\sum_{j=1}^M \tilde{w}^j}$.

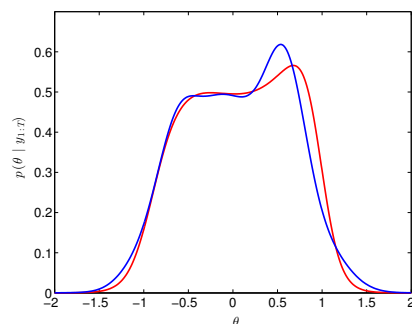
Each step is carried out for $i = 1, \dots, M$.



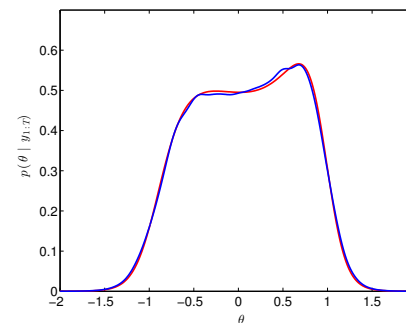
Using $T = 15$ measurements, $y_{1:15}$.



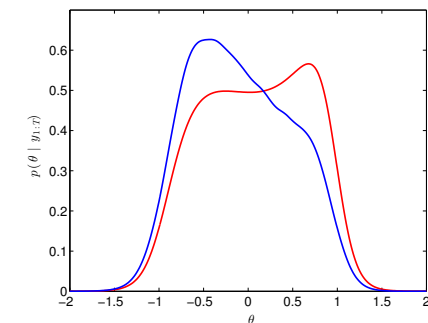
10 samples,
 $q(\theta) = \mathcal{N}(\theta | 0, 1.5^2)$



100 samples,
 $q(\theta) = \mathcal{N}(\theta | 0, 1.5^2)$



50 000 samples,
 $q(\theta) = \mathcal{N}(\theta | 0, 1.5^2)$



100 000 samples,
 $q(\theta) = \mathcal{N}(\theta | -0.5, 1^2)$

Lesson learned (again): Note the different proposal distributions used above. It is very important to be careful in selecting the importance density.

Markov Chain Monte Carlo



What's the point with the AR(1) example? (I/II)

34(48)

Task: How do we generate samples from the stationary distribution $\pi^s(x) = \mathcal{N}\left(x \mid 0, \frac{q}{1-a^2}\right)$? Put in other words, the target distribution $\pi(x)$ is given by the stationary distribution $\pi^s(x)$, i.e., $\pi(x) = \pi^s(x)$.

Two solutions for this problem:

1. Simulate sufficiently many samples from the Markov chain and discard the initial samples. The remaining samples will then be approximately distributed according to the target distribution (we just proved that x^t is distributed according to $\pi(x)$ for a large enough t).
2. We proved that the stationary distribution is Gaussian. Generate samples directly from this distribution.

Clearly a somewhat contrived example (obviously solution 2 is preferred), **but** solution 1 is a simple illustration of the strategy underlying all MCMC methods.



What's the point with the AR(1) example? (II/II)

35(48)

In the example, the Markov chain was fully specified and it was possible to explicitly compute the stationary distribution.

We are of course interested in the **reverse** situation, where we want to generate samples from a (typically rather complicated) target distribution $\pi(z)$.

The task is now to find a transition kernel such that the resulting Markov chain has the target distribution $\pi(z)$ as its stationary distribution.

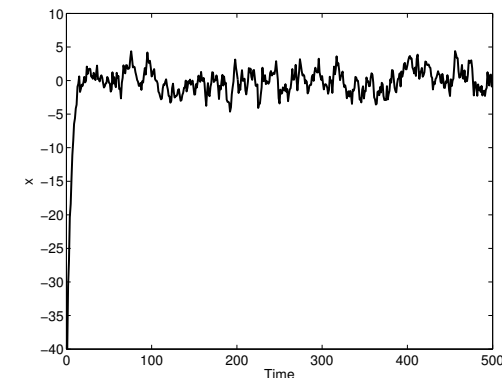
This can be done in many different ways and **constructive strategies** for doing this are provided by the Gibbs sampler and the Metropolis Hastings sampler.

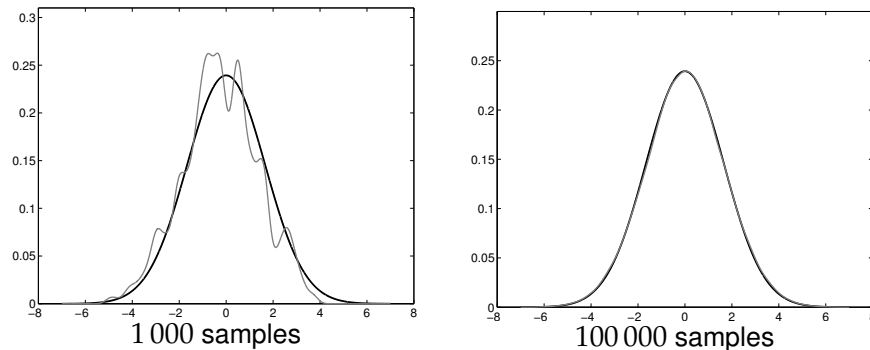


AR(1) example again (I/II)

36(48)

One realisation from $x^{t+1} = ax^t + v^t$ using $a = 0.8$, $v^t \sim \mathcal{N}(0, 1)$. The process is initialised in $x_0 = -40$.





The true stationary distribution is showed in black and the empirical histogram obtained by simulating the Markov chain $x^{t+1} = ax^t + v^t$ is plotted in gray.

The initial 1 000 samples are discarded (burn-in).

The Metropolis Hastings (MH) sampler provides a **constructive way** of producing a Markov chain that can be used to obtain samples approximately distributed according to the target distribution.

More pragmatically speaking, the MH sampler generates samples $\{z^m\}_{m=1}^M$ which can for example be used to approximately compute integrals.

The basic idea underlying the Metropolis Hastings sampler is surprisingly simple.

Starting from an initial state of the Markov chain z^1 , a new candidate sample z^* is generated using a **proposal distribution** $z^* \sim q(z | z^1)$.

This proposed sample z^* is then accepted with a certain probability, the so called **acceptance probability**

$$a(z^*, z^m) = \min \left(1, \frac{\pi(z^*)q(z^m | z^*)}{\pi(z^m)q(z^* | z^m)} \right).$$

If the sample is accepted, the new state of the Markov chain is set to the proposed sample $z^2 = z^*$, otherwise it is simply set to the previous value, $z^2 = z^1$.

Algorithm 4 Metropolis Hastings (MH) sampler

1. **Initialise:** Set the initial state of the Markov chain z^1 .

2. **For $m = 1$ to M , iterate:**

- a. Sample $z^* \sim q(z | z^m)$.
- b. Sample $u \sim \mathcal{U}[0, 1]$.
- c. Compute the acceptance probability

$$a(z^*, z^m) = \min(1, \alpha(z^*, z^m)), \text{ where } \alpha(z^*, z^m) = \frac{\pi(z^*)q(z^m | z^*)}{\pi(z^m)q(z^* | z^m)}$$

d. Set the next state z^{m+1} of the Markov chain according to

$$z^{m+1} = \begin{cases} z^* & u \leq a(z^*, z^m) \\ z^m & \text{otherwise} \end{cases}$$

Note that the MH sampler only requires two things,

1. It requires the definition of a proposal distribution $q(\cdot | \cdot)$ that can be used to generate candidate samples.
2. It must be possible to point-wise evaluate the target distribution up to a possibly unknown normalization factor.

Point-wise evaluation of the target density $\pi(\theta)$ for a specific $\theta = \bar{\theta}$

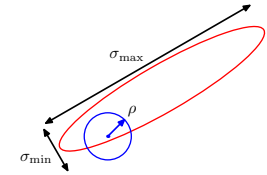
$$\pi(\bar{\theta}) = p(\bar{\theta} | y_{1:T}) = \frac{p(y_{1:T} | \bar{\theta})p(\bar{\theta})}{p(y_{1:T})}$$

- Suppose the target density over $x \in \mathbb{R}^2$ is

$$\pi(x) = \mathcal{N}\left(x; \begin{bmatrix} 4 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right)$$

- Choose the proposal density $q(x|z)$ as

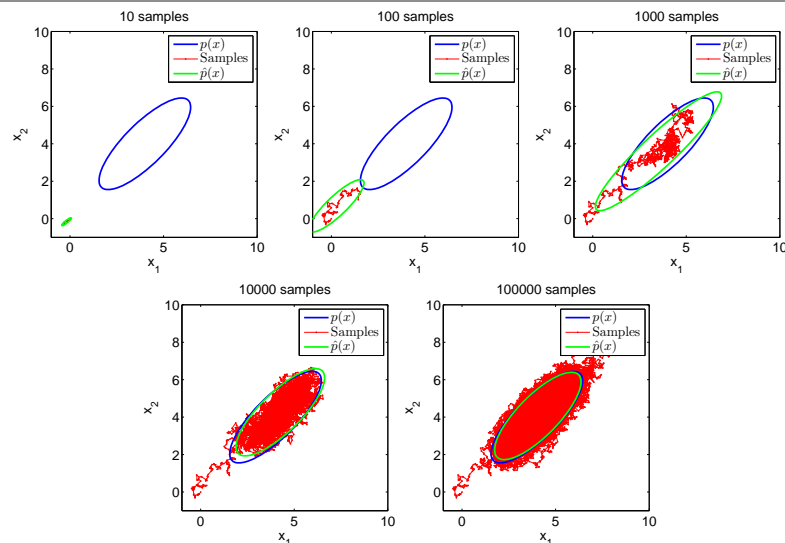
$$q(x|x^{(i-1)}) = \mathcal{N}\left(x; x^{(i-1)}, \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix}\right)$$



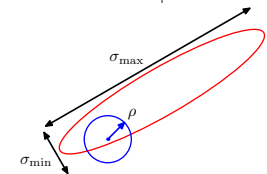
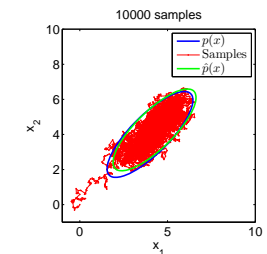
- Noticing that $q(x|x^{(i-1)}) = q(x^{(i-1)}|x)$ for all x , we have

$$\min\left(1, \frac{\pi(\bar{x})}{\pi(x^{(i-1)})} \frac{q(x^{(i-1)}|\bar{x})}{q(\bar{x}|x^{(i-1)})}\right) = \min\left(1, \frac{\pi(\bar{x})}{\pi(x^{(i-1)})}\right)$$

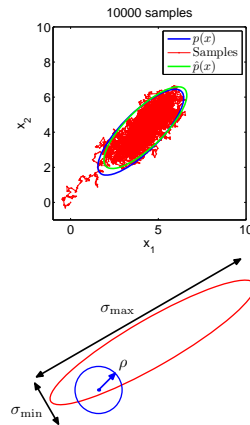
- This version of the Metropolis Hastings algorithm is called the Metropolis algorithm.



- The MH sampler generates samples that converge to the samples of a stationary distribution which is the target distribution.
- The time that passes before the samples start to represent the target density is called **burn-in period**.
- We generally have to use only the samples obtained after the burn-in period.
- Diagnosing convergence to the target distribution with MCMC algorithms is still an active area of research.
- After the burn-in period is over, the Markov chain is said to be **mixed**.



- Proposal selection is still an important problem.
- If the proposal is selected too narrow, then step-sizes get smaller and the burn-in period becomes longer.
- If the proposal is too wide, then the burn-in gets shorter, however, the acceptance rate is decreased significantly.



- Gibbs sampling is a special case of the Metropolis-Hastings algorithm where the proposal function is set to be the conditional distribution of the variables.
- It is especially useful when the dimension of the space to sample is very large e.g. images.
- Suppose, we are sampling in a two dimensional space $x = [x_1, x_2]^T$. Then the Gibbs sampler works as follows.

Gibbs Sampler for 2D

- Sample $x^{(1)} \sim q(\cdot)$.
- For $i = 2, 3, \dots$,
 - Sample $x_1^{(i)} \sim p(x_1 | x_2^{(i-1)})$.
 - Sample $x_2^{(i)} \sim p(x_2 | x_1^{(i)})$.
 - Set $x^{(i)} = [x_1^{(i)}, x_2^{(i)}]^T$.
- Note that due to the special proposal, a Gibbs sampler does not have an accept-reject step as MH.



Targeting $p(x_t | y_{1:t})$ in a nonlinear/non-Gaussian SSM using an importance sampler results in the particle filter (a member of the more general class of **Sequential Monte Carlo (SMC)** methods).

The Bayesian parameter inference problem in a general nonlinear/non-Gaussian SSM can be solved using the so called Particle Markov Chain Monte Carlo (PMCMC) methods. Here, an SMC algorithm is used as proposal to generate samples in an MCMC sampler.

Should you find this interesting I have a PhD course – *Computational inference in dynamical systems* – covering this material, see

users.isy.liu.se/rt/schon/course_CIDS.html



Monte Carlo Methods: Approximate inference tools using the samples from the target densities.

Basic Sampling Methods: The sampling methods to obtain independent samples from target densities. Though quite powerful, these would give bad results with high dimensions.

MCMC: Monte Carlo methods which produce dependent samples but more robust in high dimensions.

Metropolis Hastings Algorithm: The most well-known MCMC algorithm using arbitrary proposal densities.

Gibbs Sampler: A specific case of the MH sampler, which samples from conditionals iteratively and always accepts a new sample.

