# Probabilistic Graphical Models

Brown University CSCI 2950-P, Spring 2013 Prof. Erik Sudderth

Lecture 16: Markov Chain Monte Carlo Methods, Metropolis-Hastings Algorithm, Gibbs Sampler

> Some slides and figures courtesy lain Murray's tutorial, Markov Chain Monte Carlo, MLSS 2009

# **General Sequential Monte Carlo**



Exploit temporal structure to propose sequences recursively:

$$\begin{split} q_{0:t}(x_{0:t}|y_{0:t}) &= \overbrace{q_{0:t-1}(x_{0:t-1}|y_{0:t-1})}^{\text{Keep existing path}} \overbrace{q_{t}(x_{t}|x_{t-1},y_{t})}^{\text{extend path}} \\ \tilde{\omega}_{t}^{(i)} &= \frac{\pi_{0:t|0:t}\left(\tilde{x}_{0:t}^{(i)}|y_{0:t}\right)}{q_{0:t}\left(\tilde{x}_{0:t}^{(i)}|y_{0:t}\right)} \propto \omega_{t-1}^{(i)} \times \frac{f\left(\tilde{x}_{t}^{(i)}|\tilde{x}_{t-1}^{(i)}\right)g\left(y_{t}|\tilde{x}_{t}^{(i)}\right)}{q_{t}\left(\tilde{x}_{0:t}^{(i)}|y_{0:t}\right)} \end{split}$$

• Common choices of proposal distribution:

$$q_t(x_t \mid x_{t-1}, y_t) = f(x_t \mid x_{t-1})$$
$$q_t(x_t \mid x_{t-1}, y_t) \approx p(x_t \mid x_{t-1}, y_t)$$

Bootstrap filter simulates prior dynamical model

*If local posterior intractable, use Gaussian approximations* 

## **Particle Resampling or Selection**

$$p(x_t \mid y_{\bar{t}}) \approx \sum_{\ell=1}^{L} \omega_t^{(\ell)} \delta_{x_t^{(\ell)}}(x_t)$$

Resampling with replacement produces a random discrete distribution whose mean is the original distribution

While remaining unbiased, resampling avoids degeneracies in which most weights go to zero









- What is the probability that a state sequence, sampled from the prior model, is consistent with all observations?
- Marginal estimates degenerate on single, mediocre sample



- After each resampling step, some particles are discarded, and can never be restored in subsequent stages.
- Estimates of *smoothed* marginals/sequences typically poor

# **Beyond Temporal Particle Filters**

- Can we avoid degeneracies in estimating "smoothed" marginals based on all observations, past and future?
- Can we implement particle-based approximations to BP for tree-structured models, or arbitrary factor graphs?
- Yes! We can apply importance sampling with resampling to any sequence of distributions, such as a sequence of approximations to a complicated graphical model:



Hot Coupling, Hamze & de Freitas, NIPS 2006

• But to do this well, we need to not only be able to reweight and resample particles, we must shift their locations

Markov Chain Monte Carlo (MCMC) methods

# Markov Chain Monte Carlo (MCMC) $\xrightarrow{z^{(0)}} \xrightarrow{z^{(1)}} \xrightarrow{z^{(2)}} \xrightarrow{z^{(t+1)}} \sim q(z \mid z^{(t)})$

- At each time point, state  $z^{(t)}$  is a configuration of *all the variables in the model:* parameters, hidden variables, etc.
- We design the transition distribution  $q(z \mid z^{(t)})$  so that the chain is *irreducible* and *ergodic*, with a unique stationary distribution  $p^*(z)$

$$p^*(z) = \int_{\mathcal{Z}} q(z \mid z') p^*(z') \, dz'$$

- For learning, the target equilibrium distribution is usually the posterior distribution given data *x*:  $p^*(z) = p(z \mid x)$
- Popular recipes: *Metropolis-Hastings and Gibbs samplers*

#### **Importance Sampling for Regression**



w = 0.00548 w = 1.59e-08 w = 9.65e-06 w = 0.371 w = 0.103



w = 1.01e-08 w = 0.111 w = 1.92e-09 w = 0.0126 w = 1.1e-51

- Model: Gaussian noise around some unknown straight line
- Propose from prior on lines, weight by data likelihood

#### Metropolis Algorithm for Regression



- Perturb parameters:  $Q(\theta';\theta)$ , e.g.  $\mathcal{N}(\theta,\sigma^2)$
- Accept with probability  $\min\left(1, \frac{\tilde{P}(\theta'|\mathcal{D})}{\tilde{P}(\theta|\mathcal{D})}\right)$
- $\frac{\tilde{P}(\theta'|\mathcal{D})}{\tilde{P}(\theta|\mathcal{D})}\right)^{2.5} \\
  \frac{1.5}{1} \\
  \frac{1.5}{0.5}$



• Otherwise keep old parameters

Detail: Metropolis, as stated, requires  $Q(\theta'; \theta) = Q(\theta; \theta')$ 

#### **Markov Chain Monte Carlo** Construct a biased random walk that explores target dist $P^{\star}(x)$

Markov steps,  $x_t \sim T(x_t \leftarrow x_{t-1})$ 



MCMC gives approximate, correlated samples from  $P^{\star}(x)$ 

# Transition (Proposal) Distributions Discrete example

$$P^{\star} = \begin{pmatrix} 3/5\\1/5\\1/5 \end{pmatrix} \qquad T = \begin{pmatrix} 2/3 & 1/2 & 1/2\\1/6 & 0 & 1/2\\1/6 & 1/2 & 0 \end{pmatrix} \qquad T_{ij} = T(x_i \leftarrow x_j)$$

 $P^{\star}$  is an invariant distribution of T because  $TP^{\star} = P^{\star}$ , i.e.

$$\sum_{x} T(x' \leftarrow x) P^{\star}(x) = P^{\star}(x')$$

Also  $P^*$  is the equilibrium distribution of T: To machine precision:  $T^{100} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 3/5 \\ 1/5 \\ 1/5 \end{pmatrix} = P^*$ 

*Ergodicity* requires:  $T^{K}(x' \leftarrow x) > 0$  for all  $x' : P^{\star}(x') > 0$ , for some K

#### Sufficient: Detailed Balance

Detailed balance means  $\rightarrow x \rightarrow x'$  and  $\rightarrow x' \rightarrow x$  are equally probable:



**Detailed balance implies the invariant condition:** 

$$\sum_{x} T(x' \leftarrow x) P^{\star}(x) = P^{\star}(x') \sum_{x} T(x \leftarrow x')^{\perp}$$

Enforcing detailed balance is easy: it only involves isolated pairs

#### **Necessary: Generalized Balance**

If T satisfies stationarity, we can define a reverse operator

$$\widetilde{T}(x \leftarrow x') \propto T(x' \leftarrow x) P^{\star}(x) = \frac{T(x' \leftarrow x) P^{\star}(x)}{\sum_{x} T(x' \leftarrow x) P^{\star}(x)} = \frac{T(x' \leftarrow x) P^{\star}(x)}{P^{\star}(x')}$$

#### **Generalized balance condition:**

$$T(x'\!\leftarrow\!x)P^{\star}(x) \ = \ \widetilde{T}(x\!\leftarrow\!x')P^{\star}(x')$$

also implies the invariant condition and is necessary.

Operators satisfying detailed balance are their own reverse operator.

### **Metropolis-Hastings Algorithm**

#### **Transition operator**

- Propose a move from the current state Q(x';x) , e.g.  $\mathcal{N}(x,\sigma^2)$
- Accept with probability  $\min\left(1, \frac{P(x')Q(x;x')}{P(x)Q(x';x)}\right)$
- Otherwise next state in chain is a copy of current state

#### Notes

- Can use  $\tilde{P} \propto P(x)$ ; normalizer cancels in acceptance ratio
- Satisfies detailed balance (shown below)
- Q must be chosen to fulfill the other technical requirements  $P(x) \cdot T(x' \leftarrow x) = P(x) \cdot Q(x'; x) \min\left(1, \frac{P(x')Q(x; x')}{P(x)Q(x'; x)}\right)$   $= \min\left(P(x)Q(x'; x), P(x')Q(x; x')\right)$   $= P(x') \cdot Q(x; x') \min\left(1, \frac{P(x)Q(x'; x)}{P(x')Q(x; x')}\right) = P(x') \cdot T(x \leftarrow x')$

**Example: Gaussian Metropolis Explore**  $\mathcal{N}(0,1)$  with different step sizes  $\sigma$  $P(x) = \mathcal{N}(x \mid 0, 1)$   $Q(x'; x) = \mathcal{N}(x' \mid x, \sigma^2)$ 



## Limitations: Metropolis Random Walk



Generic proposals use  $Q(x';x) = \mathcal{N}(x,\sigma^2)$ 

 $\sigma \; {\rm large} \to {\rm many} \; {\rm rejections}$ 

 $\sigma$  small  $\rightarrow$  slow diffusion:  $\sim (L/\sigma)^2$  iterations required

#### Limitations: Metropolis Random Walk

Discrete target distribution is uniform over all states

$$Q(x';x) = \begin{cases} \frac{1}{2} & x' = x \pm 1\\ 0 & \text{otherwise} \end{cases}$$



#### **Combining MCMC Transition Proposals**

A sequence of operators, each with  $P^*$  invariant:

- $x_{0} \sim P^{\star}(x)$   $x_{1} \sim T_{a}(x_{1} \leftarrow x_{0}) \qquad P(x_{1}) = \sum_{x_{0}} T_{a}(x_{1} \leftarrow x_{0})P^{\star}(x_{0}) = P^{\star}(x_{1})$   $x_{2} \sim T_{b}(x_{2} \leftarrow x_{1}) \qquad P(x_{2}) = \sum_{x_{1}} T_{b}(x_{2} \leftarrow x_{1})P^{\star}(x_{1}) = P^{\star}(x_{2})$   $x_{3} \sim T_{c}(x_{3} \leftarrow x_{2}) \qquad P(x_{3}) = \sum_{x_{1}} T_{c}(x_{3} \leftarrow x_{2})P^{\star}(x_{2}) = P^{\star}(x_{3})$ ...
  - Combination  $T_cT_bT_a$  leaves  $P^{\star}$  invariant
  - If they can reach any x,  $T_cT_bT_a$  is a valid MCMC operator
  - Individually  $T_c$ ,  $T_b$  and  $T_a$  need not be ergodic

## **Gibbs Samplers**

A method with no rejections:

- Initialize  $\mathbf{x}$  to some value
- Pick each variable in turn or randomly and resample  $P(x_i | \mathbf{x}_{j \neq i})$



Figure from PRML, Bishop (2006)

**Proof of validity:** a) check detailed balance for component update. b) Metropolis–Hastings 'proposals'  $P(x_i|\mathbf{x}_{j\neq i}) \Rightarrow$  accept with prob. 1 Apply a series of these operators. Don't need to check acceptance.

## **Gibbs Samplers**

A method with no rejections:

- Initialize  $\mathbf{x}$  to some value
- Pick each variable in turn or randomly and resample  $P(x_i | \mathbf{x}_{j \neq i})$

#### **Alternative Justification:**

At equilibrium can assume  $\mathbf{x} \sim P(\mathbf{x})$ 

Figure from PRML, Bishop (2006)

Consistent with  $\mathbf{x}_{j\neq i} \sim P(\mathbf{x}_{j\neq i}), \ x_i \sim P(x_i | \mathbf{x}_{j\neq i})$ 

Pretend  $x_i$  was never sampled and do it again.



## **Gibbs Sampling Implementation**

Gibbs sampling benefits from few free choices and convenient features of conditional distributions:

• Conditionals with a few discrete settings can be explicitly normalized:

$$P(x_i | \mathbf{x}_{j \neq i}) \propto P(x_i, \mathbf{x}_{j \neq i})$$
  
= 
$$\frac{P(x_i, \mathbf{x}_{j \neq i})}{\sum_{x'_i} P(x'_i, \mathbf{x}_{j \neq i})} \leftarrow \text{this sum is small and easy}$$

 $Y_6$ 

Y<sub>5</sub>

 $Y_1$ 

Y2

- Continuous conditionals only univariate
  - $\Rightarrow$  amenable to standard sampling methods.
  - Inverse CDF sampling
  - Rejection sampling
  - Slice sampling

▶ ...

## **Undirected Graphical Models**



 $p(x_A, x_C \mid x_B) = p(x_A \mid x_B)p(x_C \mid x_B)$ 

• This global Markov property implies a local Markov property:

$$p(x_i \mid x_{\mathcal{V}\setminus i}) = p(x_i \mid x_{\Gamma(i)})$$

- Practical benefits of Gibbs sampling algorithm:
  - Model and algorithm have same modular structure
  - Conditionals can often be evaluated quickly, because they depend only on the neighboring nodes
  - Exponential families offer further efficiency improvements, by caching and recursively updating sufficient statistics