

# **Unsupervised Learning**

## **Sampling Methods**

**Maneesh Sahani**

`maneesh@gatsby.ucl.ac.uk`

**Gatsby Computational Neuroscience Unit, and  
MSc in Intelligent Systems, Dept Computer Science  
University College London**

**Term 1, Autumn 2006**

# Integrals in Statistical Modelling

- **Parameter estimation**

$$\hat{\theta} = \operatorname{argmax}_{\theta} \int d\mathcal{Y} P(\mathcal{Y}|\theta) P(\mathcal{X}|\mathcal{Y}, \theta)$$

(or using EM)

$$\theta^{\text{new}} = \operatorname{argmax}_{\theta} \int d\mathcal{Y} P(\mathcal{Y}|\mathcal{X}, \theta^{\text{old}}) \log P(\mathcal{X}, \mathcal{Y}|\theta)$$

- **Prediction**

$$p(x|\mathcal{D}, m) = \int d\theta p(\theta|\mathcal{D}, m) p(x|\theta, \mathcal{D}, m)$$

- **Model selection or weighting** (by marginal likelihood)

$$p(\mathcal{D}|m) = \int d\theta p(\theta|m) p(\mathcal{D}|\theta, m)$$

These integrals are often intractable:

- **Analytic intractability:** integrals may not have closed form in non-linear, non-Gaussian models  $\Rightarrow$  numerical integration.
- **Computational intractability:** Numerical integral (or sum if  $\mathcal{Y}$  or  $\theta$  are discrete) may be exponential in data or model size.

# Examples of Intractability

- Bayesian marginal likelihood/model evidence for Mixture of Gaussians: exact computations are exponential in number of data points

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N) &= \int d\theta p(\theta) \prod_{i=1}^N \sum_{s_i} p(\mathbf{x}_i | s_i, \theta) p(s_i | \theta) \\ &= \sum_{s_1} \sum_{s_2} \dots \sum_{s_N} \int d\theta p(\theta) \prod_{i=1}^N p(\mathbf{x}_i | s_i, \theta) p(s_i | \theta) \end{aligned}$$

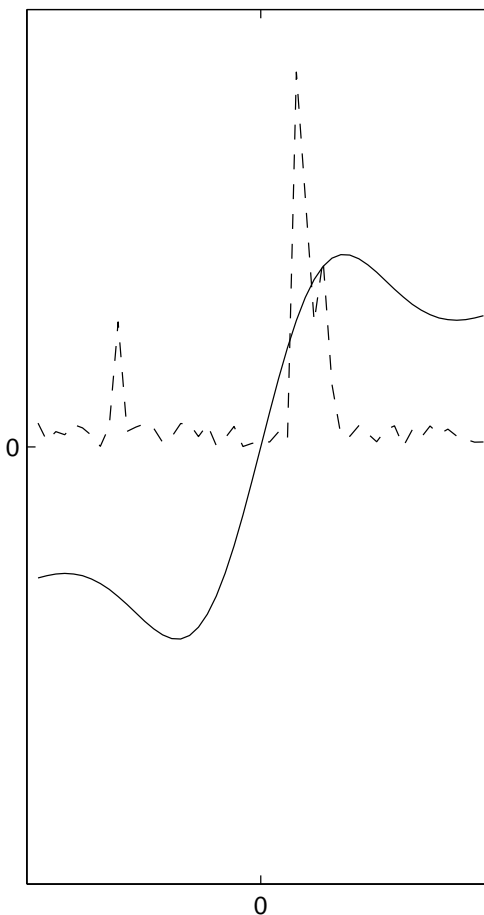
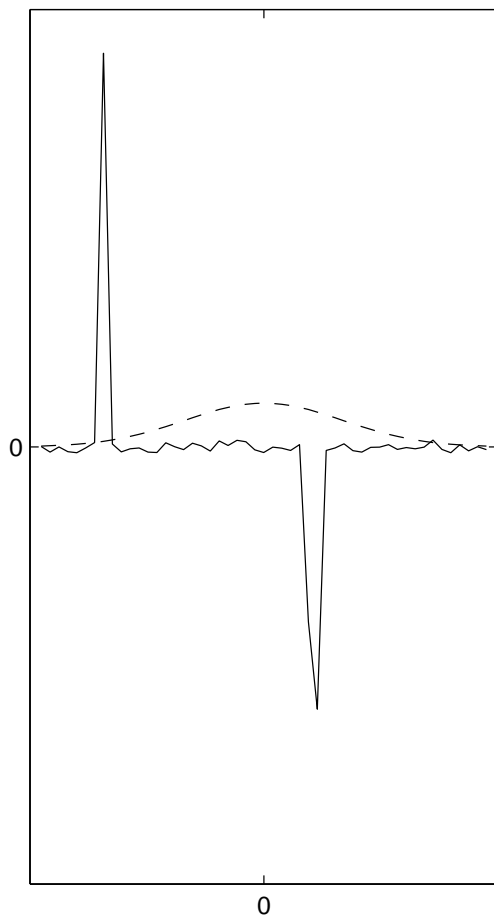
- Computing the conditional probability of a variable in a very large multiply connected directed graphical model:

$$p(x_i | X_j = a) = \sum_{\text{all settings of } \mathbf{y} \setminus \{i, j\}} p(x_i, \mathbf{y}, X_j = a) / p(X_j = a)$$

- Computing the hidden state distribution in a general nonlinear dynamical system

$$p(\mathbf{y}_t | \mathbf{x}_1, \dots, \mathbf{x}_T) \propto \int p(\mathbf{y}_t | \mathbf{y}_{t-1}) p(\mathbf{x}_t | \mathbf{y}_t) p(\mathbf{y}_{t-1} | \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) p(\mathbf{x}_{t+1}, \dots, \mathbf{x}_t | \mathbf{y}_t) d\mathbf{y}_{t-1}$$

# The integration problem



We commonly need to compute expected value integrals of the form:

$$\int F(x) p(x) dx,$$

where  $F(x)$  is some function of a random variable  $X$  which has probability density  $p(x)$ .

Three typical difficulties:

**left panel:** full line is some **complicated function**, dashed is density;

**right panel:** full line is some function and dashed is **complicated density**;

**not shown:** non-analytic integral (or sum) in **very many dimensions**

# Sampling Methods

The basic idea of sampling methods is to approximate an intractable integral or sum using **samples** from some distribution.

Monte Carlo Methods:

- Simple Monte Carlo Sampling
- Rejection Sampling
- Importance Sampling
- ...

Sequential Monte Carlo Methods:

- Particle Filtering
- ...

Markov Chain Monte Carlo Methods:

- Gibbs Sampling
- Metropolis Algorithm
- Hybrid Monte Carlo
- ...

# Simple Monte Carlo Sampling

**Idea:** Sample from  $p(x)$ , average values of  $F(x)$ .

Simple Monte Carlo:

$$\int F(x)p(x)dx \simeq \frac{1}{T} \sum_{t=1}^T F(x^{(t)}),$$

where  $x^{(t)}$  are (independent) samples drawn from  $p(x)$ .

[For example:  $x^{(t)} = G^{-1}(u^{(t)})$  with  $u \sim \text{Uniform}[0, 1]$  and  $G(x) = \int_{-\infty}^x p(x')dx'$ ]

## Attractions:

- unbiased
- variance goes as  $1/T$ , independent of dimension!

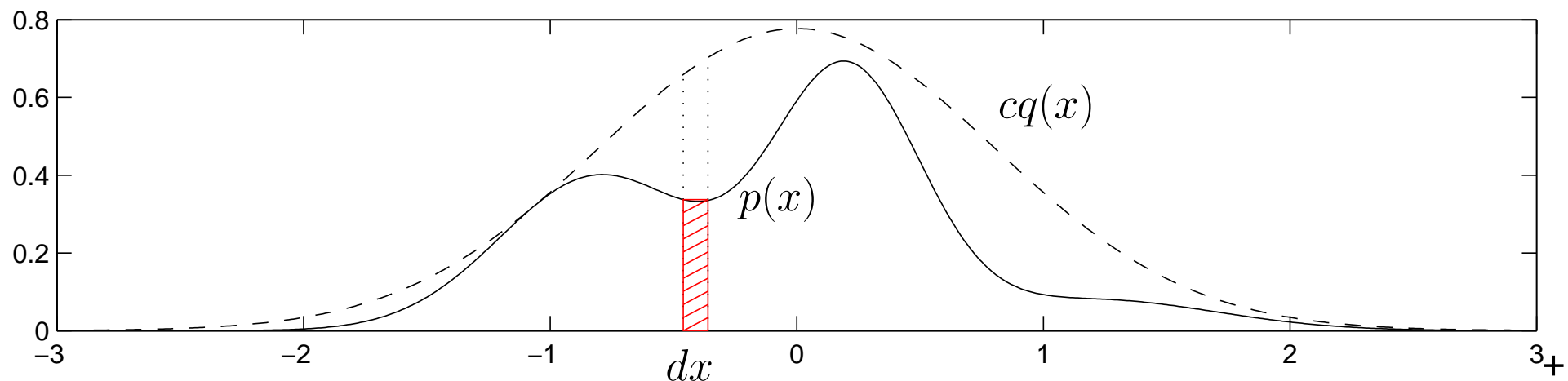
## Problems:

- it may be difficult or impossible to obtain the samples directly from  $p(x)$
- regions of high density  $p(x)$  may not correspond to regions where  $F(x)$  varies a lot (thus each evaluation might have very high variance).

# Rejection Sampling

**Idea:** sample from an upper bound on  $p(x)$ , rejecting some samples.

- Find a distribution  $q(x)$  and a constant  $c$  such that  $\forall x, p(x) \leq cq(x)$
- Sample  $x^*$  from  $q(x)$  and accept  $x^*$  with probability  $p(x^*)/(cq(x^*))$ .
- Use accepted points as in simple Monte Carlo:  $\sum_{t=1}^T F(x^{(t)})$



If  $y \sim \text{Uniform}[0, cq(x^*)]$ , we accept  $x^*$  if  $y \leq p(x^*)$ . Thus the probability of a point falling in the box =  $q(x)dx * p(x)/cq(x) = p(x)/c$ .

**Problem:** it may be difficult to find a  $q(x)$  with a small  $c$  which is easy to sample from  $\Rightarrow$  lots of wasted area.

## Examples:

- Compute  $P(X_i = b | X_j = a)$  in a directed graphical model: sample from  $P(X)$ , reject if  $X_j \neq a$ , averaging the indicator function  $I(X_i = b)$
- Compute  $E(x^2 | x > 4)$  for  $x \sim \mathcal{N}(0, 1)$

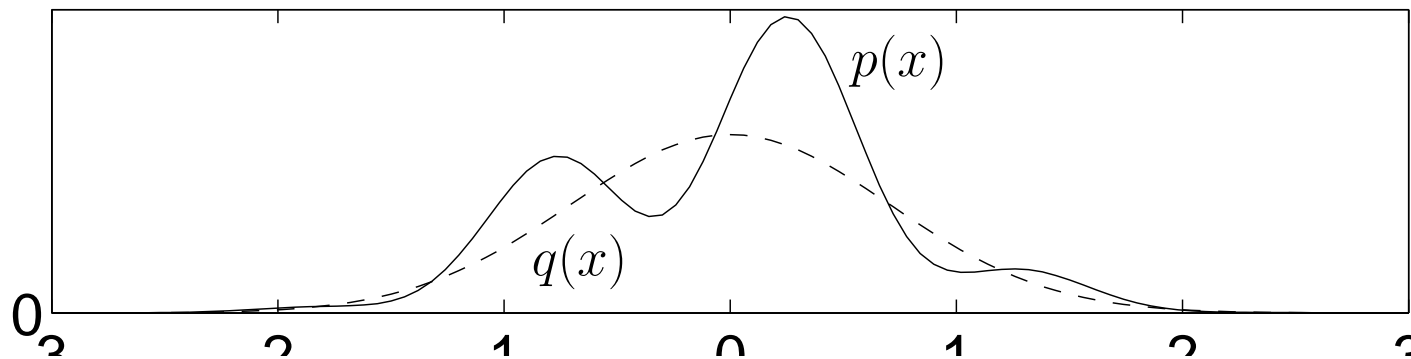
# Importance Sampling

**Idea:** Sample from a **different** distribution  $q(x)$  and weight those samples by  $p(x)/q(x)$

Sample  $x^{(t)}$  from  $q(x)$ :

$$\int F(x)p(x)dx = \int F(x)\frac{p(x)}{q(x)}q(x)dx \simeq \frac{1}{T} \sum_{t=1}^T F(x^{(t)})\frac{p(x^{(t)})}{q(x^{(t)})},$$

where  $q(x)$  is non-zero wherever  $p(x)$  is; weights  $w^{(t)} \equiv p(x^{(t)})/q(x^{(t)})$



**Attraction:** unbiased; no need for upper bound (cf rejection sampling).

**Problems:** it may be difficult to find a suitable  $q(x)$ . Monte Carlo average may be dominated by few samples (high variance); or none of the high weight samples may be found!



# Analysis of Importance Sampling

Weights:

$$w^{(t)} \equiv \frac{p(x^{(t)})}{q(x^{(t)})}$$

Define a weighting *function*  $w(x) = p(x)/q(x)$ .

Importance sample is unbiased:

$$\mathbb{E}_q [w(x)F(x)] = \int q(x)w(x)F(x)dx = \int p(x)F(x)dx$$

$$\mathbb{E}_q [w(x)] = \int q(x)w(x)dx = 1$$

The weights have variance  $\text{Var} [w(x)] = \mathbb{E}_q [w(x)^2] - 1$ , with:

$$\mathbb{E}_q [(w(x)^2)] = \int \frac{p(x)^2}{q(x)^2}q(x)dx = \int \frac{p(x)^2}{q(x)}dx$$

- How does variance effect the estimated integral?
- How does it relate to the *effective number of samples*?
- What happens if  $p(x) = \mathcal{N}(0, \sigma_p^2)$  and  $q(x) = \mathcal{N}(0, \sigma_q^2)$ ?

# Sampling - Importance Resampling (SIR)

Another (approximate) approach is to **resample** from the importance-weighted samples:

- Sample  $\xi^{(s)} \sim q(x)$ , and calculate importance weights  $w^{(s)} = p(\xi^{(s)})/q(\xi^{(s)})$ .
- Define  $\tilde{q}(x) = \sum_{s=1}^S w^{(s)} \delta(x - \xi^{(s)}) / \sum_{s=1}^S w^{(s)}$ .
- Resample  $x^{(t)} \sim \tilde{q}(x)$ .

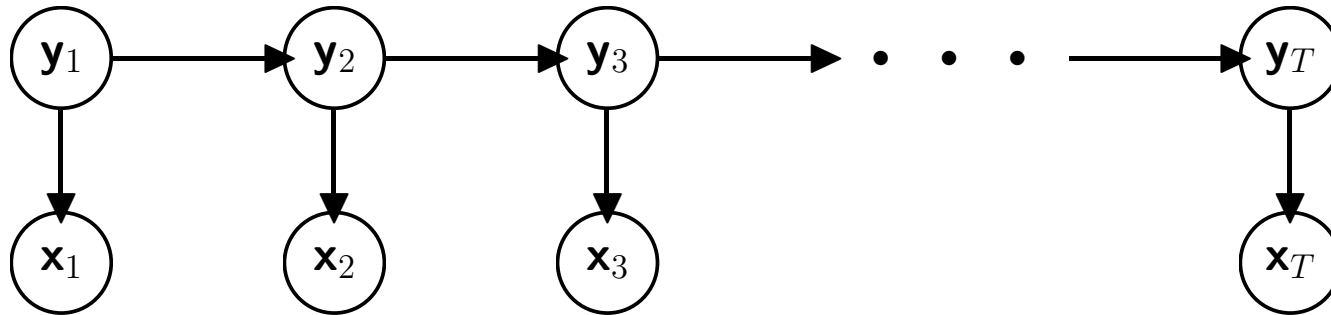
Then,

$$\begin{aligned} E_x[F(x)] &= \int dx F(x) \tilde{q}(x) \\ &= \int dx F(x) \frac{\sum_{s=1}^S w^{(s)} \delta(x - \xi^{(s)})}{\sum_{s=1}^S w^{(s)}} \\ &= \frac{\sum_{s=1}^S w^{(s)} F(\xi^{(s)})}{\sum_{s=1}^S w^{(s)}} \end{aligned}$$

but the expected value of this expression with respect to  $\xi \sim q$  is only correct as  $S \rightarrow \infty$ .

By itself, SIR looks unattractive relative to IS due to this bias. But we sometimes really do need *samples* (*i.e.* a picture of the distribution) rather than just expectations. E.g., if propagating beliefs.

# Sequential Monte Carlo



Suppose we want to compute  $p(\mathbf{y}_t | \mathbf{x}_1 \dots \mathbf{x}_t)$  in a non-linear ssm. We have

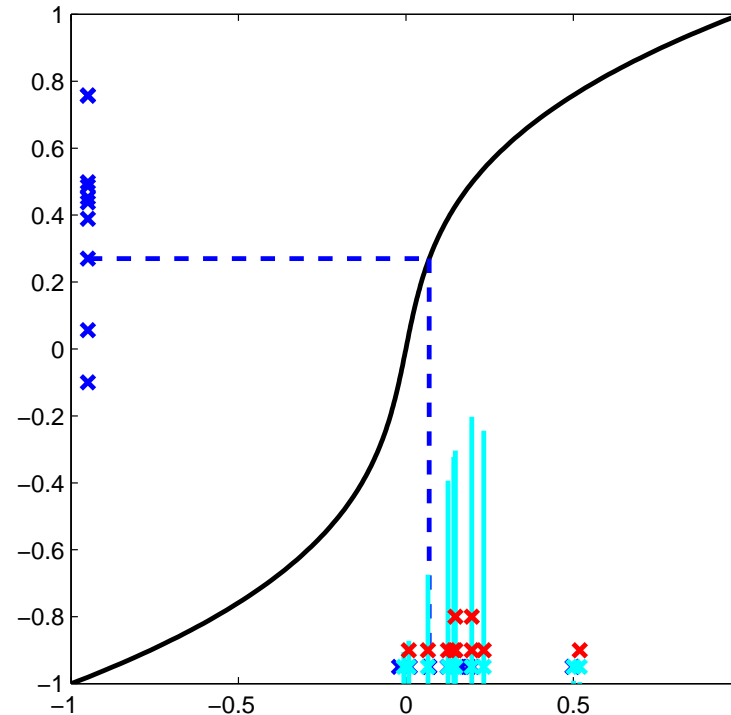
$$\begin{aligned} p(\mathbf{y}_t | \mathbf{x}_1 \dots \mathbf{x}_t) &\propto \int d\mathbf{y}_{t-1} p(\mathbf{y}_t \mathbf{y}_{t-1} \mathbf{x}_t | \mathbf{x}_1 \dots \mathbf{x}_{t-1}) \\ &= \int d\mathbf{y}_{t-1} p(\mathbf{x}_t | \mathbf{y}_t) p(\mathbf{y}_t | \mathbf{y}_{t-1}) p(\mathbf{y}_{t-1} | \mathbf{x}_1 \dots \mathbf{x}_{t-1}) \end{aligned}$$

If we have samples  $\mathbf{y}_{t-1}^{(s)} \sim p(\mathbf{y}_{t-1} | \mathbf{x}_1 \dots \mathbf{x}_{t-1})$  we can recurse (approximately):

- draw  $\mathbf{y}_t^{(s)} \sim p(\mathbf{y}_t | \mathbf{y}_{t-1}^{(s)})$
- calculate (unnormalised) weights  $w_t^{(s)} = p(\mathbf{x}_t | \mathbf{y}_t^{(s)})$ .
- resample  $\mathbf{y}_t^{(s')} \sim \sum_{s=1}^S w_t^{(s)} \delta(\mathbf{y} - \mathbf{y}_t^{(s)}) / \sum_{s=1}^S w_t^{(s)}$

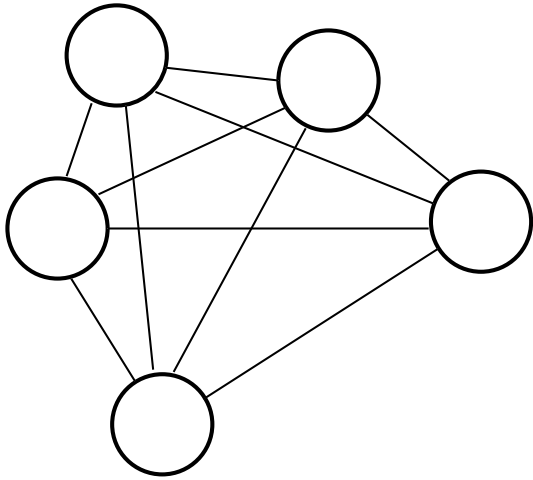
This is called Particle Filtering (this version, with  $q = p(\mathbf{y}_t | \mathbf{y}_{t-1}^{(s)})$  is also called a “bootstrap filter” or “condensation” algorithm).

# Particle Filtering



- Could avoid resampling by propagating weights. However variance in weights accumulates. Resampling helps eliminate unlikely particles.
- Can trigger resamples conditioned on variance – “stratified resampling”.
- Can use better proposal ( $q$ ) distributions (including  $p(\mathbf{y}_t | \mathbf{x}_t \mathbf{y}_{t-1}^{(s)})$  if available).
- Particle *smoothing* is possible, but often inaccurate. Difficult to create a good proposal.
- EM learning is not easy because of smoothing problems and also obtaining joint on  $(\mathbf{y}_{t-1}, \mathbf{y}_t)$ . Often use dual formulation.
- Widely used in engineering tracking applications, where filtering is most appropriate.
- Many variants ...

# Learning in Boltzmann Machines



$$\log P(\mathbf{s}^V \mathbf{s}^H | W, \mathbf{b}) = \sum_{ij} W_{ij} s_i s_j - \sum_i b_i s_i - \log Z$$

with  $Z = \sum_{\mathbf{s}} e^{\sum_{ij} W_{ij} s_i s_j - \sum_i b_i s_i}$

Generalised (gradient M-step) EM requires parameter step

$$\Delta W_{ij} \propto \frac{\partial}{\partial W_{ij}} \langle \log P(\mathbf{s}^V \mathbf{s}^H | W, \mathbf{b}) \rangle_{P(\mathbf{s}^H | \mathbf{s}^V)}$$

Write  $\langle \rangle_c$  (**clamped**) for expectations under  $P(\mathbf{s} | \mathbf{s}^V)$  (with delta function  $P(\mathbf{s}^V | \mathbf{s}^V)$ ). Then

$$\begin{aligned} \Delta W_{ij} &\propto \frac{\partial}{\partial W_{ij}} \left[ \sum_{ij} W_{ij} \langle s_i s_j \rangle_c - \sum_i b_i \langle s_i \rangle_c - \log Z \right] \\ &= \langle s_i s_j \rangle_c - \frac{\partial}{\partial W_{ij}} \log Z \\ &= \langle s_i s_j \rangle_c - \frac{1}{Z} \frac{\partial}{\partial W_{ij}} \sum_{\mathbf{s}} e^{\sum_{ij} W_{ij} s_i s_j - \sum_i b_i s_i} \\ &= \langle s_i s_j \rangle_c - \sum_{\mathbf{s}} \frac{1}{Z} e^{\sum_{ij} W_{ij} s_i s_j - \sum_i b_i s_i} s_i s_j \\ &= \langle s_i s_j \rangle_c - \sum_{\mathbf{s}} P(\mathbf{s} | W, \mathbf{b}) s_i s_j = \langle s_i s_j \rangle_c - \langle s_i s_j \rangle_u \end{aligned}$$

with  $\langle \rangle_u$  (**unclamped**) an expectation under the current joint distribution.

# Learning in Boltzmann Machines

How do we find the required expectations?

- **Junction tree** is generally intractable in all but the sparsest nets (triangulation of loops makes cliques grow very large).
- **Loopy belief propagation** fails in nets with strong correlations.
- **Rejection and Importance sampling** require proposal distributions, which are difficult to come by.
- **Mean-field methods** are possible, but approximate.

What is easy is **conditional** sampling. Given settings of all nodes in the Markov blanket of  $s_i$  can easily sample  $s_i$ . This suggests an iterative sampling algorithm:

- Choose variable settings randomly, perhaps from best available importance-like distribution. (Set any clamped nodes to clamped values).
- Cycle through (unclamped)  $s_i$ , choosing  $s_i \sim P(s_i | \mathbf{s}_{\setminus i})$ .

After enough samples, we might expect to reach the correct distribution.

This is an example of **Gibbs Sampling**.

# Markov chain Monte Carlo (MCMC) methods

Assume we are interested in drawing samples from some desired distribution  $p^*(x)$ .

We define a Markov chain:

$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow x_4 \rightarrow x_5 \dots$$

where  $x_0 \sim p_0(x)$ ,  $x_1 \sim p_1(x)$ , etc, with the property that:

$$p_t(x') = \sum_x p_{t-1}(x)T(x \rightarrow x')$$

where  $T(x \rightarrow x') = p(X_t = x' | X_{t-1} = x)$  is the **Markov chain transition probability** from  $x$  to  $x'$ .

We say that  $p^*(x)$  is an **invariant (or stationary) distribution** of the Markov chain defined by  $T$  iff:

$$p^*(x') = \sum_x p^*(x)T(x \rightarrow x') \quad \forall x'$$

# Markov chain Monte Carlo (MCMC) methods

We have a Markov chain  $x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots$  where  $x_0 \sim p_0(x)$ ,  $x_1 \sim p_1(x)$ , etc, with the property that:

$$p_t(x') = \sum_x p_{t-1}(x) T(x \rightarrow x')$$

where  $T(x \rightarrow x')$  is the Markov chain transition probability from  $x$  to  $x'$ .

A useful condition that implies invariance of  $p^*(x)$  is **detailed balance**:

$$p^*(x') T(x' \rightarrow x) = p^*(x) T(x \rightarrow x')$$

We wish to find **ergodic** Markov chains, which converge to a unique stationary distribution (also called an *equilibrium distribution*) regardless of the initial conditions  $p_0(x)$ :

$$\lim_{t \rightarrow \infty} p_t(x) = p^*(x)$$

A sufficient condition for the Markov chain to be ergodic is that

$$T^k(x \rightarrow x') > 0 \text{ for all } x \text{ and } x' \text{ where } p^*(x') > 0.$$

That is, if the equilibrium distribution gives non-zero probability to state  $x'$ , then the Markov chain should be able to reach  $x'$  from any  $x$  after some finite number of steps,  $k$ .



# Gibbs Sampling

A method for sampling from a multivariate distribution,  $p(\mathbf{x})$

**Idea:** sample from the conditional of each variable given the settings of the other variables.

Repeatedly:

- 1) pick  $i$  (either at random or in turn)
- 2) replace  $x_i$  by a sample from the conditional distribution

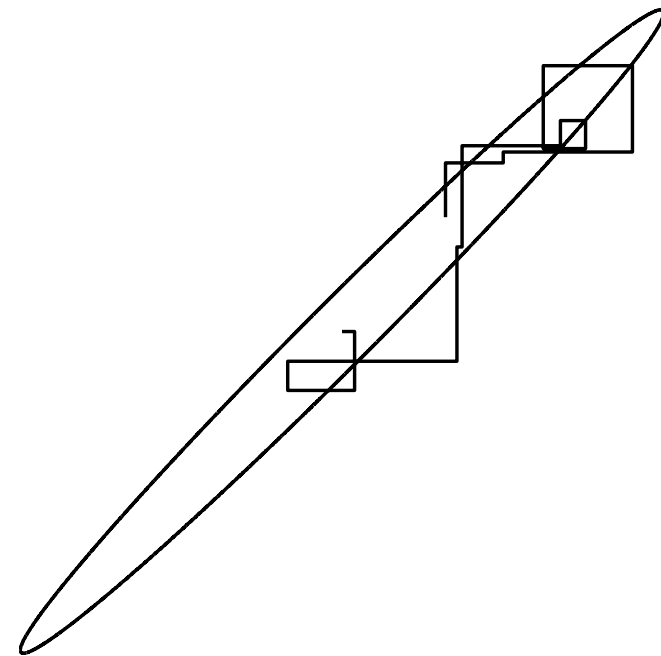
$$p(x_i | \mathbf{x}_{\setminus i}) = p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

Gibbs sampling is feasible if it is easy to sample from the conditional probabilities.

This creates a Markov chain

$$\mathbf{x}^{(1)} \rightarrow \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(3)} \rightarrow \dots$$

Under some (mild) conditions, the **equilibrium distribution**, i.e.  $p(\mathbf{x}^{(\infty)})$ , of this Markov chain is  $p(\mathbf{x})$ .



Example: 20 (half-) iterations of Gibbs sampling on a bivariate Gaussian

# Detailed balance for Gibbs sampling

We can show that Gibbs sampling has the right stationary distribution  $p(\mathbf{x})$  by showing that the **detailed balance** condition is met.

The transition probabilities are given by:

$$T(\mathbf{x} \rightarrow \mathbf{x}') = \pi_i p(x'_i | \mathbf{x}_{\setminus i})$$

where  $\pi_i$  is the probability of choosing to update the  $i$ th variable (to handle rotation updates instead of random ones, we need to consider transitions due to one full sweep).

Then we have:

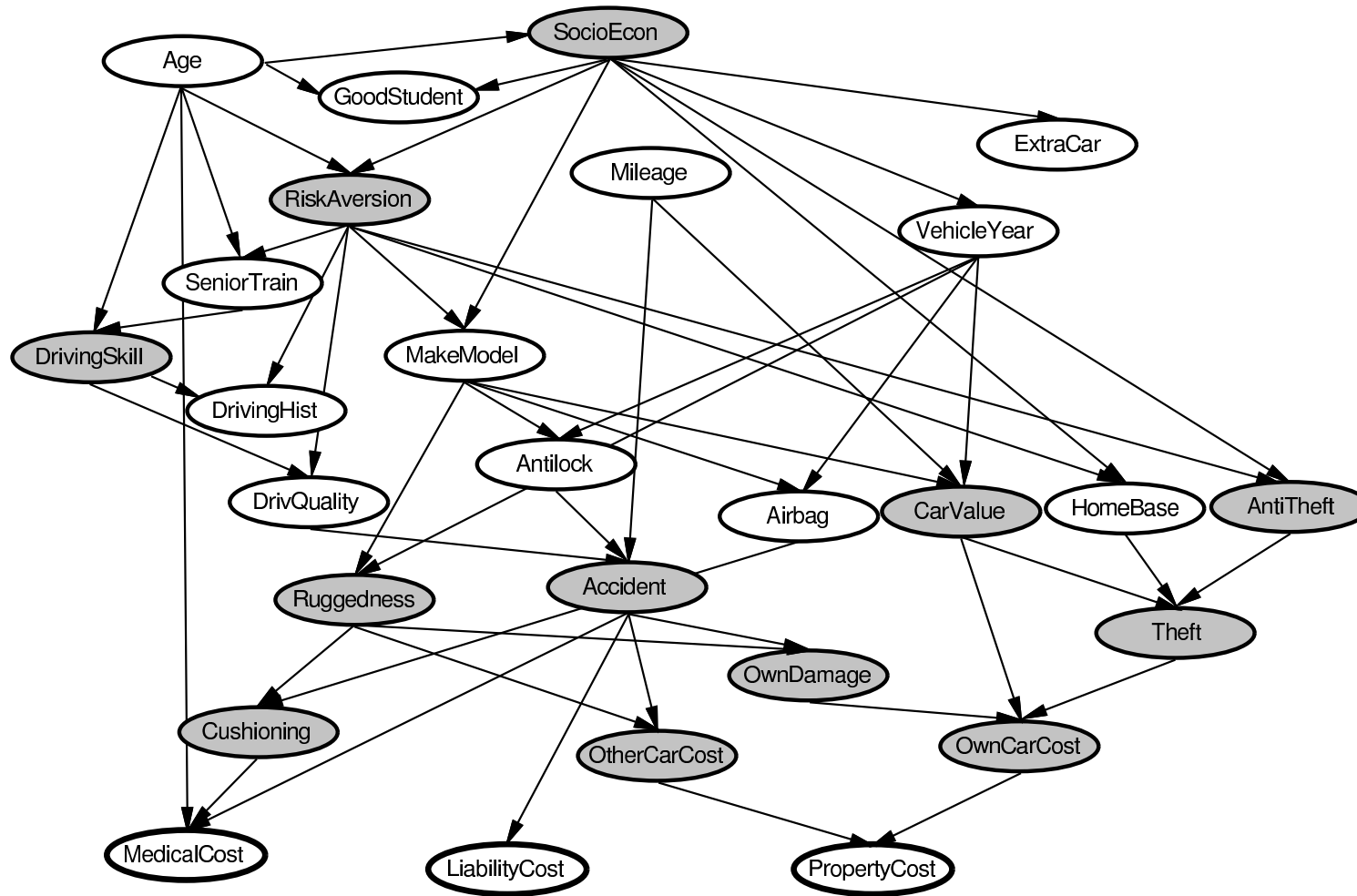
$$p(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{x}') = \pi_i p(x'_i | \mathbf{x}_{\setminus i}) \underbrace{p(x_i | \mathbf{x}_{\setminus i}) p(\mathbf{x}_{\setminus i})}_{p(\mathbf{x})}$$

and

$$p(\mathbf{x}')T(\mathbf{x}' \rightarrow \mathbf{x}) = \pi_i p(x_i | \mathbf{x}'_{\setminus i}) \underbrace{p(x'_i | \mathbf{x}'_{\setminus i}) p(\mathbf{x}'_{\setminus i})}_{p(\mathbf{x}')}$$

But  $\mathbf{x}'_{\setminus i} = \mathbf{x}_{\setminus i}$  so detailed balance holds.

# Gibbs Sampling in Graphical Models



Initialize all variables to some settings.

Sample each variable conditional on other variables.

The BUGS software implements this algorithm for a variety of graphical models.

# The Metropolis-Hastings algorithm

Gibbs sampling can be slow ( $p(x_i)$  may be well determined by  $\mathbf{x}_{\setminus i}$ ). Global transition might be better. (Also conditionals might be difficult for parameter integrals.)

**Idea:** Propose a change to current state; accept or reject.  
(A kind of rejection sampling)

**Each step:** Starting from the current state  $\mathbf{x}$ ,

1. Propose a new state  $\mathbf{x}'$  using a **proposal distribution**

$$S(\mathbf{x}'|\mathbf{x}) = S(\mathbf{x} \rightarrow \mathbf{x}').$$

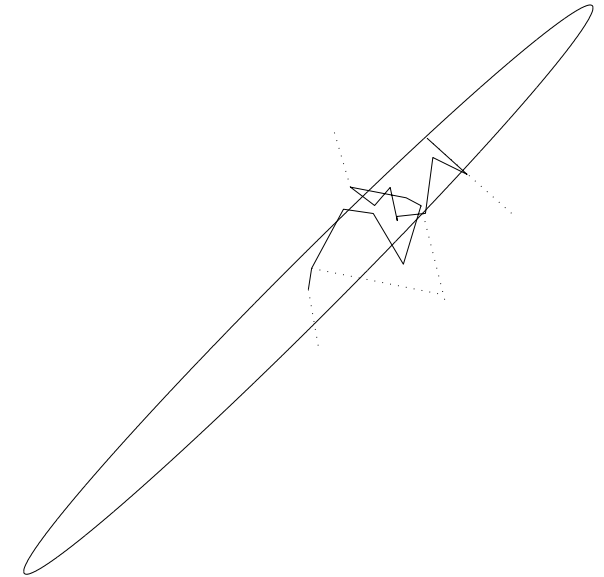
2. Accept the new state with probability:

$$\min\left(1, p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x})/p(\mathbf{x})S(\mathbf{x} \rightarrow \mathbf{x}')\right);$$

3. Otherwise retain the old state (or try again).

**Example:** 20 iterations of global metropolis sampling from bivariate Gaussian; rejected proposals are dotted.

- Metropolis algorithm was symmetric  $S(\mathbf{x}'|\mathbf{x}) = S(\mathbf{x}|\mathbf{x}')$ . Hastings generalised.
- **Local** (changing one  $x_i$ ) vs **global** (changing all  $\mathbf{x}$ ) proposal distributions.
- Efficiency dictated by rejection rate (and step size).
- May **adapt**  $S(\mathbf{x} \rightarrow \mathbf{x}')$  to balance these, but stationarity only holds once  $S$  is fixed.
- Note, we need only to compute ratios of probabilities (no normalizing constants).



# Detailed balance for Metropolis-Hastings

Analyse the case where we don't move on rejection:

$$T(\mathbf{x} \rightarrow \mathbf{x}') = S(\mathbf{x} \rightarrow \mathbf{x}') \min \left( 1, \frac{p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x})}{p(\mathbf{x})S(\mathbf{x} \rightarrow \mathbf{x}')} \right)$$

with  $T(\mathbf{x} \rightarrow \mathbf{x}')$  the expected rejection probability.

Without loss of generality we assume  $p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x}) \leq p(\mathbf{x})S(\mathbf{x} \rightarrow \mathbf{x}')$ .

Then

$$\begin{aligned} p(\mathbf{x})T(\mathbf{x} \rightarrow \mathbf{x}') &= p(\mathbf{x})S(\mathbf{x} \rightarrow \mathbf{x}') \cdot \frac{p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x})}{p(\mathbf{x})S(\mathbf{x} \rightarrow \mathbf{x}')} \\ &= p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x}) \end{aligned}$$

and

$$\begin{aligned} p(\mathbf{x}')T(\mathbf{x}' \rightarrow \mathbf{x}) &= p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x}) \cdot 1 \\ &= p(\mathbf{x}')S(\mathbf{x}' \rightarrow \mathbf{x}) \end{aligned}$$

# Gibbs distributions, temperature and annealing

Very often, need to sample from unnormalised distribution:

$$p(\mathbf{x}) = \frac{1}{Z} e^{-E(x)}$$

MCMC sampling works well in this setting (for Gibbs sampling, usually possible to normalise conditional), but may mix slowly.

Often useful to introduce **temperature**  $1/\beta$ :

$$p_\beta(\mathbf{x}) = \frac{1}{Z} e^{-\beta E(x)}$$

When  $\beta \rightarrow 0$  (temperature  $\rightarrow \infty$ ) all states are equally likely: easy to sample and mix.

As  $\beta \rightarrow 1$ ,  $p_\beta \rightarrow p$ .

**Idea:** start chain with  $\beta$  small and gradually increase to 1. **Simulated annealing** (can be used for optimisation by taking  $\beta \rightarrow \infty$ ).

For sampling, note that stationarity only holds once  $\beta$  is fixed, so need to run long enough to mix.

# Hybrid Monte Carlo: overview

The typical distance traveled by a random walk in  $n$  steps is proportional to  $\sqrt{n}$ . We want to seek regions of high probability while **avoiding random walk behavior**.

Assume that we wish to sample from  $p(\mathbf{x})$  while avoiding random walk behaviour. If we can compute derivatives of  $p(\mathbf{x})$  with respect to  $\mathbf{x}$ , this is *useful information* and we should be able to use it to draw samples better.

**Hybrid Monte Carlo:** We think of a fictitious physical system with a particle which has position  $\mathbf{x}$  and momentum  $\mathbf{v}$ . We will design a sampler which avoids random walks in  $\mathbf{x}$  by simulating a dynamical system.

We simulate the dynamical system in such a way that the marginal distribution of positions,  $p(\mathbf{x})$ , ignoring the momentum variables corresponds to the desired distribution.

# Hybrid Monte Carlo: the dynamical system

In the physical system, positions  $\mathbf{x}$  corresponding to random variables of interest are augmented by momentum variables  $\mathbf{v}$ :

$$\begin{aligned} p(\mathbf{x}, \mathbf{v}) &\propto \exp(-H(\mathbf{x}, \mathbf{v})) & H(\mathbf{x}, \mathbf{v}) &= E(\mathbf{x}) + K(\mathbf{v}) \\ E(\mathbf{x}) &= -\log p(\mathbf{x}) & K(\mathbf{v}) &= \frac{1}{2} \sum_i v_i^2 \end{aligned}$$

Importantly, note that  $\int p(\mathbf{x}, \mathbf{v}) d\mathbf{v} = p(\mathbf{x})$ , the desired distribution and  $p(\mathbf{v}) = N(0, I)$ . We think of  $E(\mathbf{x})$  as the **potential energy** of being in state  $\mathbf{x}$ , and  $K(\mathbf{v})$  as the **kinetic energy** associated with momentum  $\mathbf{v}$ . We assume “mass” = 1, so momentum = velocity.

The physical system evolves at constant **total energy**  $H$  according to Hamiltonian dynamics:

$$\frac{dx_i}{dt} = \frac{\partial H}{\partial v_i} = v_i \quad \frac{dv_i}{dt} = -\frac{\partial H}{\partial x_i} = -\frac{\partial E}{\partial x_i}.$$

The first equation says derivative of position is velocity. The second equation says that the system accelerates in the direction that decreases potential energy.

*Think of a ball rolling on a frictionless hilly surface.*



# Hybrid Monte Carlo: how to simulate the dynamical system

We can simulate the above differential equations by discretising time and running some difference equations on a computer. This introduces small (we hope) errors. (The errors we care about are errors which change the total energy—we will correct for these by occasionally rejecting moves that change the energy.)

A good way to simulate this is using **leapfrog simulation**. We take  $L$  discrete steps of size  $\epsilon$  to simulate the system evolving for  $L\epsilon$  time:

$$\begin{aligned}\hat{v}_i(t + \frac{\epsilon}{2}) &= \hat{v}_i(t) - \frac{\epsilon}{2} \frac{\partial E(\hat{x}(t))}{\partial x_i} \\ \hat{x}_i(t + \epsilon) &= \hat{x}_i(t) + \epsilon \frac{\hat{v}_i(t + \frac{\epsilon}{2})}{m_i} \\ \hat{v}_i(t + \epsilon) &= \hat{v}_i(t + \frac{\epsilon}{2}) - \frac{\epsilon}{2} \frac{\partial E(\hat{x}(t + \epsilon))}{\partial x_i}\end{aligned}$$

# Hybrid Monte Carlo: properties of the dynamical system

Hamiltonian dynamics has the following important properties:

- 1) preserves total energy,  $H$ ,
- 2) is reversible in time
- 3) preserves phase space volumes (Liouville's theorem)

The leapfrog discretisation only approximately preserves the total energy  $H$ , and

- 1) is reversible in time
- 2) preserves phase space volume

The dynamical system is simulated using the leapfrog discretisation and the new state is used as a proposal in the Metropolis algorithm to eliminate the bias caused by the leapfrog approximation

# Hybrid Monte Carlo Algorithm

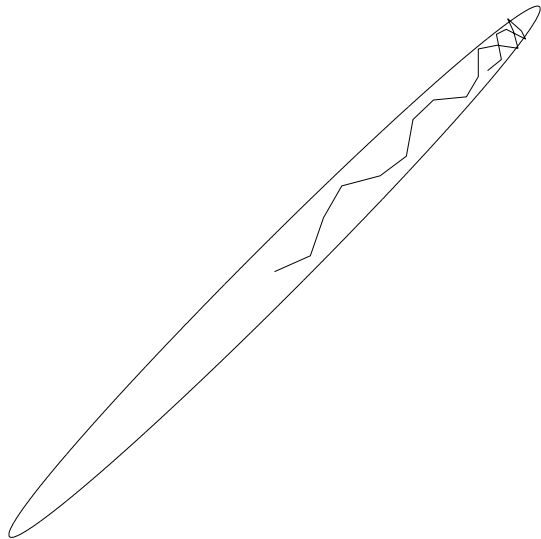
1) A new state is proposed by deterministically simulating a trajectory with  $L$  discrete steps from  $(\mathbf{x}, \mathbf{v})$  to  $(\mathbf{x}^*, \mathbf{v}^*)$ . The new state  $(\mathbf{x}^*, \mathbf{v}^*)$  is **accepted** with probability:

$$\min(1, \exp(-(H(\mathbf{v}^*, \mathbf{x}^*) - H(\mathbf{v}, \mathbf{x}))))),$$

otherwise the state remains the same.

2) Stochastically update the momenta using Gibbs sampling

$$\mathbf{v} \sim p(\mathbf{v}|\mathbf{x}) = p(\mathbf{v}) = N(0, I)$$



Example:  $L = 20$  leapfrog iterations when sampling from a bivariate Gaussian

# Other Ideas

- **MCMC importance sampling.**

- Start from known distribution ( $\mathbf{x}_0 \sim p_0$ ).
- Run MCMC  $n$  (fixed) steps and compute  $p_n$ .
- Importance sample  $w_i = p(\mathbf{x}_i^{(n)})/p_n(\mathbf{x}_i^{(n)})$ .

**Annealed importance sampling** uses annealed MCMC.

- **Coupling from the past** yields **exact** samples, by explicitly checking for mixing wrt a fixed set of random numbers (inputs to transitions).

- Applies to specific systems, where transitions can be written to agglomerate.

- **Slice sampling.**

- **Nested sampling.**

- ...